

# Lecture Notes for MAT-INF 4130

Tom Lyche

July 2, 2013



# Contents

|   |           |
|---|-----------|
| <b>Preface</b>  | <b>ix</b> |
| <b>0 A Short Review of Linear Algebra</b>                           | <b>1</b>  |
| 0.1 Notation . . . . .  | 1         |
| 0.2 Vector Spaces and Subspaces . . . . .                           | 4         |
| 0.2.1 Linear independence and bases . . . . .                       | 6         |
| 0.2.2 Subspaces . . . . .   | 8         |
| 0.2.3 The vector spaces $\mathbb{R}^n$ and $\mathbb{C}^n$ . . . . . | 11        |
| 0.3 Vector Norms . . . . .  | 11        |
| 0.4 Inner Products . . . . .  | 14        |
| 0.4.1 Real and complex inner products . . . . .                     | 14        |
| 0.4.2 Orthogonality . . . . .                                       | 17        |
| 0.5 Linear Systems . . . . .  | 20        |
| 0.5.1 Basic properties . . . . .                                    | 20        |
| 0.5.2 The inverse matrix . . . . .                                  | 22        |
| 0.6 Determinants . . . . .  | 24        |
| 0.7 Eigenpairs . . . . .  | 29        |
| 0.8 Algorithms and Numerical Stability . . . . .                    | 30        |
| <b>I Direct Methods for Linear Systems</b>                          | <b>33</b> |
| <b>1 A Special Linear System</b>                                    | <b>35</b> |
| 1.1 Gaussian Elimination Example . . . . .                          | 36        |
| 1.2 The Tridiagonal Second Derivative Matrix . . . . .              | 38        |
| 1.3 LU Factorization of a Tridiagonal System . . . . .              | 39        |
| 1.3.1 Algorithms . . . . .  | 39        |
| 1.3.2 Diagonal dominance . . . . .                                  | 41        |
| 1.4 Block Multiplication . . . . .                                  | 44        |
| 1.5 Triangular Matrices; Basic facts . . . . .                      | 48        |

---

|          |  |            |
|----------|--|------------|
| 1.6      | Review Questions . . . . .   | 50         |
| <b>2</b> | <b>LU Factorizations</b>   | <b>51</b>  |
| 2.1      | Algorithms for triangular systems . . . . .  | 51         |
| 2.2      | The LU Factorization . . . . .   | 54         |
| 2.2.1    | The LU theorem . . . . .   | 55         |
| 2.2.2    | Operation count . . . . .  | 59         |
| 2.3      | The Symmetric LU Factorization . . . . .   | 61         |
| 2.4      | Block LU factorization . . . . .   | 63         |
| 2.5      | Positive Definite and Semidefinite Matrices . . . . .                                | 64         |
| 2.5.1    | Definitions and examples . . . . .   | 64         |
| 2.5.2    | The nonsymmetric case . . . . .  | 67         |
| 2.5.3    | The symmetric case . . . . .   | 68         |
| 2.6      | The Cholesky Factorization . . . . .   | 70         |
| 2.7      | The Symmetric Positive Semidefinite Case . . . . .                                   | 72         |
| 2.8      | Semi-Cholesky factorization of a banded matrix . . . . .                             | 74         |
| 2.9      | Gaussian Elimination . . . . .   | 77         |
| 2.9.1    | Reduction to upper triangular form . . . . .   | 77         |
| 2.9.2    | Pivot strategies . . . . .   | 78         |
| 2.9.3    | Permutation matrices . . . . .   | 80         |
| 2.9.4    | Gauss transformations . . . . .  | 81         |
| 2.9.5    | PLU factorization . . . . .  | 82         |
| 2.9.6    | The LU factorization . . . . .   | 85         |
| 2.10     | Review Questions . . . . .   | 87         |
| <b>3</b> | <b>The Kronecker Product</b>   | <b>89</b>  |
| 3.1      | Test Matrices . . . . .  | 89         |
| 3.1.1    | The 2D Poisson problem . . . . .   | 90         |
| 3.1.2    | The test matrices . . . . .  | 92         |
| 3.2      | The Kronecker Product . . . . .  | 94         |
| 3.3      | Properties of the 1D and 2D Test Matrices . . . . .                                  | 97         |
| 3.4      | Review Questions . . . . .   | 102        |
| <b>4</b> | <b>Fast Direct Solution of a Large Linear System</b>                                 | <b>103</b> |
| 4.1      | Algorithms for a Banded Positive Definite System . . . . .                           | 103        |
| 4.1.1    | Cholesky factorization . . . . .   | 104        |
| 4.1.2    | Block LU factorization of a block tridiagonal matrix                                 | 104        |
| 4.1.3    | Other methods . . . . .  | 105        |
| 4.2      | A Fast Poisson Solver based on Diagonalization . . . . .                             | 105        |
| 4.3      | A Fast Poisson Solver based on the discrete sine and Fourier<br>transforms . . . . . | 107        |
| 4.3.1    | The discrete sine transform (DST) . . . . .  | 108        |

|           |  |            |
|-----------|--|------------|
| 4.3.2     | The discrete Fourier transform (DFT) . . . . .                               | 108        |
| 4.3.3     | The fast Fourier transform (FFT) . . . . .                                   | 110        |
| 4.3.4     | A poisson solver based on the FFT . . . . .                                  | 113        |
| 4.4       | Review Questions . . . . .   | 116        |
| <b>II</b> | <b>Some Matrix Theory</b>  | <b>117</b> |
| <b>5</b>  | <b>Matrix Reduction by Similarity Transformations</b>                        | <b>119</b> |
| 5.1       | Some Properties of Eigenpairs . . . . .                                      | 119        |
| 5.1.1     | Transformations of eigenpairs and trace . . . . .                            | 119        |
| 5.1.2     | Similarity transformations . . . . .   | 122        |
| 5.2       | Unitary Similarity Transformations . . . . .                                 | 123        |
| 5.2.1     | Unitary and orthonormal and matrices . . . . .                               | 123        |
| 5.2.2     | The Schur decomposition . . . . .  | 124        |
| 5.2.3     | Normal matrices . . . . .  | 127        |
| 5.3       | Minmax theorems for Hermitian Matrices . . . . .                             | 129        |
| 5.3.1     | The Rayleigh quotient . . . . .  | 129        |
| 5.3.2     | Minmax and maxmin . . . . .  | 129        |
| 5.3.3     | The Hoffman-Wielandt theorem . . . . .                                       | 132        |
| 5.4       | The Jordan Form . . . . .  | 133        |
| 5.4.1     | Diagonalizable matrices and linear independence<br>of eigenvectors . . . . . | 133        |
| 5.4.2     | Algebraic and geometric multiplicity of eigenvalues                          | 134        |
| 5.4.3     | The Jordan form . . . . .  | 136        |
| 5.5       | The Minimal Polynomial . . . . .   | 139        |
| 5.6       | Left Eigenvectors . . . . .  | 141        |
| 5.7       | Proof of the Real Schur Form . . . . .                                       | 143        |
| 5.8       | Conclusions . . . . .  | 144        |
| 5.9       | Review Questions . . . . .   | 145        |
| <b>6</b>  | <b>The Singular Value Decomposition</b>                                      | <b>147</b> |
| 6.1       | SVD and SVF . . . . .  | 147        |
| 6.1.1     | Definition and examples . . . . .  | 147        |
| 6.1.2     | Existence . . . . .  | 149        |
| 6.1.3     | The singular value factorization . . . . .                                   | 151        |
| 6.1.4     | Examples . . . . .   | 152        |
| 6.2       | SVD and the Four Fundamental Subspaces . . . . .                             | 155        |
| 6.3       | A Geometric Interpretation . . . . .   | 157        |
| 6.4       | Determining the Rank of a Matrix Numerically . . . . .                       | 158        |
| 6.4.1     | The Frobenius norm . . . . .   | 159        |
| 6.4.2     | Low rank approximation . . . . .   | 160        |

|            |   |            |
|------------|---|------------|
| 6.5        | The Minmax Theorem for Singular Values and the Hoffman-Wielandt Theorem . . . . . | 161        |
| 6.6        | Proof of the Hoffman-Wielandt Theorem for Singular Values .                       | 162        |
| 6.7        | Review Questions . . . . .  | 163        |
| <b>7</b>   | <b>Matrix Norms</b>   | <b>165</b> |
| 7.1        | Matrix Norms . . . . .  | 165        |
| 7.1.1      | Consistent and subordinate matrix norms . . . . .                                 | 166        |
| 7.1.2      | Operator norms . . . . .  | 168        |
| 7.1.3      | The operator $p$ -norms . . . . .   | 169        |
| 7.1.4      | Unitary invariant matrix norms . . . . .  | 172        |
| 7.1.5      | Absolute and monotone norms . . . . .   | 173        |
| 7.2        | The Condition Number with Respect to Inversion . . . . .                          | 173        |
| 7.3        | Proof that the $p$ -Norms are Norms . . . . .                                     | 179        |
| 7.4        | Review Questions . . . . .  | 185        |
| <b>III</b> | <b>Iterative Methods for Large Linear Systems</b>                                 | <b>187</b> |
| <b>8</b>   | <b>The Classical Iterative Methods</b>  | <b>189</b> |
| 8.1        | Classical Iterative Methods; Component Form . . . . .                             | 190        |
| 8.1.1      | The discrete Poisson system . . . . .   | 192        |
| 8.2        | Classical Iterative Methods; Matrix Form . . . . .                                | 195        |
| 8.2.1      | Fixed-point form . . . . .  | 196        |
| 8.2.2      | The preconditioning and splitting matrix . . . . .                                | 196        |
| 8.2.3      | The splitting matrices for the classical methods . .                              | 196        |
| 8.3        | Convergence . . . . .   | 198        |
| 8.3.1      | Convergence of Richardson's method. . . . .                                       | 199        |
| 8.3.2      | Convergence of SOR . . . . .  | 200        |
| 8.3.3      | Convergence of the classical methods for the discrete Poisson matrix . . . . .    | 202        |
| 8.3.4      | Number of iterations . . . . .  | 204        |
| 8.3.5      | Stopping the iteration . . . . .  | 206        |
| 8.4        | Powers of a matrix . . . . .  | 207        |
| 8.4.1      | The spectral radius . . . . .   | 207        |
| 8.4.2      | Neumann series . . . . .  | 209        |
| 8.5        | The Optimal SOR Parameter $\omega$ . . . . .                                      | 210        |
| 8.6        | Review Questions . . . . .  | 213        |
| <b>9</b>   | <b>The Conjugate Gradient Method</b>  | <b>215</b> |
| 9.1        | Quadratic Minimization and Steepest Descent . . . . .                             | 216        |
| 9.2        | The Conjugate Gradient Method . . . . .   | 219        |

|   |  |            |
|---|--|------------|
| 9.2.1   | Derivation of the method . . . . .   | 219        |
| 9.2.2   | The conjugate gradient algorithm . . . . .                                 | 222        |
| 9.2.3   | Numerical example . . . . .  | 222        |
| 9.2.4   | Implementation issues . . . . .  | 223        |
| 9.3   | Convergence . . . . .  | 225        |
| 9.3.1   | The $\mathbf{A}$ -norm . . . . .   | 225        |
| 9.3.2   | The Main Theorem . . . . .   | 225        |
| 9.3.3   | The number of iterations for the model problems .                          | 226        |
| 9.4   | Proof of the Convergence Estimates . . . . .                               | 227        |
| 9.4.1   | Convergence proof for steepest descent . . . . .                           | 227        |
| 9.4.2   | Krylov spaces and the best approximation property                          | 229        |
| 9.4.3   | Chebyshev polynomials . . . . .  | 233        |
| 9.4.4   | Monotonicity of the error . . . . .  | 236        |
| 9.5   | Preconditioning . . . . .  | 237        |
| 9.6   | Preconditioning Example . . . . .  | 240        |
| 9.6.1   | A variable coefficient problem . . . . .                                   | 240        |
| 9.6.2   | Applying preconditioning . . . . .   | 243        |
| 9.7   | Review Questions . . . . .   | 245        |
| <b>IV Orthonormal Transformations and Least Squares</b> |  | <b>247</b> |
| <b>10</b>   | <b>Orthonormal and Unitary Transformations</b>                             | <b>249</b> |
| 10.1  | The Householder Transformation . . . . .                                   | 250        |
| 10.2  | Householder Triangulation . . . . .  | 253        |
| 10.2.1  | Solving linear systems using unitary transformations                       | 255        |
| 10.2.2  | The number of arithmetic operations . . . . .                              | 256        |
| 10.3  | The QR Decomposition and QR Factorization . . . . .                        | 256        |
| 10.3.1  | Existence . . . . .  | 256        |
| 10.3.2  | QR and Gram-Schmidt . . . . .  | 259        |
| 10.4  | Givens Rotations . . . . .   | 260        |
| 10.5  | Review Questions . . . . .   | 263        |
| <b>11</b>   | <b>Least Squares</b>   | <b>265</b> |
| 11.1  | Numerical Examples . . . . .   | 266        |
| 11.2  | Curve Fitting . . . . .  | 267        |
| 11.3  | Least Squares and Singular Value Decomposition and Factorization . . . . . | 270        |
| 11.3.1  | Sum of subspaces and orthogonal projections . . .                          | 271        |
| 11.3.2  | The generalized inverse . . . . .  | 274        |
| 11.4  | Numerical Solution . . . . .   | 277        |
| 11.4.1  | Normal equations . . . . .   | 277        |

---

|           |  |            |
|-----------|--|------------|
| 11.4.2    | QR factorization . . . . .   | 278        |
| 11.4.3    | Singular value factorization . . . . .   | 279        |
| 11.5      | Perturbation Theory for Least Squares . . . . .                                    | 280        |
| 11.5.1    | Perturbing the right hand side . . . . .   | 280        |
| 11.5.2    | Perturbing the matrix . . . . .  | 282        |
| 11.6      | Perturbation Theory for Singular Values . . . . .                                  | 283        |
| 11.7      | Review Questions . . . . .   | 284        |
| <b>V</b>  | <b>Eigenvalues and Eigenvectors</b>  | <b>287</b> |
| <b>12</b> | <b>Numerical Eigenvalue Problems</b>   | <b>289</b> |
| 12.1      | Eigenpars . . . . .  | 289        |
| 12.2      | Gerschgorin's Theorem . . . . .  | 290        |
| 12.3      | Perturbation of Eigenvalues . . . . .  | 293        |
| 12.4      | Unitary Similarity Transformation of a Matrix into Upper Hessenberg Form . . . . . | 296        |
| 12.5      | Computing a Selected Eigenvalue of a Symmetric Matrix . . . . .                    | 299        |
| 12.5.1    | The inertia theorem . . . . .  | 301        |
| 12.5.2    | Approximating $\lambda_m$ . . . . .  | 303        |
| 12.6      | Review Questions . . . . .   | 304        |
| <b>13</b> | <b>The QR Algorithm</b>  | <b>307</b> |
| 13.1      | The Power Method and its variants . . . . .  | 307        |
| 13.1.1    | The power method . . . . .   | 307        |
| 13.1.2    | The inverse power method . . . . .   | 311        |
| 13.1.3    | Rayleigh quotient iteration . . . . .  | 312        |
| 13.2      | The basic QR Algorithm . . . . .   | 313        |
| 13.2.1    | Relation to the power method . . . . .   | 315        |
| 13.2.2    | Invariance of the Hessenberg form . . . . .  | 316        |
| 13.2.3    | Deflation . . . . .  | 316        |
| 13.3      | The Shifted QR Algorithms . . . . .  | 317        |
| 13.4      | A Convergence Theorem . . . . .  | 318        |
| 13.5      | Review Questions . . . . .   | 319        |
| <b>VI</b> | <b>Appendix</b>  | <b>321</b> |
| <b>A</b>  | <b>Determinants</b>  | <b>323</b> |
| A.1       | Permutations . . . . .   | 323        |
| A.2       | Basic Properties of Determinants . . . . .   | 325        |
| A.3       | The Adjoint Matrix and Cofactor Expansion . . . . .                                | 329        |
| A.4       | Computing Determinants . . . . .   | 332        |



---

|          |  |            |
|----------|--|------------|
| A.5      | Some Useful Determinant Formulas . . . . .   | 332        |
| <b>B</b> | <b>Computer Arithmetic</b>                   | <b>335</b> |
| B.1      | Absolute and Relative Errors . . . . .       | 335        |
| B.2      | Floating Point Numbers . . . . .             | 336        |
| B.3      | Rounding and Arithmetic Operations . . . . . | 339        |
| B.3.1    | Rounding . . . . .                           | 339        |
| B.3.2    | Arithmetic operations . . . . .              | 340        |
| B.4      | Backward Rounding-Error Analysis . . . . .   | 340        |
| B.4.1    | Computing a sum . . . . .                    | 340        |
| B.4.2    | Computing an inner product . . . . .         | 343        |
| B.4.3    | Computing a matrix product . . . . .         | 343        |
| <b>C</b> | <b>Differentiation of Vector Functions</b>   | <b>345</b> |
|          | <b>Bibliography</b>                          | <b>349</b> |
|          | <b>Index</b>                                 | <b>375</b> |



# Preface

These lecture notes contains the text for a course in matrix analysis and numerical linear algebra given at the beginning graduate level at the University of Oslo. Most of the chapters correspond approximately to one week of lectures. Earlier versions of this manuscript were converted to LaTeX by Are Magnus Bruaset and Njål Foldnes. A special thanks goes to Christian Schulz and Georg Muntingh who helped me with the exercise sessions and have provided solutions to all problems in this book.

Oslo, 1. July 2013

Tom Lyche



## Chapter 0

# A Short Review of Linear Algebra

In this introductory chapter we give a compact introduction to linear algebra with emphasis on  $\mathbb{R}^n$  and  $\mathbb{C}^n$ . For a more elementary introduction, see for example the book [22]. We start by introducing the notation used.

## 0.1 Notation

The following sets and notations will be used in this book.

1. The sets of natural numbers, integers, rational numbers, real numbers, and complex numbers are denoted by  $\mathbb{N}, \mathbb{Z}, \mathbb{Q}, \mathbb{R}, \mathbb{C}$ , respectively.
2. We use the “colon equal” symbol  $v := e$  to indicate that the symbol  $v$  is defined by the expression  $e$ .
3.  $\mathbb{R}^n$  is the set of  $n$ -tuples of real numbers which we will represent as column vectors. Thus  $\mathbf{x} \in \mathbb{R}^n$  means

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix},$$

where  $x_i \in \mathbb{R}$  for  $i = 1, \dots, n$ . Row vectors are normally identified using the transpose operation. Thus if  $\mathbf{x} \in \mathbb{R}^n$  then  $\mathbf{x}$  is a column vector and  $\mathbf{x}^T$  is a row vector.

4. Addition and scalar multiplication are denoted and defined by

$$\mathbf{x} + \mathbf{y} = \begin{bmatrix} x_1 + y_1 \\ \vdots \\ x_n + y_n \end{bmatrix}, \quad a\mathbf{x} = \begin{bmatrix} ax_1 \\ \vdots \\ ax_n \end{bmatrix}, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad a \in \mathbb{R}.$$

5.  $\mathbb{R}^{m \times n}$  is the set of matrices  $\mathbf{A}$  with real elements. The integers  $m$  and  $n$  are the number of rows and columns in the tableau

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}.$$

The element in the  $i$ th row and  $j$ th column of  $\mathbf{A}$  will be denoted by  $a_{i,j}$ ,  $a_{ij}$ ,  $\mathbf{A}(i, j)$  or  $(\mathbf{A})_{i,j}$ . We use the notations

$$\mathbf{a}_{:j} = \begin{bmatrix} a_{1j} \\ a_{2j} \\ \vdots \\ a_{mj} \end{bmatrix}, \quad \mathbf{a}_{i:}^T = [a_{i1}, a_{i2}, \dots, a_{in}], \quad \mathbf{A} = [\mathbf{a}_{:1}, \mathbf{a}_{:2}, \dots, \mathbf{a}_{:n}] = \begin{bmatrix} \mathbf{a}_{1:}^T \\ \mathbf{a}_{2:}^T \\ \vdots \\ \mathbf{a}_{m:}^T \end{bmatrix}$$

for the columns  $\mathbf{a}_{:j}$  and rows  $\mathbf{a}_{i:}^T$  of  $\mathbf{A}$ . We often drop the colon and write  $\mathbf{a}_j$  and  $\mathbf{a}_i^T$  with the risk of some confusion. If  $m = 1$  then  $\mathbf{A}$  is a row vector, if  $n = 1$  then  $\mathbf{A}$  is a column vector, while if  $m = n$  then  $\mathbf{A}$  is a square matrix. In this text we will denote matrices by boldface capital letters  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots$  and vectors most often by boldface lower case letters  $\mathbf{x}, \mathbf{y}, \mathbf{z}, \dots$ .

6. The imaginary unit  $\sqrt{-1}$  is denoted by  $i$ . The complex conjugate and the modulus of a complex number  $z$  is denoted by  $\bar{z}$  and  $|z|$ , respectively. Thus if  $z = x + iy = re^{i\phi} = r(\cos \phi + i \sin \phi)$ , with  $x, y \in \mathbb{R}$ , is a complex number then  $\bar{z} := x - iy = re^{-i\phi} = \cos \phi - i \sin \phi$  and  $|z| := \sqrt{\bar{z}z} = \sqrt{x^2 + y^2} = r$ .  $\text{Re}(z) := x$  and  $\text{Im}(z) := y$  denote the real and imaginary part of the complex number  $z$ .
7. For matrices and vectors with complex elements we use the notation  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and  $\mathbf{x} \in \mathbb{C}^n$ . We define complex row vectors using either the transpose  $\mathbf{x}^T$  or the conjugate transpose operation  $\mathbf{x}^* := \bar{\mathbf{x}}^T = [\bar{x}_1, \dots, \bar{x}_n]$ .
8. For  $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$  and  $a \in \mathbb{C}$  the operations of vector addition and scalar multiplication is defined by component operations as in the real case (cf. 4.).
9. The arithmetic operations on rectangular matrices are
- **matrix addition**  $\mathbf{C} = \mathbf{A} + \mathbf{B}$  if  $\mathbf{A}, \mathbf{B}, \mathbf{C}$  are matrices of the same size, i. e., with the same number of rows and columns, and  $c_{ij} = a_{ij} + b_{ij}$  for all  $i, j$ .

- **multiplication by a scalar**  $C = \alpha A$ , where  $c_{ij} = \alpha a_{ij}$  for all  $i, j$ .
  - **matrix multiplication**  $C = AB$ ,  $C = A \cdot B$  or  $C = A * B$ , where  $A \in \mathbb{C}^{m \times p}$ ,  $B \in \mathbb{C}^{p \times n}$ ,  $C \in \mathbb{C}^{m \times n}$ , and  $c_{ij} = \sum_{k=1}^p a_{ik} b_{kj}$  for  $i = 1, \dots, m$ ,  $j = 1, \dots, n$ .
  - **element-by-element matrix operations**  $C = A \times B$ ,  $D = A/B$ , and  $E = A \wedge r$  where all matrices are of the same size and  $c_{ij} = a_{ij} b_{ij}$ ,  $d_{ij} = a_{ij}/b_{ij}$  and  $e_{ij} = a_{ij}^r$  for all  $i, j$  and suitable  $r$ . The element-by-element product  $C = A \times B$  is known as the **Schur product** and also the **Hadamard product**.
10. Let  $A \in \mathbb{R}^{m \times n}$  or  $A \in \mathbb{C}^{m \times n}$ . The **transpose**  $A^T$  and **conjugate transpose**  $A^*$  are  $n \times m$  matrices with elements  $a_{ij}^T = a_{ji}$  and  $a_{ij}^* = \bar{a}_{ji}$ , respectively. If  $B$  is an  $n, p$  matrix then  $(AB)^T = B^T A^T$  and  $(AB)^* = B^* A^*$ .
11. The **unit vectors** in  $\mathbb{R}^n$  and  $\mathbb{C}^n$  are denoted by

$$e_1 := \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad e_2 := \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad e_3 := \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \quad \dots, \quad e_n := \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix},$$

while  $I_n = I := [\delta_{ij}]_{i,j=1}^n$ , where

$$\delta_{ij} := \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

is the **identity matrix** of order  $n$ . Both the columns and the transpose of the rows of  $I$  are the unit vectors  $e_1, e_2, \dots, e_n$ .

12. Some matrices with many zeros have names indicating their “shape”. Suppose  $A \in \mathbb{R}^{n \times n}$  or  $A \in \mathbb{C}^{n \times n}$ . Then  $A$  is
- **diagonal** if  $a_{ij} = 0$  for  $i \neq j$ .
  - **upper triangular** or **right triangular** if  $a_{ij} = 0$  for  $i > j$ .
  - **lower triangular** or **left triangular** if  $a_{ij} = 0$  for  $i < j$ .
  - **upper Hessenberg** if  $a_{ij} = 0$  for  $i > j + 1$ .
  - **lower Hessenberg** if  $a_{ij} = 0$  for  $i < j + 1$ .
  - **tridiagonal** if  $a_{ij} = 0$  for  $|i - j| > 1$ .
  - **$d$ -banded** if  $a_{ij} = 0$  for  $|i - j| > d$ .

13. We use the following notations for diagonal- and tridiagonal  $n \times n$  matrices

$$\text{diag}(d_i) = \text{diag}(d_1, \dots, d_n) := \begin{bmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & d_n \end{bmatrix} = \begin{bmatrix} d_1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & d_n \end{bmatrix},$$

$$\mathbf{B} = \text{tridiag}(a_i, d_i, c_i) = \text{tridiag}(\mathbf{a}, \mathbf{d}, \mathbf{c}) := \begin{bmatrix} d_1 & c_1 & & & \\ a_1 & d_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & a_{n-2} & d_{n-1} & c_{n-1} \\ & & & a_{n-1} & d_n \end{bmatrix}.$$

Here  $b_{ii} = d_i$  for  $i = 1, \dots, n$ ,  $b_{i+1,i} = a_i$ ,  $b_{i,i+1} = c_i$  for  $i = 1, \dots, n-1$ , and  $b_{ij} = 0$  otherwise.

14. Suppose  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and  $1 \leq i_1 < i_2 < \cdots < i_r \leq m$ ,  $1 \leq j_1 < j_2 < \cdots < j_c \leq n$ . The matrix  $\mathbf{A}(\mathbf{i}, \mathbf{j}) \in \mathbb{C}^{r \times c}$  is the submatrix of  $\mathbf{A}$  consisting of rows  $\mathbf{i} := [i_1, \dots, i_r]$  and columns  $\mathbf{j} := [j_1, \dots, j_c]$

$$\mathbf{A}(\mathbf{i}, \mathbf{j}) := \mathbf{A} \begin{pmatrix} i_1 & i_2 & \cdots & i_r \\ j_1 & j_2 & \cdots & j_c \end{pmatrix} = \begin{bmatrix} a_{i_1, j_1} & a_{i_1, j_2} & \cdots & a_{i_1, j_c} \\ a_{i_2, j_1} & a_{i_2, j_2} & \cdots & a_{i_2, j_c} \\ \vdots & \vdots & \ddots & \vdots \\ a_{i_r, j_1} & a_{i_r, j_2} & \cdots & a_{i_r, j_c} \end{bmatrix}.$$

For the special case of consecutive rows and columns we use the notation

$$\mathbf{A}(r_1 : r_2, c_1 : c_2) := \begin{bmatrix} a_{r_1, c_1} & a_{r_1, c_1+1} & \cdots & a_{r_1, c_2} \\ a_{r_1+1, c_1} & a_{r_1+1, c_1+1} & \cdots & a_{r_1+1, c_2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{r_2, c_1} & a_{r_2, c_1+1} & \cdots & a_{r_2, c_2} \end{bmatrix}.$$

## 0.2 Vector Spaces and Subspaces

Many mathematical systems have analogous properties to vectors in  $\mathbb{R}^2$  or  $\mathbb{R}^3$ .

### Definition 0.1 (Real vector space)

A **real vector space** is a nonempty set  $\mathcal{V}$ , whose objects are called **vectors**, together with two operations  $+$  :  $\mathcal{V} \times \mathcal{V} \rightarrow \mathcal{V}$  and  $\cdot$  :  $\mathbb{R} \times \mathcal{V} \rightarrow \mathcal{V}$ , called **addition** and **scalar multiplication**, satisfying the following axioms for all vectors  $\mathbf{u}, \mathbf{v}, \mathbf{w}$  in  $\mathcal{V}$  and scalars  $c, d$  in  $\mathbb{R}$ .

(V1) The sum  $\mathbf{u} + \mathbf{v}$  is in  $\mathcal{V}$ ,



$$\text{(V2)} \quad \mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u},$$

$$\text{(V3)} \quad \mathbf{u} + (\mathbf{v} + \mathbf{w}) = (\mathbf{u} + \mathbf{v}) + \mathbf{w},$$

$$\text{(V4)} \quad \text{There is a zero vector } \mathbf{0} \text{ such that } \mathbf{u} + \mathbf{0} = \mathbf{u},$$

$$\text{(V5)} \quad \text{For each } \mathbf{u} \text{ in } \mathcal{V} \text{ there is a vector } -\mathbf{u} \text{ in } \mathcal{V} \text{ such that } \mathbf{u} + (-\mathbf{u}) = \mathbf{0},$$

$$\text{(S1)} \quad \text{The scalar multiple } c \cdot \mathbf{u} \text{ is in } \mathcal{V},$$

$$\text{(S2)} \quad c \cdot (\mathbf{u} + \mathbf{v}) = c \cdot \mathbf{u} + c \cdot \mathbf{v},$$

$$\text{(S3)} \quad (c + d) \cdot \mathbf{u} = c \cdot \mathbf{u} + d \cdot \mathbf{u},$$

$$\text{(S4)} \quad c \cdot (d \cdot \mathbf{u}) = (cd) \cdot \mathbf{u},$$

$$\text{(S5)} \quad 1 \cdot \mathbf{u} = \mathbf{u}.$$

The scalar multiplication symbol  $\cdot$  is often omitted, writing  $c\mathbf{v}$  instead of  $c \cdot \mathbf{v}$ . We define  $\mathbf{u} - \mathbf{v} := \mathbf{u} + (-\mathbf{v})$ . We call  $\mathcal{V}$  a **complex vector space** if the scalars consist of all complex numbers  $\mathbb{C}$ . In this book a vector space is either real or complex.

From the axioms it follows that

1. The zero vector is unique.
2. For each  $\mathbf{u} \in \mathcal{V}$  the **negative**  $-\mathbf{u}$  of  $\mathbf{u}$  is unique.
3.  $0\mathbf{u} = \mathbf{0}$ ,  $c\mathbf{0} = \mathbf{0}$ , and  $-\mathbf{u} = (-1)\mathbf{u}$ .

Here are some examples

1. The space  $\mathbb{R}^n$ , where  $n \in \mathbb{N}$ , is a real vector space.
2. Similarly,  $\mathbb{C}^n$  is a complex vector space.
3. Let  $\mathcal{D}$  be a subset of  $\mathbb{R}$  and  $d \in \mathbb{N}$ . The set  $\mathcal{V}$  of all functions  $\mathbf{f}, \mathbf{g} : \mathcal{D} \rightarrow \mathbb{R}^d$  is a real vector space with

$$(\mathbf{f} + \mathbf{g})(t) := \mathbf{f}(t) + \mathbf{g}(t), \quad (c\mathbf{f})(t) := c\mathbf{f}(t), \quad t \in \mathcal{D}, \quad c \in \mathbb{R}.$$

Two functions  $\mathbf{f}, \mathbf{g}$  in  $\mathcal{V}$  are equal if  $\mathbf{f}(t) = \mathbf{g}(t)$  for all  $t \in \mathcal{D}$ . The zero element is the **zero function** given by  $\mathbf{f}(t) = \mathbf{0}$  for all  $t \in \mathcal{D}$  and the negative of  $\mathbf{f}$  is given by  $-\mathbf{f} = (-1)\mathbf{f}$ . In the following we will use boldface letters for functions only if  $d > 1$ .

4. For  $n \geq 0$  the space  $\Pi_n$  of polynomials of degree at most  $n$  consists of all polynomials  $p: \mathbb{R} \rightarrow \mathbb{R}$ ,  $p: \mathbb{R} \rightarrow \mathbb{C}$ , or  $p: \mathbb{C} \rightarrow \mathbb{C}$  of the form

$$p(t) = a_0 + a_1t + a_2t^2 + \cdots + a_nt^n, \quad (2)$$

where the coefficients  $a_0, \dots, a_n$  are real or complex numbers.  $p$  is called the **zero polynomial** if all coefficients are zero. All other polynomials are said to be **nontrivial**. The **degree** of a nontrivial polynomial  $p$  given by (2) is the smallest integer  $0 \leq k \leq n$  such that  $p(t) = a_0 + \cdots + a_kt^k$  with  $a_k \neq 0$ . The degree of the zero polynomial is not defined.  $\Pi_n$  is a vector space if we define addition and scalar multiplication as for functions.

### Definition 0.2 (Linear combination)

For  $n \geq 1$  let  $\mathcal{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be a set of vectors in a vector space  $\mathcal{V}$  and let  $c_1, \dots, c_n$  be scalars.

1. The sum  $c_1\mathbf{x}_1 + \cdots + c_n\mathbf{x}_n$  is called a **linear combination** of  $\mathbf{x}_1, \dots, \mathbf{x}_n$ .
2. The linear combination is **nontrivial** if  $c_j\mathbf{x}_j \neq \mathbf{0}$  for at least one  $j$ .
3. The set of all linear combinations of elements in  $\mathcal{X}$  is denoted  $\text{span}(\mathcal{X})$ .
4. A vector space is **finite dimensional** if it has a finite spanning set; i. e., there exists  $n \in \mathbb{N}$  and  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  in  $\mathcal{V}$  such that  $\mathcal{V} = \text{span}(\{\mathbf{x}_1, \dots, \mathbf{x}_n\})$ .

### Example 0.3 (Linear combinations)

1. Any  $\mathbf{x} = [x_1, \dots, x_m]^T$  in  $\mathbb{C}^m$  can be written as a linear combination of the unit vectors as  $\mathbf{x} = x_1\mathbf{e}_1 + x_2\mathbf{e}_2 + \cdots + x_m\mathbf{e}_m$ . Thus,  $\mathbb{C}^m = \text{span}(\{\mathbf{e}_1, \dots, \mathbf{e}_m\})$  and  $\mathbb{C}^m$  is finite dimensional. Similarly  $\mathbb{R}^m$  is finite dimensional.
2. Let  $\Pi = \cup_n \Pi_n$  be the space of all polynomials.  $\Pi$  is a vector space that is not finite dimensional. For suppose  $\Pi$  is finite dimensional. Then  $\Pi = \text{span}(\{p_1, \dots, p_m\})$  for some polynomials  $p_1, \dots, p_m$ . Let  $d$  be an integer such that the degree of  $p_j$  is less than  $d$  for  $j = 1, \dots, m$ . A polynomial of degree  $d$  cannot be written as a linear combination of  $p_1, \dots, p_m$ , a contradiction.

## 0.2.1 Linear independence and bases

### Definition 0.4 (Linear independence)

A set  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  of nonzero vectors in a vector space is **linearly dependent** if  $\mathbf{0}$  can be written as a nontrivial linear combination of  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ . Otherwise  $\mathcal{X}$  is **linearly independent**.

A set of vectors  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  is linearly independent if and only if

$$c_1\mathbf{x}_1 + \dots + c_n\mathbf{x}_n = \mathbf{0} \quad \implies \quad c_1 = \dots = c_n = 0. \quad (3)$$

Suppose  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  is linearly independent. Then

1. If  $\mathbf{x} \in \text{span}(\mathcal{X})$  then the scalars  $c_1, \dots, c_n$  in the representation  $\mathbf{x} = c_1\mathbf{x}_1 + \dots + c_n\mathbf{x}_n$  are unique.
2. Any nontrivial linear combination of  $\mathbf{x}_1, \dots, \mathbf{x}_n$  is nonzero,

**Lemma 0.5 (Linear independence and span)**

Suppose  $\mathbf{v}_1, \dots, \mathbf{v}_n$  span a vector space  $\mathcal{V}$  and that  $\mathbf{w}_1, \dots, \mathbf{w}_k$  are linearly independent vectors in  $\mathcal{V}$ . Then  $k \leq n$ .

*Proof.* Suppose  $k > n$ . Write  $\mathbf{w}_1$  as a linear combination of elements from the set  $\mathcal{X}_0 := \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ , say  $\mathbf{w}_1 = c_1\mathbf{v}_1 + \dots + c_n\mathbf{v}_n$ . Since  $\mathbf{w}_1 \neq \mathbf{0}$  not all the  $c$ 's are equal to zero. Pick a nonzero  $c$ , say  $c_{i_1}$ . Then  $\mathbf{v}_{i_1}$  can be expressed as a linear combination of  $\mathbf{w}_1$  and the remaining  $\mathbf{v}$ 's. So the set  $\mathcal{X}_1 := \{\mathbf{w}_1, \mathbf{v}_1, \dots, \mathbf{v}_{i_1-1}, \mathbf{v}_{i_1+1}, \dots, \mathbf{v}_n\}$  must also be a spanning set for  $\mathcal{V}$ . We repeat this for  $\mathbf{w}_2$  and  $\mathcal{X}_1$ . In the linear combination  $\mathbf{w}_2 = d_{i_1}\mathbf{w}_1 + \sum_{j \neq i_1} d_j\mathbf{v}_j$ , we must have  $d_{i_2} \neq 0$  for some  $i_2$  with  $i_2 \neq i_1$ . For otherwise  $\mathbf{w}_2 = d_1\mathbf{w}_1$  contradicting the linear independence of the  $\mathbf{w}$ 's. So the set  $\mathcal{X}_2$  consisting of the  $\mathbf{v}$ 's with  $\mathbf{v}_{i_1}$  replaced by  $\mathbf{w}_1$  and  $\mathbf{v}_{i_2}$  replaced by  $\mathbf{w}_2$  is again a spanning set for  $\mathcal{V}$ . Repeating this process  $n - 2$  more times we obtain a spanning set  $\mathcal{X}_n$  where  $\mathbf{v}_1, \dots, \mathbf{v}_n$  have been replaced by  $\mathbf{w}_1, \dots, \mathbf{w}_n$ . Since  $k > n$  we can then write  $\mathbf{w}_k$  as a linear combination of  $\mathbf{w}_1, \dots, \mathbf{w}_n$  contradicting the linear independence of the  $\mathbf{w}$ 's. We conclude that  $k \leq n$ .  $\square$

**Definition 0.6 (basis)**

A finite set of vectors  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  in a vector space  $\mathcal{V}$  is a **basis** for  $\mathcal{V}$  if

1.  $\text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\} = \mathcal{V}$ .
2.  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is linearly independent.

**Theorem 0.7 (Basis subset of a spanning set)**

Suppose  $\mathcal{V}$  is a vector space and that  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is a spanning set for  $\mathcal{V}$ . Then we can find a subset  $\{\mathbf{v}_{i_1}, \dots, \mathbf{v}_{i_k}\}$  that forms a basis for  $\mathcal{V}$ .

*Proof.* If  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is linearly dependent we can express one of the  $\mathbf{v}$ 's as a nontrivial linear combination of the remaining  $\mathbf{v}$ 's and drop that  $\mathbf{v}$  from the spanning set. Continue this process until the remaining  $\mathbf{v}$ 's are linearly independent. They still span the vector space and therefore form a basis.  $\square$

**Corollary 0.8 (Existence of a basis)**

A vector space is finite dimensional if and only if it has a basis.

**Proof.** Let  $\mathcal{V} = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  be a finite dimensional vector space. By Theorem 0.7,  $\mathcal{V}$  has a basis. Conversely, if  $\mathcal{V} = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  and  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is a basis then it is by definition a finite spanning set.  $\square$

**Theorem 0.9 (Dimension of a vector space)**

Every basis for a vector space  $\mathcal{V}$  has the same number of elements. This number is called the **dimension** of the vector space and denoted  $\dim \mathcal{V}$ .

**Proof.** Suppose  $\mathcal{X} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  and  $\mathcal{Y} = \{\mathbf{w}_1, \dots, \mathbf{w}_k\}$  are two bases for  $\mathcal{V}$ . By Lemma 0.5 we have  $k \leq n$ . Using the same Lemma with  $\mathcal{X}$  and  $\mathcal{Y}$  switched we obtain  $n \leq k$ . We conclude that  $n = k$ .  $\square$

The set of unit vectors  $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  form a basis for both  $\mathbb{R}^n$  and  $\mathbb{C}^n$ .

**Theorem 0.10 (Enlarging vectors to a basis)**

Every linearly independent set of vectors  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  in a finite dimensional vector space  $\mathcal{V}$  can be enlarged to a basis for  $\mathcal{V}$ .

**Proof.** If  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  does not span  $\mathcal{V}$  we can enlarge the set by one vector  $\mathbf{v}_{k+1}$  which cannot be expressed as a linear combination of  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ . The enlarged set is also linearly independent. Continue this process. Since the space is finite dimensional it must stop after a finite number of steps.  $\square$

**0.2.2 Subspaces****Definition 0.11 (Subspace)**

A nonempty subset  $\mathcal{S}$  of a real or complex vector space  $\mathcal{V}$  is called a **subspace** of  $\mathcal{V}$  if

(V1) The sum  $\mathbf{u} + \mathbf{v}$  is in  $\mathcal{S}$  for any  $\mathbf{u}, \mathbf{v} \in \mathcal{S}$ .

(S1) The scalar multiple  $c\mathbf{u}$  is in  $\mathcal{S}$  for any scalar  $c$  and any  $\mathbf{u} \in \mathcal{S}$ .

Using the operations in  $\mathcal{V}$ , any subspace  $\mathcal{S}$  of  $\mathcal{V}$  is a vector space, i. e., all 10 axioms  $V1 - V5$  and  $S1 - S5$  are satisfied for  $\mathcal{S}$ . In particular,  $\mathcal{S}$  must contain the zero element in  $\mathcal{V}$ . This follows since the operations of vector addition and scalar multiplication are inherited from  $\mathcal{V}$ .

**Example 0.12 (Examples of subspaces)**

1.  $\{\mathbf{0}\}$ , where  $\mathbf{0}$  is the zero vector is a subspace, the **trivial subspace**. The dimension of the trivial subspace is defined to be zero. All other subspaces are **nontrivial**.
2.  $\mathcal{V}$  is a subspace of itself.
3.  $\text{span}(\mathcal{X})$  is a subspace of  $\mathcal{V}$  for any  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathcal{V}$ . Indeed, it is easy to see that **(V1)** and **(S1)** hold.
4. The **sum** of two subspaces  $\mathcal{R}$  and  $\mathcal{S}$  of a vector space  $\mathcal{V}$  is defined by

$$\mathcal{R} + \mathcal{S} := \{\mathbf{r} + \mathbf{s} : \mathbf{r} \in \mathcal{R} \text{ and } \mathbf{s} \in \mathcal{S}\}. \quad (4)$$

Clearly **(V1)** and **(S1)** hold and it is a subspace of  $\mathcal{V}$ .

5. The **intersection** of two subspaces  $\mathcal{R}$  and  $\mathcal{S}$  of a vector space  $\mathcal{V}$  is defined by

$$\mathcal{R} \cap \mathcal{S} := \{\mathbf{x} : \mathbf{x} \in \mathcal{R} \text{ and } \mathbf{x} \in \mathcal{S}\}. \quad (5)$$

It is a subspace of  $\mathcal{V}$ .

6. The **union** of two subspaces  $\mathcal{R}$  and  $\mathcal{S}$  of a vector space  $\mathcal{V}$  is defined by

$$\mathcal{R} \cup \mathcal{S} := \{\mathbf{x} : \mathbf{x} \in \mathcal{R} \text{ or } \mathbf{x} \in \mathcal{S}\}. \quad (6)$$

In general it is not a subspace of  $\mathcal{V}$ .

7. A sum of two subspaces  $\mathcal{R}$  and  $\mathcal{S}$  of a vector space  $\mathcal{V}$  is called a **direct sum** and denoted  $\mathcal{R} \oplus \mathcal{S}$  if  $\mathcal{R} \cap \mathcal{S} = \{\mathbf{0}\}$ . The subspaces  $\mathcal{R}$  and  $\mathcal{S}$  are called **complementary** in the subspace  $\mathcal{R} \oplus \mathcal{S}$ .

### Theorem 0.13 (Dimension formula for sums of subspaces)

Let  $\mathcal{R}$  and  $\mathcal{S}$  be two finite dimensional subspaces of a vector space  $\mathcal{V}$ . Then

$$\dim(\mathcal{R} + \mathcal{S}) = \dim(\mathcal{R}) + \dim(\mathcal{S}) - \dim(\mathcal{R} \cap \mathcal{S}). \quad (7)$$

In particular, for a direct sum

$$\dim(\mathcal{R} \oplus \mathcal{S}) = \dim(\mathcal{R}) + \dim(\mathcal{S}). \quad (8)$$

**Proof.** Let  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  be a basis for  $\mathcal{R} \cap \mathcal{S}$ , where  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\} = \emptyset$ , the empty set, in the case  $\mathcal{R} \cap \mathcal{S} = \{\mathbf{0}\}$ . We use Theorem 0.10 to extend  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  to a basis  $\{\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{r}_1, \dots, \mathbf{r}_q\}$  for  $\mathcal{R}$  and a basis  $\{\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{s}_1, \dots, \mathbf{s}_t\}$  for  $\mathcal{S}$ . Every  $\mathbf{x} \in \mathcal{R} + \mathcal{S}$  can be written as a linear combination of  $\{\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{r}_1, \dots, \mathbf{r}_q, \mathbf{s}_1, \dots, \mathbf{s}_t\}$  so these vectors span  $\mathcal{R} + \mathcal{S}$ . We show that they are linearly independent and hence a basis. Suppose  $\mathbf{u} + \mathbf{r} + \mathbf{s} = \mathbf{0}$ , where  $\mathbf{u} := \sum_{j=1}^p \alpha_j \mathbf{u}_j$ ,  $\mathbf{r} := \sum_{j=1}^q \rho_j \mathbf{r}_j$ ,

and  $\mathbf{s} := \sum_{j=1}^t \sigma_j \mathbf{s}_j$ . Now  $\mathbf{r} = -(\mathbf{u} + \mathbf{s})$  belongs to both  $\mathcal{R}$  and to  $\mathcal{S}$  and hence  $\mathbf{r} \in \mathcal{R} \cap \mathcal{S}$ . Therefore  $\mathbf{r}$  can be written as a linear combination of  $\mathbf{u}_1, \dots, \mathbf{u}_p$  say  $\mathbf{r} := \sum_{j=1}^p \beta_j \mathbf{u}_j$ . But then  $\mathbf{0} = \sum_{j=1}^p \beta_j \mathbf{u}_j - \sum_{j=1}^q \rho_j \mathbf{r}_j$  and since  $\{\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{r}_1, \dots, \mathbf{r}_q\}$  is linearly independent we must have  $\beta_1 = \dots = \beta_p = \rho_1 = \dots = \rho_q = 0$  and hence  $\mathbf{r} = \mathbf{0}$ . We then have  $\mathbf{u} + \mathbf{s} = \mathbf{0}$  and by linear independence of  $\{\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{s}_1, \dots, \mathbf{s}_t\}$  we obtain  $\alpha_1 = \dots = \alpha_p = \sigma_1 = \dots = \sigma_t = 0$ . We have shown that the vectors  $\{\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{r}_1, \dots, \mathbf{r}_q, \mathbf{s}_1, \dots, \mathbf{s}_t\}$  constitute a basis for  $\mathcal{R} + \mathcal{S}$ . But then

$$\dim(\mathcal{R} + \mathcal{S}) = p + q + t = (p + q) + (p + t) - p = \dim(\mathcal{R}) + \dim(\mathcal{S}) - \dim(\mathcal{R} \cap \mathcal{S})$$

and (7) follows. (7) implies (8) since  $\dim\{\mathbf{0}\} = 0$ .  $\square$

It is convenient to introduce a matrix transforming a basis in a subspace into a basis for the space itself.

**Lemma 0.14 (Change of basis matrix)**

Suppose  $\mathcal{S}$  is a subspace of a finite dimensional vector space  $\mathcal{V}$  and let  $\{\mathbf{s}_1, \dots, \mathbf{s}_n\}$  be a basis for  $\mathcal{S}$  and  $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$  a basis for  $\mathcal{V}$ . Then each  $\mathbf{s}_j$  can be expressed as a linear combination of  $\mathbf{v}_1, \dots, \mathbf{v}_m$ , say

$$\mathbf{s}_j = \sum_{i=1}^m a_{ij} \mathbf{v}_i \text{ for } j = 1, \dots, n. \quad (9)$$

If  $\mathbf{x} \in \mathcal{S}$  then  $\mathbf{x} = \sum_{j=1}^n c_j \mathbf{s}_j = \sum_{i=1}^m b_i \mathbf{v}_i$  for some coefficients  $\mathbf{b} := [b_1, \dots, b_m]^T$ ,  $\mathbf{c} := [c_1, \dots, c_n]^T$ . Moreover  $\mathbf{b} = \mathbf{A}\mathbf{c}$ , where  $\mathbf{A} = [a_{ij}] \in \mathbb{C}^{m \times n}$  is given by (9). The matrix  $\mathbf{A}$  has linearly independent columns.

**Proof.** (9) holds for some  $a_{ij}$  since  $\mathbf{s}_j \in \mathcal{V}$  and  $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$  spans  $\mathcal{V}$ . Since  $\{\mathbf{s}_1, \dots, \mathbf{s}_n\}$  is a basis for  $\mathcal{S}$  and  $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$  a basis for  $\mathcal{V}$ , every  $\mathbf{x} \in \mathcal{S}$  can be written  $\mathbf{x} = \sum_{j=1}^n c_j \mathbf{s}_j = \sum_{i=1}^m b_i \mathbf{v}_i$  for some scalars  $(c_j)$  and  $(b_i)$ . But then

$$\sum_{i=1}^m b_i \mathbf{v}_i = \mathbf{x} = \sum_{j=1}^n c_j \mathbf{s}_j \stackrel{(9)}{=} \sum_{j=1}^n c_j \left( \sum_{i=1}^m a_{ij} \mathbf{v}_i \right) = \sum_{i=1}^m \left( \sum_{j=1}^n a_{ij} c_j \right) \mathbf{v}_i.$$

Since  $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$  is linearly independent it follows that  $b_i = \sum_{j=1}^n a_{ij} c_j$  for  $i = 1, \dots, m$  or  $\mathbf{b} = \mathbf{A}\mathbf{c}$ . Finally, to show that  $\mathbf{A}$  has linearly independent columns suppose  $\mathbf{b} := \mathbf{A}\mathbf{c} = \mathbf{0}$  for some  $\mathbf{c} = [c_1, \dots, c_n]^T$ . Define  $\mathbf{x} := \sum_{j=1}^n c_j \mathbf{s}_j$ . Then  $\mathbf{x} = \sum_{i=1}^m b_i \mathbf{v}_i$  and since  $\mathbf{b} = \mathbf{0}$  we have  $\mathbf{x} = \mathbf{0}$ . But since  $\{\mathbf{s}_1, \dots, \mathbf{s}_n\}$  is linearly independent it follows that  $\mathbf{c} = \mathbf{0}$ .  $\square$

The matrix  $\mathbf{A}$  in Lemma 0.14 is called a **change of basis matrix**.

### 0.2.3 The vector spaces $\mathbb{R}^n$ and $\mathbb{C}^n$

When  $\mathcal{V} = \mathbb{R}^m$  we can think of  $n$  vectors in  $\mathbb{R}^m$ , say  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , as a set  $\mathcal{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  or as the columns of a matrix  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{m \times n}$ . A linear combination can then be written as a matrix times vector  $\mathbf{X}\mathbf{c}$ , where  $\mathbf{c} = [c_1, \dots, c_n]^T$  is the vector of scalars. Thus

$$\text{span}(\mathcal{X}) = \text{span}(\mathbf{X}) = \{\mathbf{X}\mathbf{c} : \mathbf{c} \in \mathbb{R}^n\}.$$

Of course the same holds for  $\mathbb{C}^m$ .

In  $\mathbb{R}^m$  and  $\mathbb{C}^m$  each of the following statements is equivalent to linear independence of  $\mathcal{X}$ .

- (i)  $\mathbf{X}\mathbf{c} = \mathbf{0} \Rightarrow \mathbf{c} = \mathbf{0}$ ,
- (ii)  $\mathbf{X}$  has linearly independent columns,

#### Definition 0.15 (Column space and null space)

Associated with a matrix  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{m \times n}$  are the following subspaces

1. The subspace  $\text{span}(\mathbf{X})$  is called the **column space** of  $\mathbf{X}$ . It is the smallest subspace containing  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ .
2.  $\text{span}(\mathbf{X}^T)$  is called the **row space** of  $\mathbf{X}$ . It is generated by the rows of  $\mathbf{X}$  written as column vectors.
3. The subspace  $\ker(\mathbf{X}) := \{\mathbf{y} \in \mathbb{R}^n : \mathbf{X}\mathbf{y} = \mathbf{0}\}$  is called the **null space** or **kernel space** of  $\mathbf{X}$ .

Note that the subspace  $\ker(\mathbf{X})$  is nontrivial if and only if  $\mathcal{X}$  is linearly dependent.

## 0.3 Vector Norms

To measure the size of a vector we use norms.

#### Definition 0.16 (Vector norm)

A **(vector) norm** in a real (resp. complex) vector space  $\mathcal{V}$  is a function  $\|\cdot\| : \mathcal{V} \rightarrow \mathbb{R}$  that satisfies for all  $\mathbf{x}, \mathbf{y}$  in  $\mathcal{V}$  and all  $a$  in  $\mathbb{R}$  (resp.  $\mathbb{C}$ )

1.  $\|\mathbf{x}\| \geq 0$  with equality if and only if  $\mathbf{x} = \mathbf{0}$ . (positivity)
2.  $\|a\mathbf{x}\| = |a| \|\mathbf{x}\|$ . (homogeneity)
3.  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ . (subadditivity)

The triple  $(\mathcal{V}, \mathbb{R}, \|\cdot\|)$  (resp.  $(\mathcal{V}, \mathbb{C}, \|\cdot\|)$ ) is called a **normed vector space** and the inequality 3. is called the **triangle inequality**.

In this book we will use the following family of vector norms on  $\mathcal{V} = \mathbb{C}^n$  and  $\mathcal{V} = \mathbb{R}^n$ .



Otto Ludwig Hölder, 1859-1937 (left), Hermann Minkowski, 1864-1909 (right).

**Definition 0.17 (Vector p-norms)**

We define for  $p \geq 1$  and  $\mathbf{x} \in \mathbb{R}^n$  or  $\mathbf{x} \in \mathbb{C}^n$  the **p-norms** by

$$\|\mathbf{x}\|_p := \left( \sum_{j=1}^n |x_j|^p \right)^{1/p}, \quad (10)$$

$$\|\mathbf{x}\|_\infty := \max_{1 \leq j \leq n} |x_j|. \quad (11)$$

The most important cases are  $p = 1, 2, \infty$ :

1.  $\|\mathbf{x}\|_1 = \sum_{j=1}^n |x_j|$ , (the one-norm or  $l_1$ -norm)
2.  $\|\mathbf{x}\|_2 = \sqrt{\sum_{j=1}^n |x_j|^2}$ , (the two-norm,  $l_2$ -norm, or Euclidian norm)
3.  $\|\mathbf{x}\|_\infty = \max_{1 \leq j \leq n} |x_j|$ , (the infinity-norm,  $l_\infty$ -norm, or max norm)

Some remarks are in order.

1. That the Euclidian norm is a vector norm follows from Theorem 0.23. In Section 7.3, we show that the  $p$ -norms are vector norms for  $1 \leq p \leq \infty$ .
2. The triangle inequality  $\|\mathbf{x} + \mathbf{y}\|_p \leq \|\mathbf{x}\|_p + \|\mathbf{y}\|_p$  is called **Minkowski's inequality**.



3. To prove it one first establishes **Hölder's inequality**

$$\sum_{j=1}^n |x_j y_j| \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q, \quad \frac{1}{p} + \frac{1}{q} = 1, \quad \mathbf{x}, \mathbf{y} \in \mathbb{C}^n. \quad (12)$$

The relation  $\frac{1}{p} + \frac{1}{q} = 1$  means that if  $p = 1$  then  $q = \infty$  and if  $p = 2$  then  $q = 2$  and Hölder's inequality is the same as the Cauchy-Schwarz inequality (cf. Theorem 0.22) for the Euclidian norm.

4. The infinity norm is related to the other  $p$ -norms by

$$\lim_{p \rightarrow \infty} \|\mathbf{x}\|_p = \|\mathbf{x}\|_\infty \text{ for all } \mathbf{x} \in \mathbb{C}^n. \quad (13)$$

5. The equation (13) clearly holds for  $\mathbf{x} = \mathbf{0}$ . For  $\mathbf{x} \neq \mathbf{0}$  we write

$$\|\mathbf{x}\|_p := \|\mathbf{x}\|_\infty \left( \sum_{j=1}^n \left( \frac{|x_j|}{\|\mathbf{x}\|_\infty} \right)^p \right)^{1/p}.$$

Now each term in the sum is not greater than one and at least one term is equal to one, and we obtain

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_p \leq n^{1/p} \|\mathbf{x}\|_\infty, \quad p \geq 1. \quad (14)$$

Since  $\lim_{p \rightarrow \infty} n^{1/p} = 1$  for any  $n \in \mathbb{N}$  we see that (13) follows.

We return now to the general case.

**Definition 0.18 (Equivalent norms)**

We say that two norms  $\|\cdot\|$  and  $\|\cdot\|'$  on  $\mathcal{V}$  are **equivalent** if there are positive constants  $m$  and  $M$  such that for all vectors  $\mathbf{x} \in \mathcal{V}$  we have

$$m\|\mathbf{x}\|' \leq \|\mathbf{x}\| \leq M\|\mathbf{x}\|. \quad (15)$$

By (14) the  $p$ - and  $\infty$ -norms are equivalent for any  $p \geq 1$ . This result is generalized in the following theorem.

**Theorem 0.19 (Basic properties of vector norms)**

The following holds for a normed vector space  $(\mathcal{V}, \mathbb{C}, \|\cdot\|)$ .

1.  $\|\mathbf{x} - \mathbf{y}\| \geq | \|\mathbf{x}\| - \|\mathbf{y}\| |$ , for all  $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$  (inverse triangle inequality).
2. The vector norm is a continuous function  $\mathcal{V} \rightarrow \mathbb{R}$ .
3. All vector norms on  $\mathcal{V}$  are equivalent provided  $\mathcal{V}$  is finite dimensional.

*Proof.*

1. Since  $\|\mathbf{x}\| = \|\mathbf{x} - \mathbf{y} + \mathbf{y}\| \leq \|\mathbf{x} - \mathbf{y}\| + \|\mathbf{y}\|$  we obtain  $\|\mathbf{x} - \mathbf{y}\| \geq \|\mathbf{x}\| - \|\mathbf{y}\|$ . By symmetry  $\|\mathbf{x} - \mathbf{y}\| = \|\mathbf{y} - \mathbf{x}\| \geq \|\mathbf{y}\| - \|\mathbf{x}\|$  and we obtain the inverse triangle inequality.
2. This follows from the inverse triangle inequality.
3. The following proof can be skipped by those who do not have the necessary background in advanced calculus. Define the  $\|\cdot\|'$  unit sphere

$$\mathcal{S} := \{\mathbf{y} \in \mathcal{V} : \|\mathbf{y}\|' = 1\}.$$

The set  $\mathcal{S}$  is a closed and bounded set and the function  $f : \mathcal{S} \rightarrow \mathbb{R}$  given by  $f(\mathbf{y}) = \|\mathbf{y}\|$  is continuous by what we just showed. Therefore  $f$  attains its minimum and maximum value on  $\mathcal{S}$ . Thus, there are positive constants  $m$  and  $M$  such that

$$m \leq \|\mathbf{y}\| \leq M, \quad \mathbf{y} \in \mathcal{S}. \quad (16)$$

For any  $\mathbf{x} \in \mathcal{V}$  one has  $\mathbf{y} := \mathbf{x}/\|\mathbf{x}\|' \in \mathcal{S}$ , and (15) follows if we apply (16) to these  $\mathbf{y}$ .

□

## 0.4 Inner Products

An **inner product** or **scalar product** in a vector space is a function mapping pairs of vectors into a scalar.

### 0.4.1 Real and complex inner products

We consider first the real case.

#### Definition 0.20 (Real inner product)

An **inner product** in a real vector space  $\mathcal{V}$  is a function  $\langle \cdot, \cdot \rangle : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$  satisfying for all  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{V}$  and all  $a, b \in \mathbb{R}$  the following conditions:

1.  $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$  with equality if and only if  $\mathbf{x} = \mathbf{0}$ . *(positivity)*
2.  $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$  *(symmetry)*
3.  $\langle a\mathbf{x} + b\mathbf{y}, \mathbf{z} \rangle = a\langle \mathbf{x}, \mathbf{z} \rangle + b\langle \mathbf{y}, \mathbf{z} \rangle$ . *(linearity)*

The pair  $(\mathcal{V}, \langle \cdot, \cdot \rangle)$  is called a **real inner product space**. The function

$$\|\cdot\| : \mathcal{V} \longrightarrow \mathbb{R}, \quad \mathbf{x} \longmapsto \|\mathbf{x}\| := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} \quad (17)$$

is called the **inner product norm**.

The **standard inner product in**  $\mathcal{V} = \mathbb{R}^n$  is given by  $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x}$ . It is clearly an inner product in  $\mathbb{R}^n$ . The corresponding inner product norm is the Euclidian norm  $\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{x}} = \|\mathbf{x}\|_2$ .

Consider next inner products in a complex vector space. Property 2. in the definition of a real inner product is altered from symmetry to skew symmetry.

**Definition 0.21 (Complex inner product)**

An **inner product** in a complex vector space  $\mathcal{V}$  is a function  $\mathcal{V} \times \mathcal{V} \rightarrow \mathbb{C}$  satisfying for all  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{V}$  and all  $a, b \in \mathbb{C}$  the following conditions:

1.  $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$  with equality if and only if  $\mathbf{x} = \mathbf{0}$ . (positivity)
2.  $\langle \mathbf{x}, \mathbf{y} \rangle = \overline{\langle \mathbf{y}, \mathbf{x} \rangle}$  (skew symmetry)
3.  $\langle a\mathbf{x} + b\mathbf{y}, \mathbf{z} \rangle = a\langle \mathbf{x}, \mathbf{z} \rangle + b\langle \mathbf{y}, \mathbf{z} \rangle$ . (linearity)

The pair  $(\mathcal{V}, \langle \cdot, \cdot \rangle)$  is called a **complex inner product space** The function

$$\|\cdot\| : \mathcal{V} \rightarrow \mathbb{R}, \quad \mathbf{x} \mapsto \|\mathbf{x}\| := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} \quad (18)$$

is called the **inner product norm**.

Note the complex conjugate in 2. We find

$$\langle \mathbf{x}, a\mathbf{y} + b\mathbf{z} \rangle = \bar{a}\langle \mathbf{x}, \mathbf{y} \rangle + \bar{b}\langle \mathbf{x}, \mathbf{z} \rangle, \quad \langle a\mathbf{x}, a\mathbf{y} \rangle = |a|^2 \langle \mathbf{x}, \mathbf{y} \rangle. \quad (19)$$

The **standard inner product in**  $\mathbb{C}^n$  is given by

$$\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{y}^* \mathbf{x} = \mathbf{x}^T \bar{\mathbf{y}} = \sum_{j=1}^n x_j \bar{y}_j.$$

It is clearly an inner product in  $\mathbb{C}^n$ . The corresponding inner product norm is the Euclidian norm  $\|\mathbf{x}\| = \|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^* \mathbf{x}}$ .



Viktor Yakovlevich Bunyakovsky, 1804-1889 (left), Augustin-Louis Cauchy, 1789-1857 (center), Karl Hermann Amandus Schwarz, 1843-1921 (right). The name Bunyakovsky is also associated with the Cauchy-Schwarz inequality.

The following inequality holds for any inner product.

**Theorem 0.22 (Cauchy-Schwarz inequality)**

For any  $\mathbf{x}, \mathbf{y}$  in a real or complex inner product space

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|, \quad (20)$$

with equality if and only if  $\mathbf{x}$  and  $\mathbf{y}$  are linearly dependent.

**Proof.** If  $\mathbf{y} = \mathbf{0}$  then  $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{0} \rangle = 0$  and  $\|\mathbf{y}\| = 0$ . Thus the inequality holds with equality, and  $\mathbf{x}$  and  $\mathbf{y}$  are linearly dependent. So assume  $\mathbf{y} \neq \mathbf{0}$ . Define

$$\mathbf{z} := \mathbf{x} - a\mathbf{y}, \quad a := \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{y}, \mathbf{y} \rangle}.$$

Then  $\langle \mathbf{z}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle - a\langle \mathbf{y}, \mathbf{y} \rangle = 0$  so that by 2. and (19)

$$\langle a\mathbf{y}, \mathbf{z} \rangle + \langle \mathbf{z}, a\mathbf{y} \rangle = a\overline{\langle \mathbf{z}, \mathbf{y} \rangle} + \bar{a}\langle \mathbf{z}, \mathbf{y} \rangle = 0. \quad (21)$$

But then

$$\begin{aligned} \|\mathbf{x}\|^2 &= \langle \mathbf{x}, \mathbf{x} \rangle = \langle \mathbf{z} + a\mathbf{y}, \mathbf{z} + a\mathbf{y} \rangle \\ &\stackrel{(21)}{=} \langle \mathbf{z}, \mathbf{z} \rangle + \langle a\mathbf{y}, a\mathbf{y} \rangle \stackrel{(19)}{=} \|\mathbf{z}\|^2 + |a|^2 \|\mathbf{y}\|^2 \\ &\geq |a|^2 \|\mathbf{y}\|^2 = \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|^2}{\|\mathbf{y}\|^2}. \end{aligned}$$

Multiplying by  $\|\mathbf{y}\|^2$  gives (20). We have equality if and only if  $\mathbf{z} = \mathbf{0}$ , which means that  $\mathbf{x}$  and  $\mathbf{y}$  are linearly dependent.  $\square$

**Theorem 0.23 (Inner product norm)**

The inner product norm is a vector norm.

**Proof.** For all  $\mathbf{x}, \mathbf{y}$  in an inner product space and all  $a$  in  $\mathbb{C}$  we need to show

1.  $\|\mathbf{x}\| \geq 0$  with equality if and only if  $\mathbf{x} = \mathbf{0}$ . (positivity)
2.  $\|a\mathbf{x}\| = |a| \|\mathbf{x}\|$ . (homogeneity)
3.  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ . (subadditivity)

The first statement is an immediate consequence of positivity, while the second one follows from (19). Expanding  $\|\mathbf{x} + a\mathbf{y}\|^2 = \langle \mathbf{x} + a\mathbf{y}, \mathbf{x} + a\mathbf{y} \rangle$  using (19) we obtain

$$\|\mathbf{x} + a\mathbf{y}\|^2 = \|\mathbf{x}\|^2 + a\langle \mathbf{y}, \mathbf{x} \rangle + \bar{a}\langle \mathbf{x}, \mathbf{y} \rangle + |a|^2 \|\mathbf{y}\|^2, \quad a \in \mathbb{C}, \quad \mathbf{x}, \mathbf{y} \in \mathcal{V}. \quad (22)$$

Now (22) with  $a = 1$  and the Cauchy-Schwarz inequality implies

$$\|\mathbf{x} + \mathbf{y}\|^2 \leq \|\mathbf{x}\|^2 + 2\|\mathbf{x}\|\|\mathbf{y}\| + \|\mathbf{y}\|^2 = (\|\mathbf{x}\| + \|\mathbf{y}\|)^2.$$

Taking square roots completes the proof.  $\square$

In the real case the Cauchy-Schwarz inequality implies that  $-1 \leq \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\|\|\mathbf{y}\|} \leq 1$  for nonzero  $\mathbf{x}$  and  $\mathbf{y}$ , so there is a unique angle  $\theta$  in  $[0, \pi]$  such that

$$\cos \theta = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\|\|\mathbf{y}\|}. \quad (23)$$

This defines the **angle** between vectors in a real inner product space.

**Exercise 0.24 (The  $\mathbf{A}^T \mathbf{A}$  inner product)**

Suppose  $\mathbf{A} \in \mathbb{R}^{m \times n}$  has linearly independent columns. Show that  $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{y}$  defines an inner product on  $\mathbb{R}^n$ .

**Exercise 0.25 (Angle between vectors in complex case)**

Show that in the complex case there is a unique angle  $\theta$  in  $[0, \pi/2]$  such that

$$\cos \theta = \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|}{\|\mathbf{x}\|\|\mathbf{y}\|}. \quad (24)$$

## 0.4.2 Orthogonality

**Definition 0.26 (Orthogonality)**

Two vectors  $\mathbf{x}, \mathbf{y}$  in a real or complex inner product space are **orthogonal** or **perpendicular**, denoted as  $\mathbf{x} \perp \mathbf{y}$ , if  $\langle \mathbf{x}, \mathbf{y} \rangle = 0$ . The vectors are **orthonormal** if in addition  $\|\mathbf{x}\| = \|\mathbf{y}\| = 1$ .

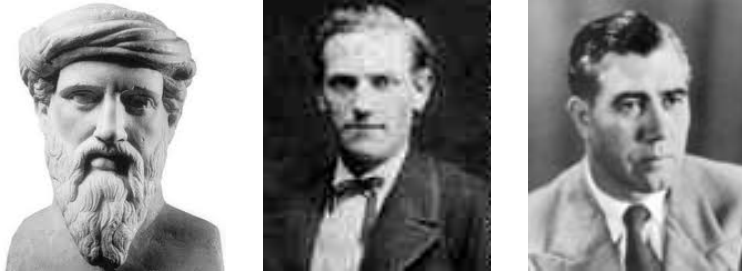
From the definitions (23), (24) of angle  $\theta$  between two vectors in  $\mathbb{R}^n$  or  $\mathbb{C}^n$  it follows that  $\mathbf{x} \perp \mathbf{y}$  if and only if  $\theta = \pi/2$ .

**Theorem 0.27 (Pythagoras)**

For a real or complex inner product space

$$\|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2, \quad \text{if } \mathbf{x} \perp \mathbf{y}. \quad (25)$$

*Proof.* We set  $a = 1$  in (22) and use the orthogonality.  $\square$



Pythagoras of Samos, BC 570-BC 495 (left), Jørgen Pedersen Gram, 1850-1916 (center), Erhard Schmidt, 1876-1959 (right).

**Definition 0.28 (Orthogonal- and orthonormal bases)**

A set of nonzero vectors  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  in a subspace  $\mathcal{S}$  of a real or complex inner product space is an **orthogonal basis** for  $\mathcal{S}$  if it is a basis for  $\mathcal{S}$  and  $\langle \mathbf{v}_i, \mathbf{v}_j \rangle = 0$  for  $i \neq j$ . It is an **orthonormal basis** for  $\mathcal{S}$  if it is a basis for  $\mathcal{S}$  and  $\langle \mathbf{v}_i, \mathbf{v}_j \rangle = \delta_{ij}$  for all  $i, j$ .

A basis for a subspace of an inner product space can be turned into an orthogonal- or orthonormal basis for the subspace by the following construction.

**Theorem 0.29 (Gram-Schmidt)**

Let  $\{\mathbf{s}_1, \dots, \mathbf{s}_k\}$  be a basis for a real or complex inner product space  $(\mathcal{S}, \langle \cdot, \cdot \rangle)$ . Define

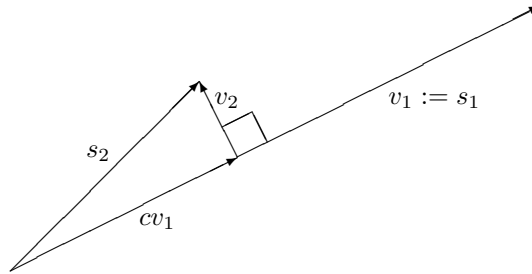
$$\mathbf{v}_1 := \mathbf{s}_1, \quad \mathbf{v}_j := \mathbf{s}_j - \sum_{i=1}^{j-1} \frac{\langle \mathbf{s}_j, \mathbf{v}_i \rangle}{\langle \mathbf{v}_i, \mathbf{v}_i \rangle} \mathbf{v}_i, \quad j = 2, \dots, k. \quad (26)$$

Then  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  is an orthogonal basis for  $\mathcal{S}$  and the normalized vectors

$$\{\mathbf{u}_1, \dots, \mathbf{u}_k\} := \left\{ \frac{\mathbf{v}_1}{\|\mathbf{v}_1\|}, \dots, \frac{\mathbf{v}_k}{\|\mathbf{v}_k\|} \right\}$$

form an orthonormal basis for  $\mathcal{S}$ .

**Proof.** To show that  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  is an orthogonal basis for  $\mathcal{S}$  we use induction on  $k$ . Define subspaces  $\mathcal{S}_j := \text{span}\{\mathbf{s}_1, \dots, \mathbf{s}_j\}$  for  $j = 1, \dots, k$ . Clearly  $\mathbf{v}_1 = \mathbf{s}_1$  is an orthogonal basis for  $\mathcal{S}_1$ . Suppose for some  $j \geq 2$  that  $\mathbf{v}_1, \dots, \mathbf{v}_{j-1}$  is an orthogonal basis for  $\mathcal{S}_{j-1}$  and let  $\mathbf{v}_j$  be given by (26) as a linear combination of  $\mathbf{s}_j$  and  $\mathbf{v}_1, \dots, \mathbf{v}_{j-1}$ . Now each of these  $\mathbf{v}_i$  is a linear combination of  $\mathbf{s}_1, \dots, \mathbf{s}_i$ , and we obtain  $\mathbf{v}_j = \sum_{i=1}^j a_i \mathbf{s}_i$  for some  $a_0, \dots, a_j$  with  $a_j = 1$ . Since  $\mathbf{s}_1, \dots, \mathbf{s}_j$



**Figure 1.** The construction of  $\mathbf{v}_1$  and  $\mathbf{v}_2$  in Gram-Schmidt. The constant  $c$  is given by  $c := \langle \mathbf{s}_2, \mathbf{v}_1 \rangle / \langle \mathbf{v}_1, \mathbf{v}_1 \rangle$ .

are linearly independent and  $a_j \neq 0$  we deduce that  $\mathbf{v}_j \neq 0$ . By the induction hypothesis

$$\langle \mathbf{v}_j, \mathbf{v}_l \rangle = \langle \mathbf{s}_j, \mathbf{v}_l \rangle - \sum_{i=1}^{j-1} \frac{\langle \mathbf{s}_j, \mathbf{v}_i \rangle}{\langle \mathbf{v}_i, \mathbf{v}_i \rangle} \langle \mathbf{v}_i, \mathbf{v}_l \rangle = \langle \mathbf{s}_j, \mathbf{v}_l \rangle - \frac{\langle \mathbf{s}_j, \mathbf{v}_l \rangle}{\langle \mathbf{v}_l, \mathbf{v}_l \rangle} \langle \mathbf{v}_l, \mathbf{v}_l \rangle = 0$$

for  $l = 1, \dots, j-1$ . Thus  $\mathbf{v}_1, \dots, \mathbf{v}_j$  is an orthogonal basis for  $\mathcal{S}_j$ .

If  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  is an orthogonal basis for  $\mathcal{S}$  then clearly  $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  is an orthonormal basis for  $\mathcal{S}$ .  $\square$

Sometimes we want to extend an orthogonal basis for a subspace to an orthogonal basis for a larger space.

### Theorem 0.30 (Orthogonal Extension of basis)

Suppose  $\mathcal{S} \subset \mathcal{T}$  are finite dimensional subspaces of a vector space  $\mathcal{V}$ . An orthogonal basis for  $\mathcal{S}$  can always be extended to an orthogonal basis for  $\mathcal{T}$ .

**Proof.** Suppose  $\dim \mathcal{S} := k < \dim \mathcal{T} = n$ . Using Theorem 0.10 we first extend an orthogonal basis  $\mathbf{s}_1, \dots, \mathbf{s}_k$  for  $\mathcal{S}$  to a basis  $\mathbf{s}_1, \dots, \mathbf{s}_k, \mathbf{s}_{k+1}, \dots, \mathbf{s}_n$  for  $\mathcal{T}$ , and then apply the Gram-Schmidt process to this basis obtaining an orthogonal basis  $\mathbf{v}_1, \dots, \mathbf{v}_n$  for  $\mathcal{T}$ . This is an extension of the basis for  $\mathcal{S}$  since  $\mathbf{v}_i = \mathbf{s}_i$  for  $i = 1, \dots, k$ . We show this by induction. Clearly  $\mathbf{v}_1 = \mathbf{s}_1$ . Suppose for some  $2 \leq r < k$  that  $\mathbf{v}_j = \mathbf{s}_j$  for  $j = 1, \dots, r-1$ . Consider (26) for  $j = r$ . Since  $\langle \mathbf{s}_r, \mathbf{v}_i \rangle = \langle \mathbf{s}_r, \mathbf{s}_i \rangle = 0$  for  $i < r$  we obtain  $\mathbf{v}_r = \mathbf{s}_r$ .  $\square$

Letting  $\mathcal{S} = \text{span}(\mathbf{s}_1, \dots, \mathbf{s}_k)$  and  $\mathcal{T}$  be  $\mathbb{R}^n$  or  $\mathbb{C}^n$  we obtain

**Corollary 0.31 (Extending orthogonal vectors to a basis)**

For  $1 \leq k < n$  a set  $\{\mathbf{s}_1, \dots, \mathbf{s}_k\}$  of nonzero orthogonal vectors in  $\mathbb{R}^n$  or  $\mathbb{C}^n$  can be extended to an orthogonal basis for the whole space.

**0.5 Linear Systems**

Consider a linear system

$$\begin{array}{cccccc} a_{11}x_1 & + & a_{12}x_2 & + & \cdots & + & a_{1n}x_n & = & b_1 \\ a_{21}x_1 & + & a_{22}x_2 & + & \cdots & + & a_{2n}x_n & = & b_2 \\ \vdots & & \vdots & & & & \vdots & & \vdots \\ a_{m1}x_1 & + & a_{m2}x_2 & + & \cdots & + & a_{mn}x_n & = & b_m \end{array}$$

of  $m$  equations in  $n$  unknowns. Here for all  $i, j$ , the coefficients  $a_{ij}$ , the unknowns  $x_j$ , and the components of the right hand sides  $b_i$ , are real or complex numbers. The system can be written as a vector equation

$$x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_n\mathbf{a}_n = \mathbf{b},$$

where  $\mathbf{a}_j = [a_{1j}, \dots, a_{mj}]^T \in \mathbb{C}^m$  for  $j = 1, \dots, n$  and  $\mathbf{b} = [b_1, \dots, b_m]^T \in \mathbb{C}^m$ . It can also be written as a matrix equation

$$\mathbf{A}\mathbf{x} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix} = \mathbf{b}.$$

The system is **homogeneous** if  $\mathbf{b} = \mathbf{0}$  and it is said to be **underdetermined**, **square**, or **overdetermined** if  $m < n$ ,  $m = n$ , or  $m > n$ , respectively.

**0.5.1 Basic properties**

A linear system has a unique solution, infinitely many solutions, or no solution. To discuss this we first consider the real case, and a homogeneous underdetermined system.

**Lemma 0.32 (Underdetermined system)**

Suppose  $\mathbf{A} \in \mathbb{R}^{m \times n}$  with  $m < n$ . Then there is a nonzero  $\mathbf{x} \in \mathbb{R}^n$  such that  $\mathbf{A}\mathbf{x} = \mathbf{0}$ .

**Proof.** Suppose  $\mathbf{A} \in \mathbb{R}^{m \times n}$  with  $m < n$ . The  $n$  columns of  $\mathbf{A}$  span a subspace of  $\mathbb{R}^m$ . Since  $\mathbb{R}^m$  has dimension  $m$  the dimension of this subspace is at most  $m$ . By



Lemma 0.5 the columns of  $\mathbf{A}$  must be linearly dependent. It follows that there is a nonzero  $\mathbf{x} \in \mathbb{R}^n$  such that  $\mathbf{A}\mathbf{x} = \mathbf{0}$ .  $\square$

A square matrix is either **nonsingular** or **singular**.

**Definition 0.33 (Real nonsingular or singular matrix)**

A square matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is said to be **nonsingular** if the only real solution of the homogeneous system  $\mathbf{A}\mathbf{x} = \mathbf{0}$  is  $\mathbf{x} = \mathbf{0}$ . The matrix is **singular** if there is a nonzero  $\mathbf{x} \in \mathbb{R}^n$  such that  $\mathbf{A}\mathbf{x} = \mathbf{0}$ .

**Theorem 0.34 (Linear systems; existence and uniqueness)**

Suppose  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . The linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  has a unique solution  $\mathbf{x} \in \mathbb{R}^n$  for any  $\mathbf{b} \in \mathbb{R}^n$  if and only if the matrix  $\mathbf{A}$  is nonsingular.

*Proof.* Suppose  $\mathbf{A}$  is nonsingular. We define  $\mathbf{B} = [\mathbf{A} \ \mathbf{b}] \in \mathbb{R}^{n \times (n+1)}$  by adding a column to  $\mathbf{A}$ . By Lemma 0.32 there is a nonzero  $\mathbf{z} \in \mathbb{R}^{n+1}$  such that  $\mathbf{B}\mathbf{z} = \mathbf{0}$ . If we write  $\mathbf{z} = \begin{bmatrix} \tilde{\mathbf{z}} \\ z_{n+1} \end{bmatrix}$  where  $\tilde{\mathbf{z}} = [z_1, \dots, z_n]^T \in \mathbb{R}^n$  and  $z_{n+1} \in \mathbb{R}$ , then

$$\mathbf{B}\mathbf{z} = [\mathbf{A} \ \mathbf{b}] \begin{bmatrix} \tilde{\mathbf{z}} \\ z_{n+1} \end{bmatrix} = \mathbf{A}\tilde{\mathbf{z}} + z_{n+1}\mathbf{b} = \mathbf{0}.$$

We cannot have  $z_{n+1} = 0$  for then  $\mathbf{A}\tilde{\mathbf{z}} = \mathbf{0}$  for a nonzero  $\tilde{\mathbf{z}}$ , contradicting the nonsingularity of  $\mathbf{A}$ . Define  $\mathbf{x} := -\tilde{\mathbf{z}}/z_{n+1}$ . Then

$$\mathbf{A}\mathbf{x} = -\mathbf{A} \begin{pmatrix} \tilde{\mathbf{z}} \\ z_{n+1} \end{pmatrix} = -\frac{1}{z_{n+1}} \mathbf{A}\tilde{\mathbf{z}} = -\frac{1}{z_{n+1}} (-z_{n+1}\mathbf{b}) = \mathbf{b},$$

so  $\mathbf{x}$  is a solution.

Suppose  $\mathbf{A}\mathbf{x} = \mathbf{b}$  and  $\mathbf{A}\mathbf{y} = \mathbf{b}$  for  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ . Then  $\mathbf{A}(\mathbf{x} - \mathbf{y}) = \mathbf{0}$  and since  $\mathbf{A}$  is nonsingular we conclude that  $\mathbf{x} - \mathbf{y} = \mathbf{0}$  or  $\mathbf{x} = \mathbf{y}$ . Thus the solution is unique.

Conversely, if  $\mathbf{A}\mathbf{x} = \mathbf{b}$  has a unique solution for any  $\mathbf{b} \in \mathbb{R}^n$  then  $\mathbf{A}\mathbf{x} = \mathbf{0}$  has a unique solution which must be  $\mathbf{x} = \mathbf{0}$ . Thus  $\mathbf{A}$  is nonsingular.  $\square$

For the complex case we have

**Lemma 0.35 (Complex underdetermined system)**

Suppose  $\mathbf{A} \in \mathbb{C}^{m \times n}$  with  $m < n$ . Then there is a nonzero  $\mathbf{x} \in \mathbb{C}^n$  such that  $\mathbf{A}\mathbf{x} = \mathbf{0}$ .

**Definition 0.36 (Complex nonsingular matrix)**

A square matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is said to be **nonsingular** if the only complex solution of the homogeneous system  $\mathbf{A}\mathbf{x} = \mathbf{0}$  is  $\mathbf{x} = \mathbf{0}$ . The matrix is **singular** if it is not nonsingular.

**Theorem 0.37 (Complex linear system; existence and uniqueness)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . The linear system  $\mathbf{Ax} = \mathbf{b}$  has a unique solution  $\mathbf{x} \in \mathbb{C}^n$  for any  $\mathbf{b} \in \mathbb{C}^n$  if and only if the matrix  $\mathbf{A}$  is nonsingular.



James Joseph Sylvester, 1814-1897. The word matrix to denote a rectangular array of numbers, was first used by Sylvester in 1850.

**0.5.2 The inverse matrix**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is a square matrix. A matrix  $\mathbf{B} \in \mathbb{C}^{n \times n}$  is called a **right inverse** of  $\mathbf{A}$  if  $\mathbf{AB} = \mathbf{I}$ . A matrix  $\mathbf{C} \in \mathbb{C}^{n \times n}$  is said to be a **left inverse** of  $\mathbf{A}$  if  $\mathbf{CA} = \mathbf{I}$ . We say that  $\mathbf{A}$  is **invertible** if it has both a left- and a right inverse. If  $\mathbf{A}$  has a right inverse  $\mathbf{B}$  and a left inverse  $\mathbf{C}$  then

$$\mathbf{C} = \mathbf{CI} = \mathbf{C}(\mathbf{AB}) = (\mathbf{CA})\mathbf{B} = \mathbf{IB} = \mathbf{B}$$

and this common inverse is called the **inverse** of  $\mathbf{A}$  and denoted by  $\mathbf{A}^{-1}$ . Thus the inverse satisfies  $\mathbf{A}^{-1}\mathbf{A} = \mathbf{AA}^{-1} = \mathbf{I}$ .

We want to characterize the class of invertible matrices and start with a lemma.

**Theorem 0.38 (Product of nonsingular matrices)**

If  $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{C}^{n \times n}$  with  $\mathbf{AB} = \mathbf{C}$  then  $\mathbf{C}$  is nonsingular if and only if both  $\mathbf{A}$  and  $\mathbf{B}$  are nonsingular. In particular, if  $\mathbf{AB} = \mathbf{I}$  or  $\mathbf{BA} = \mathbf{I}$  then  $\mathbf{A}$  is nonsingular and  $\mathbf{A}^{-1} = \mathbf{B}$ .

*Proof.* Suppose both  $\mathbf{A}$  and  $\mathbf{B}$  are nonsingular and let  $\mathbf{Cx} = \mathbf{0}$ . Then  $\mathbf{ABx} = \mathbf{0}$  and since  $\mathbf{A}$  is nonsingular we see that  $\mathbf{Bx} = \mathbf{0}$ . Since  $\mathbf{B}$  is nonsingular we have  $\mathbf{x} = \mathbf{0}$ . We conclude that  $\mathbf{C}$  is nonsingular.

For the converse suppose first that  $\mathbf{B}$  is singular and let  $\mathbf{x} \in \mathbb{C}^n$  be a nonzero vector so that  $\mathbf{Bx} = \mathbf{0}$ . But then  $\mathbf{Cx} = (\mathbf{AB})\mathbf{x} = \mathbf{A}(\mathbf{Bx}) = \mathbf{A}\mathbf{0} = \mathbf{0}$  so  $\mathbf{C}$  is singular. Finally suppose  $\mathbf{B}$  is nonsingular, but  $\mathbf{A}$  is singular. Let  $\tilde{\mathbf{x}}$  be a nonzero

vector such that  $\mathbf{A}\tilde{\mathbf{x}} = \mathbf{0}$ . By Theorem 0.37 there is a vector  $\mathbf{x}$  such that  $\mathbf{B}\mathbf{x} = \tilde{\mathbf{x}}$  and  $\mathbf{x}$  is nonzero since  $\tilde{\mathbf{x}}$  is nonzero. But then  $\mathbf{C}\mathbf{x} = (\mathbf{A}\mathbf{B})\mathbf{x} = \mathbf{A}(\mathbf{B}\mathbf{x}) = \mathbf{A}\tilde{\mathbf{x}} = \mathbf{0}$  for a nonzero vector  $\mathbf{x}$  and  $\mathbf{C}$  is singular.  $\square$

**Theorem 0.39 (When is a square matrix invertible?)**

*A square matrix is invertible if and only if it is nonsingular.*

**Proof.** Suppose first  $\mathbf{A}$  is a nonsingular matrix. By Theorem 0.37 each of the linear systems  $\mathbf{A}\mathbf{b}_i = \mathbf{e}_i$  has a unique solution  $\mathbf{b}_i$  for  $i = 1, \dots, n$ . Let  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_n]$ . Then  $\mathbf{A}\mathbf{B} = [\mathbf{A}\mathbf{b}_1, \dots, \mathbf{A}\mathbf{b}_n] = [\mathbf{e}_1, \dots, \mathbf{e}_n] = \mathbf{I}$  so that  $\mathbf{A}$  has a right inverse  $\mathbf{B}$ . By Theorem 0.38  $\mathbf{B}$  is nonsingular since  $\mathbf{I}$  is nonsingular and  $\mathbf{A}\mathbf{B} = \mathbf{I}$ . Since  $\mathbf{B}$  is nonsingular we can use what we have shown for  $\mathbf{A}$  to conclude that  $\mathbf{B}$  has a right inverse  $\mathbf{C}$ , i.e.  $\mathbf{B}\mathbf{C} = \mathbf{I}$ . But then  $\mathbf{A}\mathbf{B} = \mathbf{B}\mathbf{C} = \mathbf{I}$  so  $\mathbf{B}$  has both a right inverse and a left inverse which must be equal so  $\mathbf{A} = \mathbf{C}$ . Since  $\mathbf{B}\mathbf{C} = \mathbf{I}$  we have  $\mathbf{B}\mathbf{A} = \mathbf{I}$ , so  $\mathbf{B}$  is also a left inverse of  $\mathbf{A}$  and  $\mathbf{A}$  is invertible.

Conversely, if  $\mathbf{A}$  is invertible then it has a right inverse  $\mathbf{B}$ . Since  $\mathbf{A}\mathbf{B} = \mathbf{I}$  and  $\mathbf{I}$  is nonsingular, we again use Theorem 0.38 to conclude that  $\mathbf{A}$  is nonsingular.  $\square$

To verify that some matrix  $\mathbf{B}$  is an inverse of another matrix  $\mathbf{A}$  it is enough to show that  $\mathbf{B}$  is either a left inverse or a right inverse of  $\mathbf{A}$ . This calculation also proves that  $\mathbf{A}$  is nonsingular. We use this observation to give simple proofs of the following results.

**Corollary 0.40 (Basic properties of the inverse matrix)**

*Suppose  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$  are nonsingular and  $c$  is a nonzero constant.*

1.  $\mathbf{A}^{-1}$  is nonsingular and  $(\mathbf{A}^{-1})^{-1} = \mathbf{A}$ .
2.  $\mathbf{C} = \mathbf{A}\mathbf{B}$  is nonsingular and  $\mathbf{C}^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$ .
3.  $\mathbf{A}^T$  is nonsingular and  $(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T =: \mathbf{A}^{-T}$ .
4.  $\mathbf{A}^*$  is nonsingular and  $(\mathbf{A}^*)^{-1} = (\mathbf{A}^{-1})^* =: \mathbf{A}^{-*}$ .
5.  $c\mathbf{A}$  is nonsingular and  $(c\mathbf{A})^{-1} = \frac{1}{c}\mathbf{A}^{-1}$ .

**Proof.**

1. Since  $\mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$  the matrix  $\mathbf{A}$  is a right inverse of  $\mathbf{A}^{-1}$ . Thus  $\mathbf{A}^{-1}$  is nonsingular and  $(\mathbf{A}^{-1})^{-1} = \mathbf{A}$ .
2. We note that  $(\mathbf{B}^{-1}\mathbf{A}^{-1})(\mathbf{A}\mathbf{B}) = \mathbf{B}^{-1}(\mathbf{A}^{-1}\mathbf{A})\mathbf{B} = \mathbf{B}^{-1}\mathbf{B} = \mathbf{I}$ . Thus  $\mathbf{A}\mathbf{B}$  is invertible with the indicated inverse since it has a left inverse.
3. Now  $\mathbf{I} = \mathbf{I}^T = (\mathbf{A}^{-1}\mathbf{A})^T = \mathbf{A}^T(\mathbf{A}^{-1})^T$  showing that  $(\mathbf{A}^{-1})^T$  is a right inverse of  $\mathbf{A}^T$ . The proof of part 4 is similar.

4. The matrix  $\frac{1}{c}\mathbf{A}^{-1}$  is a one sided inverse of  $c\mathbf{A}$ .

□

**Exercise 0.41 (The inverse of a general  $2 \times 2$  matrix)**

Show that

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \alpha \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}, \quad \alpha = \frac{1}{ad - bc},$$

for any  $a, b, c, d$  such that  $ad - bc \neq 0$ .

**Exercise 0.42 (The inverse of a special  $2 \times 2$  matrix)**

Find the inverse of

$$A = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}.$$

**Exercise 0.43 (Sherman-Morrison formula)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$ , and  $\mathbf{B}, \mathbf{C} \in \mathbb{R}^{n \times m}$  for some  $n, m \in \mathbb{N}$ . If  $(\mathbf{I} + \mathbf{C}^T \mathbf{A}^{-1} \mathbf{B})^{-1}$  exists then

$$(\mathbf{A} + \mathbf{B}\mathbf{C}^T)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{I} + \mathbf{C}^T \mathbf{A}^{-1} \mathbf{B})^{-1}\mathbf{C}^T \mathbf{A}^{-1}.$$

## 0.6 Determinants

Determinants, denoted by  $\det(\cdot)$  or  $|\cdot|$ , are useful for studying eigenvalues. Recall that if  $\mathbf{A}, \mathbf{B}$  are square matrices of order  $n$  with real or complex elements, then (see Appendix A for proofs)

1.  $\det(\mathbf{A}\mathbf{B}) = \det(\mathbf{A})\det(\mathbf{B})$ .
2. If  $\mathbf{A}$  is triangular then  $\det(\mathbf{A}) = a_{11}a_{22} \cdots a_{nn}$ . In particular,  $\det(\mathbf{I}) = 1$ .
3.  $\det(\mathbf{A}^T) = \det(\mathbf{A})$ , and  $\det(\mathbf{A}^*) = \overline{\det(\mathbf{A})}$ , (complex conjugate).
4.  $\det(a\mathbf{A}) = a^n \det(\mathbf{A})$ , for  $a \in \mathbb{C}$ .
5.  $\mathbf{A}$  is singular if and only if  $\det(\mathbf{A}) = 0$ .
6. If  $\mathbf{A} = \begin{bmatrix} \mathbf{C} & \mathbf{D} \\ \mathbf{0} & \mathbf{E} \end{bmatrix}$  for some square matrices  $\mathbf{C}, \mathbf{E}$  then  $\det(\mathbf{A}) = \det(\mathbf{C})\det(\mathbf{E})$ .

7. **Cramer's rule** Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is nonsingular and  $\mathbf{b} \in \mathbb{C}^n$ . Let  $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$  be the unique solution of  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Then

$$x_j = \frac{\det(\mathbf{A}_j(\mathbf{b}))}{\det(\mathbf{A})}, \quad j = 1, 2, \dots, n,$$

where  $\mathbf{A}_j(\mathbf{b})$  denote the matrix obtained from  $\mathbf{A}$  by replacing the  $j$ th column of  $\mathbf{A}$  by  $\mathbf{b}$ .

8. **Adjoint.** Let  $\mathbf{A}_{i,j}$  denote the submatrix of  $\mathbf{A}$  obtained by deleting the  $i$ th row and  $j$ th column of  $\mathbf{A}$ . For  $\mathbf{A} \in \mathbb{C}^{n \times n}$  and  $1 \leq i, j \leq n$  the determinant  $\det(\mathbf{A}_{i,j})$  is called the **cofactor** of  $a_{ij}$ . The matrix  $\text{adj}(\mathbf{A}) \in \mathbb{C}^{n \times n}$  with elements  $\text{adj}(\mathbf{A})_{i,j} = (-1)^{i+j} \det(\mathbf{A}_{j,i})$  is called the **adjoint** of  $\mathbf{A}$ .
9. **Adjoint formula for the inverse.** If  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is nonsingular then

$$\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \text{adj}(\mathbf{A}).$$

10. **Cofactor expansion.** For any  $\mathbf{A} \in \mathbb{C}^{n \times n}$  we have

$$\det(\mathbf{A}) = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{i,j}) \quad \text{for } i = 1, \dots, n, \quad (27)$$

$$\det(\mathbf{A}) = \sum_{i=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{i,j}) \quad \text{for } j = 1, \dots, n. \quad (28)$$



Arthur Cayley, 1821-1895 (left), Gabriel Cramer 1704-1752 (center), Alexandre-Thophile Vandermonde, 1735-1796 (right). The notation  $||$  for determinants is due to Cayley 1841.

To compute the value of a determinant it is often convenient to use row- or column operations to introduce zeros in a row or column of  $\mathbf{A}$  and then use one of the cofactor expansions.

**Exercise 0.44 (Cramer's rule; special case)**

Solve the following system by Cramer's rule:

$$\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 3 \\ 6 \end{bmatrix}$$

**Exercise 0.45 (Adjoint matrix; special case)**

Show that if

$$\mathbf{A} = \begin{bmatrix} 2 & -6 & 3 \\ 3 & -2 & -6 \\ 6 & 3 & 2 \end{bmatrix},$$

then

$$\text{adj}(\mathbf{A}) = \begin{bmatrix} 14 & 21 & 42 \\ -42 & -14 & 21 \\ 21 & -42 & 14 \end{bmatrix}.$$

Moreover,

$$\text{adj}(\mathbf{A})\mathbf{A} = \begin{bmatrix} 343 & 0 & 0 \\ 0 & 343 & 0 \\ 0 & 0 & 343 \end{bmatrix} = \det(\mathbf{A})\mathbf{I}.$$

**Example 0.46 (Determinant equation for a straight line)**

The equation for a straight line through two points  $(x_1, y_1)$  and  $(x_2, y_2)$  in the plane can be written as the equation

$$\det(\mathbf{A}) := \begin{vmatrix} 1 & x & y \\ 1 & x_1 & y_1 \\ 1 & x_2 & y_2 \end{vmatrix} = 0$$

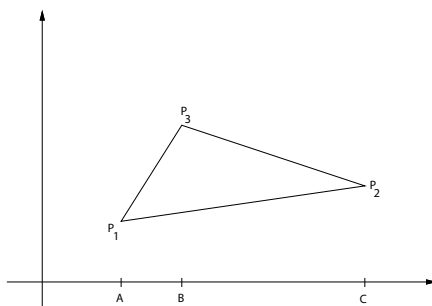
involving a determinant of order 3. We can compute this determinant using row operations of type 3. Subtracting row 2 from row 3 and then row 1 from row 2 we obtain

$$\begin{vmatrix} 1 & x & y \\ 1 & x_1 & y_1 \\ 1 & x_2 & y_2 \end{vmatrix} = \begin{vmatrix} 1 & x & y \\ 0 & x_1 - x & y_1 - y \\ 0 & x_2 - x_1 & y_2 - y_1 \end{vmatrix} = (x_1 - x)(y_2 - y_1) - (y_1 - y)(x_2 - x_1).$$

Rearranging the equation  $\det(\mathbf{A}) = 0$  we obtain

$$y - y_1 = \frac{y_2 - y_1}{x_2 - x_1}(x - x_1)$$

which is the slope form of the equation of a straight line.



**Figure 2.** The triangle  $T$  defined by the three points  $P_1$ ,  $P_2$  and  $P_3$ .

**Exercise 0.47 (Determinant equation for a plane)**

Show that

$$\begin{vmatrix} x & y & z & 1 \\ x_1 & y_1 & z_1 & 1 \\ x_2 & y_2 & z_2 & 1 \\ x_3 & y_3 & z_3 & 1 \end{vmatrix} = 0.$$

is the equation for a plane through three points  $(x_1, y_1, z_1)$ ,  $(x_2, y_2, z_2)$  and  $(x_3, y_3, z_3)$  in space.

**Exercise 0.48 (Signed area of a triangle)**

Let  $P_i = (x_i, y_i)$ ,  $i = 1, 2, 3$ , be three points in the plane defining a triangle  $T$ . Show that the area of  $T$  is<sup>1</sup>

$$A(T) = \frac{1}{2} \begin{vmatrix} 1 & 1 & 1 \\ x_1 & x_2 & x_3 \\ y_1 & y_2 & y_3 \end{vmatrix}.$$

The area is positive if we traverse the vertices in counterclockwise order.

**Exercise 0.49 (Vandermonde matrix)**

Show that

$$\begin{vmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \cdots & x_2^{n-1} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^{n-1} \end{vmatrix} = \prod_{i>j} (x_i - x_j),$$

<sup>1</sup>Hint:  $A(T) = A(ABP_3P_1) + A(P_3BCP_2) - A(P_1ACP_2)$ , c.f. Figure 2

where  $\prod_{i>j}(x_i - x_j) = \prod_{i=2}^n(x_i - x_1)(x_i - x_2) \cdots (x_i - x_{i-1})$ . This determinant is called the Vandermonde determinant.<sup>2</sup>

**Exercise 0.50 (Cauchy determinant (1842))**

Let  $\alpha = [\alpha_1, \dots, \alpha_n]^T$ ,  $\beta = [\beta_1, \dots, \beta_n]^T$  be in  $\mathbb{R}^n$ .

- a) Consider the matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  with elements  $a_{i,j} = 1/(\alpha_i + \beta_j)$ ,  $i, j = 1, 2, \dots, n$ . Show that

$$\det(\mathbf{A}) = Pg(\alpha)g(\beta)$$

where  $P = \prod_{i=1}^n \prod_{j=1}^n a_{ij}$ , and for  $\gamma = [\gamma_1, \dots, \gamma_n]^T$

$$g(\gamma) = \prod_{i=2}^n (\gamma_i - \gamma_1)(\gamma_i - \gamma_2) \cdots (\gamma_i - \gamma_{i-1})$$

*Hint: Multiply the  $i$ th row of  $\mathbf{A}$  by  $\prod_{j=1}^n (\alpha_i + \beta_j)$  for  $i = 1, 2, \dots, n$ . Call the resulting matrix  $\mathbf{C}$ . Each element of  $\mathbf{C}$  is a product of  $n-1$  factors  $\alpha_r + \beta_s$ . Hence  $\det(\mathbf{C})$  is a sum of terms where each term contain precisely  $n(n-1)$  factors  $\alpha_r + \beta_s$ . Thus  $\det(\mathbf{C}) = q(\alpha, \beta)$  where  $q$  is a polynomial of degree at most  $n(n-1)$  in  $\alpha_i$  and  $\beta_j$ . Since  $\det(\mathbf{A})$  and therefore  $\det(\mathbf{C})$  vanishes if  $\alpha_i = \alpha_j$  for some  $i \neq j$  or  $\beta_r = \beta_s$  for some  $r \neq s$ , we have that  $q(\alpha, \beta)$  must be divisible by each factor in  $g(\alpha)$  and  $g(\beta)$ . Since  $g(\alpha)$  and  $g(\beta)$  is a polynomial of degree  $n(n-1)$ , we have*

$$q(\alpha, \beta) = kg(\alpha)g(\beta)$$

for some constant  $k$  independent of  $\alpha$  and  $\beta$ . Show that  $k = 1$  by choosing  $\beta_i + \alpha_i = 0$ ,  $i = 1, 2, \dots, n$ .

- b) Notice that the cofactor of any element in the above matrix  $\mathbf{A}$  is the determinant of a matrix of similar form. Use the cofactor and determinant of  $\mathbf{A}$  to represent the elements of  $\mathbf{A}^{-1} = (b_{j,k})$ . Answer:

$$b_{j,k} = (\alpha_k + \beta_j)A_k(-\beta_j)B_j(-\alpha_k),$$

where

$$A_k(x) = \prod_{s \neq k} \left( \frac{\alpha_s - x}{\alpha_s - \alpha_k} \right), \quad B_k(x) = \prod_{s \neq k} \left( \frac{\beta_s - x}{\beta_s - \beta_k} \right).$$

<sup>2</sup>Hint: subtract  $x_n^k$  times column  $k$  from column  $k+1$  for  $k = n-1, n-2, \dots, 1$ .



**Exercise 0.51 (Inverse of the Hilbert matrix)**

Let  $\mathbf{H}_n = (h_{i,j})$  be the  $n \times n$  matrix with elements  $h_{i,j} = 1/(i+j-1)$ . Use Exercise 0.50 to show that the elements  $t_{i,j}^n$  in  $\mathbf{T}_n = \mathbf{H}_n^{-1}$  are given by

$$t_{i,j}^n = \frac{f(i)f(j)}{i+j-1},$$

where

$$f(i+1) = \left( \frac{i^2 - n^2}{i^2} \right) f(i), \quad i = 1, 2, \dots, \quad f(1) = -n.$$

## 0.7 Eigenpairs

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is a square matrix,  $\lambda \in \mathbb{C}$  and  $\mathbf{x} \in \mathbb{C}^n$ . We say that  $(\lambda, \mathbf{x})$  is an **eigenpair** for  $\mathbf{A}$  if  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  and  $\mathbf{x}$  is nonzero. The scalar  $\lambda$  is called an **eigenvalue** and  $\mathbf{x}$  is said to be an **eigenvector**.<sup>3</sup> The set of eigenvalues is called the **spectrum** of  $\mathbf{A}$  and is denoted by  $\sigma(\mathbf{A})$ . For example,  $\sigma(\mathbf{I}) = \{1, \dots, 1\} = \{1\}$ .

**Lemma 0.52 (Characteristic equation)**

For any  $\mathbf{A} \in \mathbb{C}^{n \times n}$  we have  $\lambda \in \sigma(\mathbf{A}) \iff \det(\mathbf{A} - \lambda\mathbf{I}) = 0$ .

**Proof.** Suppose  $(\lambda, \mathbf{x})$  is an eigenpair for  $\mathbf{A}$ . The equation  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  can be written  $(\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0}$ . Since  $\mathbf{x}$  is nonzero the matrix  $\mathbf{A} - \lambda\mathbf{I}$  must be singular with a zero determinant. Conversely, if  $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$  then  $\mathbf{A} - \lambda\mathbf{I}$  is singular and  $(\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0}$  for some nonzero  $\mathbf{x} \in \mathbb{C}^n$ . Thus  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  and  $(\lambda, \mathbf{x})$  is an eigenpair for  $\mathbf{A}$ .  $\square$

The expression  $\det(\mathbf{A} - \lambda\mathbf{I})$  is a polynomial of exact degree  $n$  in  $\lambda$ . For  $n = 3$  we have

$$\det(\mathbf{A} - \lambda\mathbf{I}) = \begin{vmatrix} a_{11} - \lambda & a_{12} & a_{13} \\ a_{21} & a_{22} - \lambda & a_{23} \\ a_{31} & a_{32} & a_{33} - \lambda \end{vmatrix}.$$

Expanding this determinant by the first column we find

$$\begin{aligned} \det(\mathbf{A} - \lambda\mathbf{I}) &= (a_{11} - \lambda) \begin{vmatrix} a_{22} - \lambda & a_{23} \\ a_{32} & a_{33} - \lambda \end{vmatrix} - a_{21} \begin{vmatrix} a_{12} & a_{13} \\ a_{32} & a_{33} - \lambda \end{vmatrix} \\ &\quad + a_{31} \begin{vmatrix} a_{12} & a_{13} \\ a_{22} - \lambda & a_{23} \end{vmatrix} = (a_{11} - \lambda)(a_{22} - \lambda)(a_{33} - \lambda) + r(\lambda) \end{aligned}$$

<sup>3</sup>The word “eigen” is derived from German and means “own”

for some polynomial  $r$  of degree at most one. In general

$$\det(\mathbf{A} - \lambda\mathbf{I}) = (a_{11} - \lambda)(a_{22} - \lambda) \cdots (a_{nn} - \lambda) + r(\lambda), \quad (29)$$

where each term in  $r(\lambda)$  has at most  $n - 2$  factors containing  $\lambda$ . It follows that  $r$  is a polynomial of degree at most  $n - 2$ ,  $\pi_{\mathbf{A}}$  is a polynomial of exact degree  $n$ , and the eigenvalues are the roots of this polynomial.

We observe that  $\det(\mathbf{A} - \lambda\mathbf{I}) = (-1)^n \det(\lambda\mathbf{I} - \mathbf{A})$  so  $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$  if and only if  $\det(\lambda\mathbf{I} - \mathbf{A}) = 0$ .

**Definition 0.53 (Characteristic polynomial of a matrix)**

The function  $\pi_{\mathbf{A}}: \mathbb{C} \rightarrow \mathbb{C}$  given by  $\pi_{\mathbf{A}}(\lambda) = \det(\mathbf{A} - \lambda\mathbf{I})$  is called the **characteristic polynomial** of  $\mathbf{A}$ . The equation  $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$  is called the **characteristic equation** of  $\mathbf{A}$ .

By the fundamental theorem of algebra an  $n \times n$  matrix has, counting multiplicities, precisely  $n$  eigenvalues  $\lambda_1, \dots, \lambda_n$  some of which might be complex even if  $\mathbf{A}$  is real. The complex eigenpairs of a real matrix occur in complex conjugate pairs. Indeed, taking the complex conjugate on both sides of the equation  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  with  $\mathbf{A}$  real gives  $\mathbf{A}\bar{\mathbf{x}} = \bar{\lambda}\bar{\mathbf{x}}$ .

Using Property 6. of determinants we have an additional characterization of a singular matrix.

**Theorem 0.54 (Zero eigenvalue)**

The matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is singular if and only if zero is an eigenvalue.

*Proof.* Zero is an eigenvalue if and only if  $\pi_{\mathbf{A}}(0) = \det(\mathbf{A}) = 0$  which happens if and only if  $\mathbf{A}$  is singular.  $\square$

In general it is not easy to find all eigenvalues of a matrix. One notable exception is a triangular matrix. By Property 2. of determinants we obtain

**Theorem 0.55 (Eigenvalues of a triangular matrix)**

The eigenvalues of a triangular matrix are given by its diagonal elements.

## 0.8 Algorithms and Numerical Stability

In this text we consider mathematical problems (i. e., linear algebra problems) and many detailed numerical algorithms to solve them. Complexity is discussed briefly in Section 2.2.2. As for programming issues we often vectorize the algorithms leading to shorter and more efficient programs. Stability is important both for the mathematical problems and for the numerical algorithms. Stability can be studied in terms of perturbation theory leading to condition numbers, see

---

Chapters 7, 11, 12. We will often use phrases like “the algorithm is numerically stable” or “the algorithm is not numerically stable” without saying precisely what we mean by this. Loosely speaking, an algorithm is numerically stable if the solution, computed in floating point arithmetic, is the exact solution of a slightly perturbed problem. To determine upper bounds for these perturbations is the topic of **backward error analysis**. We give a rather limited introduction to floating point arithmetic and backward error analysis in Appendix B, but in the text we will not discuss this. This does not mean that numerical stability is not an important issue. In fact, numerical stability is crucial for a good algorithm. For thorough treatments of numerical stability issues we refer to the books [12] and [26, 27].

A list of freely available software for solving linear algebra problems can be found at

<http://www.netlib.org/utk/people/JackDongarra/la-sw.html>



## **Part I**

# **Direct Methods for Linear Systems**



## Chapter 1

# A Special Linear System

Consider a system of  $n$  linear equations in  $n$  unknowns. In component form the system can be written

$$\begin{array}{ccccccc} a_{11}x_1 + a_{12}x_2 + & \cdots & + a_{1n}x_n = b_1, \\ a_{21}x_1 + a_{22}x_2 + & \cdots & + a_{2n}x_n = b_2, \\ \vdots & & & & \vdots & & \vdots \\ a_{n1}x_1 + a_{n2}x_2 + & \cdots & + a_{nn}x_n = b_n, \end{array}$$

and in matrix form

$$\mathbf{A}\mathbf{x} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} = \mathbf{b}.$$

The elements of  $\mathbf{A}$  and  $\mathbf{b}$  can be either real or complex numbers.

We recall (see Theorem 0.34) that the square system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  has a unique solution for all right hand sides  $\mathbf{b}$  if and only if  $\mathbf{A}$  is nonsingular, i. e., the homogeneous system  $\mathbf{A}\mathbf{x} = \mathbf{0}$  only has the solution  $\mathbf{x} = \mathbf{0}$ . We also recall (cf. Theorem 0.39) that a square matrix has an inverse if and only if  $\mathbf{A}$  is nonsingular, and the solution of  $\mathbf{A}\mathbf{x} = \mathbf{b}$  can be written  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ , where  $\mathbf{A}^{-1}$  is the inverse of  $\mathbf{A}$ .



Carl Friedrich Gauss, 1777-1855 (left), Myrick Hascall Doolittle, 1830-1911 (right).

Gaussian elimination with row interchanges is the classical method for solving  $n$  linear equations in  $n$  unknowns<sup>4</sup>. After an introductory and elementary discussion of Gaussian elimination we consider a problem leading to a linear system where the coefficient matrix is tridiagonal. This special matrix will occur repeatedly throughout this text. We then give an introduction to block multiplication which is an indispensable tool in matrix analysis. We end the chapter with some basic properties of triangular matrices.

## 1.1 Gaussian Elimination Example

We illustrate how Gaussian elimination works on a  $3 \times 3$  system. For a general discussion see Section 2.9.

### Example 1.1 (Gaussian elimination on a $3 \times 3$ system)

Consider a nonsingular system of three equations in three unknowns:

$$\begin{aligned} a_{11}^1 x_1 + a_{12}^1 x_2 + a_{13}^1 x_3 &= b_1^1, & \text{I} \\ a_{21}^1 x_1 + a_{22}^1 x_2 + a_{23}^1 x_3 &= b_2^1, & \text{II} \\ a_{31}^1 x_1 + a_{32}^1 x_2 + a_{33}^1 x_3 &= b_3^1. & \text{III.} \end{aligned}$$

To solve this system by Gaussian elimination suppose  $a_{11}^1 \neq 0$ . We subtract  $m_{21} := a_{21}^1/a_{11}^1$  times equation I from equation II and  $m_{31} := a_{31}^1/a_{11}^1$  times equation I

<sup>4</sup>The method was known long before Gauss used it in 1809. It was further developed by Doolittle in 1881, see [6].



from equation III. The result is

$$\begin{aligned} a_{11}^1 x_1 + a_{12}^1 x_2 + a_{13}^1 x_3 &= b_1^1, & \text{I} \\ a_{22}^2 x_2 + a_{23}^2 x_3 &= b_2^2, & \text{II}' \\ a_{32}^2 x_2 + a_{33}^2 x_3 &= b_3^2, & \text{III}', \end{aligned}$$

where  $b_i^2 = b_i^1 - m_{i1} b_1^1$  for  $i = 2, 3$  and  $a_{ij}^2 = a_{ij}^1 - m_{i1} a_{1j}^1$  for  $i, j = 2, 3$ . If  $a_{11}^1 = 0$  and  $a_{21}^1 \neq 0$  we first interchange equation I and equation II. If  $a_{11}^1 = a_{21}^1 = 0$  we interchange equation I and III. Since the system is nonsingular the first column cannot be zero and an interchange is always possible.

If  $a_{22}^2 \neq 0$  we subtract  $m_{32} := a_{32}^2/a_{22}^2$  times equation II' from equation III' to obtain

$$\begin{aligned} a_{11}^1 x_1 + a_{12}^1 x_2 + a_{13}^1 x_3 &= b_1^1, & \text{I} \\ a_{22}^2 x_2 + a_{23}^2 x_3 &= b_2^2, & \text{II}' \\ a_{33}^3 x_3 &= b_3^3, & \text{III}'', \end{aligned}$$

where  $b_3^3 = b_3^2 - m_{32} b_2^2$  and  $a_{33}^3 = a_{33}^2 - m_{32} a_{23}^2$ . If  $a_{22}^2 = 0$  then  $a_{32}^2 \neq 0$  (cf. Theorem 2.59) and we first interchange equation II' and equation III'. The reduced system is easy to solve since it is upper triangular. Starting from the bottom and moving upwards we find

$$\begin{aligned} x_3 &= b_3^3/a_{33}^3 \\ x_2 &= (b_2^2 - a_{23}^2 x_3)/a_{22}^2 \\ x_1 &= (b_1^1 - a_{12}^1 x_2 - a_{13}^1 x_3)/a_{11}^1. \end{aligned}$$

This is known as **back substitution**.

### Exercise 1.2 (Gaussian elimination example)

Solve the linear system  $\mathbf{Ax} := \begin{bmatrix} 1 & 1 & -1 \\ -1 & 1 & 3 \\ 2 & 8 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$  using Gaussian elimination.

Gaussian elimination with row interchanges can in principle be used to solve any nonsingular linear system (cf. Theorem 2.59). However, for many systems occurring in applications this method in its general form is not necessarily the method of choice. Some of the issues are:

1. **Computing time.** Solving a dense system of order  $n$  by Gaussian elimination requires  $O(n^3)$  arithmetic operations and solving large linear systems can require more time than we are willing to spend. For example, if  $n = 10^6$  and one arithmetic operation takes  $10^{-12}$  seconds then the computing time could be a staggering  $10^{-12} n^3 \approx 278$  hours.

2. **Row interchanges** is another issue in Gaussian elimination. For example, if we interchange two rows in a tridiagonal matrix then the tridiagonal structure is lost in general.
3. **Stability.** For a well conditioned problem<sup>5</sup> Gaussian elimination using floating point arithmetic will in most cases give an accurate solution. However there is no guarantee for this, see the book [12] for a thorough discussion.

In this chapter we present a problem leading to a  $n \times n$  tridiagonal linear system. We show that row interchanges are not necessary for the two problems we consider and derive stable algorithms that only requires  $O(n)$  arithmetic operations.

## 1.2 The Tridiagonal Second Derivative Matrix

Consider the simple **two point boundary value problem**

$$-u''(x) = f(x), \quad x \in [0, 1], \quad u(0) = 0, \quad u(1) = 0, \quad (1.1)$$

where  $f$  is a given continuous function on  $[0, 1]$ . This problem is also known as the **one-dimensional (1D) Poisson problem**. In principle it is easy to solve (1.1) exactly. We just integrate  $f$  twice and determine the two integration constants so that the homogeneous boundary conditions  $u(0) = u(1) = 0$  are satisfied. For example, if  $f(x) = 1$  then  $u(x) = x(x - 1)/2$  is the solution.

Suppose  $f$  cannot be integrated exactly. Problem (1.1) can then be solved approximately using the **finite difference method**. We need a difference approximation to the second derivative. If  $g$  is a function differentiable at  $x$  then

$$g'(x) = \lim_{h \rightarrow 0} \frac{g(x + \frac{h}{2}) - g(x - \frac{h}{2})}{h}$$

and applying this to a function  $u$  that is twice differentiable at  $x$

$$\begin{aligned} u''(x) &= \lim_{h \rightarrow 0} \frac{u'(x + \frac{h}{2}) - u'(x - \frac{h}{2})}{h} = \lim_{h \rightarrow 0} \frac{\frac{u(x+h) - u(x)}{h} - \frac{u(x) - u(x-h)}{h}}{h} \\ &= \lim_{h \rightarrow 0} \frac{u(x+h) - 2u(x) + u(x-h)}{h^2}. \end{aligned}$$

To define the points where this difference approximation is used we choose a positive integer  $m$ , let  $h := 1/(m+1)$  be the discretization parameter, and replace the interval  $[0, 1]$  by grid points  $x_j := jh$  for  $j = 0, 1, \dots, m+1$ . We then obtain approximations  $v_j$  to the exact solution  $u(x_j)$  for  $j = 1, \dots, m$  by replacing the differential equation by the difference equation

$$\frac{-v_{j-1} + 2v_j - v_{j+1}}{h^2} = f(jh), \quad j = 1, \dots, m, \quad v_0 = v_{m+1} = 0.$$

<sup>5</sup>see Chapter 7 for an introduction to condition numbers.



then  $\mathbf{A} = \mathbf{LU}$ , and if  $u_1, u_2, \dots, u_{n-1}$  are nonzero then (1.4) is well defined. If in addition  $u_n \neq 0$  then we can solve  $\mathbf{Ly} = \mathbf{b}$  and  $\mathbf{Ux} = \mathbf{y}$  for  $\mathbf{y}$  and  $\mathbf{x}$ .

$$\begin{aligned} y_1 &= b_1, & y_k &= b_k - l_{k-1}y_{k-1}, & k &= 2, 3, \dots, n, \\ x_n &= y_n/u_n, & x_k &= (y_k - c_k x_{k+1})/u_k, & k &= n-1, \dots, 2, 1. \end{aligned} \quad (1.5)$$

We formulate this as two algorithms. Since division by zero can occur, the algorithms will not work in general. We give sufficient conditions for success in Theorem 1.7 below.

### Algorithm 1.3 (trifactor)

Vectors  $\mathbf{l} \in \mathbb{C}^{n-1}$ ,  $\mathbf{u} \in \mathbb{C}^n$  are computed from  $\mathbf{a}, \mathbf{c} \in \mathbb{C}^{n-1}$ ,  $\mathbf{d} \in \mathbb{C}^n$ . This implements the LU factorization of a tridiagonal matrix.

```

1 function [l,u]=trifactor(a,d,c)
2 % [l,u]=trifactor(a,d,c)
3 u=d; l=a;
4 for k=1:length(a)
5     l(k)=a(k)/u(k);
6     u(k+1)=d(k+1)-l(k)*c(k);
7 end

```

### Algorithm 1.4 (trisolve)

The solution  $\mathbf{x}$  of the tridiagonal system  $\mathbf{LUx} = \mathbf{b}$  is found from (1.5). Here  $\mathbf{l}, \mathbf{c} \in \mathbb{C}^{n-1}$ ,  $\mathbf{u} \in \mathbb{C}^n$  and  $\mathbf{b} \in \mathbb{C}^{n,r}$  for some  $r \in \mathbb{N}$ . Thus we can solve a system with several righthand sides. The vectors  $\mathbf{l}, \mathbf{u}$  can be output from `trifactor`.

```

1 function x = trisolve(l,u,c,b)
2 % x = trisolve(l,u,c,b)
3 x=b;
4 n= size(b,1);
5 for k=2:n
6     x(k,:)=b(k,:)-l(k-1)*x(k-1,:);
7 end
8 x(n,:)=x(n,+)/u(n);
9 for k=n-1:-1:1
10    x(k,:)=(x(k,)-c(k)*x(k+1,))/u(k);
11 end

```

The number of arithmetic operations to compute the LU factorization of a tridiagonal matrix using Algorithm 1.3 is  $3n - 3$ , while the number of arithmetic operations for Algorithm 1.4 is  $5n - 4$ . This means that the complexity to solve a tridiagonal system is  $O(n)$ , or more precisely  $8n - 7$ , and this number only grows linearly with  $n$ , while Gaussian elimination on a full  $n \times n$  system is an  $O(n^3)$  process.

### 1.3.2 Diagonal dominance

We show that Algorithms 1.3, 1.4 are well defined for a class of tridiagonal linear systems. Moreover, these linear systems have unique solutions.

**Definition 1.5 (Diagonal dominance)**

The matrix  $\mathbf{A} = [a_{ij}] \in \mathbb{C}^{n \times n}$  is **weakly diagonally dominant** if

$$|a_{ii}| \geq \sum_{j \neq i} |a_{ij}|, \quad i = 1, \dots, n. \quad (1.6)$$

It is **strictly diagonally dominant** if strict inequality holds for  $i = 1, \dots, n$ .

The following holds for strictly diagonally dominant matrices.

**Theorem 1.6 (Strict diagonal dominance)**

A strictly diagonally dominant matrix is nonsingular. Moreover, the solution  $\mathbf{x}$  of  $\mathbf{Ax} = \mathbf{b}$  is bounded as follows:

$$\max_{1 \leq i \leq n} |x_i| \leq \max_{1 \leq i \leq n} \left( \frac{|b_i|}{\sigma_i} \right), \quad \text{where } \sigma_i := |a_{ii}| - \sum_{j \neq i} |a_{ij}|. \quad (1.7)$$

**Proof.** We first show that the bound (1.7) holds for any solution  $\mathbf{x}$ . Choose  $k$  so that  $|x_k| = \max_i |x_i|$ . Then

$$|b_k| = |a_{kk}x_k + \sum_{j \neq k} a_{kj}x_j| \geq |a_{kk}||x_k| - \sum_{j \neq k} |a_{kj}||x_j| \geq |x_k| \left( |a_{kk}| - \sum_{j \neq k} |a_{kj}| \right),$$

and this implies  $\max_{1 \leq i \leq n} |x_i| = |x_k| \leq \frac{|b_k|}{\sigma_k} \leq \max_{1 \leq i \leq n} \left( \frac{|b_i|}{\sigma_i} \right)$ . For nonsingularity, if  $\mathbf{Ax} = \mathbf{0}$ , then  $\max_{1 \leq i \leq n} |x_i| \leq 0$  by (1.7), and so  $\mathbf{x} = \mathbf{0}$ .  $\square$

The zero matrix is weakly diagonally dominant and we need additional condition to guarantee nonsingularity. Consider the 3 matrices

$$\mathbf{A}_1 = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 1 \end{bmatrix}, \quad \mathbf{A}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{A}_3 = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}.$$

They are all weakly diagonally dominant, but  $\mathbf{A}_1$  and  $\mathbf{A}_2$  are singular, while  $\mathbf{A}_3$  is nonsingular. Indeed, for  $\mathbf{A}_1$  column two is the sum of columns one and three,  $\mathbf{A}_2$  has a zero row, and  $\det(\mathbf{A}_3) = 4 \neq 0$ .

In the literature diagonal dominance is therefore most often defined by including some additional sufficient conditions. We also need conditions guaranteeing that the LU factorization (1.3) of a tridiagonal matrix is well defined.

**Theorem 1.7 (Weak diagonal dominance)**

Suppose  $\mathbf{A} = \text{tridiag}(a_i, d_i, c_i) \in \mathbb{C}^{n \times n}$  is tridiagonal and weakly diagonally dominant. If  $|d_1| > |c_1|$  and  $a_i \neq 0$  for  $i = 1, \dots, n-2$ , then  $\mathbf{A}$  has a unique LU factorization (1.3). If in addition  $d_n \neq 0$ , then  $\mathbf{A}$  is nonsingular.

**Proof.** The matrix  $\mathbf{A}$  has an LU factorization if the  $u_k$ 's in (1.4) are nonzero for  $k = 1, \dots, n-1$ . For this it is sufficient to show by induction that  $|u_k| > |c_k|$  for  $k = 1, \dots, n-1$ . By assumption  $|u_1| = |d_1| > |c_1|$ . Suppose  $|u_k| > |c_k|$  for some  $1 \leq k \leq n-2$ . Then  $|c_k|/|u_k| < 1$  and by (1.4) and since  $a_k \neq 0$

$$|u_{k+1}| = |d_{k+1} - l_k c_k| = |d_{k+1} - \frac{a_k c_k}{u_k}| \geq |d_{k+1}| - \frac{|a_k| |c_k|}{|u_k|} > |d_{k+1}| - |a_k|. \quad (1.8)$$

This also holds for  $k = n-1$  if  $a_{n-1} \neq 0$ . By (1.8) and weak diagonal dominance  $|u_{k+1}| > |d_{k+1}| - |a_k| \geq |c_{k+1}|$  and it follows by induction that an LU factorization exists. It is unique since any LU factorization must satisfy (1.4). For nonsingularity we need to show that  $u_n \neq 0$ . For then by Lemma 1.22, both  $\mathbf{L}$  and  $\mathbf{U}$  are nonsingular, and this is equivalent to  $\mathbf{A} = \mathbf{LU}$  being nonsingular. If  $a_{n-1} \neq 0$  then by (1.4)  $|u_n| > |d_n| - |a_{n-1}| \geq 0$  by weak diagonal dominance, while if  $a_{n-1} = 0$  then again by (1.8)  $|u_n| \geq |d_n| > 0$ .  $\square$

**Exercise 1.8 (Strict diagonal dominance)**

Show that a strictly diagonally dominant and tridiagonal matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  has a unique LU factorization<sup>6</sup>.

Consider now the special system  $\mathbf{T}\mathbf{v} = \mathbf{b}$  given by (1.2). The matrix  $\mathbf{T}$  is weakly diagonally dominant and satisfies the additional conditions in Theorem 1.7. Thus it is nonsingular and we can solve the system in  $O(n)$  arithmetic operations using Algorithms 1.3, 1.4.

We could use the explicit inverse of  $\mathbf{T}$ , given in Exercise 1.10, to compute the solution of  $\mathbf{T}\mathbf{v} = \mathbf{b}$  as  $\mathbf{v} = \mathbf{T}^{-1}\mathbf{b}$ . However this is not a good idea. In fact, all elements in  $\mathbf{T}^{-1}$  are nonzero and the calculation of  $\mathbf{T}^{-1}\mathbf{b}$  requires  $O(n^2)$  operations.

**Exercise 1.9 (LU factorization of 2. derivative matrix)**

<sup>6</sup>Hint, argue as in (1.8)

Show that  $\mathbf{T} = \mathbf{LU}$ , where

$$\mathbf{L} = \begin{bmatrix} 1 & 0 & \cdots & \cdots & 0 \\ -\frac{1}{2} & 1 & \ddots & & \vdots \\ 0 & -\frac{2}{3} & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & -\frac{m-1}{m} & 1 \end{bmatrix}, \mathbf{U} = \begin{bmatrix} 2 & -1 & 0 & \cdots & 0 \\ 0 & \frac{3}{2} & -1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ \vdots & & \ddots & \frac{m}{m-1} & -1 \\ 0 & \cdots & \cdots & 0 & \frac{m+1}{m} \end{bmatrix} \quad (1.9)$$

is the LU factorization of  $\mathbf{T}$ .

**Exercise 1.10 (Inverse of 2. derivative matrix)**

Let  $\mathbf{S} \in \mathbb{R}^{m \times m}$  have elements  $s_{ij}$  given by

$$s_{i,j} = s_{j,i} = \frac{1}{m+1} j(m+1-i), \quad 1 \leq j \leq i \leq m. \quad (1.10)$$

Show that  $\mathbf{ST} = \mathbf{I}$  and conclude that  $\mathbf{T}^{-1} = \mathbf{S}$ .

**Exercise 1.11 (Central difference approximation of 2. derivative)**

Consider

$$\delta^2 f(x) := \frac{f(x+h) - 2f(x) + f(x-h)}{h^2}, \quad h > 0, \quad f: [x-h, x+h] \rightarrow \mathbb{R}.$$

1. Show using Taylor expansion that if  $f \in C^2[x-h, x+h]$  then for some  $\eta_2$

$$\delta^2 f(x) = f''(\eta_2), \quad x-h < \eta_2 < x+h.$$

2. If  $f \in C^4[x-h, x+h]$  then for some  $\eta_4$

$$\delta^2 f(x) = f''(x) + \frac{h^2}{12} f^{(4)}(\eta_4), \quad x-h < \eta_4 < x+h.$$

$\delta^2 f(x)$  is known as the **central difference approximation** to the second derivative at  $x$ .

**Exercise 1.12 (Two point boundary value problem)**

We consider a finite difference method for the two point boundary value problem

$$\begin{aligned} -u''(x) + r(x)u'(x) + q(x)u(x) &= f(x), \quad \text{for } x \in [a, b], \\ u(a) &= g_0, \quad u(b) = g_1. \end{aligned} \quad (1.11)$$

We assume that the given functions  $f, q$  and  $r$  are continuous on  $[a, b]$  and that  $q(x) \geq 0$  for  $x \in [a, b]$ . It can then be shown that (1.11) has a unique solution  $u$ .

To solve (1.11) numerically we choose  $m \in \mathbb{N}$ ,  $h = (b-a)/(m+1)$ ,  $x_j = a+jh$  for  $j = 0, 1, \dots, m+1$  and solve the difference equation

$$\frac{-v_{j-1} + 2v_j - v_{j+1}}{h^2} + r(x_j) \frac{v_{j+1} - v_{j-1}}{2h} + q(x_j)v_j = f(x_j), \quad j = 1, \dots, m, \quad (1.12)$$

with  $v_0 = g_0$  and  $v_{m+1} = g_1$ .

- (a) Show that (1.12) leads to a tridiagonal linear system  $\mathbf{A}\mathbf{v} = \mathbf{b}$ , where  $\mathbf{A} = \text{tridiag}(a_j, d_j, c_j) \in \mathbb{R}^{m \times m}$  has elements

$$a_j = -1 - \frac{h}{2}r(x_j), \quad c_j = -1 + \frac{h}{2}r(x_j), \quad d_j = 2 + h^2q(x_j),$$

and

$$b_j = \begin{cases} h^2 f(x_1) - a_1 g_0, & \text{if } j = 1, \\ h^2 f(x_j), & \text{if } 2 \leq j \leq m-1, \\ h^2 f(x_m) - c_m g_1, & \text{if } j = m. \end{cases}$$

- (b) Show that the linear system satisfies the conditions in Theorem 1.7 if the spacing  $h$  is so small that  $\frac{h}{2}|r(x)| < 1$  for all  $x \in [a, b]$ .
- (c) Propose a method to find  $v_1, \dots, v_m$ .

### Exercise 1.13 (Two point boundary value problem; computation)

- (a) Consider the problem (1.11) with  $r = 0$ ,  $f = q = 1$  and boundary conditions  $u(0) = 1$ ,  $u(1) = 0$ . The exact solution is  $u(x) = 1 - \sinh x / \sinh 1$ . Write a computer program to solve (1.12) for  $h = 0.1, 0.05, 0.025, 0.0125$ , and compute the "error"  $\max_{1 \leq j \leq m} |u(x_j) - v_j|$  for each  $h$ .
- (b) Make a combined plot of the solution  $u$  and the computed points  $v_j$ ,  $j = 0, \dots, m+1$  for  $h = 0.1$ .
- (c) One can show that the error is proportional to  $h^p$  for some integer  $p$ . Estimate  $p$  based on the error for  $h = 0.1, 0.05, 0.025, 0.0125$ .

## 1.4 Block Multiplication

Block multiplication is a powerful and essential tool for dealing with matrices. It will be used extensively in this book.



A rectangular matrix  $\mathbf{A}$  can be partitioned into submatrices by drawing horizontal lines between selected rows and vertical lines between selected columns. For example, the matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$

can be partitioned as

$$(i) \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} = \left[ \begin{array}{c|cc} 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ 7 & 8 & 9 \end{array} \right], \quad (ii) [\mathbf{a}_{:1}, \mathbf{a}_{:2}, \mathbf{a}_{:3}] = \left[ \begin{array}{c|c|c} 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ \hline 7 & 8 & 9 \end{array} \right],$$

$$(iii) \begin{bmatrix} \mathbf{a}_{1:}^T \\ \mathbf{a}_{2:}^T \\ \mathbf{a}_{3:}^T \end{bmatrix} = \left[ \begin{array}{ccc} 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ \hline 7 & 8 & 9 \end{array} \right], \quad (iv) [\mathbf{A}_{11}, \mathbf{A}_{12}] = \left[ \begin{array}{c|cc} 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ 7 & 8 & 9 \end{array} \right].$$

In (i) the matrix  $\mathbf{A}$  is divided into four submatrices

$$\mathbf{A}_{11} = [1], \quad \mathbf{A}_{12} = [2, 3], \quad \mathbf{A}_{21} = \begin{bmatrix} 4 \\ 7 \end{bmatrix}, \quad \text{and} \quad \mathbf{A}_{22} = \begin{bmatrix} 5 & 6 \\ 8 & 9 \end{bmatrix},$$

while in (ii) and (iii)  $\mathbf{A}$  has been partitioned into columns and rows, respectively. The submatrices in a partition are often referred to as **blocks** and a partitioned matrix is sometimes called a **block matrix**.

In the following we assume that  $\mathbf{A} \in \mathbb{C}^{m \times p}$  and  $\mathbf{B} \in \mathbb{C}^{p \times n}$ . Here are some rules and observations for block multiplication.

1. If  $\mathbf{B} = [\mathbf{b}_{:1}, \dots, \mathbf{b}_{:n}]$  is partitioned into columns then the partition of the product  $\mathbf{AB}$  into columns is

$$\mathbf{AB} = [\mathbf{Ab}_{:1}, \mathbf{Ab}_{:2}, \dots, \mathbf{Ab}_{:n}].$$

In particular, if  $\mathbf{I}$  is the identity matrix of order  $p$  then

$$\mathbf{A} = \mathbf{AI} = \mathbf{A} [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p] = [\mathbf{Ae}_1, \mathbf{Ae}_2, \dots, \mathbf{Ae}_p]$$

and we see that column  $j$  of  $\mathbf{A}$  can be written  $\mathbf{Ae}_j$  for  $j = 1, \dots, p$ .

2. Similarly, if  $\mathbf{A}$  is partitioned into rows then

$$\mathbf{AB} = \begin{bmatrix} a_{1:}^T \\ a_{2:}^T \\ \vdots \\ a_{m:}^T \end{bmatrix} \mathbf{B} = \begin{bmatrix} a_{1:}^T \mathbf{B} \\ a_{2:}^T \mathbf{B} \\ \vdots \\ a_{m:}^T \mathbf{B} \end{bmatrix},$$

and taking  $\mathbf{A} = \mathbf{I}$  it follows that row  $i$  of  $\mathbf{B}$  can be written  $\mathbf{e}_i^T \mathbf{B}$  for  $i = 1, \dots, m$ .

3. It is often useful to write the matrix-vector product  $\mathbf{A}\mathbf{x}$  as a linear combination of the columns of  $\mathbf{A}$

$$\mathbf{A}\mathbf{x} = x_1\mathbf{a}_{:1} + x_2\mathbf{a}_{:2} + \cdots + x_p\mathbf{a}_{:p}.$$

4. If  $\mathbf{B} = [\mathbf{B}_1, \mathbf{B}_2]$ , where  $\mathbf{B}_1 \in \mathbb{C}^{p \times r}$  and  $\mathbf{B}_2 \in \mathbb{C}^{p \times (n-r)}$  then

$$\mathbf{A}[\mathbf{B}_1, \mathbf{B}_2] = [\mathbf{A}\mathbf{B}_1, \mathbf{A}\mathbf{B}_2].$$

This follows from Rule 1. by an appropriate grouping of columns.

5. If  $\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{bmatrix}$ , where  $\mathbf{A}_1 \in \mathbb{C}^{k \times p}$  and  $\mathbf{A}_2 \in \mathbb{C}^{(m-k) \times p}$  then

$$\begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \end{bmatrix} \mathbf{B} = \begin{bmatrix} \mathbf{A}_1 \mathbf{B} \\ \mathbf{A}_2 \mathbf{B} \end{bmatrix}.$$

This follows from Rule 2. by a grouping of rows.

6. If  $\mathbf{A} = [\mathbf{A}_1, \mathbf{A}_2]$  and  $\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{bmatrix}$ , where  $\mathbf{A}_1 \in \mathbb{C}^{m \times s}$ ,  $\mathbf{A}_2 \in \mathbb{C}^{m \times (p-s)}$ ,  $\mathbf{B}_1 \in \mathbb{C}^{s \times n}$  and  $\mathbf{B}_2 \in \mathbb{C}^{(p-s) \times n}$  then

$$[\mathbf{A}_1, \mathbf{A}_2] \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{bmatrix} = [\mathbf{A}_1 \mathbf{B}_1 + \mathbf{A}_2 \mathbf{B}_2].$$

Indeed,  $(\mathbf{A}\mathbf{B})_{ij} = \sum_{k=1}^p a_{ik}b_{kj} = \sum_{k=1}^s a_{ik}b_{kj} + \sum_{k=s+1}^p a_{ik}b_{kj} = (\mathbf{A}_1\mathbf{B}_1)_{ij} + (\mathbf{A}_2\mathbf{B}_2)_{ij} = (\mathbf{A}_1\mathbf{B}_1 + \mathbf{A}_2\mathbf{B}_2)_{ij}$ .

7. If  $\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$  and  $\mathbf{B} = \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix}$  then

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{11}\mathbf{B}_{11} + \mathbf{A}_{12}\mathbf{B}_{21} & \mathbf{A}_{11}\mathbf{B}_{12} + \mathbf{A}_{12}\mathbf{B}_{22} \\ \mathbf{A}_{21}\mathbf{B}_{11} + \mathbf{A}_{22}\mathbf{B}_{21} & \mathbf{A}_{21}\mathbf{B}_{12} + \mathbf{A}_{22}\mathbf{B}_{22} \end{bmatrix},$$

provided the vertical partition in  $\mathbf{A}$  matches the horizontal one in  $\mathbf{B}$ , i.e. the number of columns in  $\mathbf{A}_{11}$  and  $\mathbf{A}_{21}$  equals the number of rows in  $\mathbf{B}_{11}$  and  $\mathbf{B}_{12}$  and the number of columns in  $\mathbf{A}$  equals the number of rows in  $\mathbf{B}$ . To show this we use Rule 4. to obtain

$$\mathbf{A}\mathbf{B} = \left[ \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{11} \\ \mathbf{B}_{21} \end{bmatrix}, \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{12} \\ \mathbf{B}_{22} \end{bmatrix} \right].$$

We complete the proof using Rules 5. and 6.

8. Consider finally the general case. If all the matrix products  $\mathbf{A}_{ik}\mathbf{B}_{kj}$  in

$$\mathbf{C}_{ij} = \sum_{k=1}^s \mathbf{A}_{ik}\mathbf{B}_{kj}, \quad i = 1, \dots, p, \quad j = 1, \dots, q$$

are well defined then

$$\begin{bmatrix} \mathbf{A}_{11} & \cdots & \mathbf{A}_{1s} \\ \vdots & & \vdots \\ \mathbf{A}_{p1} & \cdots & \mathbf{A}_{ps} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{11} & \cdots & \mathbf{B}_{1q} \\ \vdots & & \vdots \\ \mathbf{B}_{s1} & \cdots & \mathbf{B}_{sq} \end{bmatrix} = \begin{bmatrix} \mathbf{C}_{11} & \cdots & \mathbf{C}_{1q} \\ \vdots & & \vdots \\ \mathbf{C}_{p1} & \cdots & \mathbf{C}_{pq} \end{bmatrix}.$$

The requirements are that

- the number of columns in  $\mathbf{A}$  is equal to the number of rows in  $\mathbf{B}$ .
- the position of the vertical partition lines in  $\mathbf{A}$  has to match the position of the horizontal partition lines in  $\mathbf{B}$ . The horizontal lines in  $\mathbf{A}$  and the vertical lines in  $\mathbf{B}$  can be anywhere.

**Exercise 1.14 (Matrix element as a quadratic form)**

For any matrix  $\mathbf{A}$  show that  $a_{ij} = \mathbf{e}_i^T \mathbf{A} \mathbf{e}_j$  for all  $i, j$ .

**Exercise 1.15 (Outer product expansion of a matrix)**

For any matrix  $\mathbf{A} \in \mathbb{C}^{m \times n}$  show that  $\mathbf{A} = \sum_{i=1}^m \sum_{j=1}^n a_{ij} \mathbf{e}_i \mathbf{e}_j^T$ .

**Exercise 1.16 (The product  $\mathbf{A}^T \mathbf{A}$ )**

Let  $\mathbf{B} = \mathbf{A}^T \mathbf{A}$ . Explain why this product is defined for any matrix  $\mathbf{A}$ . Show that  $b_{ij} = \mathbf{a}_{:i}^T \mathbf{a}_{:j}$  for all  $i, j$ .

**Exercise 1.17 (Outer product expansion)**

For  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{B} \in \mathbb{R}^{p \times n}$  show that

$$\mathbf{A} \mathbf{B}^T = \mathbf{a}_{:1} \mathbf{b}_{:1}^T + \mathbf{a}_{:2} \mathbf{b}_{:2}^T + \cdots + \mathbf{a}_{:n} \mathbf{b}_{:n}^T.$$

This is called the **outer product expansion** of the columns of  $\mathbf{A}$  and  $\mathbf{B}$ .

**Exercise 1.18 (System with many right hand sides; compact form)**

Suppose  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{m \times p}$ , and  $\mathbf{X} \in \mathbb{R}^{n \times p}$ . Show that

$$\mathbf{A} \mathbf{X} = \mathbf{B} \iff \mathbf{A} \mathbf{x}_{:j} = \mathbf{b}_{:j}, \quad j = 1, \dots, p.$$

**Exercise 1.19 (Block multiplication example)**

Suppose  $\mathbf{A} = [\mathbf{A}_1, \mathbf{A}_2]$  and  $\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{0} \end{bmatrix}$ . When is  $\mathbf{A} \mathbf{B} = \mathbf{A}_1 \mathbf{B}_1$ ?

**Exercise 1.20 (Another block multiplication example)**

Suppose  $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{n \times n}$  are given in block form by

$$\mathbf{A} := \begin{bmatrix} \lambda & \mathbf{a}^T \\ \mathbf{0} & \mathbf{A}_1 \end{bmatrix}, \quad \mathbf{B} := \begin{bmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{B}_1 \end{bmatrix}, \quad \mathbf{C} := \begin{bmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{C}_1 \end{bmatrix},$$

where  $\mathbf{A}_1, \mathbf{B}_1, \mathbf{C}_1 \in \mathbb{R}^{(n-1) \times (n-1)}$ . Show that

$$\mathbf{CAB} = \begin{bmatrix} \lambda & \mathbf{a}^T \mathbf{B}_1 \\ \mathbf{0} & \mathbf{C}_1 \mathbf{A}_1 \mathbf{B}_1 \end{bmatrix}.$$

**1.5 Triangular Matrices; Basic facts**

We need some basic facts about triangular matrices and we start with

**Lemma 1.21 (Inverse of a block triangular matrix)**

Suppose

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{0} & \mathbf{A}_{22} \end{bmatrix}$$

where  $\mathbf{A}, \mathbf{A}_{11}$  and  $\mathbf{A}_{22}$  are square matrices. Then  $\mathbf{A}$  is nonsingular if and only if both  $\mathbf{A}_{11}$  and  $\mathbf{A}_{22}$  are nonsingular. In that case

$$\mathbf{A}^{-1} = \begin{bmatrix} \mathbf{A}_{11}^{-1} & \mathbf{C} \\ \mathbf{0} & \mathbf{A}_{22}^{-1} \end{bmatrix}, \quad (1.13)$$

for some matrix  $\mathbf{C}$ .

**Proof.** Suppose  $\mathbf{A}$  is nonsingular. We partition  $\mathbf{B} := \mathbf{A}^{-1}$  conformally with  $\mathbf{A}$  and have

$$\mathbf{BA} = \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{0} & \mathbf{A}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} = \mathbf{I}$$

Using block-multiplication we find

$$\mathbf{B}_{11} \mathbf{A}_{11} = \mathbf{I}, \quad \mathbf{B}_{21} \mathbf{A}_{11} = \mathbf{0}, \quad \mathbf{B}_{21} \mathbf{A}_{12} + \mathbf{B}_{22} \mathbf{A}_{22} = \mathbf{I}, \quad \mathbf{B}_{11} \mathbf{A}_{12} + \mathbf{B}_{12} \mathbf{A}_{22} = \mathbf{0}.$$

The first equation implies that  $\mathbf{A}_{11}$  is nonsingular, this in turn implies that  $\mathbf{B}_{21} = \mathbf{0} \mathbf{A}_{11}^{-1} = \mathbf{0}$  in the second equation, and then the third equation simplifies to  $\mathbf{B}_{22} \mathbf{A}_{22} = \mathbf{I}$ . We conclude that also  $\mathbf{A}_{22}$  is nonsingular. From the fourth equation we find

$$\mathbf{B}_{12} = \mathbf{C} = -\mathbf{A}_{11}^{-1} \mathbf{A}_{12} \mathbf{A}_{22}^{-1}.$$

Conversely, if  $\mathbf{A}_{11}$  and  $\mathbf{A}_{22}$  are nonsingular then

$$\begin{bmatrix} \mathbf{A}_{11}^{-1} & -\mathbf{A}_{11}^{-1} \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \\ \mathbf{0} & \mathbf{A}_{22}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{0} & \mathbf{A}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} = \mathbf{I}$$

and  $\mathbf{A}$  is nonsingular with the indicated inverse.  $\square$

Consider now a triangular matrix.

**Lemma 1.22 (Inverse of a triangular matrix)**

An upper (lower) triangular matrix  $\mathbf{A} = [a_{ij}] \in \mathbb{C}^{n \times n}$  is nonsingular if and only if the diagonal elements  $a_{ii}$ ,  $i = 1, \dots, n$  are nonzero. In that case the inverse is upper (lower) triangular with diagonal elements  $a_{ii}^{-1}$ ,  $i = 1, \dots, n$ .

*Proof.* We use induction on  $n$ . The result holds for  $n = 1$ . The 1-by-1 matrix  $\mathbf{A} = [a_{11}]$  is nonsingular if and only if  $a_{11} \neq \mathbf{0}$  and in that case  $\mathbf{A}^{-1} = [a_{11}^{-1}]$ . Suppose the result holds for  $n = k$  and let  $\mathbf{A} \in \mathbb{C}^{(k+1) \times (k+1)}$  be upper triangular. We partition  $\mathbf{A}$  in the form

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_k & \mathbf{a}_k \\ \mathbf{0} & a_{k+1,k+1} \end{bmatrix}$$

and note that  $\mathbf{A}_k \in \mathbb{C}^{k \times k}$  is upper triangular. By Lemma 1.21  $\mathbf{A}$  is nonsingular if and only if  $\mathbf{A}_k$  and  $(a_{k+1,k+1})$  are nonsingular and in that case

$$\mathbf{A}^{-1} = \begin{bmatrix} \mathbf{A}_k^{-1} & \mathbf{c} \\ \mathbf{0} & a_{k+1,k+1}^{-1} \end{bmatrix},$$

for some  $\mathbf{c} \in \mathbb{C}^n$ . By the induction hypothesis  $\mathbf{A}_k$  is nonsingular if and only if the diagonal elements  $a_{11}, \dots, a_{kk}$  of  $\mathbf{A}_k$  are nonzero and in that case  $\mathbf{A}_k^{-1}$  is upper triangular with diagonal elements  $a_{ii}^{-1}$ ,  $i = 1, \dots, k$ . The result for  $\mathbf{A}$  follows.  $\square$

**Lemma 1.23 (Product of triangular matrices)**

The product  $\mathbf{C} = \mathbf{AB} = (c_{ij})$  of two upper (lower) triangular matrices  $\mathbf{A} = (a_{ij})$  and  $\mathbf{B} = (b_{ij})$  is upper (lower) triangular with diagonal elements  $c_{ii} = a_{ii}b_{ii}$  for all  $i$ .

*Proof.* Exercise.  $\square$

A matrix is **unit triangular** if it is triangular with 1's on the diagonal.

**Lemma 1.24 (Unit triangular matrices)**

For a unit upper (lower) triangular matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$ :

1.  $\mathbf{A}$  is nonsingular and the inverse is unit upper(lower) triangular.
2. The product of two unit upper (lower) triangular matrices is unit upper (lower) triangular.

*Proof.* 1. follows from Lemma 1.22, while Lemma 1.23 implies 2.  $\square$

## 1.6 Review Questions

**1.6.1** Define the second derivative matrix  $\mathbf{T}$ . Why is it nonsingular?

**1.6.2** Is a weakly diagonally dominant matrix nonsingular?

**1.6.3** Why do we not use the explicit inverse of  $\mathbf{T}$  to solve the linear system  $\mathbf{T}\mathbf{x} = \mathbf{b}$

**1.6.4** Show that a strictly diagonally dominant matrix is nonsingular.

**1.6.5** Does a tridiagonal matrix always have an LU factorization?

## Chapter 2

# LU Factorizations

Numerical methods for solving systems of linear equations are often based on writing a matrix as a product of simpler matrices. Such a **factorization** is useful if the corresponding matrix problem for each of the factors is simple to solve, and extra numerical stability issues are not introduced. Examples of factorizations were encountered in Chapter 1 and we saw how an LU factorization can be used to solve certain tridiagonal systems efficiently. Other factorizations based on unitary matrices will be considered later in this book.

In this chapter we consider the general theory of LU factorizations. We consider some related factorizations called symmetric LU or LDLT, and Cholesky. The latter can be used for symmetric positive matrices, and we give an introduction to positive definite and positive semidefinite matrices. We consider a matrix formulation of Gaussian elimination using Gauss transformations and permutation matrices leading to the PLU factorization of a matrix.

## 2.1 Algorithms for triangular systems

Recall that a matrix  $\mathbf{U}$  is upper triangular if  $u_{ij} = 0$  for  $i > j$ , and a matrix  $\mathbf{L}$  is lower triangular if  $l_{ij} = 0$  for  $i < j$ . If  $\mathbf{U}$  is upper triangular then  $\mathbf{U}^T$  is lower triangular.

A nonsingular triangular linear system  $\mathbf{Ax} = \mathbf{b}$  is easy to solve. By Lemma 1.22  $\mathbf{A}$  has nonzero diagonal elements. Consider first the lower triangular case. For  $n = 3$  the system is

$$\begin{bmatrix} a_{11} & 0 & 0 \\ a_{21} & a_{22} & 0 \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}.$$

$$\begin{bmatrix} a_{11} & 0 & 0 & 0 & 0 \\ a_{21} & a_{22} & 0 & 0 & 0 \\ 0 & a_{32} & a_{33} & 0 & 0 \\ 0 & 0 & a_{43} & a_{44} & 0 \\ 0 & 0 & 0 & a_{54} & a_{55} \end{bmatrix}, \quad \begin{bmatrix} a_{11} & 0 & 0 & 0 & 0 \\ a_{21} & a_{22} & 0 & 0 & 0 \\ a_{31} & a_{32} & a_{33} & 0 & 0 \\ 0 & a_{42} & a_{43} & a_{44} & 0 \\ 0 & 0 & a_{53} & a_{54} & a_{55} \end{bmatrix}$$

**Figure 2.1.** Lower triangular  $5 \times 5$  band matrices:  $d = 1$  (left) and  $d = 2$  (right).

From the first equation we find  $x_1 = b_1/a_{11}$ . Solving the second equation for  $x_2$  we obtain  $x_2 = (b_2 - a_{21}x_1)/a_{22}$ . Finally the third equation gives  $x_3 = (b_3 - a_{31}x_1 - a_{32}x_2)/a_{33}$ . This process is known as **forward substitution**. In general

$$x_k = (b_k - \sum_{j=1}^{k-1} a_{k,j}x_j)/a_{kk}, \quad k = 1, 2, \dots, n. \quad (2.1)$$

When  $\mathbf{A}$  is a lower triangular band matrix the number of arithmetic operations necessary to find  $\mathbf{x}$  can be reduced. Suppose  $\mathbf{A}$  is a lower triangular  $d$ -banded, so that  $a_{k,j} = 0$  for  $j \notin \{l_k, l_k + 1, \dots, k\}$  for  $k = 1, 2, \dots, n$ , and where  $l_k := \max(1, k - d)$ , see Figure 2.1. For a lower triangular  $d$ -band matrix the calculation in (2.1) can be simplified as follows

$$x_k = (b_k - \sum_{j=l_k}^{k-1} a_{k,j}x_j)/a_{kk}, \quad k = 1, 2, \dots, n. \quad (2.2)$$

Note that (2.2) reduces to (2.1) if  $d = n$ . Letting  $A(k, l_k : k - 1) * x(l_k : k - 1)$  denote the sum  $\sum_{j=l_k}^{k-1} a_{k,j}x_j$  we arrive at the following algorithm.

**Algorithm 2.1 (forwardsolve (row oriented))**

Given a nonsingular lower triangular  $d$ -banded matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  and  $\mathbf{b} \in \mathbb{C}^n$ .

An  $\mathbf{x} \in \mathbb{C}^n$  is computed so that  $\mathbf{Ax} = \mathbf{b}$ .

```

1 function x=forwardsolve(A,b,d)
2 n=length(b); x=b;
3 x(1)=b(1)/A(1,1);
4 for k=2:n
5     lk=max(1,k-d);
6     x(k)=(b(k)-A(k,lk:k-1)*x(lk:k-1))/A(k,k);
7 end

```

A system  $\mathbf{Ax} = \mathbf{b}$ , where  $\mathbf{A}$  is upper triangular must be solved by **back substitution** or 'bottom-up'. We first find  $x_n$  from the last equation and then



move upwards for the remaining unknowns. For an upper triangular  $d$ -banded matrix this leads to the following algorithm.

**Algorithm 2.2 (backsolve (row oriented))**

Given a nonsingular upper triangular  $d$ -banded matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  and  $\mathbf{b} \in \mathbb{C}^n$ . An  $\mathbf{x} \in \mathbb{C}^n$  is computed so that  $\mathbf{Ax} = \mathbf{b}$ .

```

1 function x=rbacksolve(A,b,d)
2 n=length(b); x=b;
3 x(n)=b(n)/A(n,n);
4 for k=n-1:-1:1
5     uk=min(n,k+d);
6     x(k)=(b(k)-A(k,k+1:uk)*x(k+1:uk))/A(k,k);
7 end

```

**Exercise 2.3 ( Column oriented forward- and backsolve)**

The initial "r" in the names of Algorithms 2.1,2.2 signals that these algorithms are row oriented. For each  $k$  we take the inner product of a part of a row with the already computed unknowns. In this exercise we develop column oriented vectorized versions of forward and backward substitution. Consider the system  $\mathbf{Ax} = \mathbf{b}$ , where  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is lower triangular. Suppose after  $k - 1$  steps of the algorithm we have a reduced system in the form

$$\begin{bmatrix} a_{k,k} & 0 & \cdots & 0 \\ a_{k+1,k} & a_{k+1,k+1} & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ a_{n,k} & & \cdots & a_{n \times n} \end{bmatrix} \begin{bmatrix} x_k \\ x_{k+1} \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_k \\ b_{k+1} \\ \vdots \\ b_n \end{bmatrix}.$$

This system is of order  $n - k + 1$ . The unknowns are  $x_k, \dots, x_n$ .

a) We see that  $x_k = b_k/a_{k,k}$  and eliminating  $x_k$  from the remaining equations we obtain a system of order  $n - k$  with unknowns  $x_{k+1}, \dots, x_n$

$$\begin{bmatrix} a_{k+1,k+1} & 0 & \cdots & 0 \\ a_{k+2,k+1} & a_{k+2,k+2} & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ a_{n,k+1} & & \cdots & a_{n \times n} \end{bmatrix} \begin{bmatrix} x_{k+1} \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_{k+1} \\ \vdots \\ b_n \end{bmatrix} - x_k \begin{bmatrix} a_{k+1,k} \\ \vdots \\ a_{n,k} \end{bmatrix}.$$

Thus at the  $k$ th step,  $k = 1, 2, \dots, n$  we set  $x_k = b_k/A(k, k)$  and update  $b$  as follows:

$$b(k+1:n) = b(k+1:n) - x(k) * A(k+1:n, k).$$

This leads to the following algorithm.

**Algorithm 2.4 (Forward solve (column oriented))**

Given a nonsingular lower triangular  $d$ -banded matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  and  $\mathbf{b} \in \mathbb{C}^n$ . An  $\mathbf{x} \in \mathbb{C}^n$  is computed so that  $\mathbf{Ax} = \mathbf{b}$ .

```

1 function x=cforwardsolve(A,b,d)
2 x=b; n=length(b);
3 for k=1:n-1
4     x(k)=b(k)/A(k,k); uk=min(n,k+d);
5     b(k+1:uk)=b(k+1:uk)-A(k+1:uk,k)*x(k);
6 end
7 x(n)=b(n)/A(n,n);
8 end

```

b) Suppose now  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is nonsingular, upper triangular,  $d$ -banded, and  $\mathbf{b} \in \mathbb{C}^n$ . Justify the following column oriented vectorized algorithms for solving  $\mathbf{Ax} = \mathbf{b}$ .

**Algorithm 2.5 (Backsolve (column oriented))**

Given a nonsingular upper triangular  $d$ -banded matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  and  $\mathbf{b} \in \mathbb{C}^n$ . An  $\mathbf{x} \in \mathbb{C}^n$  is computed so that  $\mathbf{Ax} = \mathbf{b}$ .

```

1 function x=cbacksolve(A,b,d)
2 x=b; n=length(b);
3 for k=n:-1:2
4     x(k)=b(k)/A(k,k); lk=max(1,k-d);
5     b(lk:k-1)=b(lk:k-1)-A(lk:k-1,k)*x(k);
6 end
7 x(1)=b(1)/A(1,1);
8 end

```

**Exercise 2.6 (Computing the inverse of a triangular matrix)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is a nonsingular triangular matrix with inverse  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_n]$ . The  $k$ th column  $\mathbf{b}_k$  of  $\mathbf{B}$  is the solution of the linear systems  $\mathbf{Ab}_k = \mathbf{e}_k$ . Write this system as a  $2 \times 2$  triangular block system and explain why we can find  $\mathbf{b}_k$  by solving the linear systems

$$\mathbf{A}(k:n, k:n)\mathbf{b}_k(k:n) = \mathbf{I}(k:n, k), \quad k = 1, \dots, n \quad \text{lower triangular}, \quad (2.3)$$

$$\mathbf{A}(1:k, 1:k)\mathbf{b}_k(1:k) = \mathbf{I}(1:k, k), \quad k = n, n-1, \dots, 1, \quad \text{upper triangular} \quad (2.4)$$

Is it possible to store the interesting part of  $\mathbf{b}_k$  in  $\mathbf{A}$  as soon as it is computed?

## 2.2 The LU Factorization

We say that  $\mathbf{A} = \mathbf{LU}$  is an **LU factorization** of  $\mathbf{A} \in \mathbb{C}^{n \times n}$  if  $\mathbf{L} \in \mathbb{C}^{n \times n}$  is lower triangular (**left triangular**) and  $\mathbf{U} \in \mathbb{C}^{n \times n}$  is upper triangular (**right triangular**).

### 2.2.1 The LU theorem

Consider finding  $\mathbf{L}$  and  $\mathbf{U}$ . Equating the  $i, j$  element in  $\mathbf{A}$  and the product  $\mathbf{LU}$ , and noting that  $l_{i,j} = 0$  for  $j > i$  and  $u_{i,j} = 0$  for  $i > j$ , we obtain an equation

$$a_{ij} = \sum_{k=1}^{\min(i,j)} l_{ik}u_{kj}, \quad i, j = 1, 2, \dots, n \quad (2.5)$$

involving the unknown elements in  $\mathbf{L}$  and  $\mathbf{U}$ . This is an underdetermined system of  $n^2$  equations in  $n^2 + n$  unknowns. One way to reduce the number of unknowns is to require that one of the triangular matrices should be **unit triangular**, i. e., have ones on the diagonal. Other scalings of the diagonals are also possible, see Section 2.6. Choosing  $\mathbf{U}$  to be unit triangular is sometimes known as a **Crout factorization**.



Henry Jensen, 1915-1974 (left), Prescott Durand Crout, 1907-1984. Jensen worked on LU factorizations. His name is also associated with a very useful inequality (cf. Theorem 7.37).

For our discussion we will assume that  $\mathbf{L}$  is unit triangular. Three things can happen. An LU factorization exists and is unique, it exists, but it is not unique, or it does not exist. The following  $2 \times 2$  example illustrates this.

**Example 2.7 (LU of  $2 \times 2$  matrix)**

Let  $a, b, c, d \in \mathbb{C}$ . An LU factorization of  $\mathbf{A} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$  must satisfy the equations

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ l_1 & 1 \end{bmatrix} \begin{bmatrix} u_1 & u_3 \\ 0 & u_2 \end{bmatrix} = \begin{bmatrix} u_1 & u_3 \\ l_1 u_1 & l_1 u_3 + u_2 \end{bmatrix}$$

for the unknowns  $l_1$  in  $\mathbf{L}$  and  $u_1, u_2, u_3$  in  $\mathbf{U}$ . The equations are

$$u_1 = a, \quad u_3 = b, \quad l_1 a = c, \quad u_2 = d - l_1 b.$$

These equations do not always have a solution. Indeed, the main problem is the nonlinear equation  $l_1 a = c$ . There are three cases

1.  $a \neq 0$ : The matrix has a unique LU factorization with  $l_1 = c/a$ .
2.  $a = 0, c \neq 0$ : No LU factorization exists.
3.  $a = c = 0$ : The LU factorization exists, but it is not unique. Any value for  $l_1$  can be used.

Of the four matrices

$$\mathbf{A}_1 := \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}, \quad \mathbf{A}_2 := \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{A}_3 := \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{A}_4 := \begin{bmatrix} 0 & 1 \\ 0 & 2 \end{bmatrix}.$$

$\mathbf{A}_1$  has a unique LU factorization,  $\mathbf{A}_2$  has no LU factorization,  $\mathbf{A}_3$  has a unique LU factorization even if it is singular, and  $\mathbf{A}_4$  has an LU factorization, but it is not unique.

### Example 2.8 (LU of $3 \times 3$ matrices)

The matrix

$$\mathbf{A} := \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 2 \\ 1 & 1 & 3 \end{bmatrix}$$

has an LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{U}$ , with

$$\mathbf{L} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & y & 1 \end{bmatrix}, \quad \mathbf{U} = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 2-y \end{bmatrix}.$$

It is not unique since  $\mathbf{A} = \mathbf{L}\mathbf{U}$  for any  $y \in \mathbb{C}$ .

To characterize matrices with a unique LU factorization we first give a definition.

### Definition 2.9 (Principal submatrix)

For  $k = 1, \dots, n$  the matrices  $\mathbf{A}_{[k]} \in \mathbb{C}^{k \times k}$  given by

$$\mathbf{A}_{[k]} := \mathbf{A}(1:k, 1:k) = \begin{bmatrix} a_{11} & \cdots & a_{k1} \\ \vdots & & \vdots \\ a_{k1} & \cdots & a_{kk} \end{bmatrix}$$

are called the **leading principal submatrices** of  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . More generally, a matrix  $\mathbf{B} \in \mathbb{C}^{k \times k}$  is called a **principal submatrix** of  $\mathbf{A}$  if  $\mathbf{B} = \mathbf{A}(\mathbf{r}, \mathbf{r})$ , where  $\mathbf{r} = [r_1, \dots, r_k]$  for some  $1 \leq r_1 < \dots < r_k \leq n$ . Thus,

$$b_{i,j} = a_{r_i, r_j}, \quad i, j = 1, \dots, k.$$

The determinant of a (leading) principal submatrix is called a **(leading) principal minor**.

A principal submatrix is leading if  $r_j = j$  for  $j = 1, \dots, k$ . Also a principal submatrix is special in that it uses the same rows and columns of  $\mathbf{A}$ . For  $k = 1$  The only principal submatrices of order  $k = 1$  are the diagonal elements of  $\mathbf{A}$ .

**Example 2.10 (Principal submatrices)**

The principal submatrices of  $\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$  are

$$[1], [5], [9], \begin{bmatrix} 1 & 2 \\ 4 & 5 \end{bmatrix}, \begin{bmatrix} 1 & 3 \\ 7 & 9 \end{bmatrix}, \begin{bmatrix} 5 & 6 \\ 8 & 9 \end{bmatrix}, \mathbf{A}.$$

The leading principal submatrices are

$$[1], \begin{bmatrix} 1 & 2 \\ 4 & 5 \end{bmatrix}, \mathbf{A}.$$

In preparation for the main theorem about LU factorization we prove a simple lemma.

**Lemma 2.11 (LU of leading principal sub matrices)**

Suppose  $\mathbf{A} = \mathbf{L}\mathbf{U}$  is an LU factorization of  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . For  $k = 1, \dots, n$  let  $\mathbf{A}_{[k]}, \mathbf{L}_{[k]}, \mathbf{U}_{[k]}$  be the leading principal submatrices of  $\mathbf{A}, \mathbf{L}, \mathbf{U}$ , respectively. Then  $\mathbf{A}_{[k]} = \mathbf{L}_{[k]}\mathbf{U}_{[k]}$  is an LU factorization of  $\mathbf{A}_{[k]}$  for  $k = 1, \dots, n$ .

*Proof.* For  $k = 1, \dots, n - 1$  we partition  $\mathbf{A} = \mathbf{L}\mathbf{U}$  as follows:

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{[k]} & \mathbf{B}_k \\ \mathbf{C}_k & \mathbf{F}_k \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{[k]} & \mathbf{0} \\ \mathbf{M}_k & \mathbf{N}_k \end{bmatrix} \begin{bmatrix} \mathbf{U}_{[k]} & \mathbf{S}_k \\ \mathbf{0} & \mathbf{T}_k \end{bmatrix} = \mathbf{L}\mathbf{U}, \quad (2.6)$$

where  $\mathbf{F}_k, \mathbf{N}_k, \mathbf{T}_k \in \mathbb{C}^{(n-k) \times (n-k)}$ . Using block multiplication we find  $\mathbf{A}_{[k]} = \mathbf{L}_{[k]}\mathbf{U}_{[k]}$ . Since  $\mathbf{L}_{[k]}$  is unit lower triangular and  $\mathbf{U}_{[k]}$  is upper triangular this is an LU factorization of  $\mathbf{A}_{[k]}$ .  $\square$

The following theorem give a necessary and sufficient condition for existence of a unique LU factorization.

**Theorem 2.12 (LU theorem)**

A square matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  has a unique LU factorization if and only if the leading principal submatrices  $\mathbf{A}_{[k]}$  of  $\mathbf{A}$  are nonsingular for  $k = 1, \dots, n - 1$ .

*Proof.* Suppose  $\mathbf{A}_{[k]}$  is nonsingular for  $k = 1, \dots, n - 1$ . We use induction on  $n$  to show that  $\mathbf{A}$  has a unique LU factorization. The result is clearly true for  $n = 1$ , since the unique LU factorization of a 1-by-1 matrix is  $[a_{11}] = [1][a_{11}]$ .

Suppose that  $\mathbf{A}_{[n-1]}$  has a unique LU factorization  $\mathbf{A}_{[n-1]} = \mathbf{L}_{n-1}\mathbf{U}_{n-1}$ , and that  $\mathbf{A}_{[1]}, \dots, \mathbf{A}_{[n-1]}$  are nonsingular. By block multiplication

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{[n-1]} & \mathbf{b} \\ \mathbf{c}^T & a_{nn} \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{n-1} & \mathbf{0} \\ \mathbf{m}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{U}_{n-1} & \mathbf{s} \\ 0 & t \end{bmatrix} = \mathbf{LU}, \quad (2.7)$$

if and only if  $\mathbf{m}, \mathbf{s} \in \mathbb{C}^{n-1}$  and  $t \in \mathbb{C}$  satisfy  $\mathbf{b} = \mathbf{L}_{n-1}\mathbf{s}$ ,  $\mathbf{c}^T = \mathbf{m}^T\mathbf{U}_{n-1}$ , and  $a_{nn} = \mathbf{m}^T\mathbf{s} + t$ . Since  $\mathbf{A}_{[n-1]}$  is nonsingular it follows that  $\mathbf{L}_{n-1}$  and  $\mathbf{U}_{n-1}$  are nonsingular and therefore  $\mathbf{m}, \mathbf{s}$ , and  $t$  are uniquely given. Thus (2.7) gives a unique LU factorization of  $\mathbf{A}$ .

Conversely, suppose  $\mathbf{A}$  has a unique LU factorization  $\mathbf{A} = \mathbf{LU}$ . By Lemma 2.11  $\mathbf{A}_{[k]} = \mathbf{L}_{[k]}\mathbf{U}_{[k]}$  is an LU factorization of  $\mathbf{A}_{[k]}$  for  $k = 1, \dots, n-1$ . Suppose  $\mathbf{A}_{[k]}$  is singular for some  $k \leq n-1$ . We will show that this leads to a contradiction. Let  $k$  be the smallest integer so that  $\mathbf{A}_{[k]}$  is singular. Since  $\mathbf{A}_{[j]}$  is nonsingular for  $j \leq k-1$  it follows from what we have already shown that  $\mathbf{A}_{[k]} = \mathbf{L}_{[k]}\mathbf{U}_{[k]}$  is the unique LU factorization of  $\mathbf{A}_{[k]}$ . The matrix  $\mathbf{U}_{[k]}$  is singular since  $\mathbf{A}_{[k]}$  is singular and  $\mathbf{L}_{[k]}$  is nonsingular. By block multiplication in (2.6) we have  $\mathbf{C}_k = \mathbf{M}_k\mathbf{U}_{[k]}$  or  $\mathbf{U}_{[k]}^T\mathbf{M}_k^T = \mathbf{C}_k^T$ . This can be written as  $n-k$  linear systems for the columns of  $\mathbf{M}_k^T$ . By assumption  $\mathbf{M}_k^T$  exists, but since  $\mathbf{U}_{[k]}^T$  is singular  $\mathbf{M}_k$  is not unique, a contradiction.  $\square$

A matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  can have an LU factorization even if  $\mathbf{A}_{[k]}$  is singular for some  $k < n$ . By Theorem 2.12 such an LU factorization cannot be unique.

**Remark 2.13 (LU of upper triangular matrix)**

An LU factorization of an upper triangular matrix  $\mathbf{A}$  is  $\mathbf{A} = \mathbf{IA}$  so it always exists even if  $\mathbf{A}$  has zeros somewhere on the diagonal. By Lemma 1.22, if some  $a_{kk}$  is zero then  $\mathbf{A}_{[k]}$  is singular and the LU factorization cannot be unique. In particular, for the zero matrix any unit lower triangular matrix can be used as  $\mathbf{L}$  in an LU factorization.

**Remark 2.14 (PLU factorization)**

We have shown that a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  has a unique LU factorization if and only if the leading principle submatrices  $\mathbf{A}_{[k]}$  are nonsingular for  $k = 1, \dots, n-1$ . This condition seems fairly restrictive. However, for a nonsingular matrix  $\mathbf{A}$  there always is a permutation of the rows so that the permuted matrix has an LU factorization. We obtain a factorization of the form  $\mathbf{P}^T\mathbf{A} = \mathbf{LU}$  or equivalently  $\mathbf{A} = \mathbf{PLU}$ , where  $\mathbf{P}$  is a permutation matrix,  $\mathbf{L}$  is unit lower triangular, and  $\mathbf{U}$  is upper triangular. We call this a **PLU factorization** of  $\mathbf{A}$ . (Cf. Theorem 2.60).

**Exercise 2.15 (Row interchange)**

Show that  $\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$  has a unique LU factorization. Note that we have only interchanged rows in Example 2.7.

**Exercise 2.16 (LU of singular matrix)**

Find an LU factorization of the singular matrix  $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ . Is it unique?

**Exercise 2.17 (LU and determinant)**

Suppose  $\mathbf{A}$  has an LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{U}$ . Show that  $\det(\mathbf{A}_{[k]}) = u_{11}u_{22} \cdots u_{kk}$  for  $k = 1, \dots, n$ .

**Exercise 2.18 (Diagonal elements in U)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  and  $\mathbf{A}_{[k]}$  is nonsingular for  $k = 1, \dots, n-1$ . Use Exercise 2.17 to show that the diagonal elements  $u_{kk}$  in the LU factorization are

$$u_{11} = a_{11}, \quad u_{kk} = \frac{\det(\mathbf{A}_{[k]})}{\det(\mathbf{A}_{[k-1]})}, \quad \text{for } k = 2, \dots, n. \quad (2.8)$$

**2.2.2 Operation count**

It is useful to have a number which indicates the amount of work an algorithm requires. In this book we measure this by estimating the total number of arithmetic operations. We count both additions, subtractions, multiplications and divisions, but not work on indices. As an example it is shown below that the calculation to find the LU factorization of a full matrix of order  $n$  is exactly

$$N_{LU} := \frac{2}{3}n^3 - \frac{1}{2}n^2 - \frac{1}{6}n. \quad (2.9)$$

We are only interested in this number when  $n$  is large and for such  $n$  the term  $\frac{2}{3}n^3$  dominates. We therefore regularly ignore lower order terms and use **number of arithmetic operations** both for the exact count and for the highest order term. We also say more loosely that the the number of operations is  $O(n^3)$ . We will use the number of operations counted in one of these ways as a measure of the **complexity of an algorithm** and say that the complexity of LU factorization of a full matrix is  $O(n^3)$  or more precisely  $\frac{2}{3}n^3$ .

We will compare the number of arithmetic operations of many algorithms with the number of arithmetic operations of LU factorization and define for  $n \in \mathbb{N}$  the number  $G_n$ <sup>7</sup> as follows:

**Definition 2.19** ( $G_n := \frac{2}{3}n^3$ )

We define  $G_n := \frac{2}{3}n^3$ .

The complexity of solving a system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is a full upper or lower triangular matrix is easily shown to be exactly  $n^2$ . This number is reduced to  $n^2 - n$  if  $\mathbf{A}$  has ones on the diagonal.

<sup>7</sup>It can be shown that the complexity of Gaussian elimination is also equal to  $G_n$ .

Consider now finding the total number of arithmetic operations,  $N_{LU}$  for LU factorization. Suppose in (2.7) that  $k := n - 1$  and  $\mathbf{L}_k$  and  $\mathbf{U}_k$  are already computed. To find  $\mathbf{L}$  and  $\mathbf{U}$  we have to solve the triangular systems  $\mathbf{L}_k \mathbf{s} = \mathbf{b}$ ,  $\mathbf{U}_k^T \mathbf{m} = \mathbf{c}$ , and then  $t = a_{k+1,k+1} - \mathbf{m}^T \mathbf{s}$ . Since  $\mathbf{L}_k$  is unit lower triangular and  $\mathbf{U}_k^T$  is lower triangular of order  $k$  this requires  $k(k-1)$ ,  $k^2$ , and  $2k$  operations for  $\mathbf{s}$ ,  $\mathbf{m}$ , and  $t$ , respectively, a total of  $k(2k+1)$  operations. Taking also  $\mathbf{L}_k$  and  $\mathbf{U}_k$  into consideration we obtain

$$N_{LU} = \sum_{k=1}^{n-1} k(2k+1) = 2 \sum_{k=1}^{n-1} k^2 + \sum_{k=1}^{n-1} k = \frac{2}{3}n(n-1)(n-\frac{1}{2}) + \frac{1}{2}n(n-1)$$

which equals the number in (2.9).

There is a quick way to arrive at the estimate  $2n^3/3$ . We only consider the arithmetic operations contributing to the leading term (the inner loops). Then we replace sums by integrals letting the summation indices be continuous variables and adjust limits of integration in an insightful way to simplify the calculation. We find

$$N_{LU} = \sum_{k=1}^{n-1} k(2k+1) \approx \sum_{k=1}^{n-1} 2k^2 \approx \int_1^{n-1} 2k^2 dk \approx \int_0^n 2k^2 dk = G_n.$$

We see that LU factorization is an  $O(n^3)$  process while solving a triangular system requires  $O(n^2)$  arithmetic operations. Thus, if  $n = 10^6$  and one arithmetic operation requires  $c = 10^{-12}$  seconds of computing time then  $cn^3 = 10^6$  seconds  $\approx 278$  hours and  $cn^2 = 1$  second, giving dramatic differences in computing time.

### Exercise 2.20 (Finite sums of integers)

Use induction on  $m$ , or some other method, to show that

$$1 + 2 + \cdots + m = \frac{1}{2}m(m+1), \quad (2.10)$$

$$1^2 + 2^2 + \cdots + m^2 = \frac{1}{3}m(m + \frac{1}{2})(m+1), \quad (2.11)$$

$$1 + 3 + 5 + \cdots + 2m - 1 = m^2, \quad (2.12)$$

$$1 * 2 + 2 * 3 + 3 * 4 + \cdots + (m-1)m = \frac{1}{3}(m-1)m(m+1). \quad (2.13)$$

### Exercise 2.21 (Operations)

To solve an upper triangular linear system by back substitution takes  $n^2$  arithmetic operations. Show that the number of arithmetic operations in (2.4) is  $\frac{1}{3}n(n+\frac{1}{2})(n+1) \approx \frac{1}{2}G_n$ .



**Exercise 2.22 (Multiplying triangular matrices)**

Show that the matrix multiplication  $\mathbf{AB}$  can be done in  $\frac{1}{3}n(2n^2 + 1) \approx G_n$  arithmetic operations when  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is lower triangular and  $\mathbf{B} \in \mathbb{R}^{n \times n}$  is upper triangular. What about  $\mathbf{BA}$ ?

**2.3 The Symmetric LU Factorization**

We consider next LU factorization of a real symmetric matrix.

**Definition 2.23 (Symmetric LU)**

Suppose  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . A factorization  $\mathbf{A} = \mathbf{LDL}^T$ , where  $\mathbf{L}$  is unit lower triangular and  $\mathbf{D}$  is diagonal is called a **symmetric LU factorization** or an **LDLT factorization** of  $\mathbf{A}$ .

A matrix which has a symmetric LU factorization must be symmetric since  $\mathbf{A}^T = (\mathbf{LDL}^T)^T = \mathbf{LDL}^T = \mathbf{A}$ .

**Example 2.24 ( $2 \times 2$  symmetric LU)**

Let  $a, b, c \in \mathbb{R}$ . A symmetric LU factorization of  $\mathbf{A} := \begin{bmatrix} a & b \\ b & c \end{bmatrix}$  must satisfy the equations

$$\begin{bmatrix} a & b \\ b & c \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ l_1 & 1 \end{bmatrix} \begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix} \begin{bmatrix} 1 & l_1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} d_1 & d_1 l_1 \\ l_1 d_1 & l_1^2 d_1 + d_2 \end{bmatrix}$$

for the unknowns  $l_1$  in  $\mathbf{L}$  and  $d_1, d_2$  in  $\mathbf{D}$ . The equations are

$$d_1 = a, \quad l_1 a = b, \quad d_2 = c - a l_1^2.$$

As in the nonsymmetric case the main problem is the nonlinear equation. Again there are three cases

1.  $a \neq 0$ : The matrix has a unique symmetric LU factorization with  $l_1 = b/a$ .
2.  $a = 0, b \neq 0$ : No symmetric LU factorization exists.
3.  $a = b = 0$ : The LU factorization exists, but it is not unique. Any value for  $l_1$  can be used.

Consider the four matrices

$$\mathbf{A}_1 := \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}, \quad \mathbf{A}_2 := \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{A}_3 := \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{A}_4 := \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix}.$$

Then the symmetric LU factorization is unique for  $\mathbf{A}_1$  and  $\mathbf{A}_3$ , is not unique for  $\mathbf{A}_4$  and does not exist for  $\mathbf{A}_2$ .

In view of this example it might come as no surprise that Theorem 2.12 carries over to the symmetric case. Again we start with an lemma.

**Lemma 2.25 (Symmetric LU of leading principal sub matrices)**

Suppose  $\mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{L}^T$  is a symmetric LU factorization of  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . For  $k = 1, \dots, n$  let  $\mathbf{A}_{[k]}, \mathbf{L}_{[k]}, \mathbf{D}_{[k]}$  be the leading principal submatrices of  $\mathbf{A}, \mathbf{L}, \mathbf{D}$ , respectively. Then  $\mathbf{A}_{[k]} = \mathbf{L}_{[k]}\mathbf{D}_{[k]}\mathbf{L}_{[k]}^T$  is a symmetric LU factorization of  $\mathbf{A}_{[k]}$  for  $k = 1, \dots, n$ .

*Proof.* For  $k = 1, \dots, n - 1$  we partition  $\mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{L}^T$  as follows:

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{[k]} & \mathbf{B}_k \\ \mathbf{C}_k & \mathbf{F}_k \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{[k]} & \mathbf{0} \\ \mathbf{M}_k & \mathbf{N}_k \end{bmatrix} \begin{bmatrix} \mathbf{D}_{[k]} & \mathbf{0} \\ \mathbf{0} & \mathbf{E}_k \end{bmatrix} \begin{bmatrix} \mathbf{L}_{[k]}^T & \mathbf{M}_k^T \\ \mathbf{0} & \mathbf{N}_k^T \end{bmatrix} = \mathbf{L}\mathbf{D}\mathbf{L}^T, \quad (2.14)$$

where  $\mathbf{F}_k, \mathbf{N}_k, \mathbf{E}_k \in \mathbb{R}^{n-k, n-k}$ . Block multiplication gives  $\mathbf{A}_{[k]} = \mathbf{L}_{[k]}\mathbf{D}_{[k]}\mathbf{L}_{[k]}^T$ . Since  $\mathbf{L}_{[k]}$  is unit lower triangular and  $\mathbf{D}_{[k]}$  is diagonal this is a symmetric LU factorization of  $\mathbf{A}_{[k]}$ .  $\square$

**Theorem 2.26 (Symmetric LU theorem)**

The matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  has a unique symmetric LU factorization if and only if  $\mathbf{A} = \mathbf{A}^T$  and  $\mathbf{A}_{[k]}$  is nonsingular for  $k = 1, \dots, n - 1$ .

*Proof.* If  $\mathbf{A}$  is nonsingular then  $\mathbf{D}$  is nonsingular and it can be shown that the theorem is a simple corollary of the LU theorem. To prove the general case we repeat the proof of Theorem 2.12 incorporating the necessary changes. Suppose  $\mathbf{A}^T = \mathbf{A}$  and that  $\mathbf{A}_{[k]}$  is nonsingular for  $k = 1, \dots, n - 1$ . Note that  $\mathbf{A}_{[k]}^T = \mathbf{A}_{[k]}$  for  $k = 1, \dots, n$ . We use induction on  $n$  to show that  $\mathbf{A}$  has a unique symmetric LU factorization. The result is clearly true for  $n = 1$ , since the unique symmetric LU factorization of a 1-by-1 matrix is  $[a_{11}] = [1][a_{11}][1]$ . Suppose that  $\mathbf{A}_{[n-1]}$  has a unique symmetric LU factorization  $\mathbf{A}_{[n-1]} = \mathbf{L}_{n-1}\mathbf{D}_{n-1}\mathbf{L}_{n-1}^T$ , and that  $\mathbf{A}_{[1]}, \dots, \mathbf{A}_{[n-1]}$  are nonsingular. By block multiplication

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{[n-1]} & \mathbf{b} \\ \mathbf{b}^T & a_{nn} \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{n-1} & \mathbf{0} \\ \mathbf{x}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{D}_{n-1} & \mathbf{0} \\ \mathbf{0} & d_{nn} \end{bmatrix} \begin{bmatrix} \mathbf{L}_{n-1}^T & \mathbf{x} \\ \mathbf{0}^T & 1 \end{bmatrix}, \quad (2.15)$$

if and only if  $\mathbf{b} = \mathbf{L}_{n-1}\mathbf{D}_{n-1}\mathbf{x}$  and  $a_{nn} = d_{nn} + \mathbf{x}^T\mathbf{D}_{n-1}\mathbf{x}$ . Thus we obtain a symmetric LU factorization of  $\mathbf{A}$  that is unique since  $\mathbf{L}_{n-1}$  and  $\mathbf{D}_{n-1}$  are nonsingular.

For the converse we use Lemma 2.25 in the same way as Lemma 2.11 was used to prove Theorem 2.12.  $\square$

The number of arithmetic operations for the symmetric LU factorization is approximately  $\frac{1}{2}G_n$ , half the number of operations needed for the LU factorization.

For in the LU factorization we needed to solve two triangular systems to find the vectors  $\mathbf{s}$  and  $\mathbf{m}$ , while only one such system is needed to find  $\mathbf{x}$  in the symmetric case (2.15). The work to find  $d_{nn}$  is  $O(n)$  and does not contribute to the highest order term.

## 2.4 Block LU factorization

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is a block matrix of the form

$$\mathbf{A} := \begin{bmatrix} \mathbf{A}_{11} & \cdots & \mathbf{A}_{1m} \\ \vdots & & \vdots \\ \mathbf{A}_{m1} & \cdots & \mathbf{A}_{mm} \end{bmatrix}, \quad (2.16)$$

where each diagonal block  $\mathbf{A}_{ii}$  is square. We call the factorization

$$\mathbf{A} = \mathbf{L}\mathbf{U} = \begin{bmatrix} \mathbf{I} & & & \\ \mathbf{L}_{21} & \mathbf{I} & & \\ \vdots & & \ddots & \\ \mathbf{L}_{m1} & \cdots & \mathbf{L}_{m,m-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{U}_{11} & & \cdots & \mathbf{U}_{1m} \\ & \mathbf{U}_{21} & \cdots & \mathbf{U}_{2m} \\ & & \ddots & \vdots \\ & & & \mathbf{U}_{mm} \end{bmatrix} \quad (2.17)$$

a **block LU factorization of  $\mathbf{A}$** . Here the  $i$ th diagonal blocks  $\mathbf{I}$  and  $\mathbf{U}_{ii}$  in  $\mathbf{L}$  and  $\mathbf{U}$  have the same size as  $\mathbf{A}_{ii}$ , the  $i$ th diagonal block in  $\mathbf{A}$ .

The results for elementwise LU factorization carry over to block LU factorization as follows.

### Theorem 2.27 (Block LU theorem)

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is a block matrix of the form (2.16). Then  $\mathbf{A}$  has a unique block LU factorization (2.17) if and only if the **leading principal block submatrices**

$$\mathbf{A}_{\{k\}} := \begin{bmatrix} \mathbf{A}_{11} & \cdots & \mathbf{A}_{1k} \\ \vdots & & \vdots \\ \mathbf{A}_{k1} & \cdots & \mathbf{A}_{kk} \end{bmatrix}$$

are nonsingular for  $k = 1, \dots, m-1$ .

*Proof.* Suppose  $\mathbf{A}_{\{k\}}$  is nonsingular for  $k = 1, \dots, m-1$ . Following the proof in Theorem 2.12 suppose  $\mathbf{A}_{\{m-1\}}$  has a unique block LU factorization  $\mathbf{A}_{\{m-1\}} = \mathbf{L}_{\{m-1\}}\mathbf{U}_{\{m-1\}}$ , and that  $\mathbf{A}_{\{1\}}, \dots, \mathbf{A}_{\{m-1\}}$  are nonsingular. Then  $\mathbf{L}_{\{m-1\}}$  and  $\mathbf{U}_{\{m-1\}}$  are nonsingular and

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} \mathbf{A}_{\{m-1\}} & \mathbf{B} \\ \mathbf{C}^T & \mathbf{A}_{mm} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{L}_{\{m-1\}} & \mathbf{0} \\ \mathbf{C}^T \mathbf{U}_{\{m-1\}}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{U}_{\{m-1\}} & \mathbf{L}_{\{m-1\}}^{-1} \mathbf{B} \\ 0 & \mathbf{A}_{mm} - \mathbf{C}^T \mathbf{U}_{\{m-1\}}^{-1} \mathbf{L}_{\{m-1\}}^{-1} \mathbf{B} \end{bmatrix}, \end{aligned} \quad (2.18)$$

is a block LU factorization of  $\mathbf{A}$ . It is unique by derivation. Conversely, suppose  $\mathbf{A}$  has a unique block LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{U}$ . Then as in Lemma 2.11 it is easily seen that  $\mathbf{A}_{\{k\}} = \mathbf{L}_{\{k\}}\mathbf{U}_{\{k\}}$  is the unique block LU factorization of  $\mathbf{A}_{[k]}$  for  $k = 1, \dots, m$ . The rest of the proof is similar to the proof of Theorem 2.12.  $\square$

**Remark 2.28 (Comparing LU and block LU)**

*The number of arithmetic operations for the block LU factorization is the same as for the ordinary LU factorization. An advantage of the block method is that it combines many of the operations into matrix operations.*

**Remark 2.29 (A block LU is not an LU)**

*Note that (2.17) is not an LU factorization of  $\mathbf{A}$  since the  $\mathbf{U}_{ii}$ 's are not upper triangular in general. To relate the block LU factorization to the usual LU factorization we assume that each  $\mathbf{U}_{ii}$  has an LU factorization  $\mathbf{U}_{ii} = \tilde{\mathbf{L}}_{ii}\tilde{\mathbf{U}}_{ii}$ . Then  $\mathbf{A} = \hat{\mathbf{L}}\hat{\mathbf{U}}$ , where  $\hat{\mathbf{L}} := \mathbf{L} \operatorname{diag}(\tilde{\mathbf{L}}_{ii})$  and  $\hat{\mathbf{U}} := \operatorname{diag}(\tilde{\mathbf{L}}_{ii}^{-1})\mathbf{U}$ , and this is an ordinary LU factorization of  $\mathbf{A}$ .*

**Exercise 2.30 (Making block LU into LU)**

*Show that  $\hat{\mathbf{L}}$  is unit lower triangular and  $\hat{\mathbf{U}}$  is upper triangular.*

## 2.5 Positive Definite and Semidefinite Matrices

Symmetric positive definite matrices occur often in scientific computing. In this section we study some properties of positive definite matrices. We study only real matrices, but consider both the symmetric and nonsymmetric case.

### 2.5.1 Definitions and examples

Suppose  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is a square matrix. The function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$$

is called a **quadratic form**. We say that  $\mathbf{A}$  is

- (i) **positive definite** if  $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$  for all nonzero  $\mathbf{x} \in \mathbb{R}^n$ .
- (ii) **positive semidefinite** if  $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$  for all  $\mathbf{x} \in \mathbb{R}^n$ .
- (iii) **negative (semi)definite** if  $-\mathbf{A}$  is positive (semi)definite.

- (iv) **symmetric positive (semi)definite** if  $\mathbf{A}$  is symmetric in addition to being positive (semi)definite.
- (v) **symmetric negative (semi)definite** if  $\mathbf{A}$  is symmetric in addition to being negative (semi)definite.

We observe the following.

- A matrix is positive definite if it is positive semidefinite and in addition

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = 0 \Rightarrow \mathbf{x} = \mathbf{0}. \quad (2.19)$$

- A positive definite matrix must be nonsingular. Indeed, if  $\mathbf{A} \mathbf{x} = \mathbf{0}$  for some  $\mathbf{x} \in \mathbb{R}^n$  then  $\mathbf{x}^T \mathbf{A} \mathbf{x} = 0$  which by (2.19) implies that  $\mathbf{x} = \mathbf{0}$ .

The zero-matrix is symmetric positive semidefinite, while the unit matrix is symmetric positive definite.

We considered only real valued vectors  $\mathbf{x}$  above. For symmetric matrices we have:

**Lemma 2.31 (Quadratic form with  $\mathbf{x} \in \mathbb{C}^n$ )**

If  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is symmetric positive definite then  $\mathbf{x}^* \mathbf{A} \mathbf{x} > 0$  for all nonzero  $\mathbf{x} \in \mathbb{C}^n$ .

*Proof.* Suppose  $\mathbf{x} := \mathbf{y} + i\mathbf{z}$  is nonzero with  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  and  $i := \sqrt{-1}$ . Since  $\mathbf{A}$  is symmetric we find  $\mathbf{x}^* \mathbf{A} \mathbf{x} = (\mathbf{y} - i\mathbf{z})^T \mathbf{A} (\mathbf{y} + i\mathbf{z}) = \mathbf{y}^T \mathbf{A} \mathbf{y} + \mathbf{z}^T \mathbf{A} \mathbf{z}$ . This is positive since  $\mathbf{A}$  is positive definite and at least one of the vectors  $\mathbf{y}, \mathbf{z}$  is nonzero.  $\square$

**Example 2.32 ( $2 \times 2$  positive definite)**

The family of matrices

$$\mathbf{A}[a] := \begin{bmatrix} 2 & 2-a \\ a & 1 \end{bmatrix}, \quad a \in \mathbb{R}$$

is positive definite for any  $a \in \mathbb{R}$ . Indeed for any nonzero  $\mathbf{x} = [x_1, x_2]^T \in \mathbb{R}^2$

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = 2x_1^2 + (2-a)x_1x_2 + ax_2x_1 + x_2^2 = x_1^2 + (x_1 + x_2)^2 > 0.$$

**Lemma 2.33 ( $\mathbf{T}$  is symmetric positive definite)**

The second derivative matrix  $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{n \times n}$  is symmetric positive definite.

**Proof.** Clearly  $\mathbf{T}$  is symmetric. For any  $\mathbf{x} \in \mathbb{R}^n$

$$\begin{aligned} \mathbf{x}^T \mathbf{T} \mathbf{x} &= 2 \sum_{i=1}^n x_i^2 - \sum_{i=1}^{n-1} x_i x_{i+1} - \sum_{i=2}^n x_{i-1} x_i \\ &= \sum_{i=1}^{n-1} x_i^2 - 2 \sum_{i=1}^{n-1} x_i x_{i+1} + \sum_{i=1}^{n-1} x_{i+1}^2 + x_1^2 + x_n^2 \\ &= x_1^2 + x_n^2 + \sum_{i=1}^{n-1} (x_{i+1} - x_i)^2. \end{aligned}$$

Thus  $\mathbf{x}^T \mathbf{T} \mathbf{x} \geq 0$  and if  $\mathbf{x}^T \mathbf{T} \mathbf{x} = 0$  then  $x_1 = x_n = 0$  and  $x_i = x_{i+1}$  for  $i = 1, \dots, n-1$  which implies that  $\mathbf{x} = \mathbf{0}$ . Hence  $\mathbf{T}$  is positive definite.  $\square$

Symmetric positive definite matrices is important in nonlinear optimization.

**Example 2.34 (Gradient and hessian)**

Consider (cf. (C.1)) the gradient  $\nabla f$  and hessian  $Hf$  of a function  $f : \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}$

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n, \quad Hf(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_1} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \vdots & & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_n} \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

We assume that  $f$  has continuous first and second partial derivatives on  $\Omega$ .

Under suitable conditions on the domain  $\Omega$  it is shown in advanced calculus texts that if  $\nabla f(\mathbf{x}) = \mathbf{0}$  and  $Hf(\mathbf{x})$  is symmetric positive definite then  $\mathbf{x}$  is a local minimum for  $f$ . This can be shown using the second-order Taylor expansion (C.2). Moreover,  $\mathbf{x}$  is a local maximum if  $\nabla f(\mathbf{x}) = \mathbf{0}$  and  $Hf(\mathbf{x})$  is negative definite.

**Theorem 2.35 (A general criterium)**

Let  $m, n$  be positive integers. If  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is (symmetric) positive semidefinite then  $\mathbf{B} := \mathbf{X}^T \mathbf{A} \mathbf{X} \in \mathbb{R}^{m \times m}$  is (symmetric) positive semidefinite for any  $\mathbf{X} \in \mathbb{R}^{n \times m}$ . If in addition  $\mathbf{A}$  is (symmetric) positive definite and  $\mathbf{X}$  has linearly independent columns then  $\mathbf{B}$  is (symmetric) positive definite.

**Proof.** Let  $\mathbf{y} \in \mathbb{R}^m$  and set  $\mathbf{x} := \mathbf{X} \mathbf{y} \in \mathbb{R}^n$ . Then  $\mathbf{y}^T \mathbf{B} \mathbf{y} = \mathbf{y}^T \mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{y} = \mathbf{x}^T \mathbf{A} \mathbf{x}$ . This is nonnegative if  $\mathbf{A}$  is positive semidefinite and positive if  $\mathbf{A}$  is positive definite and  $\mathbf{X}$  has linearly independent columns. For then  $\mathbf{x}$  is nonzero if  $\mathbf{y}$  is nonzero. If  $\mathbf{A}$  is symmetric then  $\mathbf{B}$  is symmetric and the statements about symmetry follows.  $\square$

**Corollary 2.36** ( $A^T A$  is symmetric positive semidefinite)

Let  $m, n$  be positive integers. If  $A \in \mathbb{R}^{m,n}$  then  $A^T A$  is symmetric positive semidefinite. It is symmetric positive definite if and only if  $A$  has linearly independent columns.

**Proof.** If  $A^T A$  is symmetric positive definite then  $\mathbf{x} A^T A \mathbf{x} = \|A\mathbf{x}\|_2^2 > 0$  for all nonzero  $\mathbf{x}$  and  $A$  has linearly independent columns. Taking  $A := I$  and  $X := A$  in Theorem 2.35 gives  $B = A^T I A = A^T A$ , and we obtain the remaining statements of the corollary.  $\square$

**2.5.2 The nonsymmetric case**

A positive definite matrix has the following properties:

**Theorem 2.37 (The nonsymmetric case)**

Suppose  $A \in \mathbb{R}^{n,n}$  is a positive definite matrix and let  $B$  be a principal submatrix. Then

1.  $B$  is positive definite,
2.  $A$  has a unique LU factorization,
3. the real eigenvalues of  $B$  are positive,
4.  $\det(B) > 0$ ,
5.  $a_{ii}a_{jj} > a_{ij}a_{ji}$ , for  $i \neq j$ .

**Proof.**

1. Suppose the submatrix  $B = A(\mathbf{r}, \mathbf{r})$  is defined by the rows and columns  $\mathbf{r} = [r_1, \dots, r_k]^T$  of  $A$ . Let  $X = [\mathbf{e}_{r_1}, \dots, \mathbf{e}_{r_k}] \in \mathbb{R}^{n \times k}$ . Then  $B := X^T A X$ , and  $B$  is positive definite by Theorem 2.35.
2. Since all leading submatrices are nonsingular this follows from the LU Theorem 2.12.
3. Suppose  $(\lambda, \mathbf{x})$  is an eigenpair of  $A$  and that  $\lambda$  is real. Since  $A$  is real we can choose  $\mathbf{x}$  to be real. Multiplying  $A\mathbf{x} = \lambda\mathbf{x}$  by  $\mathbf{x}^T$  and solving for  $\lambda$  we find  $\lambda = \frac{\mathbf{x}^T A \mathbf{x}}{\mathbf{x}^T \mathbf{x}} > 0$ .
4. The determinant of  $B$  equals the product of its eigenvalues. The eigenvalues are either real and positive or occur in complex conjugate pairs. The product of two nonzero complex conjugate numbers is positive.

5. The principal submatrix  $\begin{bmatrix} a_{ii} & a_{ij} \\ a_{ji} & a_{jj} \end{bmatrix}$  has a positive determinant.

□

We note that:

1. Part 5 of Theorem 2.37 implies that all diagonal elements of a positive definite matrix are positive. Moreover, the greatest element in absolute value is a diagonal element. This can be used to decide by inspection that a given matrix cannot be positive definite.
2. A nonsymmetric positive definite matrix can have complex eigenvalues. For example, the eigenvalues of  $\mathbf{A}[a]$  in Example 2.32 are positive for  $a \in [1 - \frac{\sqrt{5}}{2}, 1 + \frac{\sqrt{5}}{2}]$  and complex for other values of  $\mathbf{A}$ .

### 2.5.3 The symmetric case

Theorem 2.37 can be strengthened considerably when  $\mathbf{A}$  is symmetric positive definite.

#### Theorem 2.38 (Symmetric positive definite characterization)

The following statements are equivalent for a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ .

1.  $\mathbf{A}$  is symmetric positive definite.
2.  $\mathbf{A}$  has only positive eigenvalues.
3. All leading principal submatrices have a positive determinant.
4.  $\mathbf{A} = \mathbf{B}\mathbf{B}^T$  for a nonsingular  $\mathbf{B} \in \mathbb{R}^{n \times n}$ .

**Proof.**  $1 \Leftrightarrow 2$  is shown in Lemma 2.41 below. We show that  $1 \Rightarrow 3 \Rightarrow 4 \Rightarrow 1$ .

$1 \Rightarrow 3$ : This follows from Theorem 2.37

$3 \Rightarrow 4$ : By Lemma 2.42 below  $\mathbf{A}$  has a unique symmetric LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{L}^T$  with positive diagonal elements in  $\mathbf{D}$ . But then  $\mathbf{A} = \mathbf{B}\mathbf{B}^T$ , where  $\mathbf{B} := \mathbf{L}\mathbf{D}^{1/2}$ , with  $\mathbf{D}^{1/2} := \text{diag}(d_{1,1}^{1/2}, \dots, d_{n \times n}^{1/2})$ .

$4 \Rightarrow 1$ : This follows from Corollary 2.36. □

#### Exercise 2.39 (Positive definite characterizations)

Show directly that all 4 characterizations in Theorem 2.38 hold for the matrix

$$\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}.$$

Consider the eigenvalues of a real symmetric positive definite matrix. Note that such a matrix is Hermitian.



**Lemma 2.40 (Eigenvalues of a Hermitian matrix)**

All eigenvalues of a Hermitian matrix are real.

*Proof.* Suppose  $\mathbf{A}^* = \mathbf{A}$  and  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  with  $\mathbf{x} \neq 0$ . We multiply both sides of  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  from the left by  $\mathbf{x}^*$  and divide by  $\mathbf{x}^*\mathbf{x}$  to obtain  $\lambda = \frac{\mathbf{x}^*\mathbf{A}\mathbf{x}}{\mathbf{x}^*\mathbf{x}}$ . Taking complex conjugates we find  $\bar{\lambda} = \lambda^* = \frac{(\mathbf{x}^*\mathbf{A}\mathbf{x})^*}{(\mathbf{x}^*\mathbf{x})^*} = \frac{\mathbf{x}^*\mathbf{A}^*\mathbf{x}}{\mathbf{x}^*\mathbf{x}} = \frac{\mathbf{x}^*\mathbf{A}\mathbf{x}}{\mathbf{x}^*\mathbf{x}} = \lambda$ , and  $\lambda$  is real.  $\square$

**Lemma 2.41 (Symmetry and positive eigenvalues)**

A matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is symmetric positive definite if and only if it is symmetric and all eigenvalues are positive.

*Proof.* By Lemma 2.40 all eigenvalues of  $\mathbf{A}$  are real, and by Theorem 2.37 all eigenvalues are positive. Suppose conversely that  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is symmetric with positive eigenvalues  $\lambda_1, \dots, \lambda_n$ . By the spectral theorem (cf. Corollary 5.23) we have  $\mathbf{U}^T\mathbf{A}\mathbf{U} = \mathbf{D}$ , where  $\mathbf{U}^T\mathbf{U} = \mathbf{U}\mathbf{U}^T = \mathbf{I}$  and  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$ . Let  $\mathbf{x} \in \mathbb{R}^n$  be nonzero and define  $\mathbf{c} := \mathbf{U}^T\mathbf{x} = [c_1, \dots, c_n]^T$ . Then  $\mathbf{c}^T\mathbf{c} = \mathbf{x}^T\mathbf{U}\mathbf{U}^T\mathbf{x} = \mathbf{x}^T\mathbf{x}$  so  $\mathbf{c}$  is nonzero. Since  $\mathbf{x} = \mathbf{U}\mathbf{c}$  we find

$$\mathbf{x}^T\mathbf{A}\mathbf{x} = (\mathbf{U}\mathbf{c})^T\mathbf{A}\mathbf{U}\mathbf{c} = \mathbf{c}^T\mathbf{U}^T\mathbf{A}\mathbf{U}\mathbf{c} = \mathbf{c}^T\mathbf{D}\mathbf{c} = \sum_{j=1}^n \lambda_j c_j^2 > 0$$

and it follows that  $\mathbf{A}$  is positive definite.  $\square$

**Lemma 2.42 (Symmetric positive definite and symmetric LU)**

A matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is symmetric positive definite if and only if it has a symmetric LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{L}^T$  with positive diagonal elements in  $\mathbf{D}$ .

*Proof.* Suppose  $\mathbf{A}$  is symmetric positive definite. By Theorem 2.37 the leading principal submatrices  $\mathbf{A}_{[k]} \in \mathbb{R}^{k \times k}$  are nonsingular for  $k = 1, \dots, n-1$ , and  $\mathbf{A}$  has a unique symmetric LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{L}^T$  by Theorem 2.26. The  $i$ th diagonal element in  $\mathbf{D}$  is positive, since  $d_{ii} = \mathbf{e}_i^T\mathbf{D}\mathbf{e}_i = \mathbf{e}_i^T\mathbf{L}^{-1}\mathbf{A}\mathbf{L}^{-T}\mathbf{e}_i = \mathbf{x}_i^T\mathbf{A}\mathbf{x}_i > 0$ . Indeed,  $\mathbf{x}_i := \mathbf{L}^{-T}\mathbf{e}_i$  is nonzero since  $\mathbf{L}^{-T}$  is nonsingular.

Conversely, suppose  $\mathbf{A}$  has a symmetric LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{L}^T$  with positive diagonal elements in  $\mathbf{D}$ . Then  $\mathbf{A}$  is symmetric, and for any nonzero  $\mathbf{y} \in \mathbb{R}^n$  we have  $\mathbf{x}^T\mathbf{A}\mathbf{x} = \mathbf{x}^T\mathbf{L}\mathbf{D}\mathbf{L}^T\mathbf{x} = \mathbf{y}^T\mathbf{D}\mathbf{y} > 0$ , since  $\mathbf{y} := \mathbf{L}^T\mathbf{x} \neq \mathbf{0}$ .  $\square$

## 2.6 The Cholesky Factorization



André-Louis Cholesky, 1875-1918 (left), John von Neumann, 1903-1957 (right).

Lemma 2.42 implies that  $\mathbf{A}$  is symmetric positive definite if and only if it has a symmetric LU factorization, and from the proof of 3. implies 4 in that theorem we can write this in the form  $\mathbf{A} = \mathbf{B}\mathbf{B}^T$  where  $\mathbf{B}$  is lower triangular matrix with positive diagonal elements. Such a factorization has a special name.

### Definition 2.43 (Cholesky)

A factorization  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$  where  $\mathbf{L}$  is lower triangular with positive diagonal elements is called a **Cholesky factorization** of  $\mathbf{A}$ . The matrix  $\mathbf{L}$  is called a **Cholesky factor**.

From the discussion before the definition we have

### Theorem 2.44 (Cholesky)

A matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  has a Cholesky factorization if and only if it is symmetric positive definite. Moreover, the Cholesky factorization is unique.

**Proof.** We still need to show uniqueness. Suppose  $\mathbf{L}\mathbf{L}^T = \mathbf{S}\mathbf{S}^T$  are two Cholesky factorizations of the symmetric positive definite matrix  $\mathbf{A}$ . Since  $\mathbf{A}$  is nonsingular both  $\mathbf{L}$  and  $\mathbf{S}$  are nonsingular. Then  $\mathbf{S}^{-1}\mathbf{L} = \mathbf{S}^T\mathbf{L}^{-T}$ , where by Lemma 1.22  $\mathbf{S}^{-1}\mathbf{L}$  is lower triangular and  $\mathbf{S}^T\mathbf{L}^{-T}$  is upper triangular, with diagonal elements  $\ell_{ii}/s_{ii}$  and  $s_{ii}/\ell_{ii}$ , respectively. But then both matrices must be equal to the same diagonal matrix and  $\ell_{ii}^2 = s_{ii}^2$ . By positivity  $\ell_{ii} = s_{ii}$  and we conclude that  $\mathbf{S}^{-1}\mathbf{L} = \mathbf{I} = \mathbf{S}^T\mathbf{L}^{-T}$  which means that  $\mathbf{L} = \mathbf{S}$ .  $\square$

A Cholesky factorization can also be written in the equivalent form  $\mathbf{A} =$

$\mathbf{R}^T \mathbf{R}$ , where  $\mathbf{R} = \mathbf{L}^T$  is upper triangular with positive diagonal elements. The matrix  $\mathbf{A}$  must be symmetric since  $\mathbf{L}\mathbf{L}^T$  is symmetric.

**Example 2.45** ( $2 \times 2$ )

The matrix  $\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$  has a symmetric LU- and a Cholesky-factorization given by

$$\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ -\frac{1}{2} & 1 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & \frac{3}{2} \end{bmatrix} \begin{bmatrix} 1 & -\frac{1}{2} \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} \sqrt{2} & 0 \\ -1/\sqrt{2} & \sqrt{3/2} \end{bmatrix} \begin{bmatrix} \sqrt{2} & -1/\sqrt{2} \\ 0 & \sqrt{3/2} \end{bmatrix}.$$

Consider computing the Cholesky factorization directly. The equation  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$  implies that

$$a_{ik} = \sum_{j=1}^n \ell_{ij} \ell_{kj} = \sum_{j=1}^{\min(i,k)} \ell_{ij} \ell_{kj}, \quad i, k = 1, \dots, n. \quad (2.20)$$

The unknown elements in  $\mathbf{L}$  can be computed row by row or column by column. Consider the column case. Suppose we have computed the  $k-1$  first columns of  $\mathbf{L}$ . The  $k$ th column can then be computed from (2.20). Indeed, letting  $i = k$  and solving for  $\ell_{kk}$  we find

$$\ell_{kk} = \left( a_{kk} - \sum_{j=1}^{k-1} \ell_{kj}^2 \right)^{1/2}, \quad (2.21)$$

and similarly for  $i > k$

$$\ell_{ik} = \left( a_{ik} - \sum_{j=1}^{k-1} \ell_{ij} \ell_{kj} \right) / \ell_{kk} \quad i = k+1, \dots, n. \quad (2.22)$$

Since  $\mathbf{A}$  is symmetric positive definite the Cholesky factor  $\mathbf{L}$  exists, is unique, and real, and therefore the term under the square root in (2.21) must be positive. We note however that we can encounter problems in floating point computation if the term is very small.

It is easily seen that the Cholesky-factorization of an  $n$ -by- $n$  matrix based on (2.21) and (2.22) requires  $\frac{1}{2}G_n = n^3/3$  arithmetic operations. This is essentially the same as for the symmetric LU factorization. The halving of the count compared to LU factorization is due to the symmetry of  $\mathbf{A}$ .

If  $\mathbf{A}$  is  $d$ -banded then the same is true for the Cholesky factor.

**Lemma 2.46 (Banded Cholesky factor)**

Suppose  $\mathbf{A}$  is symmetric positive definite with Cholesky-factor  $\mathbf{L}$ . If  $a_{ik} = 0$  for  $i > k + d$ , then  $\ell_{ik} = 0$  for  $i > k + d$ .

**Proof.** We show that if  $\mathbf{L}$  has bandwidth  $d$  in its first  $k-1$  columns then column  $k$  also has bandwidth  $d$ . The proof then follows by induction on  $k$ . Now, if  $i > k+d$ , then  $a_{ik} = 0$ , and if  $\mathbf{L}$  has bandwidth  $d$  in its first  $k-1$  columns then  $\ell_{ij} = 0$  for  $j = 1, \dots, k-1$ . By (2.22)  $\ell_{ik} = 0$ .  $\square$

We obtain formulas for the Cholesky factorization of a symmetric positive definite band matrix by simply replacing the lower bound  $j = 1$  by  $j = \max(1, k-d)$  in (2.21) and (2.22) and letting  $i$  run from  $k+1$  to  $\min(n, k+d)$  in (2.22). The lower triangular matrix  $\mathbf{L}$  is computed in sparse form. Only the lower triangular part of  $\mathbf{A}$  is used. This leads to the following algorithm. For a different algorithm based on outer products, which can also be used for symmetric positive semidefinite matrices, see Algorithm 2.53.

**Algorithm 2.47 (bandcholesky)**

```

1 function L=bandcholesky(A,d)
2 %L=bandcholesky(A,d)
3 n=length(A);
4 L=sparse(zeros(n,n));
5 for k=1:n
6     km=max(1,k-d); kp=min(n,k+d); s=L(k,km:k-1);
7     L(k,k)=sqrt(A(k,k)-s*s');
8     L(k+1:kp,k)=(A(k+1:kp,k) - ...
9         L(k+1:kp,km:k-1)*s')/L(k,k);
10 end

```

The leading term in an operation count for a band matrix is  $O(d^2n)$ . When  $d$  is small this is a considerable saving compared to the count  $\frac{1}{2}G_n = n^3/3$  for a full matrix.

There is also a banded version of the symmetric LU factorization which requires approximately the same number of arithmetic operations as the Cholesky factorization. The choice between using a symmetric LU factorization or an  $\mathbf{LL}^T$  factorization depends on several factors. Usually an LU or a symmetric LU factorization is preferred for matrices with small bandwidth (tridiagonal, pentadiagonal), while the  $\mathbf{LL}^T$  factorization is restricted to symmetric positive definite matrices and is often used when the bandwidth is larger.

## 2.7 The Symmetric Positive Semidefinite Case

We start with the following necessary conditions for a matrix to be symmetric positive semidefinite. It shows that if a diagonal element  $a_{ii}$  is zero then all elements in row  $i$  and column  $i$  are zero.

**Lemma 2.48 (Criteria symmetric semidefinite)**

If  $\mathbf{A}$  is symmetric positive semidefinite then for all  $i, j$

1.  $|a_{ij}| \leq (a_{ii} + a_{jj})/2$ ,
2.  $|a_{ij}| \leq \sqrt{a_{ii}a_{jj}}$ ,
3.  $\max_{i,j} |a_{ij}| = \max_i a_{ii}$ ,
4.  $a_{ii} = 0 \implies a_{ij} = a_{ji} = 0$ , fixed  $i$ , all  $j$ .

**Proof.** Part 3 follows from part 1 and part 4 from part 2. Now

$$0 \leq (\alpha \mathbf{e}_i + \beta \mathbf{e}_j)^T \mathbf{A} (\alpha \mathbf{e}_i + \beta \mathbf{e}_j) = \alpha^2 a_{ii} + \beta^2 a_{jj} + 2\alpha\beta a_{ij}, \quad \text{all } i, j, \alpha, \beta \in \mathbb{R}, \quad (2.23)$$

since  $\mathbf{A}$  is symmetric positive semidefinite. Taking  $\alpha = 1, \beta = \pm 1$  we obtain  $a_{ii} + a_{jj} \pm 2a_{ij} \geq 0$  and this implies part 1. Part 2 follows trivially from part 1 if  $a_{ii} = a_{jj} = 0$ . Suppose one of them, say  $a_{ii}$  is nonzero. Note that  $a_{ii} = \mathbf{e}_i^T \mathbf{A} \mathbf{e}_i > 0$ . Taking  $\alpha = -a_{ij}, \beta = a_{ii}$  in (2.23) we find

$$0 \leq a_{ij}^2 a_{ii} + a_{ii}^2 a_{jj} - 2a_{ij}^2 a_{ii} = a_{ii}(a_{ii} a_{jj} - a_{ij}^2).$$

But then  $a_{ii} a_{jj} - a_{ij}^2 \geq 0$  and part 2 follows.  $\square$

As an illustration consider the matrices

$$\mathbf{A}_1 = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{A}_2 = \begin{bmatrix} 1 & 2 \\ 2 & 2 \end{bmatrix}, \quad \mathbf{A}_3 = \begin{bmatrix} -2 & 1 \\ 1 & 2 \end{bmatrix}.$$

None of them is positive semidefinite, since neither part 1 nor part 2 hold.

**Theorem 2.49 (Positive symmetric semidefinite characterization)**

The following is equivalent for a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ .

1.  $\mathbf{A}$  is positive semidefinite.
2.  $\mathbf{A}$  has only nonnegative eigenvalues.
3.  $\mathbf{A} = \mathbf{B}\mathbf{B}^T$  for some  $\mathbf{B} \in \mathbb{R}^{n \times n}$ .
4. All principal submatrices have a nonnegative determinant.

**Proof.** The proof of  $1 \Leftrightarrow 2$  follows as in the proof of Theorem 2.38.  $1 \Leftrightarrow 3$  follows from Theorem 2.51 while  $1 \Rightarrow 4$  is a consequence of Theorem 2.37. To prove  $4 \Rightarrow 1$ , one first shows that  $\epsilon \mathbf{I} + \mathbf{A}$  is symmetric positive definite for all  $\epsilon > 0$  (Cf. page 567 of [24]). But then  $\mathbf{x}^T \mathbf{A} \mathbf{x} = \lim_{\epsilon \rightarrow 0} \mathbf{x}^T (\epsilon \mathbf{I} + \mathbf{A}) \mathbf{x} \geq 0$  for all  $\mathbf{x} \in \mathbb{R}^n$ .  $\square$

In part 4 of Theorem 2.49 we require nonnegativity of all principal minors, while only positivity of leading principal minors was required for positive definite matrices (cf. Theorem 2.38). To see that nonnegativity of the leading principal minors is not enough consider the matrix  $\mathbf{A} := \begin{bmatrix} 0 & 0 \\ 0 & -1 \end{bmatrix}$ . The leading principal minors are nonnegative, but  $\mathbf{A}$  is not positive semidefinite.

## 2.8 Semi-Cholesky factorization of a banded matrix

A symmetric positive semidefinite matrix has a factorization that is similar to the Cholesky factorization.

### Definition 2.50 (Semi-Cholesky factorization)

A factorization  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$  where  $\mathbf{L}$  is lower triangular with nonnegative diagonal elements is called a **semi-Cholesky factorization**.

Note that a semi-Cholesky factorization of a symmetric positive definite matrix is necessarily a Cholesky factorization. For if  $\mathbf{A}$  is positive definite then it is nonsingular and then  $\mathbf{L}$  must be nonsingular. Thus the diagonal elements of  $\mathbf{L}$  cannot be zero.

### Theorem 2.51 (Characterization, semi-Cholesky factorization)

A matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  has a semi-Cholesky factorization  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$  if and only if it is symmetric positive semidefinite.

*Proof.* If  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$  is a semi-Cholesky factorization then  $\mathbf{A}$  is symmetric positive semidefinite by Corollary 2.36. For the converse we use induction on  $n$ . A positive semidefinite matrix of order one has a semi-Cholesky factorization since the one and only element in  $\mathbf{A}$  is nonnegative. Suppose any symmetric positive semidefinite matrix of order  $n - 1$  has a semi-Cholesky factorization and suppose  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is symmetric positive semidefinite. We partition  $\mathbf{A}$  as follows

$$\mathbf{A} = \begin{bmatrix} \alpha & \mathbf{v}^T \\ \mathbf{v} & \mathbf{B} \end{bmatrix}, \quad \alpha \in \mathbb{R}, \mathbf{v} \in \mathbb{R}^{n-1}, \mathbf{B} \in \mathbb{R}^{(n-1) \times (n-1)}. \quad (2.24)$$

There are two cases. Suppose first  $\alpha = \mathbf{e}_1^T \mathbf{A} \mathbf{e}_1 > 0$ . We claim that  $\mathbf{C} := \mathbf{B} - \mathbf{v}\mathbf{v}^T/\alpha$  is symmetric positive semidefinite.  $\mathbf{C}$  is symmetric. To show that  $\mathbf{C}$  is positive semidefinite we consider any  $\mathbf{y} \in \mathbb{R}^{n-1}$  and define  $\mathbf{x}^T := [-\mathbf{v}^T \mathbf{y}/\alpha, \mathbf{y}^T] \in \mathbb{R}^n$ . Then

$$\begin{aligned} 0 \leq \mathbf{x}^T \mathbf{A} \mathbf{x} &= [-\mathbf{v}^T \mathbf{y}/\alpha, \mathbf{y}^T] \begin{bmatrix} \alpha & \mathbf{v}^T \\ \mathbf{v} & \mathbf{B} \end{bmatrix} \begin{bmatrix} -\mathbf{v}^T \mathbf{y}/\alpha \\ \mathbf{y} \end{bmatrix} \\ &= [0, -(\mathbf{v}^T \mathbf{y})\mathbf{v}^T/\alpha + \mathbf{y}^T \mathbf{B}] \begin{bmatrix} -\mathbf{v}^T \mathbf{y}/\alpha \\ \mathbf{y} \end{bmatrix} \\ &= -(\mathbf{v}^T \mathbf{y})(\mathbf{v}^T \mathbf{y})/\alpha + \mathbf{y}^T \mathbf{B} \mathbf{y} = \mathbf{y}^T \mathbf{C} \mathbf{y}, \end{aligned} \quad (2.25)$$

since  $(\mathbf{v}^T \mathbf{y}) \mathbf{v}^T \mathbf{y} = (\mathbf{v}^T \mathbf{y})^T \mathbf{v}^T \mathbf{y} = \mathbf{y}^T \mathbf{v} \mathbf{v}^T \mathbf{y}$ . So  $\mathbf{C} \in \mathbb{R}^{(n-1) \times (n-1)}$  is symmetric positive semidefinite and by the induction hypothesis it has a semi-Cholesky factorization  $\mathbf{C} = \mathbf{L}_1 \mathbf{L}_1^T$ . The matrix

$$\mathbf{L}^T := \begin{bmatrix} \beta & \mathbf{v}^T / \beta \\ \mathbf{0} & \mathbf{L}_1^T \end{bmatrix}, \quad \beta := \sqrt{\alpha}, \quad (2.26)$$

is upper triangular with nonnegative diagonal elements and

$$\mathbf{L} \mathbf{L}^T = \begin{bmatrix} \beta & \mathbf{0} \\ \mathbf{v} / \beta & \mathbf{L}_1 \end{bmatrix} \begin{bmatrix} \beta & \mathbf{v}^T / \beta \\ \mathbf{0} & \mathbf{L}_1^T \end{bmatrix} = \begin{bmatrix} \alpha & \mathbf{v}^T \\ \mathbf{v} & \mathbf{B} \end{bmatrix} = \mathbf{A}$$

is a semi-Cholesky factorization of  $\mathbf{A}$ .

If  $\alpha = 0$  then it follows from 4. in Lemma 2.48 that  $\mathbf{v} = \mathbf{0}$ . Moreover,  $\mathbf{B} \in \mathbb{R}^{(n-1) \times (n-1)}$  in (2.24) is positive semidefinite and therefore has a semi-Cholesky factorization  $\mathbf{B} = \mathbf{L}_1 \mathbf{L}_1^T$ . But then  $\mathbf{L} \mathbf{L}^T$ , where  $\mathbf{L} = \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{L}_1 \end{bmatrix}$  is a semi-Cholesky factorization of  $\mathbf{A}$ . Indeed,  $\mathbf{L}$  is lower triangular and

$$\mathbf{L} \mathbf{L}^T = \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{L}_1 \end{bmatrix} \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{L}_1^T \end{bmatrix} = \begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{B} \end{bmatrix} = \mathbf{A}.$$

□

Recall that a matrix  $\mathbf{A}$  is  $d$ -banded if  $a_{ij} = 0$  for  $|i - j| > d$ . A (semi-)Cholesky factorization preserves bandwidth.

### Theorem 2.52 (Bandwidth semi-Cholesky factor)

The semi-Cholesky factor  $\mathbf{L}$  given by (2.26) has the same bandwidth as  $\mathbf{A}$ .

*Proof.* Suppose  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is  $d$ -banded. Then  $\mathbf{v}^T = [\mathbf{u}^T, \mathbf{0}^T]$  in (2.24), where  $\mathbf{u} \in \mathbb{R}^d$ , and therefore  $\mathbf{C} := \mathbf{B} - \mathbf{v} \mathbf{v}^T / \alpha$  differs from  $\mathbf{B}$  only in the upper left  $d \times d$  corner. It follows that  $\mathbf{C}$  has the same bandwidth as  $\mathbf{B}$  and  $\mathbf{A}$ . By induction on  $n$ ,  $\mathbf{C} = \mathbf{L}_1 \mathbf{L}_1^T$ , where  $\mathbf{L}_1^T$  has the same bandwidth as  $\mathbf{C}$ . But then  $\mathbf{L}$  in (2.26) has the same bandwidth as  $\mathbf{A}$ . □

**Algorithm 2.53 (bandsemi-cholesky)**

Suppose  $\mathbf{A}$  is symmetric positive semidefinite. A lower triangular matrix  $\mathbf{L}$  is computed so that  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$ . This is the Cholesky factorization of  $\mathbf{A}$  if  $\mathbf{A}$  is symmetric positive definite and a semi-Cholesky factorization of  $\mathbf{A}$  otherwise. The algorithm uses the Matlab command `tril`.

```

1 function L=bandsemicholeskyL(A,d)
2 %L=bandsemicholeskyL(A,d)
3 n=length(A);
4 for k=1:n
5     if A(k,k)>0
6         kp=min(n,k+d);
7         A(k,k)=sqrt(A(k,k));
8         A(k+1:kp,k)=A(k+1:kp,k)/A(k,k);
9         for j=k+1:kp
10            A(j:kp,j)=A(j:kp,j)-A(j,k)*A(j:kp,k);
11        end
12    else
13        A(k:kp,k)=zeros(kp-k+1,1);
14    end
15 end
16 L=tril(A);

```

Consider now implementing an algorithm based on the previous discussion. Since  $\mathbf{A}$  is symmetric we only need to use the lower part of  $\mathbf{A}$ . The first column of  $\mathbf{L}$  is  $[\beta, \mathbf{v}^T/\beta]^T$  if  $\alpha > 0$ . If  $\alpha = 0$  then by 4 in Lemma 2.48 the first column of  $\mathbf{A}$  is zero and this is also the first column of  $\mathbf{L}$ . We obtain

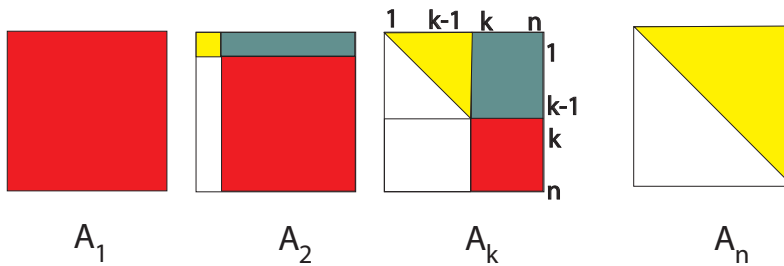
$$\begin{array}{l}
 \text{if } A(1,1) > 0 \\
 A(1,1) = \sqrt{A(1,1)} \\
 A(2:n,1) = A(2:n,1)/A(1,1) \\
 \text{for } j = 2:n \\
 A(j:n,j) = A(j:n,j) - A(j,1) * A(j:n,1)
 \end{array} \tag{2.27}$$

Here we store the first column of  $\mathbf{L}$  in the first column of  $\mathbf{A}$  and the lower part of  $\mathbf{C} = \mathbf{B} - \mathbf{v}\mathbf{v}^T/\alpha$  in the lower part of  $A(2:n,2:n)$ .

The code can be made more efficient when  $\mathbf{A}$  is a  $d$ -banded matrix. We simply replace all occurrences of  $n$  by  $\min(i+d, n)$ . Continuing the reduction we arrive at Algorithm 2.53.

In the algorithm we overwrite the lower triangle of  $\mathbf{A}$  with the elements of  $\mathbf{L}$ . Column  $k$  of  $\mathbf{L}$  is zero for those  $k$  where  $\ell_{kk} = 0$ . We reduce round-off noise by forcing those rows to be zero. In the semidefinite case no update is necessary and we “do nothing”.





**Figure 2.2.** *Gaussian elimination*

There are many versions of Cholesky factorizations, see [5]. Algorithm 2.47 is based on outer products  $\mathbf{v}\mathbf{v}^T$ . An advantage of this formulation is that it can be extended to symmetric positive semidefinite matrices. However deciding when a diagonal element is zero is a problem in floating point arithmetic.

## 2.9 Gaussian Elimination

In this section we take a closer look at Gaussian elimination. We show that if the conditions of the LU theorem is satisfied then Gaussian elimination without row interchanges is just a way of computing the LU factorization of the coefficient matrix. If row interchanges are incorporated then we need to introduce a matrix permutation matrix  $\mathbf{P}$ , and obtain a factorization of the form  $\mathbf{A} = \mathbf{P}\mathbf{L}\mathbf{U}$ . We obtain this **PLU factorization** by using a matrix formulation of Gaussian elimination.

### 2.9.1 Reduction to upper triangular form

Consider the general  $n \times n$  case (see Example 1.1 for the 3 by 3 case.). We start with a nonsingular linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  and generate a sequence of equivalent systems  $\mathbf{A}_k\mathbf{x} = \mathbf{b}_k$  for  $k = 1, \dots, n$ , where  $\mathbf{A}_1 := \mathbf{A}$ ,  $\mathbf{b}_1 := \mathbf{b}$ , and  $\mathbf{A}_k$  has zeros under the diagonal in its first  $k-1$  columns. Thus  $\mathbf{A}_n$  is upper triangular and the system  $\mathbf{A}_n\mathbf{x} = \mathbf{b}_n$  can be solved using one of Algorithms 2.2 or 2.5. The process is illustrated in Figure 2.2.

The matrix  $\mathbf{A}_k$  has the form

$$\mathbf{A}_k = \left[ \begin{array}{ccc|ccc} a_{1,1}^1 & \cdots & a_{1,k-1}^1 & a_{1,k}^1 & \cdots & a_{1,j}^1 & \cdots & a_{1,n}^1 \\ & \ddots & \vdots & \vdots & & \vdots & & \vdots \\ & & a_{k-1,k-1}^{k-1} & a_{k-1,k}^{k-1} & \cdots & a_{k-1,j}^{k-1} & \cdots & a_{k-1,n}^{k-1} \\ \hline & & & a_{k,k}^k & \cdots & a_{k,j}^k & \cdots & a_{k,n}^k \\ & & & \vdots & & \vdots & & \vdots \\ & & & a_{i,k}^k & \cdots & a_{i,j}^k & \cdots & a_{i,n}^k \\ & & & \vdots & & \vdots & & \vdots \\ & & & a_{n,k}^k & \cdots & a_{n,j}^k & \cdots & a_{n \times n}^k \end{array} \right]. \quad (2.28)$$

The process transforming  $\mathbf{A}_k$  into  $\mathbf{A}_{k+1}$  for  $k = 1, \dots, n-1$  can be described as follows:

|   |        |
|---|--------|
| <p>for <math>k = 1 : n - 1</math><br/>         Find <math>r_k \geq k</math> such that <math>a_{r_k, k} \neq 0</math>;<br/>         Interchange row <math>k</math> and <math>r_k</math> of <math>\mathbf{A}_k</math>;<br/>         for <math>i = k + 1 : n</math><br/> <math>m_{ik} = a_{ik}^k / a_{kk}^k</math><br/>         for <math>j = k : n</math><br/> <math>a_{ij}^{k+1} = a_{ij}^k - m_{ik} a_{kj}^k</math></p> | (2.29) |
|---|--------|

Since  $a_{ik}^{k+1} = a_{ik}^k - m_{ik} a_{kk}^k = 0$  for  $i = k + 1, \dots, n$  it follows that  $\mathbf{A}_{k+1}$  will have zeros under the diagonal in its first  $k$  columns and the elimination is carried one step further. The numbers  $m_{ik}$  in (2.29) are called **multipliers**. Interchanging two rows (and/or two columns) during Gaussian elimination is known as **pivoting**. The element which is moved to the diagonal position  $(k, k)$  is called the **pivot element** or **pivot** for short, and the row containing the pivot is called the **pivot row**.

## 2.9.2 Pivot strategies

The most common pivoting strategy used in the Gaussian elimination process (2.29) is

$$|a_{r_k, k}^k| := \max\{|a_{i, k}^k| : k \leq i \leq n\}$$

with  $r_k$  the smallest such index in case of a tie. This is known as **partial pivoting**. The following example illustrating that small pivots should be avoided.

**Example 2.54 (Row pivoting)**

Applying Gaussian elimination without row interchanges to the linear system

$$\begin{aligned}10^{-4}x_1 + 2x_2 &= 4 \\ x_1 + x_2 &= 3\end{aligned}$$

we obtain the upper triangular system

$$\begin{aligned}10^{-4}x_1 + 2x_2 &= 4 \\ (1 - 2 \times 10^4)x_2 &= 3 - 4 \times 10^4\end{aligned}$$

The exact solution is

$$x_2 = \frac{-39997}{-19999} \approx 2, \quad x_1 = \frac{4 - 2x_2}{10^{-4}} = \frac{20000}{19999} \approx 1.$$

Suppose we round the result of each arithmetic operation to three digits. The solutions  $\text{fl}(x_1)$  and  $\text{fl}(x_2)$  computed in this way is

$$\text{fl}(x_2) = 2, \quad \text{fl}(x_1) = 0.$$

The computed value 0 of  $x_1$  is completely wrong. Suppose instead we apply Gaussian elimination to the same system, but where we have interchanged the equations. The system is

$$\begin{aligned}x_1 + x_2 &= 3 \\ 10^{-4}x_1 + 2x_2 &= 4\end{aligned}$$

and we obtain the upper triangular system

$$\begin{aligned}x_1 + x_2 &= 3 \\ (2 - 10^{-4})x_2 &= 4 - 3 \times 10^{-4}\end{aligned}$$

Now the solution is computed as follows

$$x_2 = \frac{3.9997}{1.9999} \approx 2, \quad x_1 = 3 - x_2 \approx 1.$$

In this case rounding each calculation to three digits produces  $\text{fl}(x_1) = 1$  and  $\text{fl}(x_2) = 2$  which is quite satisfactory since it is the exact solution rounded to three digits.

Related to partial pivoting is **scaled partial pivoting**. Here  $r_k$  is the smallest index such that

$$\frac{|a_{r_k, k}^k|}{s_k} := \max\left\{\frac{|a_{i, k}^k|}{s_k} : k \leq i \leq n\right\}, \quad s_k := \max_{1 \leq j \leq n} |a_{kj}|.$$

This can sometimes give more accurate results if the coefficient matrix have coefficients of wildly different sizes. Note that the scaling factors  $s_k$  are computed using the initial matrix.

It also is possible to interchange both rows and columns. The choice

$$a_{r_k, s_k}^k := \max\{|a_{i,j}^k| : k \leq i, j \leq n\}$$

with  $r_k, s_k$  the smallest such indices in case of a tie, is known as **complete pivoting**. Complete pivoting is known to be more numerically stable than partial pivoting, but requires a lot of search and is seldom used in practice.

### 2.9.3 Permutation matrices

Pivoting can be described in terms of permutation matrices.

**Definition 2.55** *Let the components of  $\mathbf{p} = [k_1, \dots, k_n]^T$  be a permutation of the components of  $[1, 2, \dots, n]^T$ . We call  $\mathbf{P} := \mathbf{I}(:, \mathbf{p}) = [\mathbf{e}_{k_1}, \mathbf{e}_{k_2}, \dots, \mathbf{e}_{k_n}] \in \mathbb{R}^{n \times n}$  a **permutation matrix**. When discussing Gaussian elimination a permutation  $\mathbf{p}$  is sometimes called a **pivot vector**.*

Since  $\mathbf{P}^T = \mathbf{I}(\mathbf{p}, :)$  it follows that  $(\mathbf{P}^T \mathbf{P})_{i,j} = \mathbf{e}_{k_i}^T \mathbf{e}_{k_j} = \delta_{ij}$ . Thus  $\mathbf{P}^T \mathbf{P} = \mathbf{I}$ , the inverse of  $\mathbf{P}$  is equal to its transpose, and  $\mathbf{P} \mathbf{P}^T = \mathbf{I}$  as well. If  $\mathbf{p}$  and  $\mathbf{P}$  are as in Definition 2.55 and  $\mathbf{A} \in \mathbb{C}^{n \times n}$  then

$$\mathbf{A} \mathbf{P} = \mathbf{A}(:, \mathbf{p}), \quad \mathbf{P}^T \mathbf{A} = \mathbf{A}(\mathbf{p}, :), \quad \mathbf{P}^T \mathbf{A} \mathbf{P} = \mathbf{A}(\mathbf{p}, \mathbf{p}). \quad (2.30)$$

Thus, post-multiplying a matrix  $\mathbf{A}$  by a permutation matrix results in a permutation of the columns, pre-multiplying by the transpose of a permutation matrix gives a permutation of the rows, while the transformation  $\mathbf{P}^T \mathbf{A} \mathbf{P}$  permutes both the rows and columns using the same permutation  $\mathbf{p}$ . In particular, the diagonal of  $\mathbf{P}^T \mathbf{A} \mathbf{P}$  is a permutation of the diagonal of  $\mathbf{A}$ :

$$\text{diag}(\mathbf{P}^T \mathbf{A} \mathbf{P}) = \text{diag}(\mathbf{A})(\mathbf{p}). \quad (2.31)$$

We will use a particularly simple permutation matrix.

**Definition 2.56 (Interchange matrix)**

*We define a **(j,k)-Interchange Matrix**  $\mathbf{I}_{j,k}$  by interchanging column  $j$  and  $k$  of the identity matrix.*

Since  $\mathbf{I}_{j,k} = \mathbf{I}_{k,j}$ , and we obtain the identity by applying  $\mathbf{I}_{j,k}$  twice, we see that  $\mathbf{I}_{j,k}^2 = \mathbf{I}$  and an interchange matrix is symmetric and equal to its own inverse. Pre-multiplying a matrix by an interchange matrix interchanges two rows of the matrix, post-multiplication interchanges two columns. By (2.31) the diagonal of  $\mathbf{I}_{j,k} \mathbf{A} \mathbf{I}_{j,k}$  is almost conserved; only the diagonal elements  $a_{j,j}$  and  $a_{k,k}$  are interchanged.

### 2.9.4 Gauss transformations

The elimination process in (2.29) can be interpreted in matrix terms using interchange matrices and Gauss transformations.

**Definition 2.57 (Gauss transformation)**

Suppose for some  $1 \leq k < n$  that  $\mathbf{g}_k = [0, \dots, 0, g_{k+1,k}, \dots, g_{n,k}]^T \in \mathbb{R}^n$  has its first  $k$  components equal to zero and let  $\mathbf{e}_k$  be the  $k$ 'th unit vector in  $\mathbb{R}^n$ . The matrix

$$\mathbf{G}_k := \mathbf{I} - \mathbf{g}_k \mathbf{e}_k^T$$

is called a **Gauss transformation**. The name **elementary lower triangular matrix** is also used.

In the 3 by 3 case the Gauss transform takes the form

$$\mathbf{G}_1 = \begin{bmatrix} 1 & 0 & 0 \\ -g_{2,1} & 1 & 0 \\ -g_{3,1} & 0 & 1 \end{bmatrix}, \quad \mathbf{G}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -g_{3,2} & 1 \end{bmatrix},$$

where  $\mathbf{g}_1 = [0 \ g_{21} \ g_{31}]^T$  and  $\mathbf{g}_3 = [0 \ 0 \ g_{32}]^T$ .

For general  $n$  and  $1 \leq k < n$  a Gauss transformation can be used to zero out column  $k$  under the diagonal in the matrix  $\mathbf{A}_k$ . Column  $k$  in  $\mathbf{A}_k$  is transformed into column  $k$  of  $\mathbf{A}_{k+1}$  using a matrix  $\mathbf{M}_k^-$  as follows:

$$\begin{bmatrix} 1 & \cdots & 0 & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & & \vdots \\ 0 & \cdots & 1 & 0 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 1 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & -m_{k+1,k} & 1 & \cdots & 0 \\ \vdots & & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & -m_{n,k} & 0 & \cdots & 1 \end{bmatrix} \begin{bmatrix} a_{1,k-1}^1 \\ \vdots \\ a_{k,k}^k \\ a_{k+1,k}^k \\ \vdots \\ a_{n,k}^k \end{bmatrix} = \begin{bmatrix} a_{1,k-1}^1 \\ \vdots \\ a_{k,k}^k \\ 0 \\ \vdots \\ 0 \end{bmatrix}. \quad (2.32)$$

The transformation matrix  $\mathbf{M}_k^-$  in (2.32) is a Gauss transformation:

$$\mathbf{M}_k^- := \mathbf{I} - \mathbf{m}_k \mathbf{e}_k^T, \quad \mathbf{m}_k = [0, \dots, 0, m_{k+1,k}, \dots, m_{n,k}]^T. \quad (2.33)$$

We collect some properties of Gauss transformations that we will need.

**Lemma 2.58 (Gausstransformations)**

For  $1 \leq k < n$  let  $\mathbf{g}_k = [0, \dots, 0, g_{k+1,k}, \dots, g_{n,k}]^T \in \mathbb{R}^n$ . Then

1.  $(\mathbf{I} - \mathbf{g}_k \mathbf{e}_k^T)^{-1} = \mathbf{I} + \mathbf{g}_k \mathbf{e}_k^T,$
2.  $(\mathbf{I} - \mathbf{g}_1 \mathbf{e}_1^T) \cdots (\mathbf{I} - \mathbf{g}_k \mathbf{e}_k^T) = \mathbf{I} - \sum_{j=1}^k \mathbf{g}_j \mathbf{e}_j^T,$

3.  $\mathbf{I}_{i,j}(\mathbf{I} - \mathbf{g}_k \mathbf{e}_k^T) \mathbf{I}_{i,j} = \mathbf{I} - (\mathbf{I}_{i,j} \mathbf{g}_k) \mathbf{e}_k^T$ , for  $k < i, j \leq n$ .

**Proof.** We note that

$$\mathbf{e}_j^T \mathbf{g}_k = 0 \text{ for } j = 1, 2, \dots, k. \quad (2.34)$$

1. By direct multiplication using (2.34)

$$(\mathbf{I} - \mathbf{g}_k \mathbf{e}_k^T)(\mathbf{I} + \mathbf{g}_k \mathbf{e}_k^T) = \mathbf{I} - \mathbf{g}_k \mathbf{e}_k^T + \mathbf{g}_k \mathbf{e}_k^T - \mathbf{g}_k (\mathbf{e}_k^T \mathbf{g}_k) \mathbf{e}_k^T = \mathbf{I}.$$

2. Part 2. clearly holds for  $k = 1$ . Assuming by induction the result for  $k - 1$  we obtain by (2.34)

$$\begin{aligned} (\mathbf{I} - \mathbf{g}_1 \mathbf{e}_1^T) \cdots (\mathbf{I} - \mathbf{g}_k \mathbf{e}_k^T) &= (\mathbf{I} - \sum_{j=1}^{k-1} \mathbf{g}_j \mathbf{e}_j^T)(\mathbf{I} - \mathbf{g}_k \mathbf{e}_k^T) \\ &= \mathbf{I} - \sum_{j=1}^{k-1} \mathbf{g}_j \mathbf{e}_j^T - \mathbf{g}_k \mathbf{e}_k^T + \sum_{j=1}^{k-1} \mathbf{g}_j (\mathbf{e}_j^T \mathbf{g}_k) \mathbf{e}_k^T = \mathbf{I} - \sum_{j=1}^k \mathbf{g}_j \mathbf{e}_j^T. \end{aligned}$$

3. We find  $\mathbf{I}_{i,j}(\mathbf{I} - \mathbf{g}_k \mathbf{e}_k^T) \mathbf{I}_{i,j} = \mathbf{I}_{i,j}^2 - (\mathbf{I}_{i,j} \mathbf{g}_k)(\mathbf{e}_k^T \mathbf{I}_{i,j})$ . Now  $\mathbf{I}_{i,j}^2 = \mathbf{I}$  and  $\mathbf{e}_k^T \mathbf{I}_{i,j} = \mathbf{e}_k^T$  in view of  $i, j > k$ . Thus Part 3 follows.

□

It should be noted that the order of the factors in the product in Part 2 of Lemma 2.58 is important. For example for  $n = 3$

$$(\mathbf{I} + \mathbf{g}_1 \mathbf{e}_1^T)(\mathbf{I} + \mathbf{g}_2 \mathbf{e}_2^T) = \begin{bmatrix} 1 & 0 & 0 \\ g_{21} & 1 & 0 \\ g_{31} & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & g_{32} & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ g_{21} & 1 & 0 \\ g_{31} & g_{3,2} & 1 \end{bmatrix}.$$

In general,

$$(\mathbf{I} + \mathbf{g}_1 \mathbf{e}_1^T) \cdots (\mathbf{I} + \mathbf{g}_{n-1} \mathbf{e}_{n-1}^T) = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ g_{2,1} & 1 & 0 & \cdots & 0 \\ g_{3,1} & g_{3,2} & 1 & \cdots & 0 \\ \vdots & & \ddots & \ddots & \vdots \\ g_{n,1} & g_{n,2} & \cdots & g_{n,n-1} & 1 \end{bmatrix}. \quad (2.35)$$

Multiplying the factors in any other order does not give such a nice result.

### 2.9.5 PLU factorization

We can reformulate the Gaussian elimination process (2.29) in matrix terms. This leads to a **factorization** of the coefficient matrix  $\mathbf{A} = \mathbf{PLU}$ , where  $\mathbf{P}$  is a

permutation matrix,  $\mathbf{L}$  is a lower triangular, and  $\mathbf{U}$  is upper triangular. Without interchanges the matrix  $\mathbf{P}$  is the identity and we obtain an LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{U}$ .

The transformation  $\mathbf{A}_k \rightarrow \mathbf{A}_{k+1}$  in (2.29) can be written in terms of an interchange matrix and a Gauss transformation as follows:

$$\mathbf{A}_{k+1} = \mathbf{M}_k^- \mathbf{J}_k \mathbf{A}_k, \quad \mathbf{J}_k := \mathbf{I}_{k,r_k},$$

where we first interchanged row  $k$  and  $r_k \geq k$  in  $\mathbf{A}_k$  and then introduced zeros under the diagonal in column  $k$  using

$$\mathbf{M}_k^- := \mathbf{I} - \mathbf{m}_k \mathbf{e}_k^T, \quad \mathbf{m}_k = [0, \dots, 0, m_{k+1,k}, \dots, m_{n,k}]^T, \quad m_{ik} := a_{ik}^k / a_{kk}^k.$$

Applying this repeatedly we obtain

$$\mathbf{A}_k = \mathbf{M}_{k-1}^- \mathbf{J}_{k-1} \cdots \mathbf{M}_2^- \mathbf{J}_2 \mathbf{M}_1^- \mathbf{J}_1 \mathbf{A}, \quad k = 2, \dots, n. \quad (2.36)$$

Since  $\mathbf{J}_k^{-1} = \mathbf{J}_k$  and  $(\mathbf{M}_k^-)^{-1} = \mathbf{M}_k^+$ , we find the factorizations

$$\mathbf{A} = \mathbf{J}_1 \mathbf{M}_1^+ \mathbf{J}_2 \mathbf{M}_2^+ \cdots \mathbf{J}_{k-1} \mathbf{M}_{k-1}^+ \mathbf{A}_k, \quad k = 2, \dots, n. \quad (2.37)$$

Gaussian elimination with row pivoting is mathematically well defined on a nonsingular matrix.

**Theorem 2.59 (Gaussian elimination is well defined)**

*Suppose  $\mathbf{A} \in \mathbb{C}^{n,n}$  is nonsingular. Then for  $k = 1, 2, \dots, n-1$  we can in (2.29) find  $r_k \geq k$  such that  $a_{r_k,k} \neq 0$ .*

**Proof.** The result holds for  $k = 1$  since  $\mathbf{A}$  is nonsingular and therefore cannot have a zero column. Thus  $\mathbf{A}_2$  is well defined. Suppose for some  $k \geq 2$  that  $a_{i,i}^k \neq 0$  for  $i = 1, 2, \dots, k-1$ . We partition  $\mathbf{A}_k$  given by (2.28) in upper block triangular form

$$\mathbf{A}_k = \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \\ \mathbf{0} & \mathbf{E}_k \end{bmatrix}.$$

The matrix  $\mathbf{B}_k$  is upper triangular of order  $k-1$  with diagonal elements  $a_{1,1}^1 \cdots a_{k-1,k-1}^{k-1}$ . Therefore  $\mathbf{B}_k$  is nonsingular, and  $\mathbf{A}_k$  is nonsingular, since by (2.36) it is a product of nonsingular matrices. By Lemma 1.21  $\mathbf{E}_k$  is nonsingular and cannot have a zero first column. Thus  $\mathbf{A}_{k+1}$  is well defined and the result follows by induction.  $\square$

**Theorem 2.60 (PLU theorem)**

*Gaussian elimination on a nonsingular matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$ , with row pivoting as described in (2.29), leads to a factorization  $\mathbf{A} = \mathbf{P}\mathbf{L}\mathbf{U}$ , where  $\mathbf{P}$  is a permutation*

matrix,  $\mathbf{L}$  is lower triangular with ones on the diagonal, and  $\mathbf{U}$  is upper triangular. More explicitly,

$$\begin{aligned} \mathbf{P} &= \mathbf{J}_1 \mathbf{J}_2 \cdots \mathbf{J}_{n-1}, & \mathbf{J}_k &:= \mathbf{I}_{r_k, k}, \\ \mathbf{L} &= \mathbf{L}_1 \mathbf{L}_2 \cdots \mathbf{L}_{n-1}, \\ \mathbf{L}_k &= \mathbf{I} + \tilde{\mathbf{m}}_k \mathbf{e}_k^T, & \tilde{\mathbf{m}}_k &:= \mathbf{J}_{n-1} \cdots \mathbf{J}_{k+1} \mathbf{m}_k, & k &= 1, 2, \dots, n-1, \\ \mathbf{U} &= \mathbf{A}_n. \end{aligned} \tag{2.38}$$

**Proof.** By repeated use of Part 3 of Lemma 2.58 we have

$$\mathbf{L}_k = \mathbf{J}_{n-1} \cdots \mathbf{J}_{k+1} \mathbf{M}_k^+ \mathbf{J}_{k+1} \cdots \mathbf{J}_{n-1}, \quad k = 1, \dots, n-1.$$

Using  $\mathbf{J}_k^2 = \mathbf{I}$  repeatedly gives for  $n = 4$

$$\begin{aligned} \mathbf{PLU} &= (\mathbf{J}_1 \mathbf{J}_2 \mathbf{J}_3) (\mathbf{L}_1 \mathbf{L}_2 \mathbf{L}_3) \mathbf{A}_4 \\ &= (\mathbf{J}_1 \mathbf{J}_2 \mathbf{J}_3) (\mathbf{J}_3 \mathbf{J}_2 \mathbf{M}_1^+ \mathbf{J}_2 \mathbf{J}_3) (\mathbf{J}_3 \mathbf{M}_2^+ \mathbf{J}_3) (\mathbf{M}_3^+) \mathbf{A}_4 \\ &= \mathbf{J}_1 \mathbf{M}_1^+ \mathbf{J}_2 \mathbf{M}_2^+ \mathbf{J}_3 \mathbf{M}_3^+ \mathbf{A}_4. \end{aligned}$$

But then  $\mathbf{PLU} = \mathbf{A}$  by (2.37). Using the same cancellation effect for general  $n$  proves the theorem.  $\square$

Using Part 2 of Lemma 2.58 we see that the matrix  $\mathbf{L}$  in Theorem 2.60 has the form (2.35) with  $\mathbf{g}_k = \tilde{\mathbf{m}}_k$ ,  $k = 1, \dots, n$ .

Once we have a PLU factorization of  $\mathbf{A}$  the system  $\mathbf{Ax} = \mathbf{b}$  is solved easily in three steps. Since  $\mathbf{PLUx} = \mathbf{b}$  we have  $\mathbf{Pz} = \mathbf{b}$ ,  $\mathbf{Ly} = \mathbf{z}$ , and  $\mathbf{Ux} = \mathbf{y}$ . With  $\mathbf{P} = \mathbf{I}(:, \mathbf{p})$  the solution  $\mathbf{x}$  can be found from Algorithms 2.1 and 2.2 in two steps.

1. `y=rforwardsolve(L,b(p),n);`
2. `x=rbacksolve(U,y,n);`

**Exercise 2.61 (Using PLU of  $\mathbf{A}$  to solve  $\mathbf{A}^T \mathbf{x} = \mathbf{b}$ )**

Suppose we know the PLU factors  $\mathbf{P}, \mathbf{L}, \mathbf{U}$  in a PLU factorization  $\mathbf{A} = \mathbf{PLU}$  of  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . Explain how we can solve the system  $\mathbf{A}^T \mathbf{x} = \mathbf{b}$  economically.

**Exercise 2.62 (Using PLU to compute the determinant)**

Suppose we know the PLU factors  $\mathbf{P}, \mathbf{L}, \mathbf{U}$  in a PLU factorization  $\mathbf{A} = \mathbf{PLU}$  of  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . Explain how we can use this to compute the determinant of  $\mathbf{A}$ .

**Exercise 2.63 (Using PLU to compute the inverse)**

Suppose the factors  $\mathbf{P}, \mathbf{L}, \mathbf{U}$  in a PLU factorization of  $\mathbf{A} \in \mathbb{R}^{n \times n}$  are known. Use Exercises 2.21, 2.22 to show that it takes approximately  $2G_n$  arithmetic operations to compute  $\mathbf{A}^{-1} = \mathbf{U}^{-1} \mathbf{L}^{-1} \mathbf{P}^T$ . Here we have not counted the final multiplication with  $\mathbf{P}^T$  which amounts to  $n$  row interchanges.



## 2.9.6 The LU factorization

Consider now the lucky situation where no row interchanges is necessary in Gaussian elimination. In this case (2.29) simplifies to

$$\begin{array}{l}
 \text{for } k = 1 : n - 1 \\
 \text{for } i = k + 1 : n \\
 \quad m_{ik} = a_{ik}^k / a_{kk}^k \\
 \text{for } j = k : n \\
 \quad a_{ij}^{k+1} = a_{ij}^k - m_{ik} a_{kj}^k
 \end{array} \tag{2.39}$$

Gaussian elimination without row interchanges is sometimes referred to as **naive Gaussian elimination**, and we then have  $\mathbf{P} = \mathbf{I}$  in the PLU factorization. The PLU theorem then gives:

### Theorem 2.64 (LU factorization)

If Gaussian elimination without row interchanges is well defined then we obtain in (2.39) the LU factorization

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ m_{2,1} & 1 & 0 & \cdots & 0 \\ m_{3,1} & m_{3,2} & 1 & \cdots & 0 \\ \vdots & & \ddots & \ddots & \vdots \\ m_{n,1} & m_{n,2} & \cdots & m_{n,n-1} & 1 \end{bmatrix} \begin{bmatrix} a_{1,1}^1 & a_{1,2}^1 & a_{1,3}^1 & \cdots & a_{1,n}^1 \\ 0 & a_{2,2}^2 & a_{2,3}^2 & \cdots & a_{2,n}^2 \\ 0 & 0 & a_{3,3}^3 & \cdots & a_{3,n}^3 \\ \vdots & & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & a_{n,n}^n \end{bmatrix} = \mathbf{L}\mathbf{U}. \tag{2.40}$$

Since we get division by zero in (2.39) if  $a_{kk}^k = 0$  for some  $k \leq n - 1$  it is important to know when this can happen. We first show:

### Theorem 2.65 (Nonzero pivots)

Let  $\mathbf{A} \in \mathbb{C}^{n,n}$  and let  $m_{i,k}$  and  $a_{i,j}^k$  be defined by (2.39).

1. If  $a_{r,r}^r \neq 0$ ,  $r = 1, 2, \dots, k-1$ , then

$$\begin{vmatrix} a_{11} & \cdots & a_{k1} \\ \vdots & & \vdots \\ a_{k1} & \cdots & a_{kk} \end{vmatrix} = a_{1,1}^1 a_{2,2}^2 \cdots a_{k,k}^k, \quad k = 1, 2, \dots, n.$$

2. Let  $1 \leq k \leq n$ . Then  $a_{r,r}^r \neq 0$ ,  $r = 1, 2, \dots, k$ , if and only if the leading principal submatrix  $\mathbf{A}_{[r]}$  is nonsingular for  $r = 1, 2, \dots, k$ .

**Proof.**

1. Since  $a_{r,r}^r \neq 0$ ,  $r = 1, 2, \dots, k-1$  and  $\mathbf{J}_r = \mathbf{I}$  in (2.36) we obtain  $\mathbf{A}_k = \mathbf{M}\mathbf{A}$ , where  $\mathbf{M} := \mathbf{M}_{k-1}^- \cdots \mathbf{M}_1^-$ . By (2.28) we have

$$\mathbf{A}_k = \begin{bmatrix} \mathbf{F}_k & \mathbf{G}_k \\ \mathbf{0} & \mathbf{H}_k \end{bmatrix}, \quad \mathbf{F}_k = \begin{bmatrix} a_{1,1}^1 & \cdots & a_{1,k}^1 \\ & \ddots & \vdots \\ & & a_{k,k}^k \end{bmatrix} \in \mathbb{C}^{k,k}.$$

Since  $\mathbf{A}_k = \mathbf{M}\mathbf{A}$  and  $\mathbf{M}$  is lower triangular we obtain  $\mathbf{F}_k = \mathbf{M}_{[k]} \mathbf{A}_{[k]}$  and since  $\mathbf{M}_{[k]}$  has ones on the diagonal

$$a_{1,1}^1 a_{2,2}^2 \cdots a_{k,k}^k = \det(\mathbf{F}_k) = \det(\mathbf{M}_{[k]}) \det(\mathbf{A}_{[k]}) = \det(\mathbf{A}_{[k]}).$$

2. Suppose  $a_{r,r}^r \neq 0$ ,  $r = 1, 2, \dots, k$ . By Part 1  $\det(\mathbf{A}_{[r]}) = a_{1,1}^1 \cdots a_{r,r}^r$ ,  $r = 1, 2, \dots, k$ , so that  $\mathbf{A}_{[r]}$  is nonsingular for  $r = 1, 2, \dots, k$ . Conversely, suppose  $a_{i,i}^i = 0$  for some  $i \leq k$ . Let  $i$  be the smallest integer such that  $a_{i,i}^i = 0$ . We can then do Gaussian elimination without row interchanges on  $\mathbf{A}$  to obtain  $\mathbf{A}_i$ . By Part 1  $\det(\mathbf{A}_{[i]}) = a_{1,1}^1 \cdots a_{i,i}^i = 0$  so that  $(\mathbf{A})_{[i]}$  is singular.  $\square$

$\square$

The theorem implies:

**Corollary 2.66 (When is naive Gaussian elimination possible?)**

In (2.39) we have  $a_{k,k}^k \neq 0$  for  $k = 1, \dots, n-1$  if and only if the leading principal submatrices

$$\mathbf{A}_{[k]} := \begin{bmatrix} a_{11} & \cdots & a_{1k} \\ \vdots & & \vdots \\ a_{k1} & \cdots & a_{kk} \end{bmatrix}$$

of  $\mathbf{A}$  are nonsingular for  $k = 1, \dots, n-1$ .

We note that

1. The PLU factorization can alternatively be written  $\mathbf{P}^T \mathbf{A} = \mathbf{L}\mathbf{U}$ . Thus, if  $\mathbf{A}$  is nonsingular then there exists a permutation of the rows of  $\mathbf{A}$  so that the matrix with the rows permuted has an LU factorization. This means that if we knew the row pivots in advance then we can carry out Gaussian elimination without row pivoting on the matrix  $\mathbf{P}^T \mathbf{A}$ , where  $\mathbf{P}^T = \mathbf{I}_{r_{n-1}, n-1} \cdots \mathbf{I}_{r_1, 1}$ .
2. If the leading principal submatrices  $\mathbf{A}_{[k]}$  are nonsingular for  $k = 1, \dots, n-1$  then the LU factorization is unique and Gaussian elimination is just one particular way of computing the LU factorization.

3. The calculation in (2.39) requires  $\sum_{k=1}^{n-1} (n-k)^2$  multiplications, the same number of subtractions, and  $\sum_{k=1}^{n-1} k$  divisions. So the complexity of Gaussian elimination is  $\frac{2}{3}n^3 - \frac{1}{2}n^2 - \frac{1}{6}n$ . This is exactly the complexity of LU factorization (cf. (2.9)).
4. Corollary 2.66 holds even if  $\mathbf{A}$  is singular. Since  $\mathbf{L}$  is nonsingular the matrix  $\mathbf{U}$  is then singular, and since  $a_{kk}^k \neq 0$  for  $k = 1, \dots, n-1$  we must have  $a_{nn}^n = 0$  when  $\mathbf{A}$  is singular.
5. To verify the nonsingularity of the leading principal submatrices can be difficult in practice. We have show that this condition holds for a class of diagonally dominant matrices and for positive definite matrices.

**Exercise 2.67 (Direct proof of Theorem 2.64)**

Equation (2.39) implies that  $m_{ik} = (a_{ij}^k - a_{ij}^{k+1})/a_{kj}^k$  for  $k \leq \min(i-1, j-1)$ . Use this to give a proof of Theorem 2.64 by directly showing that  $(\mathbf{LU})_{ij} = \sum_{k=1}^n l_{ik}u_{kj} = a_{ij}$ . Consider separately the two cases  $i \leq j$  and  $i > j$ .

## 2.10 Review Questions

**2.9.1** When is a triangular matrix nonsingular?

**2.9.2** Approximately how many arithmetic operations are needed for

- the multiplication of two square matrices?
- The LU factorization of a matrix?
- the solution of  $\mathbf{Ax} = \mathbf{b}$ , when  $\mathbf{A}$  is triangular?

**2.9.3** What is the content of

- the LU theorem?
- the symmetric LU theorem?

**2.9.4** Is  $\mathbf{A}^T \mathbf{A}$  symmetric positive definite?

- 2.9.5**
- What class of matrices has a Cholesky factorization?
  - What is the bandwidth of the Cholesky factor of a band matrix?

**2.9.6** For a symmetric matrix give 3 conditions that are equivalent to positive definiteness.

**2.9.7** What class of matrices has a semi-Cholesky factorization?

**2.9.8** What is the general condition for Gaussian elimination without row interchanges to be well defined?

**2.9.9** What is a PLU factorization? When does it exist?

**2.9.10** What is complete pivoting?

## Chapter 3

# The Kronecker Product



Leopold Kronecker, 1823-1891 (left), Siméon Denis Poisson, 1781-1840 (right).

Matrices arising from 2D and 3D problems sometimes have a Kronecker product structure. Identifying a Kronecker structure can be very rewarding since it simplifies the study of such matrices.

### 3.1 Test Matrices

In this section we introduce some matrices which we will use to compare various algorithms in later chapters.

### 3.1.1 The 2D Poisson problem

Let  $\Omega := (0, 1)^2 = \{(x, y) : 0 < x, y < 1\}$  be the open unit square with boundary  $\partial\Omega$ . Consider the problem

$$-\Delta u := -\frac{\partial^2 u}{\partial x^2} - \frac{\partial^2 u}{\partial y^2} = f \text{ on } \Omega, \quad (3.1)$$

$$u := 0 \text{ on } \partial\Omega.$$

Here the function  $f$  is given and continuous on  $\Omega$ , and we seek a function  $u = u(x, y)$  such that (3.1) holds and which is zero on  $\partial\Omega$ .

Let  $m$  be a positive integer. We solve the problem numerically by finding approximations  $v_{j,k} \approx u(jh, kh)$  on a grid of points given by

$$\bar{\Omega}_h := \{(jh, kh) : j, k = 0, 1, \dots, m+1\}, \quad \text{where } h = 1/(m+1).$$

The points  $\Omega_h := \{(jh, kh) : j, k = 1, \dots, m\}$  are the interior points, while  $\bar{\Omega}_h \setminus \Omega_h$  are the boundary points. The solution is zero at the boundary points. Using the difference approximation from Chapter 1 for the second derivative we obtain the following approximations for the partial derivatives

$$\frac{\partial^2 u(jh, kh)}{\partial x^2} \approx \frac{v_{j-1,k} - 2v_{j,k} + v_{j+1,k}}{h^2}, \quad \frac{\partial^2 u(jh, kh)}{\partial y^2} \approx \frac{v_{j,k-1} - 2v_{j,k} + v_{j,k+1}}{h^2}.$$

Inserting this in (3.1) we get the following discrete analog of (3.1)

$$\begin{aligned} -\Delta_h v_{j,k} &= f_{j,k}, & (jh, kh) \in \Omega_h, \\ v_{j,k} &= 0, & (jh, kh) \in \partial\Omega_h, \end{aligned} \quad (3.2)$$

where  $f_{j,k} := f(jh, kh)$  and

$$-\Delta_h v_{j,k} := \frac{-v_{j-1,k} + 2v_{j,k} - v_{j+1,k}}{h^2} + \frac{-v_{j,k-1} + 2v_{j,k} - v_{j,k+1}}{h^2}. \quad (3.3)$$

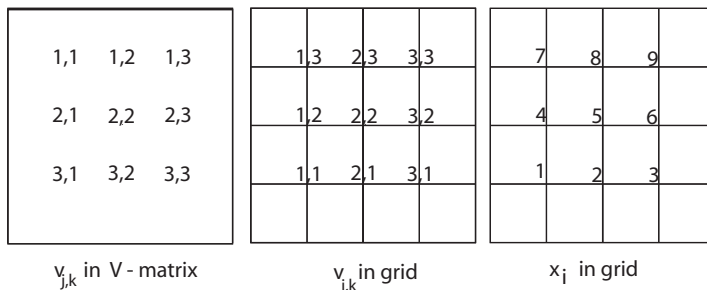
Multiplying both sides of (3.2) by  $h^2$  we obtain

$$\begin{aligned} 4v_{j,k} - v_{j-1,k} - v_{j+1,k} - v_{j,k-1} - v_{j,k+1} &= h^2 f_{j,k}, & (jh, kh) \in \Omega_h, \\ v_{0,k} = v_{m+1,k} = v_{j,0} = v_{j,m+1} &= 0, & j, k = 0, 1, \dots, m+1. \end{aligned} \quad (3.4)$$

The equations in (3.4) define a set of linear equations for the unknowns  $\mathbf{V} := [v_{jk}] \in \mathbb{R}^{m \times m}$ .

Observe that (3.4) can be written as a matrix equation in the form

$$\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2 \mathbf{F} \quad \text{with } h = 1/(m+1), \quad (3.5)$$



**Figure 3.1.** *Numbering of grid points*

where  $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m \times m}$  is the second derivative matrix given by (1.2) and  $\mathbf{F} = (f_{jk}) = (f(jh, kh)) \in \mathbb{R}^{m \times m}$ . Indeed, the  $(j, k)$  element in  $\mathbf{TV} + \mathbf{VT}$  is given by

$$\sum_{i=1}^m \mathbf{T}_{j,i} v_{i,k} + \sum_{i=1}^m v_{j,i} \mathbf{T}_{i,k},$$

and this is precisely the left hand side of (3.4).

To write (3.4) in standard form  $\mathbf{Ax} = \mathbf{b}$  we need to order the unknowns  $v_{j,k}$  in some way. The following operation of **vectorization** of a matrix gives one possible ordering.

**Definition 3.1 (vec operation)**

For any  $\mathbf{B} \in \mathbb{R}^{m \times n}$  we define the vector

$$\text{vec}(\mathbf{B}) := [b_{11}, \dots, b_{m1}, b_{12}, \dots, b_{m2}, \dots, b_{1n}, \dots, b_{mn}]^T \in \mathbb{R}^{mn}$$

by stacking the columns of  $\mathbf{B}$  on top of each other.

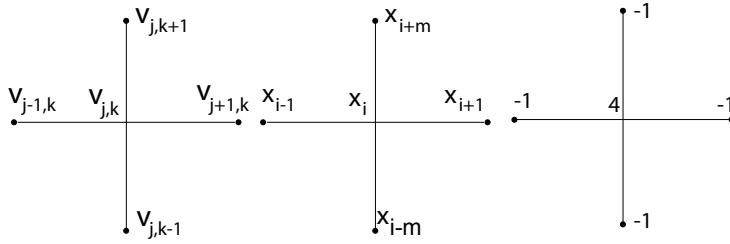
Let  $n = m^2$  and  $\mathbf{x} := \text{vec}(\mathbf{V}) \in \mathbb{R}^n$ . Note that forming  $\mathbf{x}$  by stacking the columns of  $\mathbf{V}$  on top of each other means an ordering of the grid points which for  $m = 3$  is illustrated in Figure 3.1. We call this the **natural ordering**. The elements in (3.4) form a 5-point stencil, as shown in Figure 3.2.

To find the matrix  $\mathbf{A}$  we note that for values of  $j, k$  where the 5-point stencil does not touch the boundary, (3.4) takes the form

$$4x_i - x_{i-1} - x_{i+1} - x_{i-m} - x_{i+m} = b_i,$$

where  $x_i = v_{jk}$  and  $b_i = h^2 f_{jk}$ . This must be modified close to the boundary. We obtain the linear system

$$\mathbf{Ax} = \mathbf{b}, \quad \mathbf{A} \in \mathbb{R}^{n \times n}, \quad \mathbf{b} \in \mathbb{R}^n, \quad n = m^2, \quad (3.6)$$



**Figure 3.2.** *The 5-point stencil*

where  $\mathbf{x} = \text{vec}(\mathbf{V})$ ,  $\mathbf{b} = h^2 \text{vec}(\mathbf{F})$  with  $\mathbf{F} = (f_{jk}) \in \mathbb{R}^{m \times m}$ , and  $\mathbf{A}$  is the **Poisson matrix** given by

$$\begin{aligned} a_{ii} &= 4, & i &= 1, \dots, n, \\ a_{i+1,i} &= a_{i,i+1} = -1, & i &= 1, \dots, n-1, & i &\neq m, 2m, \dots, (m-1)m, \\ a_{i+m,i} &= a_{i,i+m} = -1, & i &= 1, \dots, n-m, \\ a_{ij} &= 0, & & \text{otherwise.} \end{aligned} \quad (3.7)$$

For  $m = 3$  we have the following matrix

$$\mathbf{A} = \begin{bmatrix} 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 & 0 \\ -1 & 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & -1 & 4 & 0 & 0 & -1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 4 & -1 & 0 & -1 & 0 & 0 \\ 0 & -1 & 0 & -1 & 4 & -1 & 0 & -1 & 0 \\ 0 & 0 & -1 & 0 & -1 & 4 & 0 & 0 & -1 \\ 0 & 0 & 0 & -1 & 0 & 0 & 4 & -1 & 0 \\ 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 & -1 \\ 0 & 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 \end{bmatrix}.$$

**Exercise 3.2** ( $4 \times 4$  Poisson matrix)

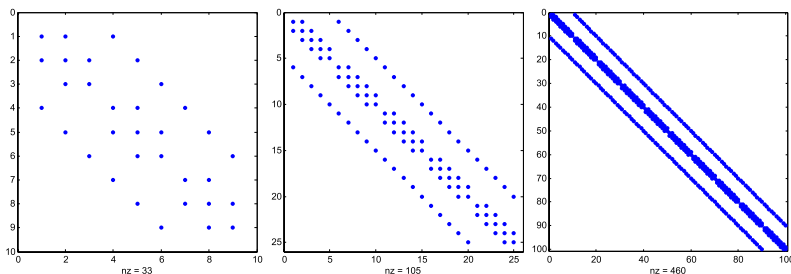
Write down the Poisson matrix for  $m = 2$  and show that it is strictly diagonally dominant.

### 3.1.2 The test matrices

The second derivative matrix  $\mathbf{T} = \text{tridiag}(-1, 2, -1)$  is a special case of the tridiagonal matrix

$$\mathbf{T}_1 := \text{tridiag}(a, d, a) \quad (3.8)$$





**Figure 3.3.** Band structure of the 2D test matrix,  $n = 9$ ,  $n = 25$ ,  $n = 100$

where  $a, d \in \mathbb{R}$ . We call this the **1D test matrix**. It is symmetric and strictly diagonally dominant if  $|d| > 2|a|$ .

The (2-dimensional) Poisson matrix is a special case of the matrix  $\mathbf{T}_2 = [a_{ij}] \in \mathbb{R}^{n \times n}$  with elements

$$\begin{aligned} a_{ii} &= 2d, & i = 1, \dots, n, \\ a_{i,i+1} = a_{i+1,i} &= a, & i = 1, \dots, n-1, \quad i \neq m, 2m, \dots, (m-1)m, \\ a_{i,i+m} = a_{i+m,i} &= a, & i = 1, \dots, n-m, \\ a_{ij} &= 0, & \text{otherwise,} \end{aligned} \quad (3.9)$$

and where  $a, d$  are real numbers. We will refer to this matrix as simply the **2D test matrix**. For  $m = 3$  the 2D test matrix looks as follows

$$\mathbf{T}_2 = \left[ \begin{array}{ccc|ccc|ccc} 2d & a & 0 & a & 0 & 0 & 0 & 0 & 0 \\ a & 2d & a & 0 & a & 0 & 0 & 0 & 0 \\ 0 & a & 2d & 0 & 0 & a & 0 & 0 & 0 \\ \hline a & 0 & 0 & 2d & a & 0 & a & 0 & 0 \\ 0 & a & 0 & a & 2d & a & 0 & a & 0 \\ 0 & 0 & a & 0 & a & 2d & 0 & 0 & a \\ \hline 0 & 0 & 0 & a & 0 & 0 & 2d & a & 0 \\ 0 & 0 & 0 & 0 & a & 0 & a & 2d & a \\ 0 & 0 & 0 & 0 & 0 & a & 0 & a & 2d \end{array} \right]. \quad (3.10)$$

The partition into  $3 \times 3$  sub matrices shows that  $\mathbf{T}_2$  is block tridiagonal.

Properties of  $\mathbf{T}_2$  can be derived from properties of  $\mathbf{T}_1$  by using properties of the Kronecker product.

## 3.2 The Kronecker Product

### Definition 3.3 (Kronecker product)

For any positive integers  $p, q, r, s$  we define the **Kronecker product** of two matrices  $\mathbf{A} \in \mathbb{R}^{p \times q}$  and  $\mathbf{B} \in \mathbb{R}^{r \times s}$  as a matrix  $\mathbf{C} \in \mathbb{R}^{pr \times qs}$  given in block form as

$$\mathbf{C} = \begin{bmatrix} \mathbf{A}b_{1,1} & \mathbf{A}b_{1,2} & \cdots & \mathbf{A}b_{1,s} \\ \mathbf{A}b_{2,1} & \mathbf{A}b_{2,2} & \cdots & \mathbf{A}b_{2,s} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{A}b_{r,1} & \mathbf{A}b_{r,2} & \cdots & \mathbf{A}b_{r,s} \end{bmatrix}.$$

We denote the Kronecker product of  $\mathbf{A}$  and  $\mathbf{B}$  by  $\mathbf{C} = \mathbf{A} \otimes \mathbf{B}$ .

This definition of the Kronecker product is known more precisely as the **left Kronecker product**. In the literature one often finds the **right Kronecker product** which in our notation is given by  $\mathbf{B} \otimes \mathbf{A}$ .

The Kronecker product  $\mathbf{u} \otimes \mathbf{v} = [\mathbf{u}^T v_1, \dots, \mathbf{u}^T v_r]^T$  of two column vectors  $\mathbf{u} \in \mathbb{R}^p$  and  $\mathbf{v} \in \mathbb{R}^r$  is a column vector of length  $p \times r$ .

As examples of Kronecker products which are relevant for our discussion, if

$$\mathbf{T}_1 = \begin{bmatrix} d & a & 0 \\ a & d & a \\ 0 & a & d \end{bmatrix} \quad \text{and} \quad \mathbf{I} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

then

$$\mathbf{T}_1 \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}_1 = \begin{bmatrix} \mathbf{T}_1 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{T}_1 \end{bmatrix} + \begin{bmatrix} d\mathbf{I} & a\mathbf{I} & \mathbf{0} \\ a\mathbf{I} & d\mathbf{I} & a\mathbf{I} \\ \mathbf{0} & a\mathbf{I} & d\mathbf{I} \end{bmatrix} = \mathbf{T}_2$$

given by (3.10). The same equation holds for any integer  $m \geq 2$

$$\mathbf{T}_1 \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}_1 = \mathbf{T}_2, \quad \mathbf{T}_1, \mathbf{I} \in \mathbb{R}^{m \times m}, \quad \mathbf{T}_2 \in \mathbb{R}^{(m^2) \times (m^2)}. \quad (3.11)$$

The sum of two Kronecker products involving the identity matrix is worthy of a special name.

### Definition 3.4 (Kronecker sum)

For positive integers  $r, s, k$ , let  $\mathbf{A} \in \mathbb{R}^{r \times r}$ ,  $\mathbf{B} \in \mathbb{R}^{s \times s}$ , and  $\mathbf{I}_k$  be the identity matrix of order  $k$ . The sum  $\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}$  is known as the **Kronecker sum** of  $\mathbf{A}$  and  $\mathbf{B}$ .

In other words, the 2D test matrix  $\mathbf{T}_2$  is the Kronecker sum involving the 1D test matrix  $\mathbf{T}_1$ .

The following simple arithmetic rules hold for Kronecker products. For scalars  $\lambda, \mu$  and matrices  $\mathbf{A}, \mathbf{A}_1, \mathbf{A}_2, \mathbf{B}, \mathbf{B}_1, \mathbf{B}_2, \mathbf{C}$  of dimensions such that the operations are defined, we have

$$\begin{aligned}(\lambda \mathbf{A}) \otimes (\mu \mathbf{B}) &= \lambda \mu (\mathbf{A} \otimes \mathbf{B}), \\(\mathbf{A}_1 + \mathbf{A}_2) \otimes \mathbf{B} &= \mathbf{A}_1 \otimes \mathbf{B} + \mathbf{A}_2 \otimes \mathbf{B}, \\ \mathbf{A} \otimes (\mathbf{B}_1 + \mathbf{B}_2) &= \mathbf{A} \otimes \mathbf{B}_1 + \mathbf{A} \otimes \mathbf{B}_2, \\(\mathbf{A} \otimes \mathbf{B}) \otimes \mathbf{C} &= \mathbf{A} \otimes (\mathbf{B} \otimes \mathbf{C}).\end{aligned}\tag{3.12}$$

Note however that in general we have  $\mathbf{A} \otimes \mathbf{B} \neq \mathbf{B} \otimes \mathbf{A}$ , but it can be shown that there are permutation matrices  $\mathbf{P}, \mathbf{Q}$  such that  $\mathbf{B} \otimes \mathbf{A} = \mathbf{P}(\mathbf{A} \otimes \mathbf{B})\mathbf{Q}$ , see [14].

### Exercise 3.5 (Properties of Kronecker products)

*Prove (3.12).*

The following **mixed product rule** is an essential tool for dealing with Kronecker products and sums.

### Lemma 3.6 (Mixed product rule)

*Suppose  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$  are rectangular matrices with dimensions so that the products  $\mathbf{AC}$  and  $\mathbf{BD}$  are defined. Then the product  $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D})$  is defined and*

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD}).\tag{3.13}$$

*Proof.* If  $\mathbf{B} \in \mathbb{R}^{r,t}$  and  $\mathbf{D} \in \mathbb{R}^{t,s}$  for some integers  $r, s, t$ , then

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \begin{bmatrix} \mathbf{A}b_{1,1} & \cdots & \mathbf{A}b_{1,t} \\ \vdots & & \vdots \\ \mathbf{A}b_{r,1} & \cdots & \mathbf{A}b_{r,t} \end{bmatrix} \begin{bmatrix} \mathbf{C}d_{1,1} & \cdots & \mathbf{C}d_{1,s} \\ \vdots & & \vdots \\ \mathbf{C}d_{t,1} & \cdots & \mathbf{C}d_{t,s} \end{bmatrix}.$$

Thus for all  $i, j$

$$((\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}))_{i,j} = \mathbf{AC} \sum_{k=1}^t b_{i,k} d_{k,j} = (\mathbf{AC})(\mathbf{BD})_{i,j} = ((\mathbf{AC}) \otimes (\mathbf{BD}))_{i,j}.$$

□

Using the mixed product rule we obtain the following properties of Kronecker products and sums.

### Theorem 3.7 (Properties of Kronecker products)

*Suppose for  $r, s \in \mathbb{N}$  that  $\mathbf{A} \in \mathbb{R}^{r,r}$  and  $\mathbf{B} \in \mathbb{R}^{s,s}$  are square matrices with eigenpairs  $(\lambda_i, \mathbf{u}_i)$   $i = 1, \dots, r$  and  $(\mu_j, \mathbf{v}_j)$ ,  $j = 1, \dots, s$ . Moreover, let  $\mathbf{F}, \mathbf{V} \in \mathbb{R}^{r \times s}$ . Then*

1.  $(\mathbf{A} \otimes \mathbf{B})^T = \mathbf{A}^T \otimes \mathbf{B}^T$ , (this also holds for rectangular matrices).
2. If  $\mathbf{A}$  and  $\mathbf{B}$  are nonsingular then  $\mathbf{A} \otimes \mathbf{B}$  is nonsingular. with  $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$ .
3. If  $\mathbf{A}$  and  $\mathbf{B}$  are symmetric then  $\mathbf{A} \otimes \mathbf{B}$  and  $\mathbf{A} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{B}$  are symmetric.
4.  $(\mathbf{A} \otimes \mathbf{B})(\mathbf{u}_i \otimes \mathbf{v}_j) = \lambda_i \mu_j (\mathbf{u}_i \otimes \mathbf{v}_j)$ ,  $i = 1, \dots, r$ ,  $j = 1, \dots, s$ ,
5.  $(\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B})(\mathbf{u}_i \otimes \mathbf{v}_j) = (\lambda_i + \mu_j)(\mathbf{u}_i \otimes \mathbf{v}_j)$ ,  $i = 1, \dots, r$ ,  $j = 1, \dots, s$ ,
6. If one of  $\mathbf{A}$ ,  $\mathbf{B}$  is symmetric positive definite and the other is symmetric positive semidefinite then  $\mathbf{A} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{B}$  is symmetric positive definite.
7.  $\mathbf{A}\mathbf{V}\mathbf{B}^T = \mathbf{F} \Leftrightarrow (\mathbf{A} \otimes \mathbf{B}) \text{vec}(\mathbf{V}) = \text{vec}(\mathbf{F})$ ,
8.  $\mathbf{A}\mathbf{V} + \mathbf{V}\mathbf{B}^T = \mathbf{F} \Leftrightarrow (\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}) \text{vec}(\mathbf{V}) = \text{vec}(\mathbf{F})$ .

Before giving the simple proofs of this theorem we present some comments.

1. The transpose (or the inverse) of an ordinary matrix product equals the transpose (or the inverse) of the matrices in reverse order. For Kronecker products the order is kept.
2. The eigenvalues of the Kronecker product (or sum) are the product (or sum) of the eigenvalues of the factors. The eigenvectors are the Kronecker products of the eigenvectors of the factors. In particular, the eigenvalues of the test matrix  $\mathbf{T}_2$  are sums of eigenvalues of  $\mathbf{T}_1$ . We will find these eigenvalues in the next section.
3. Since we already know that  $\mathbf{T} = \text{tridiag}(-1, 2, -1)$  is positive definite the 2D Poisson matrix  $\mathbf{A} = \mathbf{T} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}$  is also positive definite.
4. The system  $\mathbf{A}\mathbf{V}\mathbf{B}^T = \mathbf{F}$  in part 7 can be solved by first finding  $\mathbf{W}$  from  $\mathbf{A}\mathbf{W} = \mathbf{F}$ , and then finding  $\mathbf{V}$  from  $\mathbf{B}\mathbf{V}^T = \mathbf{W}^T$ . This is preferable to solving the much larger linear system  $(\mathbf{A} \otimes \mathbf{B}) \text{vec}(\mathbf{V}) = \text{vec}(\mathbf{F})$ .
5. A fast way to solve the 2D Poisson problem in the form  $\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = \mathbf{F}$  will be considered in the next chapter.

**Proof.**

1. Exercise.
2. By the mixed product rule  $(\mathbf{A} \otimes \mathbf{B})(\mathbf{A}^{-1} \otimes \mathbf{B}^{-1}) = (\mathbf{A}\mathbf{A}^{-1}) \otimes (\mathbf{B}\mathbf{B}^{-1}) = \mathbf{I}_r \otimes \mathbf{I}_s = \mathbf{I}_{rs}$ . Thus  $(\mathbf{A} \otimes \mathbf{B})$  is nonsingular with the indicated inverse.

3. By 1,  $(\mathbf{A} \otimes \mathbf{B})^T = \mathbf{A}^T \otimes \mathbf{B}^T = \mathbf{A} \otimes \mathbf{B}$ . Moreover, since then  $\mathbf{A} \otimes \mathbf{I}$  and  $\mathbf{I} \otimes \mathbf{B}$  are symmetric, their sum is symmetric.
4.  $(\mathbf{A} \otimes \mathbf{B})(\mathbf{u}_i \otimes \mathbf{v}_j) = (\mathbf{A}\mathbf{u}_i) \otimes (\mathbf{B}\mathbf{v}_j) = (\lambda_i \mathbf{u}_i) \otimes (\mu_j \mathbf{v}_j) = (\lambda_i \mu_j)(\mathbf{u}_i \otimes \mathbf{v}_j)$ , for all  $i, j$ , where we used the mixed product rule.
5.  $(\mathbf{A} \otimes \mathbf{I}_s)(\mathbf{u}_i \otimes \mathbf{v}_j) = \lambda_i(\mathbf{u}_i \otimes \mathbf{v}_j)$ , and  $(\mathbf{I}_r \otimes \mathbf{B})(\mathbf{u}_i \otimes \mathbf{v}_j) = \mu_j(\mathbf{u}_i \otimes \mathbf{v}_j)$ . The result now follows by summing these relations.
6. By 1,  $\mathbf{A} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{B}$  is symmetric. Moreover, the eigenvalues  $\lambda_i + \mu_j$  are positive since for all  $i, j$ , both  $\lambda_i$  and  $\mu_j$  are nonnegative and one of them is positive. It follows that  $\mathbf{A} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{B}$  is symmetric positive definite.
7. We partition  $\mathbf{V}$ ,  $\mathbf{F}$ , and  $\mathbf{B}^T$  by columns as  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_s]$ ,  $\mathbf{F} = [\mathbf{f}_1, \dots, \mathbf{f}_s]$  and  $\mathbf{B}^T = [\mathbf{b}_1, \dots, \mathbf{b}_s]$ . Then we have

$$\begin{aligned}
(\mathbf{A} \otimes \mathbf{B}) \text{vec}(\mathbf{V}) &= \text{vec}(\mathbf{F}) \\
\Leftrightarrow \begin{bmatrix} \mathbf{A}b_{11} & \cdots & \mathbf{A}b_{1s} \\ \vdots & & \vdots \\ \mathbf{A}b_{s1} & \cdots & \mathbf{A}b_{ss} \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_s \end{bmatrix} &= \begin{bmatrix} \mathbf{f}_1 \\ \vdots \\ \mathbf{f}_s \end{bmatrix} \\
\Leftrightarrow \mathbf{A} \left[ \sum_j b_{1j} \mathbf{v}_j, \dots, \sum_j b_{sj} \mathbf{v}_j \right] &= [\mathbf{f}_1, \dots, \mathbf{f}_s] \\
\Leftrightarrow \mathbf{A}[\mathbf{V}b_1, \dots, \mathbf{V}b_s] = \mathbf{F} &\Leftrightarrow \mathbf{A}\mathbf{V}\mathbf{B}^T = \mathbf{F}.
\end{aligned}$$

8. This follows immediately from (7) as follows

$$\begin{aligned}
(\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}) \text{vec}(\mathbf{V}) &= \text{vec}(\mathbf{F}) \\
\Leftrightarrow (\mathbf{A}\mathbf{V}\mathbf{I}_s^T + \mathbf{I}_r \mathbf{V}\mathbf{B}^T) &= \mathbf{F} \Leftrightarrow \mathbf{A}\mathbf{V} + \mathbf{V}\mathbf{B}^T = \mathbf{F}.
\end{aligned}$$

□

For more on Kronecker products see [14].

### 3.3 Properties of the 1D and 2D Test Matrices

Using Theorem 3.7 we can derive properties of the 2D test matrix  $\mathbf{T}_2$  from those of  $\mathbf{T}_1$ . We need to determine the eigenpairs of  $\mathbf{T}_1$ . We show that the eigenvectors are the columns of the **sine matrix** defined by

$$\mathbf{S} = \left[ \sin \frac{jk\pi}{m+1} \right]_{j,k=1}^m \in \mathbb{R}^{m \times m}. \quad (3.14)$$

For  $m = 3$ ,

$$\mathbf{S} = [\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3] = \begin{bmatrix} \sin \frac{\pi}{4} & \sin \frac{2\pi}{4} & \sin \frac{3\pi}{4} \\ \sin \frac{2\pi}{4} & \sin \frac{4\pi}{4} & \sin \frac{6\pi}{4} \\ \sin \frac{3\pi}{4} & \sin \frac{6\pi}{4} & \sin \frac{9\pi}{4} \end{bmatrix} = \begin{bmatrix} t & 1 & t \\ 1 & 0 & -1 \\ t & -1 & t \end{bmatrix}, \quad t := \frac{1}{\sqrt{2}}.$$

**Lemma 3.8 (Eigenpairs of 1D test matrix)**

Suppose  $\mathbf{T}_1 = (t_{kj})_{k,j} = \text{tridiag}(a, d, a) \in \mathbb{R}^{m \times m}$  with  $m \geq 2$ ,  $a, d \in \mathbb{R}$ , and let  $h = 1/(m+1)$ .

1. We have  $\mathbf{T}_1 \mathbf{s}_j = \lambda_j \mathbf{s}_j$  for  $j = 1, \dots, m$ , where

$$\mathbf{s}_j = [\sin(j\pi h), \sin(2j\pi h), \dots, \sin(mj\pi h)]^T, \quad (3.15)$$

$$\lambda_j = d + 2a \cos(j\pi h). \quad (3.16)$$

2. The eigenvalues are distinct and the eigenvectors are orthogonal

$$\mathbf{s}_j^T \mathbf{s}_k = \frac{m+1}{2} \delta_{j,k} = \frac{1}{2h} \delta_{j,k}, \quad j, k = 1, \dots, m. \quad (3.17)$$

**Proof.** We find for  $1 < k < m$

$$\begin{aligned} (\mathbf{T}_1 \mathbf{s}_j)_k &= \sum_{l=1}^m t_{k,l} \sin(lj\pi h) = a [\sin((k-1)j\pi h) + \sin((k+1)j\pi h)] + d \sin(kj\pi h) \\ &= 2a \cos(j\pi h) \sin(kj\pi h) + d \sin(kj\pi h) = \lambda_j \mathbf{s}_{k,j}. \end{aligned}$$

This also holds for  $k = 1, m$ , and part 1 follows. Since  $j\pi h = j\pi/(m+1) \in (0, \pi)$  for  $j = 1, \dots, m$  and the cosine function is strictly monotone decreasing on  $(0, \pi)$  the eigenvalues are distinct, and since  $\mathbf{T}_1$  is symmetric it follows from Lemma 3.9 below that the eigenvectors  $\mathbf{s}_j$  are orthogonal. To finish the proof of (3.17) we compute the square of the Euclidian norm of each  $\mathbf{s}_j$  as follows:

$$\begin{aligned} \mathbf{s}_j^T \mathbf{s}_j &= \sum_{k=1}^m \sin^2(kj\pi h) = \sum_{k=0}^m \sin^2(kj\pi h) = \frac{1}{2} \sum_{k=0}^m (1 - \cos(2kj\pi h)) \\ &= \frac{m+1}{2} - \frac{1}{2} \sum_{k=0}^m \cos(2kj\pi h) = \frac{m+1}{2}, \end{aligned}$$

since the last cosine sum is zero. We show this by summing a geometric series of complex exponentials. With  $i = \sqrt{-1}$  we find

$$\sum_{k=0}^m \cos(2kj\pi h) + i \sum_{k=0}^m \sin(2kj\pi h) = \sum_{k=0}^m e^{2ikj\pi h} = \frac{e^{2i(m+1)j\pi h} - 1}{e^{2ij\pi h} - 1} = 0,$$

and (3.17) follows.  $\square$

**Lemma 3.9 (Eigenpairs of a Hermitian matrix)**

The eigenvalues of a Hermitian matrix are real. Moreover, eigenvectors corresponding to distinct eigenvalues are orthogonal.

*Proof.* The first part was shown in Lemma 2.40. Suppose that  $(\lambda, \mathbf{x})$  and  $(\mu, \mathbf{y})$  are two eigenpairs for  $\mathbf{A}$  with  $\mu \neq \lambda$ . Multiplying  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  by  $\mathbf{y}^*$  gives

$$\lambda \mathbf{y}^* \mathbf{x} = \mathbf{y}^* \mathbf{A} \mathbf{x} = (\mathbf{x}^* \mathbf{A}^* \mathbf{y})^* = (\mathbf{x}^* \mathbf{A} \mathbf{y})^* = (\mu \mathbf{x}^* \mathbf{y})^* = \mu \mathbf{y}^* \mathbf{x},$$

using that  $\mu$  is real. Since  $\lambda \neq \mu$  it follows that  $\mathbf{y}^* \mathbf{x} = 0$ , which means that  $\mathbf{x}$  and  $\mathbf{y}$  are orthogonal.  $\square$

It is now easy to find the eigenpairs of the 2D test matrix and determine when it is positive definite.

**Theorem 3.10 (Eigenpairs of 2D test matrix)**

For fixed  $m \geq 2$  let  $\mathbf{T}_2$  be the matrix given by (3.9) and let  $h = 1/(m+1)$ .

1. We have  $\mathbf{T}_2 \mathbf{x}_{j,k} = \lambda_{j,k} \mathbf{x}_{j,k}$  for  $j, k = 1, \dots, m$ , where

$$\mathbf{x}_{j,k} = \mathbf{s}_j \otimes \mathbf{s}_k, \quad (3.18)$$

$$\mathbf{s}_j = [\sin(j\pi h), \sin(2j\pi h), \dots, \sin(mj\pi h)]^T, \quad (3.19)$$

$$\lambda_{j,k} = 2d + 2a \cos(j\pi h) + 2a \cos(k\pi h). \quad (3.20)$$

2. The eigenvectors are orthogonal

$$\mathbf{x}_{j,k}^T \mathbf{x}_{p,q} = \frac{1}{4h^2} \delta_{j,p} \delta_{k,q}, \quad j, k, p, q = 1, \dots, m. \quad (3.21)$$

3.  $\mathbf{T}_2$  is symmetric positive definite if  $d > 0$  and  $d \geq 2|a|$ .

*Proof.* By Theorem 3.7 the eigenvalues of  $\mathbf{T}_2 = \mathbf{T}_1 \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}_1$  are sums of eigenvalues of  $\mathbf{T}_1$  and the eigenvectors are Kronecker products of the eigenvectors of  $\mathbf{T}_1$ . Part 1 now follows from Lemma 3.8. Using the transpose rule, the mixed product rule and (3.17) we find for  $j, k, p, q = 1, \dots, m$

$$(\mathbf{s}_j \otimes \mathbf{s}_k)^T (\mathbf{s}_p \otimes \mathbf{s}_q) = (\mathbf{s}_j^T \otimes \mathbf{s}_k^T) (\mathbf{s}_p \otimes \mathbf{s}_q) = (\mathbf{s}_j^T \mathbf{s}_p) \otimes (\mathbf{s}_k^T \mathbf{s}_q) = \frac{1}{4h^2} \delta_{j,p} \delta_{k,q}$$

and part 2 follows. Since  $\mathbf{T}_2$  is symmetric, part 3 will follow if the eigenvalues are positive. But this is true if  $d > 0$  and  $d \geq 2|a|$ . Thus  $\mathbf{T}_2$  is positive definite.  $\square$

**Exercise 3.11 (2. derivative matrix is positive definite)**

Write down the eigenvalues of  $\mathbf{T} = \text{tridiag}(-1, 2, -1)$  using Lemma 3.8 and conclude that  $\mathbf{T}$  is symmetric positive definite.

**Exercise 3.12 (1D test matrix is positive definite?)**

Use Lemma 3.8 to show that the matrix  $\mathbf{T}_1 := \text{tridiag}(a, d, a) \in \mathbb{R}^{n \times n}$  is symmetric positive definite if  $d > 0$  and  $d \geq 2|a|$ .

**Exercise 3.13 (Eigenvalues for 2D test matrix of order 4)**

For  $m = 2$  the matrix (3.9) is given by

$$\mathbf{A} = \begin{bmatrix} 2d & a & a & 0 \\ a & 2d & 0 & a \\ a & 0 & 2d & a \\ 0 & a & a & 2d \end{bmatrix}.$$

Show that  $\lambda = 2a + 2d$  is an eigenvalue corresponding to the eigenvector  $\mathbf{x} = [1, 1, 1, 1]^T$ . Verify that apart from a scaling of the eigenvector this agrees with (3.20) and (3.19) for  $j = k = 1$  and  $m = 2$ .

**Exercise 3.14 (Nine point scheme for Poisson problem)**

Consider the following 9 point difference approximation to the Poisson problem  $-\Delta u = f$ ,  $u = 0$  on the boundary of the unit square (cf. (3.1))

$$\begin{aligned} \text{(a)} \quad & -(\square_h v)_{j,k} = (\mu f)_{j,k} & j, k = 1, \dots, m \\ \text{(b)} \quad & 0 = v_{0,k} = v_{m+1,k} = v_{j,0} = v_{j,m+1}, & j, k = 0, 1, \dots, m+1, \\ \text{(c)} \quad & -(\square_h v)_{j,k} = [20v_{j,k} - 4v_{j-1,k} - 4v_{j,k-1} - 4v_{j+1,k} - 4v_{j,k+1} & (3.22) \\ & \quad - v_{j-1,k-1} - v_{j+1,k-1} - v_{j-1,k+1} - v_{j+1,k+1}]/(6h^2), \\ \text{(d)} \quad & (\mu f)_{j,k} = [8f_{j,k} + f_{j-1,k} + f_{j,k-1} + f_{j+1,k} + f_{j,k+1}]/12. \end{aligned}$$

a) Write down the 4-by-4 system we obtain for  $m = 2$ .

b) Find  $v_{j,k}$  for  $j, k = 1, 2$ , if  $f(x, y) = 2\pi^2 \sin(\pi x) \sin(\pi y)$  and  $m = 2$ . Answer:  $v_{j,k} = 5\pi^2/66$ .

It can be shown that (3.22) defines an  $O(h^4)$  approximation to (3.1).

**Exercise 3.15 (Matrix equation for nine point scheme)**

Consider the nine point difference approximation to (3.1) given by (3.22) in Problem 3.14.

a) Show that (3.22) is equivalent to the matrix equation

$$\mathbf{TV} + \mathbf{VT} - \frac{1}{6}\mathbf{TVT} = h^2\mu\mathbf{F}. \quad (3.23)$$

Here  $\mu\mathbf{F}$  has elements  $(\mu f)_{j,k}$  given by (3.22d).



- b) Show that the standard form of the matrix equation (3.23) is  $\mathbf{Ax} = \mathbf{b}$ , where  $\mathbf{A} = \mathbf{T} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T} - \frac{1}{6}\mathbf{T} \otimes \mathbf{T}$ ,  $\mathbf{x} = \text{vec}(\mathbf{V})$ , and  $\mathbf{b} = h^2 \text{vec}(\mu \mathbf{F})$ .

### Exercise 3.16 (Biharmonic equation)

Consider the biharmonic equation

$$\begin{aligned} \Delta^2 u(s, t) &:= \Delta(\Delta u(s, t)) = f(s, t) & (s, t) \in \Omega, \\ u(s, t) = 0, \quad \Delta u(s, t) &= 0 & (s, t) \in \partial\Omega. \end{aligned} \quad (3.24)$$

Here  $\Omega$  is the open unit square. The condition  $\Delta u = 0$  is called the Navier boundary condition. Moreover,  $\Delta^2 u = u_{xxxx} + 2u_{xxyy} + u_{yyyy}$ .

- a) Let  $v = -\Delta u$ . Show that (3.24) can be written as a system

$$\begin{aligned} -\Delta v(s, t) &= f(s, t) & (s, t) \in \Omega \\ -\Delta u(s, t) &= v(s, t) & (s, t) \in \Omega \\ u(s, t) &= v(s, t) = 0 & (s, t) \in \partial\Omega. \end{aligned} \quad (3.25)$$

- b) Discretizing, using (3.3), with  $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m \times m}$ ,  $h = 1/(m+1)$ , and  $\mathbf{F} = (f(jh, kh))_{j,k=1}^m$  we get two matrix equations

$$\mathbf{TV} + \mathbf{VT} = h^2 \mathbf{F}, \quad \mathbf{TU} + \mathbf{UT} = h^2 \mathbf{V}.$$

Show that

$$(\mathbf{T} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}) \text{vec}(\mathbf{V}) = h^2 \text{vec}(\mathbf{F}), \quad (\mathbf{T} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}) \text{vec}(\mathbf{U}) = h^2 \text{vec}(\mathbf{V}).$$

and hence  $\mathbf{A} = (\mathbf{T} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T})^2$  is the matrix for the standard form of the discrete biharmonic equation.

- c) Show that with  $n = m^2$  the vector form and standard form of the systems in b) can be written

$$\mathbf{T}^2 \mathbf{U} + 2\mathbf{TUT} + \mathbf{UT}^2 = h^4 \mathbf{F} \quad \text{and} \quad \mathbf{Ax} = \mathbf{b}, \quad (3.26)$$

where  $\mathbf{A} = \mathbf{T}^2 \otimes \mathbf{I} + 2\mathbf{T} \otimes \mathbf{T} + \mathbf{I} \otimes \mathbf{T}^2 \in \mathbb{R}^{n \times n}$ ,  $\mathbf{x} = \text{vec}(\mathbf{U})$ , and  $\mathbf{b} = h^4 \text{vec}(\mathbf{F})$ .

- d) Determine the eigenvalues and eigenvectors of the matrix  $\mathbf{A}$  in c) and show that it is symmetric positive definite. Also determine the bandwidth of  $\mathbf{A}$ .
- e) Suppose we want to solve the standard form equation  $\mathbf{Ax} = \mathbf{b}$ . We have two representations for the matrix  $\mathbf{A}$ , the product one in b) and the one in c). Which one would you prefer for the basis of an algorithm? Why?

## 3.4 Review Questions

3.4.1 Consider the Poisson matrix.

- Write this matrix as a Kronecker sum,
- how are its eigenvalues and eigenvectors related to the second derivative matrix?
- is it symmetric? positive definite?

3.4.2 What are the eigenpairs of  $T_1 := \text{tridiagonal}(a, d, a)$ ?

3.4.3 What are the inverse and transpose of a Kronecker product?

- 3.4.4
- give an economical general way to solve the linear system  $(\mathbf{A} \otimes \mathbf{B}) \text{vec}(\mathbf{V}) = \text{vec}(\mathbf{F})$ ?
  - Same for  $(\mathbf{A} \otimes \mathbf{I}_s + \mathbf{I}_r \otimes \mathbf{B}) \text{vec}(\mathbf{V}) = \text{vec}(\mathbf{F})$ .

## Chapter 4

# Fast Direct Solution of a Large Linear System

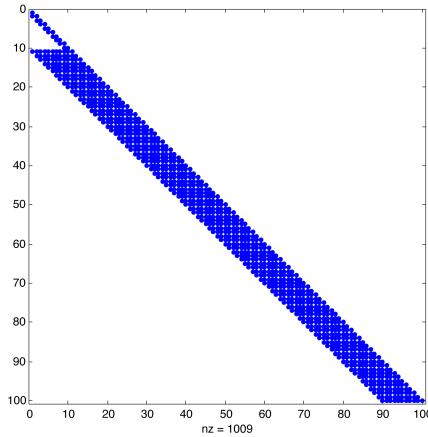
### 4.1 Algorithms for a Banded Positive Definite System

In this chapter we present a fast method for solving  $\mathbf{Ax} = \mathbf{b}$ , where  $\mathbf{A}$  is the Poisson matrix (3.7). Thus, for  $n = 9$

$$\mathbf{A} = \left[ \begin{array}{ccc|ccc|ccc} 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 & 0 \\ -1 & 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & -1 & 4 & 0 & 0 & -1 & 0 & 0 & 0 \\ \hline -1 & 0 & 0 & 4 & -1 & 0 & -1 & 0 & 0 \\ 0 & -1 & 0 & -1 & 4 & -1 & 0 & -1 & 0 \\ 0 & 0 & -1 & 0 & -1 & 4 & 0 & 0 & -1 \\ \hline 0 & 0 & 0 & -1 & 0 & 0 & 4 & -1 & 0 \\ 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 & -1 \\ 0 & 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 \end{array} \right]$$
$$= \begin{bmatrix} \mathbf{T} + 2\mathbf{I} & -\mathbf{I} & \mathbf{0} \\ -\mathbf{I} & \mathbf{T} + 2\mathbf{I} & -\mathbf{I} \\ \mathbf{0} & -\mathbf{I} & \mathbf{T} + 2\mathbf{I} \end{bmatrix},$$

where  $\mathbf{T} = \text{tridiag}(-1, 2, -1)$ . For the matrix  $\mathbf{A}$  we know by now that

1. It is symmetric positive definite.
2. It is banded.
3. It is block-tridiagonal.



**Figure 4.1.** Fill-in in the Cholesky factor of the Poisson matrix ( $n = 100$ ).

4. We know the eigenvalues and eigenvectors of  $\mathbf{A}$ .
5. The eigenvectors are orthogonal.

#### 4.1.1 Cholesky factorization

Since  $\mathbf{A}$  is symmetric positive definite we can use the Cholesky factorization  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$ , with  $\mathbf{L}$  lower triangular, to solve  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Since  $\mathbf{A}$  is banded with bandwidth  $d = \sqrt{n}$  the matrix  $\mathbf{L}$  has bandwidth  $d = \sqrt{n}$  (cf. Lemma 2.46) and the complexity of this factorization is  $O(nd^2) = O(n^2)$ . We need to store  $\mathbf{A}$ , and this can be done in sparse form.

The nonzero elements in  $\mathbf{L}$  are shown in Figure 4.1 for  $n = 100$ . Note that most of the zeros between the diagonals in  $\mathbf{A}$  have become nonzero in  $\mathbf{L}$ . This is known as **fill-in**.

#### 4.1.2 Block LU factorization of a block tridiagonal matrix

The Poisson matrix has a block tridiagonal structure. Consider finding the block LU factorization of a block tridiagonal matrix. We are looking for a factorization of the form

$$\begin{bmatrix} D_1 & C_1 & & & \\ A_1 & D_2 & C_2 & & \\ & \ddots & \ddots & \ddots & \\ & & A_{m-2} & D_{m-1} & C_{m-1} \\ & & & A_{m-1} & D_m \end{bmatrix} = \begin{bmatrix} I & & & & \\ L_1 & I & & & \\ & \ddots & \ddots & \ddots & \\ & & L_{m-1} & I & \\ & & & & I \end{bmatrix} \begin{bmatrix} U_1 & C_1 & & & \\ & \ddots & \ddots & \ddots & \\ & & U_{m-1} & C_{m-1} & \\ & & & & U_m \end{bmatrix}. \quad (4.1)$$

Here  $\mathbf{D}_1, \dots, \mathbf{D}_m$  and  $\mathbf{U}_1, \dots, \mathbf{U}_m$  are square matrices while  $\mathbf{A}_1, \dots, \mathbf{A}_{m-1}, \mathbf{L}_1, \dots, \mathbf{L}_{m-1}$  and  $\mathbf{C}_1, \dots, \mathbf{C}_{m-1}$  can be rectangular.

Using block multiplication the formulas (1.4) generalize to

$$\mathbf{U}_1 = \mathbf{D}_1, \quad \mathbf{L}_k = \mathbf{A}_k \mathbf{U}_k^{-1}, \quad \mathbf{U}_{k+1} = \mathbf{D}_{k+1} - \mathbf{L}_k \mathbf{C}_k, \quad k = 1, 2, \dots, m-1. \quad (4.2)$$

To solve the system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  we partition  $\mathbf{b}$  conformally with  $\mathbf{A}$  in the form  $\mathbf{b}^T = [\mathbf{b}_1^T, \dots, \mathbf{b}_m^T]$ . The formulas for solving  $\mathbf{L}\mathbf{y} = \mathbf{b}$  and  $\mathbf{U}\mathbf{x} = \mathbf{y}$  are as follows:

$$\begin{aligned} \mathbf{y}_1 &= \mathbf{b}_1, & \mathbf{y}_k &= \mathbf{b}_k - \mathbf{L}_{k-1} \mathbf{y}_{k-1}, & k &= 2, 3, \dots, m, \\ \mathbf{x}_m &= \mathbf{U}_m^{-1} \mathbf{y}_m, & \mathbf{x}_k &= \mathbf{U}_k^{-1} (\mathbf{y}_k - \mathbf{C}_k \mathbf{x}_{k+1}), & k &= m-1, \dots, 2, 1. \end{aligned} \quad (4.3)$$

The solution is then  $\mathbf{x}^T = [\mathbf{x}_1^T, \dots, \mathbf{x}_m^T]$ . To find  $\mathbf{L}_k$  in (4.2) we solve the linear systems  $\mathbf{L}_k \mathbf{U}_k = \mathbf{A}_k$ . Similarly we need to solve a linear system to find  $\mathbf{x}_k$  in (4.3).

The number of arithmetic operations using block factorizations is  $O(n^2)$ , asymptotically the same as for Cholesky factorization. However we only need to store the  $m \times m$  blocks and using matrix operations can be an advantage.

### 4.1.3 Other methods

Other methods include

- Iterative methods, (we study this in Chapters 8 and 9),
- multigrid. See [8],
- fast solvers based on diagonalization and the fast Fourier transform. See Sections 4.2, 4.3.

## 4.2 A Fast Poisson Solver based on Diagonalization

The algorithm we now derive will only require  $O(n^{3/2})$  arithmetic operations and we only need to work with matrices of order  $m$ . Using the fast Fourier transform the number of arithmetic operations can be reduced further to  $O(n \log n)$ .

To start we recall that  $\mathbf{A}\mathbf{x} = \mathbf{b}$  can be written as a matrix equation in the form (cf. (3.5))

$$\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2 \mathbf{F} \quad \text{with} \quad h = 1/(m+1),$$

where  $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m \times m}$  is the second derivative matrix,  $\mathbf{V} = (v_{jk}) \in \mathbb{R}^{m \times m}$  are the unknowns, and  $\mathbf{F} = (f_{jk}) = (f(jh, kh)) \in \mathbb{R}^{m \times m}$  contains function values.

Recall that the eigenpairs of  $T$  are given by

$$\begin{aligned} T\mathbf{s}_j &= \lambda_j \mathbf{s}_j, \quad j = 1, \dots, m, \\ \mathbf{s}_j &= [\sin(j\pi h), \sin(2j\pi h), \dots, \sin(mj\pi h)]^T, \\ \lambda_j &= 2 - 2\cos(j\pi h) = 4\sin^2(j\pi h/2), \quad h = 1/(m+1), \\ \mathbf{s}_j^T \mathbf{s}_k &= \delta_{jk}/(2h) \text{ for all } j, k. \end{aligned}$$

Let

$$\mathbf{S} := [\mathbf{s}_1, \dots, \mathbf{s}_m] = [\sin(jk\pi h)]_{j,k=1}^m \in \mathbb{R}^{m \times m}, \quad \mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_m). \quad (4.4)$$

Then

$$T\mathbf{S} = [T\mathbf{s}_1, \dots, T\mathbf{s}_m] = [\lambda_1 \mathbf{s}_1, \dots, \lambda_m \mathbf{s}_m] = \mathbf{S}\mathbf{D}, \quad \mathbf{S}^2 = \mathbf{S}^T \mathbf{S} = \frac{1}{2h} \mathbf{I}.$$

Define  $\mathbf{X} \in \mathbb{R}^{m \times m}$  by  $\mathbf{V} = \mathbf{S}\mathbf{X}\mathbf{S}$ , where  $\mathbf{V}$  is the solution of  $T\mathbf{V} + \mathbf{V}T = h^2\mathbf{F}$ . Then

$$\begin{aligned} T\mathbf{V} + \mathbf{V}T &= h^2\mathbf{F} \\ \stackrel{\mathbf{V}=\mathbf{S}\mathbf{X}\mathbf{S}}{\iff} T\mathbf{S}\mathbf{X}\mathbf{S} + \mathbf{S}\mathbf{X}\mathbf{S}T &= h^2\mathbf{F} \\ \stackrel{\mathbf{S}^{(\cdot)}\mathbf{S}}{\iff} \mathbf{S}T\mathbf{S}\mathbf{X}\mathbf{S}^2 + \mathbf{S}^2\mathbf{X}\mathbf{S}T\mathbf{S} &= h^2\mathbf{S}\mathbf{F}\mathbf{S} = h^2\mathbf{G} \\ \stackrel{T\mathbf{S}=\mathbf{S}\mathbf{D}}{\iff} \mathbf{S}^2\mathbf{D}\mathbf{X}\mathbf{S}^2 + \mathbf{S}^2\mathbf{X}\mathbf{S}^2\mathbf{D} &= h^2\mathbf{G} \\ \stackrel{\mathbf{S}^2=\mathbf{I}/(2h)}{\iff} \mathbf{D}\mathbf{X} + \mathbf{X}\mathbf{D} &= 4h^4\mathbf{G}. \end{aligned}$$

Since  $\mathbf{D}$  is diagonal, the equation  $\mathbf{D}\mathbf{X} + \mathbf{X}\mathbf{D} = 4h^4\mathbf{G}$ , is easy to solve. For the  $j, k$  element we find

$$(\mathbf{D}\mathbf{X} + \mathbf{X}\mathbf{D})_{j,k} = \sum_{\ell=1}^m d_{j,\ell} x_{\ell,k} + \sum_{\ell=1}^m x_{j,\ell} d_{\ell,k} = \lambda_j x_{j,k} + \lambda_k x_{j,k}$$

so that for all  $j, k$

$$x_{jk} = 4h^4 g_{jk} / (\lambda_j + \lambda_k) = h^4 g_{jk} / (\sigma_j + \sigma_k), \quad \sigma_j := \lambda_j/4 = \sin^2(j\pi h/2).$$

Thus to find  $\mathbf{V}$  we compute

1.  $\mathbf{G} = \mathbf{S}\mathbf{F}\mathbf{S}$ ,
2.  $x_{j,k} = h^4 g_{j,k} / (\sigma_j + \sigma_k)$ ,  $j, k = 1, \dots, m$ ,
3.  $\mathbf{V} = \mathbf{S}\mathbf{X}\mathbf{S}$ .

We can compute  $m\mathbf{X}$ ,  $\mathbf{S}$  and the  $\sigma$ 's without using loops. Using outer products, element by element division, and raising a matrix element by element to a power we find

$$\mathbf{X} = h^4 \mathbf{G} / \mathbf{M}, \text{ where } \mathbf{M} := \begin{bmatrix} \sigma_1 \\ \vdots \\ \sigma_m \end{bmatrix} [1, \dots, 1] + \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} [\sigma_1, \dots, \sigma_m],$$

$$\mathbf{S} = \sin\left(\pi h \begin{bmatrix} \frac{1}{2} \\ \vdots \\ m \end{bmatrix} [1 \ 2 \ \dots \ m]\right), \quad \boldsymbol{\sigma} = \sin\left(\frac{\pi h}{2} \begin{bmatrix} \frac{1}{2} \\ \vdots \\ m \end{bmatrix}\right) \wedge 2.$$

We now get the following algorithm to solve numerically the Poisson problem  $-\Delta u = f$  on  $\Omega = (0, 1)^2$  and  $u = 0$  on  $\partial\Omega$  using the 5-point scheme, i. e., let  $m \in \mathbb{N}$ ,  $h = 1/(m+1)$ , and  $\mathbf{F} = (f(jh, kh)) \in \mathbb{R}^{m \times m}$ . We compute  $\mathbf{V} \in \mathbb{R}^{(m+2) \times (m+2)}$  using diagonalization of  $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m \times m}$ .

**Algorithm 4.1 (Fast Poisson solver)**

```

1 function V=fastpoisson(F)
2 %function V=fastpoisson(F)
3 m=length(F); h=1/(m+1); hv=pi*h*(1:m)';
4 sigma=sin(hv/2).^2;
5 S=sin(hv*(1:m));
6 G=S*F*S;
7 X=h^4*G./(sigma*ones(1,m)+ones(m,1)*sigma)';
8 V=zeros(m+2,m+2);
9 V(2:m+1,2:m+1)=S*X*S;

```

The formulas are fully vectorized. Since the 6th line in Algorithm 4.1 only requires  $O(m^2)$  arithmetic operations the complexity of this algorithm is for large  $m$  determined by the 4  $m$ -by- $m$  matrix multiplications and is given by  $O(4 \times 2m^3) = O(8n^{3/2})$ .<sup>8</sup> The method is very fast and will be used as a preconditioner for a more complicated problem in Chapter 9. In 2012 it took about 0.2 seconds on a laptop to find the  $10^6$  unknowns  $v_{j,k}$  on a  $1000 \times 1000$  grid.

### 4.3 A Fast Poisson Solver based on the discrete sine and Fourier transforms

In Algorithm 4.1 we need to compute the product of the sine matrix  $\mathbf{S} \in \mathbb{R}^{m \times m}$  given by (4.4) and a matrix  $\mathbf{A} \in \mathbb{R}^{m \times m}$ . Since the matrices are  $m$ -by- $m$  this will normally require  $O(m^3)$  operations. In this section we show that it is possible to calculate the products  $\mathbf{SA}$  and  $\mathbf{AS}$  in  $O(m^2 \log_2 m)$  operations.

<sup>8</sup>It is possible to compute  $\mathbf{V}$  using only two matrix multiplications and hence reduce the complexity to  $O(4n^{3/2})$ . This is detailed in Problem 4.8.

We need to discuss certain transforms known as the **discrete sine transform**, the **discrete Fourier transform** and the **fast Fourier transform**. In addition we have the **discrete cosine transform** which will not be discussed here. These transforms are of independent interest. They have applications to signal processing and image analysis, and are often used when one is dealing with discrete samples of data on a computer.

### 4.3.1 The discrete sine transform (DST)

Given  $\mathbf{v} = [v_1, \dots, v_m]^T \in \mathbb{R}^m$  we say that the vector  $\mathbf{w} = [w_1, \dots, w_m]^T$  given by

$$w_j = \sum_{k=1}^m \sin\left(\frac{jk\pi}{m+1}\right) v_k, \quad j = 1, \dots, m$$

is the **discrete sine transform** (DST) of  $\mathbf{v}$ . In matrix form we can write the DST as the matrix times vector  $\mathbf{w} = \mathbf{S}\mathbf{v}$ , where  $\mathbf{S}$  is the sine matrix given by (4.4). We can then identify the matrix  $\mathbf{B} = \mathbf{S}\mathbf{A}$  as the DST of  $\mathbf{A} \in \mathbb{R}^{m,n}$ , i.e. as the DST of the columns of  $\mathbf{A}$ . The product  $\mathbf{B} = \mathbf{A}\mathbf{S}$  can also be interpreted as a DST. Indeed, since  $\mathbf{S}$  is symmetric we have  $\mathbf{B} = (\mathbf{S}\mathbf{A}^T)^T$  which means that  $\mathbf{B}$  is the transpose of the DST of the rows of  $\mathbf{A}$ . It follows that we can compute the unknowns  $\mathbf{V}$  in Algorithm 4.1 by carrying out discrete sine transforms on 4  $m$ -by- $m$  matrices in addition to the computation of  $\mathbf{X}$ .

### 4.3.2 The discrete Fourier transform (DFT)



Jean Baptiste Joseph Fourier, 1768 - 1830.

The fast computation of the DST is based on its relation to the discrete Fourier transform (DFT) and the fact that the DFT can be computed by a technique known as the fast Fourier transform (FFT). To define the DFT let for  $N \in \mathbb{N}$

$$\omega_N = \exp^{-2\pi i/N} = \cos(2\pi/N) - i \sin(2\pi/N), \quad (4.5)$$

where  $i = \sqrt{-1}$  is the imaginary unit. Given  $\mathbf{y} = [y_1, \dots, y_N]^T \in \mathbb{R}^N$  we say that



$\mathbf{z} = [z_1, \dots, z_N]^T$  given by

$$\mathbf{z} = \mathbf{F}_N \mathbf{y}, \quad z_{j+1} = \sum_{k=0}^{N-1} \omega_N^{jk} y_{k+1}, \quad j = 0, \dots, N-1$$

is the **discrete Fourier transform** (DFT) of  $\mathbf{y}$ . We can write this as a matrix times vector product  $\mathbf{z} = \mathbf{F}_N \mathbf{y}$ , where the **Fourier matrix**  $\mathbf{F}_N \in \mathbb{C}^{N \times N}$  has elements  $\omega_N^{jk}$ ,  $j, k = 0, 1, \dots, N-1$ . For a matrix we say that  $\mathbf{B} = \mathbf{F}_N \mathbf{A}$  is the DFT of  $\mathbf{A}$ .

As an example, since

$$\omega_4 = \exp^{-2\pi i/4} = \cos(\pi/2) - i \sin(\pi/2) = -i$$

we find  $\omega_4^2 = (-i)^2 = -1$ ,  $\omega_4^3 = (-i)(-1) = i$ ,  $\omega_4^4 = (-1)^2 = 1$ ,  $\omega_4^6 = i^2 = -1$ ,  $\omega_4^9 = i^3 = -i$ , and so

$$\mathbf{F}_4 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & \omega_4 & \omega_4^2 & \omega_4^3 \\ 1 & \omega_4^2 & \omega_4^4 & \omega_4^6 \\ 1 & \omega_4^3 & \omega_4^6 & \omega_4^9 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -i & -1 & i \\ 1 & -1 & 1 & -1 \\ 1 & i & -1 & -i \end{bmatrix}. \quad (4.6)$$

The following lemma shows how the discrete sine transform of order  $m$  can be computed from the discrete Fourier transform of order  $2m+2$ . We recall that for any complex number  $w$

$$\sin w = \frac{e^{iw} - e^{-iw}}{2i}.$$

#### Lemma 4.2 (Sine transform as Fourier transform)

Given a positive integer  $m$  and a vector  $\mathbf{x} \in \mathbb{R}^m$ . Component  $k$  of  $\mathbf{Sx}$  is equal to  $i/2$  times component  $k+1$  of  $\mathbf{F}_{2m+2}\mathbf{z}$  where

$$\mathbf{z}^T = [0, \mathbf{x}^T, 0, -\mathbf{x}_B^T] \in \mathbb{R}^{2m+2}, \quad \mathbf{x}_B^T := [x_m, \dots, x_2, x_1].$$

In symbols

$$(\mathbf{Sx})_k = \frac{i}{2} (\mathbf{F}_{2m+2}\mathbf{z})_{k+1}, \quad k = 1, \dots, m.$$

**Proof.** Let  $\omega = \omega_{2m+2} = e^{-2\pi i/(2m+2)} = e^{-\pi i/(m+1)}$ . We note that

$$\omega^{jk} = e^{-\pi ijk/(m+1)}, \quad \omega^{(2m+2-j)k} = e^{-2\pi i e^{\pi ijk/(m+1)}} = e^{\pi ijk/(m+1)}.$$

Component  $k + 1$  of  $\mathbf{F}_{2m+2}\mathbf{z}$  is then given by

$$\begin{aligned} (\mathbf{F}_{2m+2}\mathbf{z})_{k+1} &= \sum_{j=0}^{2m-1} \omega^{jk} z_{j+1} = \sum_{j=1}^m x_j \omega^{jk} - \sum_{j=1}^m x_j \omega^{(2m+2-j)k} \\ &= \sum_{j=1}^m x_j (e^{-\pi i j k / (m+1)} - e^{\pi i j k / (m+1)}) \\ &= -2i \sum_{j=1}^m x_j \sin\left(\frac{jk\pi}{m+1}\right) = -2i(\mathbf{S}_m \mathbf{x})_k. \end{aligned}$$

Dividing both sides by  $-2i$  and noting  $-1/(2i) = -i/(2i^2) = i/2$ , proves the lemma.  $\square$

It follows that we can compute the DST of length  $m$  by extracting  $m$  components from the DFT of length  $N = 2m + 2$ .

### 4.3.3 The fast Fourier transform (FFT)

From a linear algebra viewpoint the fast Fourier transform is a quick way to compute the matrix-vector product  $\mathbf{F}_N \mathbf{y}$ . Suppose  $N$  is even. The key to the FFT is a connection between  $\mathbf{F}_N$  and  $\mathbf{F}_{N/2}$  which makes it possible to compute the FFT of order  $N$  as two FFT's of order  $N/2$ . By repeating this process we can reduce the number of arithmetic operations to compute a DFT from  $O(N^2)$  to  $O(N \log_2 N)$ .

Suppose  $N$  is even. The connection between  $\mathbf{F}_N$  and  $\mathbf{F}_{N/2}$  involves a permutation matrix  $\mathbf{P}_N \in \mathbb{R}^{N \times N}$  given by

$$\mathbf{P}_N = [\mathbf{e}_1, \mathbf{e}_3, \dots, \mathbf{e}_{N-1}, \mathbf{e}_2, \mathbf{e}_4, \dots, \mathbf{e}_N],$$

where the  $\mathbf{e}_k = (\delta_{j,k})$  are unit vectors. If  $\mathbf{A}$  is a matrix with  $N$  columns  $[\mathbf{a}_1, \dots, \mathbf{a}_N]$  then

$$\mathbf{A}\mathbf{P}_N = [\mathbf{a}_1, \mathbf{a}_3, \dots, \mathbf{a}_{N-1}, \mathbf{a}_2, \mathbf{a}_4, \dots, \mathbf{a}_N],$$

i.e. post multiplying  $\mathbf{A}$  by  $\mathbf{P}_N$  permutes the columns of  $\mathbf{A}$  so that all the odd-indexed columns are followed by all the even-indexed columns. For example we have from (4.6)

$$\mathbf{P}_4 = [\mathbf{e}_1 \ \mathbf{e}_3 \ \mathbf{e}_2 \ \mathbf{e}_4] = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad \mathbf{F}_4 \mathbf{P}_4 = \left[ \begin{array}{cc|cc} 1 & 1 & 1 & 1 \\ 1 & -1 & -i & i \\ \hline 1 & 1 & -1 & -1 \\ 1 & -1 & i & -i \end{array} \right],$$

where we have indicated a certain block structure of  $\mathbf{F}_4\mathbf{P}_4$ . These blocks can be related to the 2-by-2 matrix  $\mathbf{F}_2$ . We define the diagonal scaling matrix  $\mathbf{D}_2$  by

$$\mathbf{D}_2 = \text{diag}(1, \omega_4) = \begin{bmatrix} 1 & 0 \\ 0 & -i \end{bmatrix}.$$

Since  $\omega_2 = \exp^{-2\pi i/2} = -1$  we find

$$\mathbf{F}_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}, \quad \mathbf{D}_2\mathbf{F}_2 = \begin{bmatrix} 1 & 1 \\ -i & i \end{bmatrix},$$

and we see that

$$\mathbf{F}_4\mathbf{P}_4 = \left[ \begin{array}{c|c} \mathbf{F}_2 & \mathbf{D}_2\mathbf{F}_2 \\ \hline \mathbf{F}_2 & -\mathbf{D}_2\mathbf{F}_2 \end{array} \right].$$

This result holds in general.

### Theorem 4.3 (Fast Fourier transform)

If  $N = 2m$  is even then

$$\mathbf{F}_{2m}\mathbf{P}_{2m} = \left[ \begin{array}{c|c} \mathbf{F}_m & \mathbf{D}_m\mathbf{F}_m \\ \hline \mathbf{F}_m & -\mathbf{D}_m\mathbf{F}_m \end{array} \right], \quad (4.7)$$

where

$$\mathbf{D}_m = \text{diag}(1, \omega_N, \omega_N^2, \dots, \omega_N^{m-1}). \quad (4.8)$$

**Proof.** Fix integers  $p, q$  with  $1 \leq p, q \leq m$  and set  $j := p - 1$  and  $k := q - 1$ . Since

$$\omega_m^m = 1, \quad \omega_{2m}^{2k} = \omega_m^k, \quad \omega_{2m}^m = -1, \quad (\mathbf{F}_m)_{p,q} = \omega_m^{jk}, \quad (\mathbf{D}_m\mathbf{F}_m)_{p,q} = \omega_{2m}^j \omega_m^{jk},$$

we find by considering elements in the four sub-blocks in turn

$$\begin{aligned} (\mathbf{F}_{2m}\mathbf{P}_{2m})_{p,q} &= \omega_{2m}^{j(2k)} &= \omega_m^{jk}, \\ (\mathbf{F}_{2m}\mathbf{P}_{2m})_{p+m,q} &= \omega_{2m}^{(j+m)(2k)} &= \omega_m^{(j+m)k} &= \omega_m^{jk}, \\ (\mathbf{F}_{2m}\mathbf{P}_{2m})_{p,q+m} &= \omega_{2m}^{j(2k+1)} &= \omega_{2m}^j \omega_m^{jk}, \\ (\mathbf{F}_{2m}\mathbf{P}_{2m})_{p+m,q+m} &= \omega_{2m}^{(j+m)(2k+1)} &= \omega_{2m}^{j+m} \omega_m^{(j+m)k} &= -\omega_{2m}^j \omega_m^{jk}. \end{aligned}$$

It follows that the four  $m$ -by- $m$  blocks of  $\mathbf{F}_{2m}\mathbf{P}_{2m}$  have the required structure.  $\square$

Using Theorem 4.3 we can carry out the DFT as a block multiplication. Let  $\mathbf{y} \in \mathbb{R}^{2m}$  and set  $\mathbf{w} = \mathbf{P}_{2m}^T \mathbf{y} = [\mathbf{w}_1^T, \mathbf{w}_2^T]^T$ , where

$$\mathbf{w}_1^T = [y_1, y_3, \dots, y_{2m-1}], \quad \mathbf{w}_2^T = [y_2, y_4, \dots, y_{2m}].$$

Then

$$\begin{aligned} \mathbf{F}_{2m}\mathbf{y} &= \mathbf{F}_{2m}\mathbf{P}_{2m}\mathbf{P}_{2m}^T\mathbf{y} = \mathbf{F}_{2m}\mathbf{P}_{2m}\mathbf{w} \\ &= \left[ \begin{array}{c|c} \mathbf{F}_m & \mathbf{D}_m\mathbf{F}_m \\ \hline \mathbf{F}_m & -\mathbf{D}_m\mathbf{F}_m \end{array} \right] \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{q}_1 + \mathbf{q}_2 \\ \mathbf{q}_1 - \mathbf{q}_2 \end{bmatrix}, \end{aligned}$$

where

$$\mathbf{q}_1 = \mathbf{F}_m\mathbf{w}_1, \quad \text{and} \quad \mathbf{q}_2 = \mathbf{D}_m(\mathbf{F}_m\mathbf{w}_2).$$

In order to compute  $\mathbf{F}_{2m}\mathbf{y}$  we need to compute  $\mathbf{F}_m\mathbf{w}_1$  and  $\mathbf{F}_m\mathbf{w}_2$ . Thus, by combining two FFT's of order  $m$  we obtain an FFT of order  $2m$ . If  $n = 2^k$  then this process can be applied recursively as in the following Matlab function:

**Algorithm 4.4 (Recursive FFT)**

```

1 function z=fftrecur(y)
2 %function z=fftrecur(y)
3 y=y(:);
4 n=length(y);
5 if n==1
6     z=y;
7 else
8     q1=fftrecur(y(1:2:n-1))
9     q2=exp(-2*pi*i/n).^((0:n/2-1)') .* fftrec(y(2:2:n))
10    z=[q1+q2; q1-q2];
11 end

```

Statement 3 is included so that the input  $\mathbf{y} \in \mathbb{R}^n$  can be either a row or column vector, while the output  $\mathbf{z}$  is a column vector.

Such a recursive version of FFT is useful for testing purposes, but is much too slow for large problems. A challenge for FFT code writers is to develop nonrecursive versions and also to handle efficiently the case where  $N$  is not a power of two. We refer to [32] for further details.

The complexity of the FFT is given by  $\gamma N \log_2 N$  for some constant  $\gamma$  independent of  $N$ . To show this for the special case when  $N$  is a power of two let  $x_k$  be the complexity (the number of arithmetic operations) when  $N = 2^k$ . Since we need two FFT's of order  $N/2 = 2^{k-1}$  and a multiplication with the diagonal matrix  $\mathbf{D}_{N/2}$ , it is reasonable to assume that  $x_k = 2x_{k-1} + \gamma 2^k$  for some constant  $\gamma$  independent of  $k$ . Since  $x_0 = 0$  we obtain by induction on  $k$  that  $x_k = \gamma k 2^k$ . Indeed, this holds for  $k = 0$  and if  $x_{k-1} = \gamma(k-1)2^{k-1}$  then  $x_k = 2x_{k-1} + \gamma 2^k = 2\gamma(k-1)2^{k-1} + \gamma 2^k = \gamma k 2^k$ . Reasonable implementations of FFT typically have  $\gamma \approx 5$ , see [32].

The efficiency improvement using the FFT to compute the DFT is spectacular for large  $N$ . The direct multiplication  $\mathbf{F}_N\mathbf{y}$  requires  $O(8n^2)$  arithmetic

operations since complex arithmetic is involved. Assuming that the FFT uses  $5N \log_2 N$  arithmetic operations we find for  $N = 2^{20} \approx 10^6$  the ratio

$$\frac{8N^2}{5N \log_2 N} \approx 84000.$$

Thus if the FFT takes one second of computing time and the computing time is proportional to the number of arithmetic operations then the direct multiplication would take something like 84000 seconds or 23 hours.

#### 4.3.4 A poisson solver based on the FFT

We now have all the ingredients to compute the matrix products  $\mathbf{SA}$  and  $\mathbf{AS}$  using FFT's of order  $2m + 2$  where  $m$  is the order of  $\mathbf{S}$  and  $\mathbf{A}$ . This can then be used for quick computation of the exact solution  $\mathbf{V}$  of the discrete Poisson problem in Algorithm 4.1. We first compute  $\mathbf{H} = \mathbf{SF}$  using Lemma 4.2 and  $m$  FFT's, one for each of the  $m$  columns of  $\mathbf{F}$ . We then compute  $\mathbf{G} = \mathbf{HS}$  by  $m$  FFT's, one for each of the rows of  $\mathbf{H}$ . After  $\mathbf{X}$  is determined we compute  $\mathbf{Z} = \mathbf{SX}$  and  $\mathbf{V} = \mathbf{ZS}$  by another  $2m$  FFT's. In total the work amounts to  $4m$  FFT's of order  $2m + 2$ . Since one FFT requires  $O(\gamma(2m + 2) \log_2(2m + 2))$  arithmetic operations the  $4m$  FFT's amount to

$$8\gamma m(m + 1) \log_2(2m + 2) \approx 8\gamma m^2 \log_2 m = 4\gamma n \log_2 n,$$

where  $n = m^2$  is the size of the linear system  $\mathbf{Ax} = \mathbf{b}$  we would be solving if Cholesky factorization was used. This should be compared to the  $O(8n^{3/2})$  arithmetic operations used in Algorithm 4.1 requiring 4 straightforward matrix multiplications with  $\mathbf{S}$ . What is faster will depend heavily on the programming of the FFT and the size of the problem. We refer to [32] for other efficient ways to implement the DST.

##### Exercise 4.5 (Fourier matrix)

Show that the Fourier matrix  $\mathbf{F}_4$  is symmetric, but not Hermitian.

##### Exercise 4.6 (Sine transform as Fourier transform)

Verify Lemma 4.2 directly when  $m = 1$ .

##### Exercise 4.7 (Explicit solution of the discrete Poisson equation)

Show that the exact solution of the discrete Poisson equation (3.4) can be written  $\mathbf{V} = (v_{i,j})_{i,j=1}^m$ , where

$$v_{ij} = \frac{1}{(m + 1)^4} \sum_{p=1}^m \sum_{r=1}^m \sum_{k=1}^m \sum_{l=1}^m \frac{\sin\left(\frac{ip\pi}{m+1}\right) \sin\left(\frac{jr\pi}{m+1}\right) \sin\left(\frac{kp\pi}{m+1}\right) \sin\left(\frac{lr\pi}{m+1}\right)}{\left[\sin\left(\frac{p\pi}{2(m+1)}\right)\right]^2 + \left[\sin\left(\frac{r\pi}{2(m+1)}\right)\right]^2} f_{p,r}.$$

**Exercise 4.8 (Improved version of Algorithm 4.1)**

Algorithm 4.1 involves multiplying a matrix by  $\mathbf{S}$  four times. In this problem we show that it is enough to multiply by  $\mathbf{S}$  two times. We achieve this by diagonalizing only the second  $\mathbf{T}$  in  $\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2\mathbf{F}$ . Let  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_m)$ , where  $\lambda_j = 4\sin^2(j\pi h/2)$ ,  $j = 1, \dots, m$ .

(a) Show that

$$\mathbf{T}\mathbf{X} + \mathbf{X}\mathbf{D} = \mathbf{C}, \text{ where } \mathbf{X} = \mathbf{V}\mathbf{S}, \text{ and } \mathbf{C} = h^2\mathbf{F}\mathbf{S}.$$

(b) Show that

$$(\mathbf{T} + \lambda_j\mathbf{I})\mathbf{x}_j = \mathbf{c}_j \quad j = 1, \dots, m, \quad (4.9)$$

where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]$  and  $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_m]$ . Thus we can find  $\mathbf{X}$  by solving  $m$  linear systems, one for each of the columns of  $\mathbf{X}$ . Recall that a tridiagonal  $m \times m$  system can be solved by (1.4) and (1.5) in  $8m - 7$  arithmetic operations. Give an algorithm to find  $\mathbf{X}$  which only requires  $O(\delta m^2)$  arithmetic operations for some constant  $\delta$  independent of  $m$ .

(c) Describe a method to compute  $\mathbf{V}$  which only requires  $O(4m^3) = O(4n^{3/2})$  arithmetic operations.

(d) Describe a method based on the fast Fourier transform which requires  $O(2\gamma n \log_2 n)$  where  $\gamma$  is the same constant as mentioned at the end of the last section.

**Exercise 4.9 (Fast solution of 9 point scheme)**

Consider the equation

$$\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} - \frac{1}{6}\mathbf{T}\mathbf{V}\mathbf{T} = h^2\mu\mathbf{F},$$

that was derived in Exercise 3.15 for the 9-point scheme. Define the matrix  $\mathbf{X}$  by  $\mathbf{V} = \mathbf{S}\mathbf{X}\mathbf{S} = (x_{j,k})$  where  $\mathbf{V}$  is the solution of (3.23). Show that

$$\mathbf{D}\mathbf{X} + \mathbf{X}\mathbf{D} - \frac{1}{6}\mathbf{D}\mathbf{X}\mathbf{D} = 4h^4\mathbf{G}, \text{ where } \mathbf{G} = \mathbf{S}\mu\mathbf{F}\mathbf{S},$$

and that

$$x_{j,k} = \frac{h^4 g_{j,k}}{\sigma_j + \sigma_k - \frac{2}{3}\sigma_j\sigma_k}, \text{ where } \sigma_j = \sin^2((j\pi h)/2) \text{ for } j, k = 1, 2, \dots, m.$$

Show that  $\sigma_j + \sigma_k - \frac{2}{3}\sigma_j\sigma_k > 0$  for  $j, k = 1, 2, \dots, m$ . Conclude that the matrix  $\mathbf{A}$  in Exercise 3.15 b) is symmetric positive definite and that (3.22) always has a solution  $\mathbf{V}$ .

**Exercise 4.10 (Algorithm for fast solution of 9 point scheme)**

Derive an algorithm for solving (3.22) which for large  $m$  requires essentially the same number of operations as in Algorithm 4.1. (We assume that  $\mu\mathbf{F}$  already has been formed).

**Exercise 4.11 (Fast solution of biharmonic equation)**

For the biharmonic problem we derived in Exercise 3.16 the equation

$$\mathbf{T}^2\mathbf{U} + 2\mathbf{T}\mathbf{U}\mathbf{T} + \mathbf{U}\mathbf{T}^2 = h^4\mathbf{F}.$$

Define the matrix  $\mathbf{X} = (x_{j,k})$  by  $\mathbf{U} = \mathbf{S}\mathbf{X}\mathbf{S}$  where  $\mathbf{U}$  is the solution of (3.26). Show that

$$\mathbf{D}^2\mathbf{X} + 2\mathbf{D}\mathbf{X}\mathbf{D} + \mathbf{X}\mathbf{D}^2 = 4h^6\mathbf{G}, \text{ where } \mathbf{G} = \mathbf{S}\mathbf{F}\mathbf{S},$$

and that

$$x_{j,k} = \frac{h^6 g_{j,k}}{4(\sigma_j + \sigma_k)^2}, \text{ where } \sigma_j = \sin^2((j\pi h)/2) \text{ for } j, k = 1, 2, \dots, m.$$

**Exercise 4.12 (Algorithm for fast solution of biharmonic equation)**

Use Exercise 4.11 to derive an algorithm

```
function U=simplefastbiharmonic(F)
```

which requires only  $O(\delta n^{3/2})$  operations to find  $\mathbf{U}$  in Problem 3.16. Here  $\delta$  is some constant independent of  $n$ .

**Exercise 4.13 (Check algorithm for fast solution of biharmonic equation)**

In Exercise 4.12 compute the solution  $\mathbf{U}$  corresponding to  $\mathbf{F} = \text{ones}(m, m)$ . For some small  $m$ 's check that you get the same solution obtained by solving the standard form  $\mathbf{A}\mathbf{x} = \mathbf{b}$  in (3.26). You can use  $\mathbf{x} = \mathbf{A}\backslash\mathbf{b}$  for solving  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Use  $\mathbf{F}(\cdot)$  to vectorize a matrix and  $\text{reshape}(\mathbf{x}, m, m)$  to turn a vector  $\mathbf{x} \in \mathbb{R}^{m^2}$  into an  $m \times m$  matrix. Use the Matlab command  $\text{surf}(\mathbf{U})$  for plotting  $\mathbf{U}$  for, say,  $m = 50$ . Compare the result with Exercise 4.12 by plotting the difference between both matrices.

**Exercise 4.14 (Fast solution of biharmonic equation using 9 point rule)**

Repeat Exercises 3.16, 4.12 and 4.13 using the nine point rule (3.22) to solve the system (3.25).

## 4.4 Review Questions

4.4.1 Consider the Poisson matrix.

- What is the bandwidth of its Cholesky factor?
- approximately how many arithmetic operations does it take to find the Cholesky factor?
- same question for block LU,
- same question for the fast Poisson solver with and without FFT.

4.4.2 What is the discrete sine transform and discrete Fourier transform of a vector?



**Part II**

**Some Matrix Theory**



## Chapter 5

# Matrix Reduction by Similarity Transformations

A basic problem in numerical linear algebra is to compute eigenvalues and eigenvectors of a matrix  $\mathbf{A}$ . Before attempting to find eigenvalues and eigenvectors of  $\mathbf{A}$  (exceptions are made for certain sparse matrices), it should be reduced by similarity transformations to a simpler form. The contents of this chapter is mainly theoretical, but the results are useful in numerical analysis.

## 5.1 Some Properties of Eigenpairs

We recall (cf. Section 0.7) that  $(\lambda, \mathbf{x})$  is an **eigenpair** for  $\mathbf{A}$  if  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  and  $\mathbf{x}$  is nonzero. The scalar  $\lambda$  is called an **eigenvalue** and  $\mathbf{x}$  is said to be an **eigenvector**. The set of eigenvalues is called the **spectrum** of  $\mathbf{A}$  and is denoted by  $\sigma(\mathbf{A})$ . For example,  $\sigma(\mathbf{I}) = \{1, \dots, 1\} = \{1\}$ . The eigenvalues are the roots of the **characteristic polynomial** given by  $\pi_{\mathbf{A}}(\lambda) := \det(\mathbf{A} - \lambda\mathbf{I})$  for  $\lambda \in \mathbb{C}$ .

### 5.1.1 Transformations of eigenpairs and trace

The following results will be useful.

#### **Theorem 5.1 (Transformations of eigenpairs)**

*Suppose  $(\lambda, \mathbf{x})$  is an eigenpair for  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . Then*

1. *If  $\mathbf{A}$  is nonsingular then  $(\lambda^{-1}, \mathbf{x})$  is an eigenpair for  $\mathbf{A}^{-1}$ .*
2.  *$(\lambda^k, \mathbf{x})$  is an eigenpair for  $\mathbf{A}^k$  for  $k \in \mathbb{N}$ .*
3. *If  $p$  is a polynomial given by  $p(t) = a_0 + a_1t + a_2t^2 + \dots + a_kt^k$  then  $(p(\lambda), \mathbf{x})$  is an eigenpair for the matrix  $p(\mathbf{A}) := a_0\mathbf{I} + a_1\mathbf{A} + a_2\mathbf{A}^2 + \dots + a_k\mathbf{A}^k$ .*
4.  *$\lambda$  is an eigenvalue for  $\mathbf{A}^T$ , in fact  $\pi_{\mathbf{A}^T} = \pi_{\mathbf{A}}$ .*

5.  $\bar{\lambda}$  is an eigenvalue for  $\mathbf{A}^*$ , in fact  $\pi_{\mathbf{A}^*}(\bar{\lambda}) = \overline{\pi_{\mathbf{A}}(\lambda)}$  for all  $\lambda \in \mathbb{C}$ .
6. If  $\mathbf{A} = \begin{bmatrix} \mathbf{B} & \mathbf{C} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}$  is block triangular then  $\pi_{\mathbf{A}} = \pi_{\mathbf{B}} \cdot \pi_{\mathbf{D}}$ .

**Proof.**

1.  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \implies \mathbf{A}^{-1}\mathbf{x} = \lambda^{-1}\mathbf{x}$ .
2. We use induction on  $k$ . The case  $k = 1$  is trivial, and if  $\mathbf{A}^{k-1}\mathbf{x} = \lambda^{k-1}\mathbf{x}$  then  $\mathbf{A}^k\mathbf{x} = \mathbf{A}\mathbf{A}^{k-1}\mathbf{x} = \lambda^{k-1}\mathbf{A}\mathbf{x} = \lambda^k\mathbf{x}$ .
3.  $p(\mathbf{A})\mathbf{x} = \sum_{j=0}^k a_j \mathbf{A}^j \mathbf{x} \stackrel{2.}{=} \sum_{j=0}^k a_j \lambda^j \mathbf{x} = p(\lambda)\mathbf{x}$ .
4. Using Property 3. of determinants we find for any  $\lambda \in \mathbb{C}$

$$\pi_{\mathbf{A}^T}(\lambda) = \det(\mathbf{A}^T - \lambda\mathbf{I}) = \det((\mathbf{A} - \lambda\mathbf{I})^T) = \det(\mathbf{A} - \lambda\mathbf{I}) = \pi_{\mathbf{A}}(\lambda).$$

Thus  $\mathbf{A}^T$  and  $\mathbf{A}$  have the same characteristic polynomial and hence the same eigenvalues.

5. We have  $\pi_{\mathbf{A}^*}(\bar{\lambda}) \stackrel{4.}{=} \pi_{\overline{\mathbf{A}}}(\bar{\lambda}) = \det(\overline{\mathbf{A}} - \bar{\lambda}\mathbf{I}) = \overline{\det(\mathbf{A} - \lambda\mathbf{I})} = \overline{\pi_{\mathbf{A}}(\lambda)}$ . Thus  $\pi_{\mathbf{A}}(\lambda) = 0 \Leftrightarrow \pi_{\mathbf{A}^*}(\bar{\lambda}) = 0$  and the result follows.
6. By Property 6. of determinants

$$\pi_{\mathbf{A}}(\lambda) = \begin{vmatrix} \mathbf{B} - \lambda\mathbf{I} & \mathbf{C} \\ \mathbf{0} & \mathbf{D} - \lambda\mathbf{I} \end{vmatrix} = \det(\mathbf{B} - \lambda\mathbf{I}) \det(\mathbf{D} - \lambda\mathbf{I}) = \pi_{\mathbf{B}}(\lambda) \cdot \pi_{\mathbf{D}}(\lambda).$$

□

There are two important relations between the elements of a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  and its eigenvalues  $\lambda_1, \dots, \lambda_n$ .

**Theorem 5.2 (Sums and products of eigenvalues; trace)**

For any  $\mathbf{A} \in \mathbb{C}^{n \times n}$

$$\text{trace}(\mathbf{A}) = \lambda_1 + \lambda_2 + \dots + \lambda_n, \quad \det(\mathbf{A}) = \lambda_1 \lambda_2 \dots \lambda_n, \quad (5.1)$$

where the **trace** of  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is the sum of its diagonal elements

$$\text{trace}(\mathbf{A}) := a_{11} + a_{22} + \dots + a_{nn}. \quad (5.2)$$

**Proof.** We compare two different expansion of  $\pi_{\mathbf{A}}$ . On the one hand from (29) we find

$$\pi_{\mathbf{A}}(\lambda) = (-1)^n \lambda^n + c_{n-1} \lambda^{n-1} + \dots + c_0,$$

where  $c_{n-1} = (-1)^{n-1} \text{trace}(\mathbf{A})$  and  $c_0 = \pi_{\mathbf{A}}(0) = \det(\mathbf{A})$ . On the other hand

$$\pi_{\mathbf{A}}(\lambda) = (\lambda_1 - \lambda) \cdots (\lambda_n - \lambda) = (-1)^n \lambda^n + d_{n-1} \lambda^{n-1} + \cdots + d_0,$$

where  $d_{n-1} = (-1)^{n-1}(\lambda_1 + \cdots + \lambda_n)$  and  $d_0 = \lambda_1 \cdots \lambda_n$ . Since  $c_j = d_j$  for all  $j$  we obtain (5.1).  $\square$

For a  $2 \times 2$  matrix the characteristic equation takes the convenient form

$$\lambda^2 - \text{trace}(\mathbf{A})\lambda + \det(\mathbf{A}) = 0. \quad (5.3)$$

Thus, if  $\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$  then  $\text{trace}(\mathbf{A}) = 4$ ,  $\det(\mathbf{A}) = 3$  so that  $\pi_{\mathbf{A}}(\lambda) = \lambda^2 - 4\lambda + 3$ .

**Exercise 5.3 (Eigenvalues of an idempotent matrix)**

Let  $\lambda \in \sigma(\mathbf{A})$  where  $\mathbf{A}^2 = \mathbf{A} \in \mathbb{C}^{n \times n}$ . Show that  $\lambda = 0$  or  $\lambda = 1$ . (A matrix is called **idempotent** if  $\mathbf{A}^2 = \mathbf{A}$ ).

**Exercise 5.4 (Eigenvalues of a nilpotent matrix)**

Let  $\lambda \in \sigma(\mathbf{A})$  where  $\mathbf{A}^k = 0$  for some  $k \in \mathbb{N}$ . Show that  $\lambda = 0$ . (A matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  such that  $\mathbf{A}^k = 0$  for some  $k \in \mathbb{N}$  is called **nilpotent**).

**Exercise 5.5 (Eigenvalues of a unitary matrix)**

Let  $\lambda \in \sigma(\mathbf{A})$ , where  $\mathbf{A}^* \mathbf{A} = \mathbf{I}$ . Show that  $|\lambda| = 1$ .

**Exercise 5.6 (Nonsingular approximation of a singular matrix)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is singular. Then we can find  $\epsilon_0 > 0$  such that  $\mathbf{A} + \epsilon \mathbf{I}$  is nonsingular for all  $\epsilon \in \mathbb{C}$  with  $|\epsilon| < \epsilon_0$ . Hint:  $\det(\mathbf{A}) = \lambda_1 \lambda_2 \cdots \lambda_n$ , where  $\lambda_i$  are the eigenvalues of  $\mathbf{A}$ .

**Exercise 5.7 (Companion matrix)**

For  $q_0, \dots, q_{n-1} \in \mathbb{C}$  let  $p(\lambda) = \lambda^n + q_{n-1} \lambda^{n-1} + \cdots + q_0$  be a polynomial of degree  $n$  in  $\lambda$ . We derive two matrices that have  $(-1)^n p$  as its characteristic polynomial.

a) Show that  $p = (-1)^n \pi_{\mathbf{A}}$  where

$$\mathbf{A} = \begin{bmatrix} -q_{n-1} & -q_{n-2} & \cdots & -q_1 & -q_0 \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix}.$$

$\mathbf{A}$  is called the **companion matrix** of  $f$ .

b) Show that  $p = (-1)^n \pi_B$  where

$$\mathbf{B} = \begin{bmatrix} 0 & 0 & \cdots & 0 & -q_0 \\ 1 & 0 & \cdots & 0 & -q_1 \\ 0 & 1 & \cdots & 0 & -q_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & -q_{n-1} \end{bmatrix}.$$

Thus  $\mathbf{B}$  can also be regarded as a companion matrix for  $p$ .

### 5.1.2 Similarity transformations

Row operations are used in Gaussian elimination to reduce a matrix to triangular form, but row operations change the eigenvalues of a matrix. We need a transformation which can be used to simplify a matrix without changing the eigenvalues.

#### Definition 5.8 (Similar matrices)

Two matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$  are said to be **similar** if there is a nonsingular matrix  $\mathbf{S} \in \mathbb{C}^{n \times n}$  such that  $\mathbf{B} = \mathbf{S}^{-1}\mathbf{A}\mathbf{S}$ . The transformation  $\mathbf{A} \rightarrow \mathbf{B}$  is called a **similarity transformation**. It is called a **unitary similarity transformation** if  $\mathbf{S}^* \mathbf{S} = \mathbf{I}$  and an **orthonormal similarity transformation** if  $\mathbf{S} \in \mathbb{R}^{n \times n}$  and  $\mathbf{S}^T \mathbf{S} = \mathbf{I}$ .

#### Theorem 5.9 (Eigenpairs of similar matrices)

Let  $\mathbf{B} = \mathbf{S}^{-1}\mathbf{A}\mathbf{S}$ , where  $\mathbf{S} \in \mathbb{C}^{n \times n}$  is nonsingular with columns  $\mathbf{s}_1, \dots, \mathbf{s}_n$ . Then  $\mathbf{B}$  and  $\mathbf{A}$  have the same characteristic polynomial. Moreover,  $(\lambda, \mathbf{v})$  is an eigenpair for  $\mathbf{B}$  if and only if  $(\lambda, \mathbf{S}\mathbf{v})$  is an eigenpair for  $\mathbf{A}$ .

**Proof.** By properties of determinants

$$\begin{aligned} \pi_B(\lambda) &= \det(\mathbf{S}^{-1}\mathbf{A}\mathbf{S} - \lambda\mathbf{I}) = \det(\mathbf{S}^{-1}(\mathbf{A} - \lambda\mathbf{I})\mathbf{S}) \\ &= \det(\mathbf{S}^{-1}) \det(\mathbf{A} - \lambda\mathbf{I}) \det(\mathbf{S}) = \det(\mathbf{S}^{-1}\mathbf{S}) \det(\mathbf{A} - \lambda\mathbf{I}) = \pi_A(\lambda). \end{aligned}$$

But then  $\mathbf{A}$  and  $\mathbf{B}$  have the same characteristic polynomial. Moreover,  $(\mathbf{S}^{-1}\mathbf{A}\mathbf{S})\mathbf{v} = \lambda\mathbf{v}$  if and only if  $\mathbf{A}(\mathbf{S}\mathbf{v}) = \lambda(\mathbf{S}\mathbf{v})$ .  $\square$

As a corollary we have the following useful result.

#### Corollary 5.10 (Spectra of $\mathbf{AB}$ and $\mathbf{BA}$ )

For any  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and  $\mathbf{B} \in \mathbb{C}^{n \times m}$  the matrices  $\mathbf{AB}$  and  $\mathbf{BA}$  have the same spectrum apart from some extra zero eigenvalues. More precisely,

$$\lambda^n \pi_{\mathbf{AB}}(\lambda) = \lambda^m \pi_{\mathbf{BA}}(\lambda), \quad \lambda \in \mathbb{C}.$$

If  $m = n$  then  $\pi_{\mathbf{AB}} = \pi_{\mathbf{BA}}$ .

*Proof.* Define block matrices of order  $n + m$  by

$$\mathbf{E} = \begin{bmatrix} \mathbf{AB} & \mathbf{0} \\ \mathbf{B} & \mathbf{0} \end{bmatrix}, \quad \mathbf{F} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{B} & \mathbf{BA} \end{bmatrix}, \quad \mathbf{S} = \begin{bmatrix} \mathbf{I} & \mathbf{A} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

The matrix  $\mathbf{S}$  is nonsingular with  $\mathbf{S}^{-1} = \begin{bmatrix} \mathbf{I} & -\mathbf{A} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$ . Moreover,  $\mathbf{ES} = \mathbf{SF}$  so  $\mathbf{E}$  and  $\mathbf{F}$  are similar and have the same characteristic polynomial by Theorem 5.9. By Property 6. of Theorem 5.1 we have  $\pi_{\mathbf{E}}(\lambda) = \lambda^n \pi_{\mathbf{AB}}(\lambda) = \pi_{\mathbf{F}}(\lambda) = \lambda^m \pi_{\mathbf{BA}}(\lambda)$ . If  $m = n$  we can cancel the  $\lambda$  factors.  $\square$

## 5.2 Unitary Similarity Transformations

In this section we consider the reduction of a matrix to triangular, or almost triangular form using unitary similarity transformations, and characterize matrices with orthonormal eigenvectors.

We start by reviewing some basic facts about matrices with orthonormal columns.

### 5.2.1 Unitary and orthonormal matrices

labelsec:orthmat

#### Definition 5.11 (Unitary matrix)

A matrix  $\mathbf{U} \in \mathbb{C}^{n \times n}$  is said to be **unitary** if  $\mathbf{U}^* \mathbf{U} = \mathbf{I}$ . A real unitary matrix is called **orthonormal**.

**Warning:** An orthonormal matrix is often called an “orthogonal matrix” in the literature.

In the following we consider only the complex case. The real case follows by replacing conjugate transpose “\*” by transpose “T” and  $\mathbb{C}$  by  $\mathbb{R}$ . We use the **standard inner product in  $\mathbb{C}^n$**  given by  $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{y}^* \mathbf{x}$ . In the real case we have  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x}$ . Orthogonality and orthonormality is with respect to the standard inner product.

Since  $\mathbf{U}^* \mathbf{U} = \mathbf{I}$  the matrix  $\mathbf{U}$  is nonsingular,  $\mathbf{U}^{-1} = \mathbf{U}^*$  and  $\mathbf{U} \mathbf{U}^* = \mathbf{I}$  as well. Moreover, both the columns and rows of a unitary matrix of order  $n$  form orthonormal bases for  $\mathbb{C}^n$ . We also note that the product of two unitary matrices is unitary. Indeed if  $\mathbf{U}_1^* \mathbf{U}_1 = \mathbf{I}$  and  $\mathbf{U}_2^* \mathbf{U}_2 = \mathbf{I}$  then  $(\mathbf{U}_1 \mathbf{U}_2)^* (\mathbf{U}_1 \mathbf{U}_2) = \mathbf{U}_2^* \mathbf{U}_1^* \mathbf{U}_1 \mathbf{U}_2 = \mathbf{I}$ .

#### Theorem 5.12 (Unitary matrix)

The matrix  $\mathbf{U} \in \mathbb{C}^{n \times n}$  is unitary if and only if  $\langle \mathbf{U} \mathbf{x}, \mathbf{U} \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$ . In particular, if  $\mathbf{U}$  is unitary then  $\|\mathbf{U} \mathbf{x}\|_2 = \|\mathbf{x}\|_2$  for all  $\mathbf{x} \in \mathbb{C}^n$ .

**Proof.** If  $U^*U = I$  and  $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$  then

$$\langle U\mathbf{x}, U\mathbf{y} \rangle = (U\mathbf{y})^*(U\mathbf{x}) = \mathbf{y}^*U^*U\mathbf{x} = \mathbf{y}^*\mathbf{x} = \langle \mathbf{x}, \mathbf{y} \rangle.$$

Conversely, if  $\langle U\mathbf{x}, U\mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$  then  $U^*U = I$  since for  $i, j = 1, \dots, n$

$$(U^*U)_{i,j} = \mathbf{e}_i^T U^*U\mathbf{e}_j = (U\mathbf{e}_i)^*(U\mathbf{e}_j) = \langle U\mathbf{e}_j, U\mathbf{e}_i \rangle = \langle \mathbf{e}_j, \mathbf{e}_i \rangle = \mathbf{e}_i^*\mathbf{e}_j = \delta_{i,j}.$$

The last part of the theorem follows immediately by taking  $\mathbf{y} = \mathbf{x}$ :  $\square$



Issai Schur, 1875-1941 (left), John William Strutt (Lord Rayleigh), 1842-1919 (right).

## 5.2.2 The Schur decomposition

Although not every matrix can be diagonalized it can be brought into triangular form by a *unitary* similarity transformation.

### Theorem 5.13 (Schur decomposition)

For each  $A \in \mathbb{C}^{n \times n}$  there exists a unitary matrix  $U \in \mathbb{C}^{n \times n}$  such that  $R := U^*AU$  is upper triangular.

The matrices  $U$  and  $R$  in the Schur decomposition are called **Schur factors**.

**Proof.** We use induction on  $n$ . For  $n = 1$  the matrix  $U$  is the  $1 \times 1$  identity matrix. Assume that the theorem is true for matrices of order  $k$  and suppose  $A \in \mathbb{C}^{n \times n}$ , where  $n := k + 1$ . Let  $(\lambda_1, \mathbf{v}_1)$  be an eigenpair for  $A$  with  $\|\mathbf{v}_1\|_2 = 1$ . By Theorem 0.30 we can extend  $\mathbf{v}_1$  to an orthonormal basis  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  for  $\mathbb{C}^n$ . The matrix  $V := [\mathbf{v}_1, \dots, \mathbf{v}_n] \in \mathbb{C}^{n \times n}$  is unitary, and for the first column of the product  $V^*AV$  we find

$$V^*AV\mathbf{e}_1 = V^*A\mathbf{v}_1 = \lambda_1 V^*\mathbf{v}_1 = \lambda_1 \mathbf{e}_1.$$



It follows that

$$\mathbf{V}^* \mathbf{A} \mathbf{V} = \left[ \begin{array}{c|c} \lambda_1 & \mathbf{x}^* \\ \mathbf{0} & \mathbf{M} \end{array} \right], \text{ for some } \mathbf{M} \in \mathbb{C}^{k \times k} \text{ and } \mathbf{x} \in \mathbb{C}^k. \quad (5.4)$$

By the induction hypothesis there is a unitary matrix  $\mathbf{W}_1 \in \mathbb{C}^{k \times k}$  such that  $\mathbf{W}_1^* \mathbf{M} \mathbf{W}_1$  is upper triangular. Define

$$\mathbf{W} = \left[ \begin{array}{c|c} 1 & \mathbf{0}^* \\ \mathbf{0} & \mathbf{W}_1 \end{array} \right] \text{ and } \mathbf{U} = \mathbf{V} \mathbf{W}.$$

Then  $\mathbf{W}$  and  $\mathbf{U}$  are unitary and

$$\begin{aligned} \mathbf{U}^* \mathbf{A} \mathbf{U} &= \mathbf{W}^* (\mathbf{V}^* \mathbf{A} \mathbf{V}) \mathbf{W} = \left[ \begin{array}{c|c} 1 & \mathbf{0}^* \\ \mathbf{0} & \mathbf{W}_1^* \end{array} \right] \left[ \begin{array}{c|c} \lambda_1 & \mathbf{x}^* \\ \mathbf{0} & \mathbf{M} \end{array} \right] \left[ \begin{array}{c|c} 1 & \mathbf{0}^* \\ \mathbf{0} & \mathbf{W}_1 \end{array} \right] \\ &= \left[ \begin{array}{c|c} \lambda_1 & \mathbf{x}^* \mathbf{W}_1 \\ \mathbf{0} & \mathbf{W}_1^* \mathbf{M} \mathbf{W}_1 \end{array} \right] \end{aligned}$$

is upper triangular.  $\square$

If  $\mathbf{A}$  has complex eigenvalues then  $\mathbf{U}$  will be complex even if  $\mathbf{A}$  is real. The following is a real version of Theorem 5.13.

**Theorem 5.14 (Schur form, real eigenvalues)**

For each  $\mathbf{A} \in \mathbb{R}^{n \times n}$  with real eigenvalues there exists a matrix  $\mathbf{U} \in \mathbb{R}^{n \times n}$  with  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ , such that  $\mathbf{U}^T \mathbf{A} \mathbf{U}$  is upper triangular.

*Proof.* Consider the proof of Theorem 5.13. Since  $\mathbf{A}$  and  $\lambda_1$  are real the eigenvector  $\mathbf{v}_1$  is real and the matrix  $\mathbf{W}$  is real and  $\mathbf{W}^T \mathbf{W} = \mathbf{I}$ . By the induction hypothesis  $\mathbf{V}$  is real and  $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ . But then also  $\mathbf{U} = \mathbf{V} \mathbf{W}$  is real and  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ .  $\square$

By using the unitary transformation  $\mathbf{V}$  on the  $n \times n$  matrix  $\mathbf{A}$ , we obtain a matrix  $\mathbf{M}$  of order  $n - 1$ .  $\mathbf{M}$  has the same eigenvalues as  $\mathbf{A}$  except  $\lambda$ . Thus we can find another eigenvalue of  $\mathbf{A}$  by working with a smaller matrix  $\mathbf{M}$ . This is an example of a **deflation** technique which is very useful in numerical work.

**Example 5.15 (Deflation example)**

The matrix  $\mathbf{T} := \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$  has an eigenpair  $(2, \mathbf{x}_1)$ , where  $\mathbf{x}_1 = [-1, 0, 1]^T$ .

We can extend  $\mathbf{x}_1$  to a basis  $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$  for  $\mathbb{R}^3$  by defining  $\mathbf{x}_2 = [0, 1, 0]^T$ ,  $\mathbf{x}_3 = [1, 0, 1]^T$ . This is already an orthogonal basis and normalizing we obtain the orthonormal matrix

$$\mathbf{V} = \begin{bmatrix} -\frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} \\ 0 & 1 & 0 \\ \frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} \end{bmatrix}.$$

We obtain (5.4) with  $\lambda = 2$  and

$$\mathbf{M} = \begin{bmatrix} 2 & -\sqrt{2} \\ -\sqrt{2} & 2 \end{bmatrix}.$$

We can now find the remaining eigenvalues of  $\mathbf{A}$  from the  $2 \times 2$  matrix  $\mathbf{M}$ .

**Exercise 5.16 (Schur decomposition example)**

Show that a Schur decomposition of  $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 2 \end{bmatrix}$  is  $\mathbf{U}^T \mathbf{A} \mathbf{U} = \begin{bmatrix} -1 & -1 \\ 0 & 4 \end{bmatrix}$ , where  $\mathbf{U} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$ .

How far can we reduce a real matrix  $\mathbf{A}$  with some complex eigenvalues by a real unitary similarity transformation? To study this we note that the complex eigenvalues of a real matrix occur in conjugate pairs,  $\lambda = \mu + i\nu$ ,  $\bar{\lambda} = \mu - i\nu$ , where  $\mu, \nu$  are real. The real  $2 \times 2$  matrix

$$\mathbf{M} = \begin{bmatrix} \mu & \nu \\ -\nu & \mu \end{bmatrix} \quad (5.5)$$

has eigenvalues  $\lambda = \mu + i\nu$  and  $\bar{\lambda} = \mu - i\nu$ .

**Definition 5.17 (Quasi-triangular matrix)**

We say that a matrix is **quasi-triangular** if it is block triangular with only  $1 \times 1$  and  $2 \times 2$  blocks on the diagonal. Moreover, no  $2 \times 2$  block should have real eigenvalues.

As an example consider

$$\mathbf{R} = \left[ \begin{array}{cc|cc|cc} 2 & 1 & 3 & 4 & 5 & \\ -1 & 2 & 4 & 3 & 2 & \\ \hline 0 & 0 & 1 & 2 & 3 & \\ 0 & 0 & 0 & 3 & 2 & \\ 0 & 0 & 0 & -1 & 1 & \end{array} \right].$$

$\mathbf{R}$  has three diagonal blocks:

$$\mathbf{D}_1 = \begin{bmatrix} 2 & 1 \\ -1 & 2 \end{bmatrix}, \quad \mathbf{D}_2 = [1], \quad \mathbf{D}_3 = \begin{bmatrix} 3 & 2 \\ -1 & 1 \end{bmatrix}.$$

By Theorem 5.1 the eigenvalues of  $\mathbf{R}$  are the union of the eigenvalues of  $\mathbf{D}_1$ ,  $\mathbf{D}_2$  and  $\mathbf{D}_3$ . The eigenvalues of  $\mathbf{D}_1$  are  $2+i$  and  $2-i$ , while  $\mathbf{D}_2$  has eigenvalue 1, and  $\mathbf{D}_3$  has the same eigenvalues as  $\mathbf{D}_1$ . Thus  $\mathbf{R}$  has one real eigenvalue 1 corresponding to the  $1 \times 1$  block and complex eigenvalues  $2+i$ ,  $2-i$  with multiplicity 2 corresponding to the two  $2 \times 2$  blocks.

Any  $\mathbf{A} \in \mathbb{R}^{n \times n}$  can be reduced to quasi-triangular form by a real orthonormal similarity transformation. A proof is given in Section 5.7.

### 5.2.3 Normal matrices

It is possible to characterize matrices that have a diagonal Schur factorization.

**Definition 5.18 (Normal matrix)**

A matrix  $A \in \mathbb{C}^{n \times n}$  is said to be **normal** if  $AA^* = A^*A$ .

Examples of normal matrices are

1.  $A^* = A$ , (Hermitian)
2.  $A^* = -A$ , (Skew-Hermitian)
3.  $A^* = A^{-1}$ , (Unitary)
4.  $A = D$ . (Diagonal)

The 2. derivative matrix  $T$  in (1.2) and the discrete Poisson matrix (cf. Lemma 3.10) are examples of normal matrices.

**Exercise 5.19 (Skew-Hermitian matrix)**

Suppose  $C = A + iB$ , where  $A, B \in \mathbb{R}^{n \times n}$ . Show that  $C$  is skew-Hermitian if and only if  $A^T = -A$  and  $B^T = B$ .

**Exercise 5.20 (Eigenvalues of a skew-Hermitian matrix)**

Show that any eigenvalue of a skew-Hermitian matrix is purely imaginary.

The following theorem says that a matrix has orthonormal eigenpairs if and only if it is normal.

**Theorem 5.21 (Orthonormal eigenpairs characterization)**

A matrix  $A \in \mathbb{C}^{n \times n}$  is unitary similar with a diagonal matrix if and only if it is normal.

*Proof.* If  $B = U^*AU$ , with  $B$  diagonal, and  $U^*U = I$ , then

$$\begin{aligned} AA^* &= (UBU^*)(UB^*U^*) = UBB^*U^* \text{ and} \\ A^*A &= (UB^*U^*)(UBU^*) = UB^*BU^*. \end{aligned}$$

Now  $BB^* = B^*B$  since  $B$  is diagonal, and  $A$  is normal.

Suppose  $A^*A = AA^*$ . By Theorem 5.13 we can find  $U$  with  $U^*U = I$  such that  $B := U^*AU$  is upper triangular. Since  $A$  is normal  $B$  is normal. Indeed,

$$BB^* = U^*AUU^*A^*U = U^*AA^*U = U^*A^*AU = B^*B.$$

The proof is complete if we can show that an upper triangular normal matrix  $\mathbf{B}$  must be diagonal. The diagonal elements in  $\mathbf{E} := \mathbf{B}^* \mathbf{B}$  and  $\mathbf{F} := \mathbf{B} \mathbf{B}^*$  are given by

$$e_{ii} = \sum_{k=1}^n \bar{b}_{ki} b_{ki} = \sum_{k=1}^i |b_{ki}|^2 \quad \text{and} \quad f_{ii} = \sum_{k=1}^n b_{ik} \bar{b}_{ik} = \sum_{k=i}^n |b_{ik}|^2.$$

The result now follows by equating  $e_{ii}$  and  $f_{ii}$  for  $i = 1, 2, \dots, n$ . In particular for  $i = 1$  we have  $|b_{11}|^2 = |b_{11}|^2 + |b_{12}|^2 + \dots + |b_{1n}|^2$ , so  $b_{1k} = 0$  for  $k = 2, 3, \dots, n$ . Suppose  $b_{jk} = 0$  for  $j = 1, \dots, i-1$ ,  $k = j+1, \dots, n$ . Then

$$e_{ii} = \sum_{k=1}^i |b_{ki}|^2 = |b_{ii}|^2 = \sum_{k=i}^n |b_{ik}|^2 = f_{ii}$$

so  $b_{ik} = 0$ ,  $k = i+1, \dots, n$ . By induction on the rows we see that  $\mathbf{B}$  is diagonal.  $\square$

The special cases where  $\mathbf{A}$  is Hermitian, or real and symmetric, occur often in applications and deserve special attention.

**Corollary 5.22 (Spectral theorem, complex form)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is Hermitian. Then  $\mathbf{A}$  has real eigenvalues  $\lambda_1, \dots, \lambda_n$ . Moreover, there is a unitary matrix  $\mathbf{U} \in \mathbb{C}^{n \times n}$  such that  $\mathbf{U}^* \mathbf{A} \mathbf{U} = \text{diag}(\lambda_1, \dots, \lambda_n)$ . For any such  $\mathbf{U}$  the columns  $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  of  $\mathbf{U}$  are orthonormal eigenvectors of  $\mathbf{A}$  and  $\mathbf{A} \mathbf{u}_j = \lambda_j \mathbf{u}_j$  for  $j = 1, \dots, n$ .

*Proof.* That the eigenvalues are real was shown in Lemma 3.9. The rest follows from Theorem 5.21.  $\square$

There is also a real version.

**Corollary 5.23 (Spectral theorem (real form))**

Suppose  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\mathbf{A}^T = \mathbf{A}$ . Then  $\mathbf{A}$  has real eigenvalues  $\lambda_1, \dots, \lambda_n$ . Moreover, there is an orthonormal matrix  $\mathbf{U} \in \mathbb{R}^{n \times n}$  such that  $\mathbf{U}^T \mathbf{A} \mathbf{U} = \text{diag}(\lambda_1, \dots, \lambda_n)$ . For any such  $\mathbf{U}$  the columns  $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  of  $\mathbf{U}$  are orthonormal eigenvectors of  $\mathbf{A}$  and  $\mathbf{A} \mathbf{u}_j = \lambda_j \mathbf{u}_j$  for  $j = 1, \dots, n$ .

*Proof.* Since a real symmetric matrix has real eigenvalues and eigenvectors this follows from Theorem 5.22.  $\square$

**Example 5.24** The orthonormal diagonalization of  $\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$  is  $\mathbf{U}^T \mathbf{A} \mathbf{U} = \text{diag}(1, 3)$ , where  $\mathbf{U} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$ .

## 5.3 Minmax theorems for Hermitian Matrices

There are some useful characterizations of the eigenvalues of a Hermitian matrix. They are based on the Rayleigh quotient that is a useful tool when studying eigenpairs.

### 5.3.1 The Rayleigh quotient

#### Definition 5.25 (Rayleigh quotient)

For  $\mathbf{A} \in \mathbb{C}^{n \times n}$  and a nonzero  $\mathbf{x} \in \mathbb{C}^n$  the number

$$R(\mathbf{x}) = R_{\mathbf{A}}(\mathbf{x}) := \frac{\mathbf{x}^* \mathbf{A} \mathbf{x}}{\mathbf{x}^* \mathbf{x}}$$

is called a **Rayleigh quotient**.

If  $(\lambda, \mathbf{x})$  is an eigenpair for  $\mathbf{A}$  then  $R(\mathbf{x}) = \frac{\mathbf{x}^* \mathbf{A} \mathbf{x}}{\mathbf{x}^* \mathbf{x}} = \lambda$ .

Equation (5.6) in the following lemma shows that the Rayleigh quotient of a normal matrix is a **convex combination** of its eigenvalues.

#### Lemma 5.26 (Convex combination of the eigenvalues)

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is normal with orthonormal eigenpairs  $(\lambda_j, \mathbf{u}_j)$ ,  $j = 1, 2, \dots, n$  and let  $\mathbf{x} \in \mathbb{C}^n$ . Then

$$R_{\mathbf{A}}(\mathbf{x}) = \sum_{j=1}^n b_j \lambda_j, \quad b_j \geq 0, \quad \sum_{j=1}^n b_j = 1, \quad (5.6)$$

where  $b_j = |c_j|^2 / \sum_{i=1}^n |c_i|^2$ ,  $j = 1, \dots, n$ , and  $\mathbf{x} = \sum_{j=1}^n c_j \mathbf{u}_j$  is the eigenvector expansion of  $\mathbf{x}$ .

**Proof.** By orthonormality of the eigenvectors  $\mathbf{x}^* \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n \bar{c}_i \bar{\mathbf{u}}_i c_j \mathbf{u}_j = \sum_{j=1}^n |c_j|^2$ . Similarly,  $\mathbf{x}^* \mathbf{A} \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n \bar{c}_i \bar{\mathbf{u}}_i c_j \lambda_j \mathbf{u}_j = \sum_{j=1}^n \lambda_j |c_j|^2$ . and (5.6) follows with  $b_j = |c_j|^2 / \sum_{i=1}^n |c_i|^2$ ,  $j = 1, \dots, n$ . This is clearly a combination of nonnegative quantities and a convex combination since  $\sum_{j=1}^n |c_j|^2 / \sum_{i=1}^n |c_i|^2 = 1$ .  $\square$

### 5.3.2 Minmax and maxmin

First we show

#### Theorem 5.27 (Minmax)

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is Hermitian with eigenvalues  $\lambda_1, \dots, \lambda_n$ , ordered so that  $\lambda_1 \geq \dots \geq \lambda_n$ . Let  $1 \leq k \leq n$ . For any subspace  $\mathcal{S}$  of  $\mathbb{C}^n$  of dimension  $n - k + 1$

$$\lambda_k \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}), \quad (5.7)$$

with equality for  $\mathcal{S} = \tilde{\mathcal{S}} := \text{span}(\mathbf{u}_k, \dots, \mathbf{u}_n)$  and  $\mathbf{x} = \mathbf{u}_k$ . Here  $(\lambda_j, \mathbf{u}_j)$ ,  $1 \leq j \leq n$  are orthonormal eigenpairs for  $\mathbf{A}$ .

**Proof.** Let  $\mathcal{S}$  be any subspace of  $\mathbb{C}^n$  of dimension  $n - k + 1$  and define  $\mathcal{S}' := \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$ . We need to find  $\mathbf{y} \in \mathcal{S}$  so that  $R(\mathbf{y}) \geq \lambda_k$ . Now  $\mathcal{S} + \mathcal{S}' := \{\mathbf{s} + \mathbf{s}' : \mathbf{s} \in \mathcal{S}, \mathbf{s}' \in \mathcal{S}'\}$  is a subspace of  $\mathbb{C}^n$  and by (7)

$$\dim(\mathcal{S} \cap \mathcal{S}') = \dim(\mathcal{S}) + \dim(\mathcal{S}') - \dim(\mathcal{S} + \mathcal{S}') \geq (n - k + 1) + k - n = 1.$$

It follows that  $\mathcal{S} \cap \mathcal{S}'$  is nonempty. Let  $\mathbf{y} \in \mathcal{S} \cap \mathcal{S}' = \sum_{j=1}^k c_j \mathbf{u}_j$  with  $\sum_{j=1}^k |c_j|^2 = 1$ . Defining  $c_j = 0$  for  $k + 1 \leq j \leq n$ , we obtain by Lemma 5.26

$$\max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}) \geq R(\mathbf{y}) = \sum_{j=1}^n \lambda_j |c_j|^2 = \sum_{j=1}^k \lambda_j |c_j|^2 \geq \sum_{j=1}^k \lambda_k |c_j|^2 = \lambda_k,$$

and (5.7) follows. If  $\mathbf{y} \in \tilde{\mathcal{S}}$ , say  $\mathbf{y} = \sum_{j=k}^n d_j \mathbf{u}_j$  with  $\sum_{j=k}^n |d_j|^2 = 1$  then again by Lemma 5.26  $R(\mathbf{y}) = \sum_{j=k}^n \lambda_j |d_j|^2 \leq \lambda_k$ , and since  $\mathbf{y} \in \tilde{\mathcal{S}}$  is arbitrary we have  $\max_{\substack{\mathbf{x} \in \tilde{\mathcal{S}} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}) \leq \lambda_k$  and equality in (5.7) follows for  $\mathcal{S} = \tilde{\mathcal{S}}$ . Moreover,  $R(\mathbf{u}_k) = \lambda_k$ .  $\square$

There is also a maxmin version of this result.

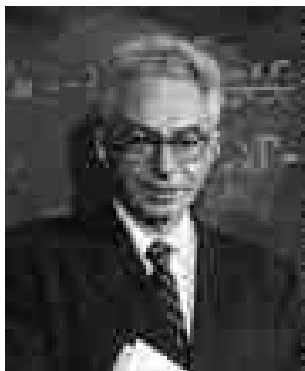
### Theorem 5.28 (Maxmin)

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is Hermitian with eigenvalues  $\lambda_1, \dots, \lambda_n$ , ordered so that  $\lambda_1 \geq \dots \geq \lambda_n$ . Let  $1 \leq k \leq n$ . For any subspace  $\mathcal{S}$  of  $\mathbb{C}^n$  of dimension  $k$

$$\lambda_k \geq \min_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}), \quad (5.8)$$

with equality for  $\mathcal{S} = \tilde{\mathcal{S}} := \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$  and  $\mathbf{x} = \mathbf{u}_k$ . Here  $(\lambda_j, \mathbf{u}_j)$ ,  $1 \leq j \leq n$  are orthonormal eigenpairs for  $\mathbf{A}$ .

**Proof.** The proof is very similar to the proof of Theorem 5.27. We define  $\mathcal{S}' := \text{span}(\mathbf{u}_k, \dots, \mathbf{u}_n)$  and show that  $R(\mathbf{y}) \leq \lambda_k$  for some  $\mathbf{y} \in \mathcal{S} \cap \mathcal{S}'$ . It is easy to see that  $R(\mathbf{y}) \geq \lambda_k$  for any  $\mathbf{y} \in \tilde{\mathcal{S}}$ .  $\square$



Richard Courant, 1888-1972 (left), Ernst Sigismund Fischer, 1875-1954 (right).

These theorems immediately lead to classical minmax and maxmin characterizations.

**Corollary 5.29 (The Courant-Fischer theorem)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is Hermitian with eigenvalues  $\lambda_1, \dots, \lambda_n$ , ordered so that  $\lambda_1 \geq \dots \geq \lambda_n$ . Then

$$\lambda_k = \min_{\dim(\mathcal{S})=n-k+1} \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}) = \max_{\dim(\mathcal{S})=k} \min_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R(\mathbf{x}), \quad k = 1, \dots, n. \quad (5.9)$$

Using the maxmin Theorem 5.27 we can prove inequalities of eigenvalues without knowing the eigenvectors and we can get both upper and lower bounds.

**Theorem 5.30 (Eigenvalue perturbation for Hermitian matrices)**

Let  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$  be Hermitian with eigenvalues  $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n$  and  $\beta_1 \geq \beta_2 \geq \dots \geq \beta_n$ . Then

$$\alpha_k + \varepsilon_n \leq \beta_k \leq \alpha_k + \varepsilon_1, \quad \text{for } k = 1, \dots, n, \quad (5.10)$$

where  $\varepsilon_1 \geq \varepsilon_2 \geq \dots \geq \varepsilon_n$  are the eigenvalues of  $\mathbf{E} := \mathbf{B} - \mathbf{A}$ .

**Proof.** Since  $\mathbf{E}$  is a sum of Hermitian matrices it is Hermitian and the eigenvalues are real. Let  $(\alpha_j, \mathbf{u}_j)$ ,  $j = 1, \dots, n$  be orthonormal eigenpairs for  $\mathbf{A}$  and let  $\mathcal{S} := \text{span}\{\mathbf{u}_k, \dots, \mathbf{u}_n\}$ . By Theorem 5.27 we obtain

$$\beta_k \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{B}}(\mathbf{x}) \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{A}}(\mathbf{x}) + \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{E}}(\mathbf{x}) \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{A}}(\mathbf{x}) + \max_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \mathbf{x} \neq \mathbf{0}}} R_{\mathbf{E}}(\mathbf{x}) = \alpha_k + \varepsilon_1,$$

and this proves the upper inequality. For the lower one we define  $\mathbf{D} := -\mathbf{E}$  and observe that  $-\varepsilon_n$  is the largest eigenvalue of  $\mathbf{D}$ . Since  $\mathbf{A} = \mathbf{B} + \mathbf{D}$  it follows from

the result just proved that  $\alpha_k \leq \beta_k - \varepsilon_n$ , which is the same as the lower inequality.  $\square$

In many applications of this result the eigenvalues of the matrix  $E$  will be small and then the theorem states that the eigenvalues of  $B$  are close to those of  $A$ . Moreover, it associates a unique eigenvalue of  $A$  with each eigenvalue of  $B$ .

**Exercise 5.31 (Eigenvalue perturbation for Hermitian matrices)**

Show that in Theorem 5.30, if  $E$  is symmetric positive semidefinite then  $\beta_i \geq \alpha_i$ .



Alan Jerome Hoffman, 1924- (left), Helmut Wielandt, 1910-2001 (right).

### 5.3.3 The Hoffman-Wielandt theorem

We can also give a bound involving all eigenvalues. The following theorem shows that the eigenvalue problem for a normal matrix is well conditioned.

**Theorem 5.32 (Hoffman-Wielandt theorem)**

Suppose  $A, B \in \mathbb{C}^{n \times n}$  are both normal matrices with eigenvalues  $\lambda_1, \dots, \lambda_n$  and  $\mu_1, \dots, \mu_n$ , respectively. Then there is a permutation  $i_1, \dots, i_n$  of  $1, 2, \dots, n$  such that

$$\sum_{j=1}^n |\mu_{i_j} - \lambda_j|^2 \leq \sum_{i=1}^n \sum_{j=1}^n |a_{ij} - b_{ij}|^2. \quad (5.11)$$

For a proof of this theorem see [[28], p. 190]. For a Hermitian matrix we can use the identity permutation if we order both set of eigenvalues in nonincreasing or nondecreasing order.

**Exercise 5.33 (Hoffman-Wielandt)**

Show that (5.11) does not hold for the matrices  $A := \begin{bmatrix} 0 & 0 \\ 0 & 4 \end{bmatrix}$  and  $B := \begin{bmatrix} -1 & -1 \\ 1 & 1 \end{bmatrix}$ . Why does this not contradict the Hoffman-Wielandt theorem?



## 5.4 The Jordan Form

For any  $\mathbf{A} \in \mathbb{C}^{n \times n}$  there is a unitary matrix  $\mathbf{U}$  such that  $\mathbf{U}^* \mathbf{A} \mathbf{U} = \mathbf{R}$  is upper triangular. Moreover  $\mathbf{R}$  is diagonal if  $\mathbf{A}$  is normal. The following question arises. How close to a diagonal matrix can we reduce a general matrix by a similarity transformation? The main result is the Jordan form in Theorem 5.44. For a proof, see for example [13].

### 5.4.1 Diagonalizable matrices and linear independence of eigenvectors

We start by giving a characterization of matrices that are similar to a diagonal matrix.

**Definition 5.34 (Diagonalizable matrix)**

A matrix of order  $n$  is **diagonalizable** if it is similar to a diagonal matrix, and **defective** if this is not the case.

We have  $\mathbf{S}^{-1} \mathbf{A} \mathbf{S} = \text{diag}(\lambda_1, \dots, \lambda_n)$  if and only if  $\mathbf{S}^* \mathbf{A}^* \mathbf{S}^{-*} = \text{diag}(\bar{\lambda}_1, \dots, \bar{\lambda}_n)$ , where  $\mathbf{S}^{-*} := (\mathbf{S}^*)^{-1} = (\mathbf{S}^{-1})^*$ . Thus  $\mathbf{A}$  is diagonalizable if and only if  $\mathbf{A}^*$  is diagonalizable.

**Theorem 5.35 (Eigenvectors of diagonalizable matrices)**

A matrix of order  $n$  is diagonalizable if and only if its eigenvectors form a basis for  $\mathbb{C}^n$ .

*Proof.* Let  $\mathbf{S} \in \mathbb{C}^{n \times n}$  be nonsingular with columns  $\mathbf{s}_1, \dots, \mathbf{s}_n$ . Then

$$\begin{aligned} \mathbf{S}^{-1} \mathbf{A} \mathbf{S} = \text{diag}(\lambda_1, \dots, \lambda_n) &\Leftrightarrow \mathbf{A} \mathbf{S} = \mathbf{S} \text{diag}(\lambda_1, \dots, \lambda_n) \\ &\Leftrightarrow \mathbf{A} \mathbf{s}_i = \lambda_i \mathbf{s}_i, \quad i = 1, \dots, n. \end{aligned}$$

Since  $\mathbf{S}$  is nonsingular the  $n$  columns of are linearly independent and therefore constitute a basis for  $\mathbb{C}^n$ .  $\square$

If the eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_n$  of a matrix of order  $n$  are linearly independent then any  $\mathbf{x} \in \mathbb{C}^n$  can be written  $\mathbf{x} = \sum_{j=1}^n c_j \mathbf{v}_j$  for some scalars  $c_1, \dots, c_n$ . We call this an **eigenvector expansion** of  $\mathbf{x}$ .

For distinct eigenvalues we have the following result.

**Theorem 5.36 (Distinct eigenvalues)**

*Eigenvectors corresponding to distinct eigenvalues are linearly independent.*

*Proof.* Suppose  $(\lambda_k, \mathbf{x}_k)$ ,  $k = 1, \dots, m$  are eigenpairs of  $\mathbf{A}$  and that  $\lambda_1, \dots, \lambda_m$  are distinct, but  $\mathbf{x}_1, \dots, \mathbf{x}_m$  are linearly dependent. With  $m$  the smallest such positive

integer we will obtain a contradiction. For some nonzero vector  $\mathbf{c} := [c_1, \dots, c_m]$  we have

$$\sum_{j=1}^m c_j \mathbf{x}_j = \mathbf{0}. \quad (5.12)$$

Clearly  $m \geq 2$  since the eigenvectors and  $\mathbf{c}$  are nonzero. Applying  $\mathbf{A}$  to (5.12) we obtain by linearity  $\sum_{j=1}^m c_j \lambda_j \mathbf{x}_j = \mathbf{0}$ . From this relation we subtract  $\lambda_m$  times (5.12) and find  $\sum_{j=1}^{m-1} c_j (\lambda_j - \lambda_m) \mathbf{x}_j = \mathbf{0}$ . But since  $\lambda_j - \lambda_m \neq 0$  for  $j = 1, \dots, m-1$  and at least one  $c_j \neq 0$  for  $j < m$  we see that  $\{\mathbf{x}_1, \dots, \mathbf{x}_{m-1}\}$  is linearly dependent, contradicting the minimality of  $m$ .  $\square$

### Corollary 5.37 (Diagonalizable matrix)

*A matrix with distinct eigenvalues is diagonalizable.*

**Proof.** By the previous theorem the eigenvectors are linearly independent.  $\square$

## 5.4.2 Algebraic and geometric multiplicity of eigenvalues

A defective matrix must necessarily have one or more multiple eigenvalues, but as the following example shows this is not sufficient.

### Example 5.38 (Two upper triangular matrices)

*Consider the 2 matrices of order 3*

$$\mathbf{A}_1 := \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{A}_2 := \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}.$$

*Both matrices are upper triangular and have an eigenvalue  $\lambda = 1$  of multiplicity 3.*

1. *The eigenvectors of  $\mathbf{A}_1$  are the linearly independent unit vectors  $\mathbf{x}_i = \mathbf{e}_i$ ,  $i = 1, 2, 3$ . Thus  $\mathbf{A}_1$  is diagonalizable.*
2. *An eigenvector  $\mathbf{x} = [x_1, x_2, x_3]^T$  of  $\mathbf{A}_2$  must be a solution of the homogenous triangular linear system*

$$(\mathbf{A} - \mathbf{I})\mathbf{x} = \mathbf{0} \text{ or } \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

*But then  $x_2 = x_3 = 0$  and any eigenvector must be a multiple of  $\mathbf{e}_1$ . We conclude that  $\mathbf{A}_2$  is defective.*

Linear independence of eigenvectors depends on the multiplicity of the eigenvalues in a nontrivial way. For multiple eigenvalues we need to distinguish between two kinds of multiplicities.

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  and

$$\pi_{\mathbf{A}}(\lambda) := \det(\mathbf{A} - \lambda \mathbf{I}) = (\mu_1 - \lambda)^{a_1} \cdots (\mu_r - \lambda)^{a_r}, \quad \mu_i \neq \mu_j, \quad i \neq j, \quad \sum_{i=1}^r a_i = n. \quad (5.13)$$

The positive integer  $a_i = a(\mu_i) = a_{\mathbf{A}}(\mu_i)$  is called the **multiplicity**, or more precisely the **algebraic multiplicity** of the eigenvalue  $\mu_i$ . The multiplicity of an eigenvalue is simple (double, triple) if  $a_i$  is equal to one (two, three).

To define a second kind of multiplicity we consider for each  $\lambda \in \sigma(\mathbf{A})$  the nullspace

$$\ker(\mathbf{A} - \lambda \mathbf{I}) := \{\mathbf{x} \in \mathbb{C}^n : (\mathbf{A} - \lambda \mathbf{I})\mathbf{x} = \mathbf{0}\} \quad (5.14)$$

of  $\mathbf{A} - \lambda \mathbf{I}$ . The nullspace is a subspace of  $\mathbb{C}^n$  consisting of all eigenvectors of  $\mathbf{A}$  corresponding to the eigenvalue  $\lambda$ . The dimension of the subspace must be at least one since  $\mathbf{A} - \lambda \mathbf{I}$  is singular.

#### Definition 5.39 (Geometric multiplicity)

The **geometric multiplicity**  $g = g(\lambda) = g_{\mathbf{A}}(\lambda)$  of an eigenvalue  $\lambda$  of  $\mathbf{A}$  is the dimension of the nullspace  $\ker(\mathbf{A} - \lambda \mathbf{I})$ .

#### Example 5.40 (Geometric multiplicity)

The  $n \times n$  identity matrix  $\mathbf{I}$  has the eigenvalue  $\lambda = 1$  with  $\pi_{\mathbf{I}}(\lambda) = (1 - \lambda)^n$ . Since  $\mathbf{I} - \lambda \mathbf{I}$  is the zero matrix when  $\lambda = 1$ , the nullspace of  $\mathbf{I} - \lambda \mathbf{I}$  is all of  $n$ -space and it follows that  $a = g = n$ . On the other hand we saw in Example 5.38 that the  $3 \times 3$  matrix  $\mathbf{A} = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$  has the eigenvalue  $\lambda = 1$  with  $a = 3$  and only one eigenvector. Thus  $g = 1$ .

#### Theorem 5.41 (Geometric multiplicity of similar matrices)

Similar matrices have the same eigenvalues with the same algebraic and geometric multiplicities.

**Proof.** Similar matrices have the same characteristic polynomials and only the invariance of geometric multiplicity needs to be shown. Suppose  $\lambda \in \sigma(\mathbf{A})$ ,  $\dim \ker(\mathbf{S}^{-1}\mathbf{A}\mathbf{S} - \lambda \mathbf{I}) = k$ , and  $\dim \ker(\mathbf{A} - \lambda \mathbf{I}) = \ell$ . We need to show that  $k = \ell$ . Suppose  $\mathbf{v}_1, \dots, \mathbf{v}_k$  is a basis for  $\ker(\mathbf{S}^{-1}\mathbf{A}\mathbf{S} - \lambda \mathbf{I})$ . Then  $\mathbf{S}^{-1}\mathbf{A}\mathbf{S}\mathbf{v}_i = \lambda \mathbf{v}_i$  or  $\mathbf{A}\mathbf{S}\mathbf{v}_i = \lambda \mathbf{S}\mathbf{v}_i$ ,  $i = 1, \dots, k$ . But then  $\{\mathbf{S}\mathbf{v}_1, \dots, \mathbf{S}\mathbf{v}_k\} \subset \ker(\mathbf{A} - \lambda \mathbf{I})$ , which implies that  $k \leq \ell$ . Similarly, if  $\mathbf{w}_1, \dots, \mathbf{w}_\ell$  is a basis for  $\ker(\mathbf{A} - \lambda \mathbf{I})$  then  $\{\mathbf{S}^{-1}\mathbf{w}_1, \dots, \mathbf{S}^{-1}\mathbf{w}_\ell\} \subset \ker(\mathbf{S}^{-1}\mathbf{A}\mathbf{S} - \lambda \mathbf{I})$ , which implies that  $k \geq \ell$ . We conclude that  $k = \ell$ .  $\square$

**Exercise 5.42 (Find eigenpair example)**

Find eigenvalues and eigenvectors of  $A = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 2 & 3 \\ 0 & 0 & 2 \end{bmatrix}$ .

**5.4.3 The Jordan form**

Marie Ennemond Camille Jordan, 1838-1922 (left), William Rowan Hamilton, 1805-1865 (right).

**Definition 5.43 (Jordan block)**

A **Jordan block** of order  $m$ , denoted  $\mathbf{J}_m(\lambda)$  is an  $m \times m$  matrix of the form

$$\mathbf{J}_m(\lambda) := \begin{bmatrix} \lambda & 1 & 0 & \cdots & 0 & 0 \\ 0 & \lambda & 1 & \cdots & 0 & 0 \\ 0 & 0 & \lambda & \cdots & 0 & 0 \\ \vdots & & & & \vdots & \\ 0 & 0 & 0 & \cdots & \lambda & 1 \\ 0 & 0 & 0 & \cdots & 0 & \lambda \end{bmatrix}.$$

A  $3 \times 3$  Jordan block has the form  $\mathbf{J}_3(\lambda) = \begin{bmatrix} \lambda & 1 & 0 \\ 0 & \lambda & 1 \\ 0 & 0 & \lambda \end{bmatrix}$ . Since a Jordan block is upper triangular  $\lambda$  is an eigenvalue of  $\mathbf{J}_m(\lambda)$  and any eigenvector must be a multiple of  $\mathbf{e}_1$ . Indeed, if  $\mathbf{J}_m(\lambda)\mathbf{v} = \lambda\mathbf{v}$  for some  $\mathbf{v} = [v_1, \dots, v_m]$  then  $v_2 = \dots = v_m = 0$ . Thus, the eigenvectors of  $\mathbf{J}_m(\lambda)$  have algebraic multiplicity  $m$  and geometric multiplicity one.

The Jordan form is a decomposition of a matrix into Jordan blocks.

**Theorem 5.44 (The Jordan form of a matrix)**

Suppose  $A \in \mathbb{C}^{n \times n}$  has  $k$  distinct eigenvalues  $\lambda_1, \dots, \lambda_k$  of algebraic multiplicities  $a_1, \dots, a_k$  and geometric multiplicities  $g_1, \dots, g_k$ . There is a nonsingular matrix



We see that the first principal vector in each Jordan block is an eigenvector of  $\mathbf{A}$ . The remaining principal vectors are not eigenvectors.

**Exercise 5.45 (Jordan example)**

For the Jordan form of the matrix  $\mathbf{A} = \begin{bmatrix} 3 & 0 & 1 \\ -4 & 1 & -2 \\ -4 & 0 & -1 \end{bmatrix}$  we have  $\mathbf{J} = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ . Find  $\mathbf{S}$ .

**Exercise 5.46 (Big Jordan example)**

Find the Jordan form of the matrix

$$\mathbf{A} = \frac{1}{9} \begin{bmatrix} 10 & 16 & -8 & -5 & 6 & 1 & -3 & 4 \\ -7 & 32 & -7 & -10 & 12 & 2 & -6 & 8 \\ -6 & 12 & 12 & -15 & 18 & 3 & -9 & 12 \\ -5 & 10 & -5 & -2 & 24 & 4 & -12 & 16 \\ -4 & 8 & -4 & -16 & 30 & 14 & -15 & 20 \\ -3 & 6 & -3 & -12 & 9 & 24 & -9 & 24 \\ -2 & 4 & -2 & -8 & 6 & -2 & 15 & 28 \\ -1 & 2 & -1 & -4 & 3 & -1 & -6 & 41 \end{bmatrix}. \quad (5.18)$$

The Jordan form implies

**Corollary 5.47 (Geometric multiplicity)**

We have

1. The geometric multiplicity of an eigenvalue is always bounded above by the algebraic multiplicity of the eigenvalue.
2. The number of linearly independent eigenvectors of a matrix equals the sum of the geometric multiplicities of the eigenvalues.
3. A matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  has  $n$  linearly independent eigenvectors if and only if the algebraic and geometric multiplicity of all eigenvalues are the same.

**Proof.**

1. The algebraic multiplicity  $a_i$  of an eigenvalue  $\lambda_i$  is equal to the size of the corresponding  $\mathbf{U}_i$ . Moreover each  $\mathbf{U}_i$  contains  $g_i$  Jordan blocks of size  $m_{i,j} \geq 1$ . Thus  $g_i \leq a_i$ .
2. Since  $\mathbf{A}$  and  $\mathbf{J}$  are similar the geometric multiplicities of the eigenvalues of these matrices are the same, and it is enough to prove statement 2 for the Jordan factor  $\mathbf{J}$ . We show this only for the matrix  $\mathbf{J}$  given by (5.17). The general case should then be clear. There are only 4 eigenvectors of  $\mathbf{J}$ , namely  $\mathbf{e}_1, \mathbf{e}_4, \mathbf{e}_6, \mathbf{e}_7$  corresponding to the 4 Jordan blocks. These 4 vectors are clearly linearly independent. Moreover there are  $k = 2$  distinct eigenvalues and  $g_1 + g_2 = 3 + 1 = 4$ .
3. Since  $g_i \leq a_i$  for all  $i$  and  $\sum_i a_i = n$  we have  $\sum_i g_i = n$  if and only if  $a_i = g_i$  for  $i = 1, \dots, k$ .

□

The following lemma and the following exercise is useful when studying powers of matrices.

**Lemma 5.48 (A nilpotent matrix)**

We have

$$(\mathbf{J}_m(\lambda) - \lambda \mathbf{I})^m = \mathbf{0}, \quad \lambda \in \mathbb{C}; \quad m \in \mathbb{N}.$$

*Proof.* Let  $\mathbf{E}_m := \mathbf{J}_m(\lambda) - \lambda \mathbf{I}$ . For  $m = 4$  we find

$$\mathbf{E}_4 = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \quad \mathbf{E}_4^2 = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \quad \mathbf{E}_4^3 = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \quad \mathbf{E}_4^4 = \mathbf{0}.$$

In general,  $\mathbf{E}_m^r = \begin{bmatrix} \mathbf{0} & \mathbf{I}_{m-r} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$  for  $1 \leq r \leq m-1$ , and it follows that  $\mathbf{E}_m^m = \mathbf{0}$ . □

**Exercise 5.49 (Properties of the Jordan form)**

Let  $\mathbf{J}$  be the Jordan form of a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  as given in Theorem 5.44. Then for  $r = 0, 1, 2, \dots, m = 2, 3, \dots$ , and any  $\lambda \in \mathbb{C}$

1.  $\mathbf{A}^r = \mathbf{S} \mathbf{J}^r \mathbf{S}^{-1}$ ,
2.  $\mathbf{J}^r = \text{diag}(\mathbf{U}_1^r, \dots, \mathbf{U}_k^r)$ ,
3.  $\mathbf{U}_i^r = \text{diag}(\mathbf{J}_{m_i,1}(\lambda_i)^r, \dots, \mathbf{J}_{m_i,g_i}(\lambda_i)^r)$ ,
4.  $\mathbf{J}_m(\lambda)^r = (\mathbf{E}_m + \lambda \mathbf{I}_m)^r = \sum_{k=0}^{\min\{r, m-1\}} \binom{r}{k} \lambda^{r-k} \mathbf{E}_m^k$ .

**Exercise 5.50 (Powers of a Jordan block)**

Find  $\mathbf{J}^{100}$  and  $\mathbf{A}^{100}$  for the matrix in Exercise 5.45.

## 5.5 The Minimal Polynomial

Let  $\mathbf{J}$  be the Jordan form of  $\mathbf{A}$  given in Theorem 5.44. Since  $\mathbf{A}$  and  $\mathbf{J}$  are similar they have the same characteristic polynomial, and since the Jordan form of  $\mathbf{A}$  is upper triangular with the eigenvalues of  $\mathbf{A}$  on the diagonal we have

$$\pi_{\mathbf{A}}(\lambda) = \pi_{\mathbf{J}}(\lambda) = \prod_{i=1}^k \prod_{j=1}^{g_i} (\lambda_i - \lambda)^{m_{i,j}}.$$

The polynomials  $p_{ij}(\lambda) := (\lambda_i - \lambda)^{m_{i,j}}$  are called the **elementary divisors** of  $\mathbf{A}$ . They divide the characteristic polynomial.

**Definition 5.51 (Minimal polynomial of a matrix)**

Suppose  $\mathbf{A} = \mathbf{SJS}^{-1}$  is the Jordan canonical form of  $\mathbf{A}$ . The polynomial

$$\mu(\lambda) := \prod_{i=1}^k (\lambda_i - \lambda)^{m_i} \text{ where } m_i := \max_{1 \leq j \leq g_i} m_{i,j},$$

is called the **minimal polynomial** of  $\mathbf{A}$ .

Since each factor in  $\mu(z)$  is also a factor in  $\pi_{\mathbf{A}}(z)$ , we have the factorization  $\pi_{\mathbf{A}}(z) = \mu(z)\nu(z)$  for some polynomial  $\nu(z)$ .

**Exercise 5.52 (Minimal polynomial example)**

What is the characteristic polynomial and the minimal polynomial of the matrix  $\mathbf{J}$  in (5.17)?

To see in what way the minimal polynomial is minimal, we consider two matrices defined from the characteristic polynomial  $\pi_{\mathbf{A}}$  and the minimal polynomial. We substitute a matrix for the independent variable in these polynomials and define

$$\pi_{\mathbf{A}}(\mathbf{A}) := \prod_{i=1}^k \prod_{j=1}^{g_i} (\lambda_i \mathbf{I} - \mathbf{A})^{m_{i,j}}, \quad \mu(\mathbf{A}) := \prod_{i=1}^k (\lambda_i \mathbf{I} - \mathbf{A})^{m_i}. \quad (5.19)$$

By induction it is easy to see that  $\mu(\mathbf{A})$  and  $\pi_{\mathbf{A}}(\mathbf{A})$  are polynomials in the matrix  $\mathbf{A}$ . Moreover,  $\mu(\mathbf{A}) = \prod_{i=1}^k (\lambda_i \mathbf{I} - \mathbf{SJS}^{-1})^{m_i} = \mathbf{S}\mu(\mathbf{J})\mathbf{S}^{-1}$ , so that  $\mu(\mathbf{A}) = \mathbf{0}$  if and only if  $\mu(\mathbf{J}) = \mathbf{0}$ . Since  $\mathbf{J}$  and  $\mathbf{U}_1, \dots, \mathbf{U}_k$  are block diagonal we find

$$\begin{aligned} \mu(\mathbf{J}) &= \prod_{i=1}^k (\lambda_i \mathbf{I} - \mathbf{J})^{m_i} = \prod_{i=1}^k \text{diag}((\lambda_i \mathbf{I} - \mathbf{U}_1)^{m_i}, \dots, (\lambda_i \mathbf{I} - \mathbf{U}_k)^{m_i}) \\ &= \text{diag}\left(\prod_{i=1}^k (\lambda_i \mathbf{I} - \mathbf{U}_1)^{m_i}, \dots, \prod_{i=1}^k (\lambda_i \mathbf{I} - \mathbf{U}_k)^{m_i}\right) = \mathbf{0}, \end{aligned}$$

since by Lemma 5.48 and the maximality of  $m_r$

$$(\lambda_r \mathbf{I} - \mathbf{U}_r)^{m_r} = \text{diag}((\lambda_r \mathbf{I} - \mathbf{J}_{m_{r,1}})^{m_r}, \dots, (\lambda_r \mathbf{I} - \mathbf{J}_{m_{r,g_r}})^{m_r}) = \mathbf{0}, \quad r = 1, \dots, k.$$

We have shown that a matrix satisfies its minimal polynomial equation  $\mu(\mathbf{A}) = \mathbf{0}$ . Moreover, the degree of any polynomial  $p$  such that  $p(\mathbf{A}) = \mathbf{0}$  is at least as large as the degree  $d = \sum_{i=1}^k m_i$  of the minimal polynomial  $\mu$ . This follows from the proof since any such polynomial must contain the elementary divisors  $(\lambda_i - \lambda)^{m_i}$  for  $i = 1, \dots, k$ . Since the minimal polynomial divides the characteristic polynomial we obtain as a corollary the **Cayley-Hamilton Theorem** which says that a matrix satisfies its characteristic equation  $\pi_{\mathbf{A}}(\mathbf{A}) = \mathbf{0}$ .



**Exercise 5.53 (Similar matrix polynomials)**

Show that  $p(B) = S^{-1}p(A)S$  for any polynomial  $p$  and any similar matrices  $B = S^{-1}AS$ .

**Exercise 5.54 (Minimal polynomial of a diagonalizable matrix)**

What is the minimal polynomial of the unit matrix and more generally of a diagonalizable matrix?

## 5.6 Left Eigenvectors

**Definition 5.55 (Left eigenpair)**

Suppose  $A \in \mathbb{C}^{n \times n}$  is a square matrix,  $\lambda \in \mathbb{C}$  and  $\mathbf{y} \in \mathbb{C}^n$ . We say that  $(\lambda, \mathbf{y})$  is a **left eigenpair** for  $A$  if  $\mathbf{y}^*A = \lambda\mathbf{y}^*$  and  $\mathbf{y}$  is nonzero.

Since  $A^*\mathbf{y} = \bar{\lambda}\mathbf{y}$  Theorem 5.1 implies that  $\lambda$  is an eigenvalue of  $A$ , while a **left eigenvector** is an eigenvector of  $A^*$ . Thus left and right eigenvalues are identical, but left and right eigenvectors are in general different. For an Hermitian matrix the right and left eigenpairs are the same

Left- and right eigenvectors corresponding to distinct eigenvalues are orthogonal.

**Theorem 5.56 (Biorthogonality)**

Suppose  $(\mu, \mathbf{y})$  and  $(\lambda, \mathbf{x})$  are left and right eigenpairs of  $A \in \mathbb{C}^{n \times n}$ . If  $\lambda \neq \mu$  then  $\mathbf{y}^*\mathbf{x} = 0$ .

*Proof.* Using the eigenpair relation in two ways we obtain  $\mathbf{y}^*A\mathbf{x} = \lambda\mathbf{y}^*\mathbf{x} = \mu\mathbf{y}^*\mathbf{x}$  and we conclude that  $\mathbf{y}^*\mathbf{x} = 0$ .  $\square$

Right and left eigenvectors corresponding to the same eigenvalue are sometimes orthogonal, sometimes not.

**Theorem 5.57 (Simple eigenvalue)**

Suppose  $(\lambda, \mathbf{x})$  and  $(\lambda, \mathbf{y})$  are right and left eigenpairs of  $A \in \mathbb{C}^{n \times n}$ . If  $\lambda$  has algebraic multiplicity one then  $\mathbf{y}^*\mathbf{x} \neq 0$ .

*Proof.* Assume that  $\|\mathbf{x}\|_2 = 1$ . We have (cf. (5.4))

$$V^*AV = \left[ \begin{array}{c|c} \lambda & \mathbf{z}^* \\ \hline \mathbf{0} & M \end{array} \right],$$

where  $V$  is unitary and  $V\mathbf{e}_1 = \mathbf{x}$ . We show that if  $\mathbf{y}^*\mathbf{x} = 0$  then  $\lambda$  is also an eigenvalue of  $M$  contradicting the multiplicity assumption of  $\lambda$ . Let  $\mathbf{u} := V^*\mathbf{y}$ . Then

$$(V^*A^*V)\mathbf{u} = V^*A^*\mathbf{y} = \bar{\lambda}V^*\mathbf{y} = \bar{\lambda}\mathbf{u},$$

so  $(\bar{\lambda}, \mathbf{u})$  is an eigenpair of  $\mathbf{V}^* \mathbf{A}^* \mathbf{V}$ . But then  $\mathbf{y}^* \mathbf{x} = \mathbf{u}^* \mathbf{V}^* \mathbf{V} \mathbf{e}_1$ . Suppose that  $\mathbf{u}^* \mathbf{e}_1 = 0$ , i. e.,  $\mathbf{u} = \begin{bmatrix} 0 \\ \mathbf{v} \end{bmatrix}$  for some nonzero  $\mathbf{v} \in \mathbb{C}^{n-1}$ . Then

$$\mathbf{V}^* \mathbf{A}^* \mathbf{V} \mathbf{u} = \begin{bmatrix} \bar{\lambda} & \mathbf{0}^* \\ \mathbf{z} & \mathbf{M}^* \end{bmatrix} \begin{bmatrix} 0 \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} 0 \\ \mathbf{M}^* \mathbf{v} \end{bmatrix} = \bar{\lambda} \begin{bmatrix} 0 \\ \mathbf{v} \end{bmatrix}$$

and by Theorem 5.1 it follows that  $\lambda$  is an eigenvalue of  $\mathbf{M}$ .  $\square$

The case with multiple eigenvalues is more complicated. For example, the matrix  $\mathbf{A} := \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$  has one eigenvalue  $\lambda = 1$  of algebraic multiplicity two, one right eigenvector  $\mathbf{x} = \mathbf{e}_1$  and one left eigenvector  $\mathbf{y} = \mathbf{e}_2$ . Thus  $\mathbf{x}$  and  $\mathbf{y}$  are orthogonal.

**Theorem 5.58 (Biorthogonal eigenvector expansion)**

If  $\mathbf{A} \in \mathbb{C}^{n \times n}$  has linearly independent right eigenvectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  then there exists a set of left eigenvectors  $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$  with  $\mathbf{y}_i^* \mathbf{x}_j = \delta_{i,j}$ . Conversely, if  $\mathbf{A} \in \mathbb{C}^{n \times n}$  has linearly independent left eigenvectors  $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$  then there exists a set of right eigenvectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  with  $\mathbf{y}_i^* \mathbf{x}_j = \delta_{i,j}$ . For any scaling of these sets we have the eigenvector expansions

$$\mathbf{v} = \sum_{j=1}^n \frac{\mathbf{y}_j^* \mathbf{v}}{\mathbf{y}_j^* \mathbf{x}_j} \mathbf{x}_j = \sum_{k=1}^n \frac{\mathbf{x}_k^* \mathbf{v}}{\mathbf{y}_k^* \mathbf{x}_k} \mathbf{y}_k. \quad (5.20)$$

**Proof.** For any right eigenpairs  $(\lambda_1, \mathbf{x}_1), \dots, (\lambda_n, \mathbf{x}_n)$  and left eigenpairs  $(\lambda_1, \mathbf{y}_1), \dots, (\lambda_n, \mathbf{y}_n)$  of  $\mathbf{A}$  we have  $\mathbf{A} \mathbf{X} = \mathbf{X} \mathbf{D}$ ,  $\mathbf{Y}^* \mathbf{A} = \mathbf{D} \mathbf{Y}^*$ , where

$$\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_n], \quad \mathbf{Y} := [\mathbf{y}_1, \dots, \mathbf{y}_n], \quad \mathbf{D} := \text{diag}(\lambda_1, \dots, \lambda_n).$$

If  $\mathbf{X}$  is nonsingular then  $\mathbf{X}^{-1} \mathbf{A} = \mathbf{D} \mathbf{X}^{-1}$  and it follows that  $\mathbf{Y}^* := \mathbf{X}^{-1}$  contains a collection of left eigenvectors such that  $\mathbf{Y}^* \mathbf{X} = \mathbf{I}$ . Thus the columns of  $\mathbf{Y}$  are linearly independent and  $\mathbf{y}_i^* \mathbf{x}_j = \delta_{i,j}$ . Similarly, if  $\mathbf{Y}$  is nonsingular then  $\mathbf{A} \mathbf{Y}^{-*} = \mathbf{Y}^{-*} \mathbf{D}$  and it follows that  $\mathbf{X} := \mathbf{Y}^{-*}$  contains a collection of linearly independent right eigenvectors such that  $\mathbf{Y}^* \mathbf{X} = \mathbf{I}$ . If  $\mathbf{v} = \sum_{j=1}^n c_j \mathbf{x}_j$  then  $\mathbf{y}_i^* \mathbf{v} = \sum_{j=1}^n c_j \mathbf{y}_i^* \mathbf{x}_j = c_i \mathbf{y}_i^* \mathbf{x}_i$ , so  $c_i = \mathbf{y}_i^* \mathbf{v} / \mathbf{y}_i^* \mathbf{x}_i$  for  $i = 1, \dots, n$  and the first expansion in (5.20) follows. The second expansion follows similarly.  $\square$

For an Hermitian matrix the right eigenvectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  are also left eigenvectors and (5.20) takes the form

$$\mathbf{v} = \sum_{j=1}^n \frac{\mathbf{x}_j^* \mathbf{v}}{\mathbf{x}_j^* \mathbf{x}_j} \mathbf{x}_j. \quad (5.21)$$

**Exercise 5.59 (Biorthogonal expansion)**

Determine right and left eigenpairs for the matrix  $\mathbf{A} := \begin{bmatrix} 3 & 1 \\ 2 & 1 \end{bmatrix}$  and the two expansions in (5.20) for any  $\mathbf{v} \in \mathbb{R}^2$ .

**Exercise 5.60 (Generalized Rayleigh quotient)**

For  $\mathbf{A} \in \mathbb{C}^{n \times n}$  and any  $\mathbf{y}, \mathbf{x} \in \mathbb{C}^n$  with  $\mathbf{y}^* \mathbf{x} \neq 0$  the quantity  $R(\mathbf{y}, \mathbf{x}) = R_{\mathbf{A}}(\mathbf{y}, \mathbf{x}) := \frac{\mathbf{y}^* \mathbf{A} \mathbf{x}}{\mathbf{y}^* \mathbf{x}}$  is called a **generalized Rayleigh quotient** for  $\mathbf{A}$ . Show that if  $(\lambda, \mathbf{x})$  is a right eigenpair for  $\mathbf{A}$  then  $R(\mathbf{y}, \mathbf{x}) = \lambda$  for any  $\mathbf{y}$  with  $\mathbf{y}^* \mathbf{x} \neq 0$ . Also show that if  $(\lambda, \mathbf{y})$  is a left eigenpair for  $\mathbf{A}$  then  $R(\mathbf{y}, \mathbf{x}) = \lambda$  for any  $\mathbf{x}$  with  $\mathbf{y}^* \mathbf{x} \neq 0$ .

**5.7 Proof of the Real Schur Form**

In this section we prove the following theorem.

**Theorem 5.61 (The real Schur form)**

Suppose  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Then we can find  $\mathbf{U} \in \mathbb{R}^{n \times n}$  with  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$  such that  $\mathbf{U}^T \mathbf{A} \mathbf{U}$  is quasi-triangular.

*Proof.* If  $\mathbf{A}$  has only real eigenvalues, Theorem 5.14 gives the result. Suppose  $\lambda = \mu + i\nu$ ,  $\mu, \nu \in \mathbb{R}$ , is an eigenvalue of  $\mathbf{A}$  with  $\nu \neq 0$ . Let  $\mathbf{z} = \mathbf{x} + i\mathbf{y}$ ,  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , be an eigenvector of  $\mathbf{A}$  corresponding to  $\lambda$ . Since

$$\mathbf{A} \mathbf{z} = \mathbf{A}(\mathbf{x} + i\mathbf{y}) = (\mu + i\nu)(\mathbf{x} + i\mathbf{y}) = \mu \mathbf{x} - \nu \mathbf{y} + i(\nu \mathbf{x} + \mu \mathbf{y}),$$

we find by comparing real and imaginary parts

$$\mathbf{A} \mathbf{x} = \mu \mathbf{x} - \nu \mathbf{y}, \quad \mathbf{A} \mathbf{y} = \nu \mathbf{x} + \mu \mathbf{y}. \quad (5.22)$$

We claim that  $\mathbf{x}$  and  $\mathbf{y}$  are linearly independent. First we note that  $\nu \neq 0$  implies  $\mathbf{x} \neq \mathbf{0}$ ,  $\mathbf{y} \neq \mathbf{0}$ . For if  $\mathbf{x} = \mathbf{0}$  then (5.22) implies that  $\mathbf{0} = -\nu \mathbf{y}$ , and hence  $\mathbf{y} = \mathbf{0}$  as well, contradicting the nonzeroness of the eigenvector. Similarly, if  $\mathbf{y} = \mathbf{0}$  then  $\mathbf{0} = \nu \mathbf{x}$ , again resulting in a zero eigenvector. Suppose  $\mathbf{y} = \alpha \mathbf{x}$  for some  $\alpha$ . Replacing  $\mathbf{y}$  by  $\alpha \mathbf{x}$  in (5.22), we find  $\mathbf{A} \mathbf{x} = (\mu - \alpha \nu) \mathbf{x}$  and  $\mathbf{A} \mathbf{x} = \mathbf{A} \mathbf{y} / \alpha = (\mu + \nu / \alpha) \mathbf{x}$ . But then  $\mu - \alpha \nu = \mu + \nu / \alpha$  or  $\alpha^2 = -1$ . Since  $\mathbf{x}$  and  $\mathbf{y}$  are real, we cannot have both  $\mathbf{y} = \alpha \mathbf{x}$  and  $\alpha^2 = -1$ . We conclude that  $\mathbf{x}$  and  $\mathbf{y}$  are linearly independent.

(5.22) can be written in matrix form as

$$\mathbf{A} \mathbf{X}_1 = \mathbf{X}_1 \mathbf{M}, \quad \mathbf{X}_1 = [\mathbf{x}, \mathbf{y}] \in \mathbb{R}^{n,2}, \quad \mathbf{M} = \begin{bmatrix} \mu & \nu \\ -\nu & \mu \end{bmatrix}. \quad (5.23)$$

By Theorem 10.12 we can find an orthonormal matrix  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  such that

$$\mathbf{Q} \mathbf{X}_1 = \begin{bmatrix} \mathbf{R} \\ \mathbf{0} \end{bmatrix}$$

where  $\mathbf{R} \in \mathbb{R}^{2,2}$  is upper triangular. Since  $\mathbf{X}_1$  has linearly independent columns,  $\mathbf{R}$  is nonsingular. Let  $\mathbf{Q} = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n]$  and define

$$\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2] = [\mathbf{x}, \mathbf{y}, \mathbf{q}_3, \dots, \mathbf{q}_n].$$

We find

$$QX = [QX_1, Qq_3, \dots, Qq_n] = \begin{bmatrix} R & \mathbf{0} \\ \mathbf{0} & I_{n-2} \end{bmatrix}.$$

Since  $R$  is nonsingular,  $QX$  and  $X$  are nonsingular. Moreover, using (5.23)

$$X^{-1}AX = [X^{-1}AX_1, X^{-1}AX_2] = [X^{-1}X_1M, X^{-1}AX_2] = \begin{bmatrix} M & B \\ \mathbf{0} & C \end{bmatrix}$$

for some matrices  $B \in \mathbb{R}^{2, n-2}$ ,  $C \in \mathbb{R}^{n-2, n-2}$ . Now

$$QAQ^T = (QX)X^{-1}AX(QX)^{-1} = \begin{bmatrix} R & \mathbf{0} \\ \mathbf{0} & I_{n-2} \end{bmatrix} \begin{bmatrix} M & B \\ \mathbf{0} & C \end{bmatrix} \begin{bmatrix} R^{-1} & \mathbf{0} \\ \mathbf{0} & I_{n-2} \end{bmatrix},$$

or

$$QAQ^T = \begin{bmatrix} RMR^{-1} & RB \\ \mathbf{0} & C \end{bmatrix}. \quad (5.24)$$

By Theorem 5.9 the  $2 \times 2$  matrix  $RMR^{-1}$  has the same eigenvalues  $\lambda$  and  $\bar{\lambda}$  as  $M$ . The remaining  $n-2$  eigenvalues of  $A$  are the eigenvalues of  $C$ .

To complete the proof we use induction on  $n$ . The theorem is trivially true for  $n = 1$  and  $n = 2$ . Suppose  $n \geq 3$  and it holds for matrices of order  $\leq n-1$ . Let

$$V = \begin{bmatrix} I_2 & \mathbf{0} \\ \mathbf{0} & \hat{V} \end{bmatrix}$$

where  $\hat{V} \in \mathbb{R}^{n-2, n-2}$ ,  $\hat{V}^T \hat{V} = I_{n-2}$  and  $\hat{V}^T C \hat{V}$  is quasi-triangular. Let  $U = QV$ . Then  $U \in \mathbb{R}^{n \times n}$ ,  $U^T U = I$  and  $U^T A U$  is quasi-triangular.  $\square$

## 5.8 Conclusions

Consider the eigenpair problem for some classes of matrices  $A \in \mathbb{C}^{n \times n}$ .

**Diagonal Matrices.** The eigenpairs are easily determined. Since  $Ae_i = a_{ii}e_i$  the eigenpairs are  $(\lambda_i, e_i)$ , where  $\lambda_i = a_{ii}$  for  $i = 1, \dots, n$ . Moreover,  $a(\lambda_i) = g(\lambda_i)$  for all  $i$ , since the eigenvectors of  $A$  are linearly independent.

**Triangular Matrices** Suppose  $A$  is upper or lower triangular. Since  $\det(A - \lambda I) = \prod_{i=1}^n (a_{ii} - \lambda)$  the eigenvalues are  $\lambda_i = a_{ii}$  for  $i = 1, \dots, n$ , the diagonal elements of  $A$ . To determine the eigenvectors can be challenging as Example 5.40 indicates.

**Block Diagonal Matrices** Suppose

$$A = \text{diag}(A_1, A_2, \dots, A_r), \quad A_i \in \mathbb{C}^{m_i \times m_i}.$$

Here the eigenpair problem reduces to  $r$  smaller problems. Let  $\mathbf{A}_i \mathbf{X}_i = \mathbf{X}_i \mathbf{D}_i$  define the eigenpairs of  $\mathbf{A}_i$  for  $i = 1, \dots, r$  and let  $\mathbf{X} := \text{diag}(\mathbf{X}_1, \dots, \mathbf{X}_r)$ ,  $\mathbf{D} := \text{diag}(\mathbf{D}_1, \dots, \mathbf{D}_r)$ . Then the eigenpairs for  $\mathbf{A}$  are given by

$$\begin{aligned} \mathbf{A}\mathbf{D} &= \text{diag}(\mathbf{A}_1, \dots, \mathbf{A}_r) \text{diag}(\mathbf{X}_1, \dots, \mathbf{X}_r) = \text{diag}(\mathbf{A}_1 \mathbf{X}_1, \dots, \mathbf{A}_r \mathbf{X}_r) \\ &= \text{diag}(\mathbf{X}_1 \mathbf{D}_1, \dots, \mathbf{X}_r \mathbf{D}_r) = \mathbf{X}\mathbf{D}. \end{aligned}$$

**Block Triangular matrices** Let  $\mathbf{A}_{11}, \mathbf{A}_{22}, \dots, \mathbf{A}_{rr}$  be the diagonal blocks of  $\mathbf{A}$ . By Property 8. of determinants

$$\det(\mathbf{A} - \lambda \mathbf{I}) = \prod_{i=1}^r \det(\mathbf{A}_{ii} - \lambda \mathbf{I})$$

and the eigenvalues are found from the eigenvalues of the diagonal blocks.

## 5.9 Review Questions

- 5.9.1** Does  $\mathbf{A}$  and  $\mathbf{A}^T$ ,  $\mathbf{A}$  and  $\mathbf{A}^*$  have the same eigenvalues? What about  $\mathbf{A}^* \mathbf{A}$  and  $\mathbf{A} \mathbf{A}^*$ ?
- 5.9.2** Can the geometric multiplicity of an eigenvalue be bigger than the algebraic multiplicity?
- 5.9.3** What are the eigenvalues of a diagonal matrix?
- 5.9.4** What are the Schur factors of a matrix?
- 5.9.5** What is a quasi-triangular matrix?
- 5.9.6** Give some classes of normal matrices. Why are normal matrices important?
- 5.9.7** State the Courant-Fischer theorem.
- 5.9.8** State the Hoffman-Wieland theorem for Hermitian matrices.
- 5.9.9** What is a left eigenvector of a matrix.



## Chapter 6

# The Singular Value Decomposition

The singular value decomposition is useful both for theory and practice. Some of its applications include solving over-determined equations, principal component analysis in statistics, numerical determination of the rank of a matrix, algorithms used in search engines, and the theory of matrices.

We know from Theorem 5.21 that a square matrix  $\mathbf{A}$  can be diagonalized by a unitary similarity transformation if and only if it is normal, that is  $\mathbf{A}^* \mathbf{A} = \mathbf{A} \mathbf{A}^*$ . In particular, if  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is normal then it has a set of orthonormal eigenpairs  $(\lambda_1, \mathbf{u}_1), \dots, (\lambda_n, \mathbf{u}_n)$ . Letting  $\mathbf{U} := [\mathbf{u}_1, \dots, \mathbf{u}_n] \in \mathbb{C}^{n \times n}$  and  $\mathbf{D} := \text{diag}(\lambda_1, \dots, \lambda_n)$  we have the spectral decomposition

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{U}^*, \text{ where } \mathbf{U}^* \mathbf{U} = \mathbf{I}. \quad (6.1)$$

## 6.1 SVD and SVF

The singular value decomposition (SVD) is a generalization of the spectral decomposition, to any matrix, even a rectangular one. For any  $m, n \in \mathbb{N}$  we say that  $\mathbf{D} \in \mathbb{C}^{m \times n}$  is a diagonal matrix if  $d_{i,j} = 0$  for all  $i \neq j$ . A diagonal matrix is a **nonnegative (positive) diagonal matrix** if all the diagonal elements  $d_{i,i}$ ,  $i = 1, \dots, \min(m, n)$  are nonnegative (positive).

### 6.1.1 Definition and examples

#### Definition 6.1 (SVD)

A decomposition of  $\mathbf{A} \in \mathbb{C}^{m \times n}$  of the form  $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*$ , where  $\mathbf{U} \in \mathbb{C}^{m \times m}$  and  $\mathbf{V} \in \mathbb{C}^{n \times n}$  are unitary, and  $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$  is a nonnegative diagonal matrix, is called a **singular value decomposition (SVD)** of  $\mathbf{A}$ . If  $\mathbf{A}$  is real, then  $\mathbf{U}$  and  $\mathbf{V}$

are real and orthogonal and an SVD takes the form  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$ . The diagonal elements of  $\Sigma$  are denoted  $\sigma_1, \dots, \sigma_{\min(m,n)}$ , and are called **singular values**. The columns  $\mathbf{u}_1, \dots, \mathbf{u}_m$  of  $\mathbf{U}$  and  $\mathbf{v}_1, \dots, \mathbf{v}_n$  of  $\mathbf{V}$  are called **left singular vectors** and **right singular vectors**, respectively. The SVD is **ordered** if  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min(m,n)}$ . For a fixed  $\mathbf{A} \in \mathbb{C}^{m \times n}$  we define  $\sigma_j := 0$  for all integers  $j > \min(m, n)$ .

### Example 6.2 (SVD1)

The decomposition

$$\mathbf{A} := \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} = \mathbf{U}\mathbf{D}\mathbf{U}^T \quad (6.2)$$

is both a spectral decomposition and a singular value decomposition. Indeed,  $\mathbf{A}$  has eigenpairs  $(2, [1, 1]^T)$  and  $(0, [1, -1]^T)$  and normalizing the eigenvectors we obtain a spectral decomposition. Since the elements of the diagonal matrix are nonnegative this is also a singular value decomposition  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$  with  $\Sigma = \mathbf{D}$  and  $\mathbf{V} = \mathbf{U}$ .

### Example 6.3 (SVD2)

The matrix  $\mathbf{A} = \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix}$  is not normal and therefore does not have a spectral decomposition. Since it has an eigenvalue zero of algebraic multiplicity 2, but only one eigenvector  $[1, 1]^T$  it is defective and cannot be diagonalized by any similarity transformation. But

$$\mathbf{A} := \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} =: \mathbf{U}\Sigma\mathbf{V}^T \quad (6.3)$$

is a singular value decomposition.

### Example 6.4 (SVD3)

The matrix  $\mathbf{A} := \frac{1}{25} \begin{bmatrix} 11 & 48 \\ 48 & 39 \end{bmatrix}$  is symmetric with the eigenpairs  $(3, [3, 4]^T)$  and  $(-1, [-4, 3]^T)$ . Normalizing the eigenvectors we obtain the spectral decomposition

$$\mathbf{A} = \frac{1}{25} \begin{bmatrix} 11 & 48 \\ 48 & 39 \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 3 & -4 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & -1 \end{bmatrix} \frac{1}{5} \begin{bmatrix} 3 & 4 \\ -4 & 3 \end{bmatrix} = \mathbf{U}\mathbf{D}\mathbf{U}^T.$$

This is not a singular value decomposition since one of the elements of the diagonal matrix is negative. A singular value decomposition is given by

$$\mathbf{A} = \frac{1}{25} \begin{bmatrix} 11 & 48 \\ 48 & 39 \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 3 & -4 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix} \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix} = \mathbf{U}\Sigma\mathbf{V}^T. \quad (6.4)$$



### 6.1.2 Existence

Every matrix has a singular value decomposition. To show this we first consider the matrices  $\mathbf{A}^* \mathbf{A}$  and  $\mathbf{A} \mathbf{A}^*$ .

**Theorem 6.5 (The matrices  $\mathbf{A}^* \mathbf{A}$ ,  $\mathbf{A} \mathbf{A}^*$ )**

Suppose  $m, n \in \mathbb{N}$  and  $\mathbf{A} \in \mathbb{C}^{m \times n}$ . Then

1. The matrices  $\mathbf{A}^* \mathbf{A} \in \mathbb{C}^{n \times n}$  and  $\mathbf{A} \mathbf{A}^* \in \mathbb{C}^{m \times m}$  are Hermitian with nonnegative eigenvalues.
2. The characteristic polynomials of these matrices are closely related:

$$\lambda^m \pi_{\mathbf{A}^* \mathbf{A}}(\lambda) = \lambda^n \pi_{\mathbf{A} \mathbf{A}^*}(\lambda), \quad \lambda \in \mathbb{C}.$$

3. Let  $(\lambda_j, \mathbf{v}_j)$  be orthonormal eigenpairs for  $\mathbf{A}^* \mathbf{A}$ . If  $\lambda_j > 0$ ,  $j = 1, \dots, r$  and  $\lambda_j = 0$ ,  $j = r + 1, \dots, n$  then  $\{\mathbf{A} \mathbf{v}_1, \dots, \mathbf{A} \mathbf{v}_r\}$  is an orthogonal basis for the column space  $\text{span}(\mathbf{A}) := \{\mathbf{A} \mathbf{y} \in \mathbb{C}^m : \mathbf{y} \in \mathbb{C}^n\}$  and  $\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$  is an orthonormal basis for the nullspace  $\ker(\mathbf{A}) := \{\mathbf{y} \in \mathbb{C}^n : \mathbf{A} \mathbf{y} = \mathbf{0}\}$ .
4. Let  $(\lambda_j, \mathbf{u}_j)$  be orthonormal eigenpairs for  $\mathbf{A} \mathbf{A}^*$ . If  $\lambda_j > 0$ ,  $j = 1, \dots, r$  and  $\lambda_j = 0$ ,  $j = r + 1, \dots, m$  then  $\{\mathbf{A}^* \mathbf{u}_1, \dots, \mathbf{A}^* \mathbf{u}_r\}$  is an orthogonal basis for the column space  $\text{span}(\mathbf{A}^*)$  and  $\{\mathbf{u}_{r+1}, \dots, \mathbf{u}_m\}$  is an orthonormal basis for the nullspace  $\ker(\mathbf{A}^*)$ .
5. The rank of  $\mathbf{A}$  equals the number of positive eigenvalues of  $\mathbf{A}^* \mathbf{A}$  and  $\mathbf{A} \mathbf{A}^*$ .

**Proof.**

1. Clearly  $\mathbf{B}_1 := \mathbf{A}^* \mathbf{A}$  and  $\mathbf{B}_2 := \mathbf{A} \mathbf{A}^*$  are Hermitian. If  $\mathbf{A}^* \mathbf{A} \mathbf{v} = \lambda \mathbf{v}$  with  $\mathbf{v} \neq \mathbf{0}$ , then

$$\lambda = \frac{\mathbf{v}^* \mathbf{A}^* \mathbf{A} \mathbf{v}}{\mathbf{v}^* \mathbf{v}} = \frac{\|\mathbf{A} \mathbf{v}\|_2^2}{\|\mathbf{v}\|_2^2} \geq 0 \quad (6.5)$$

and the eigenvalues of  $\mathbf{B}_1$  are nonnegative. Similarly,  $\mathbf{B}_2$  has nonnegative eigenvalues.

2. This follows from Corollary 5.10.
3. By orthonormality of  $\mathbf{v}_1, \dots, \mathbf{v}_n$  we have  $(\mathbf{A} \mathbf{v}_j)^* \mathbf{A} \mathbf{v}_k = \mathbf{v}_j^* \mathbf{A}^* \mathbf{A} \mathbf{v}_k = \lambda_k \mathbf{v}_j^* \mathbf{v}_k = 0$ ,  $j \neq k$ , showing that  $\mathbf{A} \mathbf{v}_1, \dots, \mathbf{A} \mathbf{v}_n$  are orthogonal vectors. Moreover, (6.5) implies that  $\mathbf{A} \mathbf{v}_1, \dots, \mathbf{A} \mathbf{v}_r$  are nonzero and  $\mathbf{A} \mathbf{v}_j = \mathbf{0}$  for  $j = r + 1, \dots, n$ . In particular, the elements of  $\{\mathbf{A} \mathbf{v}_1, \dots, \mathbf{A} \mathbf{v}_r\}$  and  $\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$  are linearly independent vectors in  $\text{span}(\mathbf{A})$  and  $\ker(\mathbf{A})$ , respectively. The proof will be complete once it is shown that  $\text{span}(\mathbf{A}) \subset \text{span}(\mathbf{A} \mathbf{v}_1, \dots, \mathbf{A} \mathbf{v}_r)$  and  $\ker(\mathbf{A}) \subset \text{span}(\mathbf{v}_{r+1}, \dots, \mathbf{v}_n)$ . Suppose  $\mathbf{x} \in \text{span}(\mathbf{A})$ . Then  $\mathbf{x} = \mathbf{A} \mathbf{y}$

for some  $\mathbf{y} \in \mathbb{C}^n$ . Let  $\mathbf{y} = \sum_{j=1}^n c_j \mathbf{v}_j$  be an eigenvector expansion of  $\mathbf{y}$ . Since  $\mathbf{A}\mathbf{v}_j = \mathbf{0}$  for  $j = r+1, \dots, n$  we obtain  $\mathbf{x} = \mathbf{A}\mathbf{y} = \sum_{j=1}^n c_j \mathbf{A}\mathbf{v}_j = \sum_{j=1}^r c_j \mathbf{A}\mathbf{v}_j \in \text{span}(\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r)$ . Finally, if  $\mathbf{y} = \sum_{j=1}^n c_j \mathbf{v}_j \in \ker(\mathbf{A})$ , then we have  $\mathbf{A}\mathbf{y} = \sum_{j=1}^r c_j \mathbf{A}\mathbf{v}_j = \mathbf{0}$ , and  $c_1 = \dots = c_r = 0$  since  $\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r$  are linearly independent. But then  $\mathbf{y} = \sum_{j=r+1}^n c_j \mathbf{v}_j \in \text{span}(\mathbf{v}_{r+1}, \dots, \mathbf{v}_n)$ .

4. Since  $\mathbf{A}\mathbf{A}^* = \mathbf{B}^*\mathbf{B}$  with  $\mathbf{B} := \mathbf{A}^*$  this follows from part 3 with  $\mathbf{A} = \mathbf{B}$ .
5. By part 1 and 2  $\mathbf{A}^*\mathbf{A}$  and  $\mathbf{A}\mathbf{A}^*$  have the same number  $r$  of positive eigenvalues and by part 3 and 4  $r$  is the rank of  $\mathbf{A}$ .

□

**Theorem 6.6 (The matrices  $\mathbf{A}^*\mathbf{A}$ ,  $\mathbf{A}\mathbf{A}^*$  and SVD)**

If  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}$  is a singular value decomposition of  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and  $\sigma_j := 0$  for  $j > \min(m, n)$  then

1.  $\mathbf{A}^*\mathbf{A} = \mathbf{V} \text{diag}(\sigma_1^2, \dots, \sigma_n^2) \mathbf{V}^*$  is a spectral decomposition of  $\mathbf{A}^*\mathbf{A}$ .
2.  $\mathbf{A}\mathbf{A}^* = \mathbf{U} \text{diag}(\sigma_1^2, \dots, \sigma_m^2) \mathbf{U}^*$  is a spectral decomposition of  $\mathbf{A}\mathbf{A}^*$ .
3. The columns of  $\mathbf{U}$  are orthonormal eigenvectors of  $\mathbf{A}\mathbf{A}^*$ .
4. The columns of  $\mathbf{V}$  are orthonormal eigenvectors of  $\mathbf{A}^*\mathbf{A}$ .
5. The rank of  $\mathbf{A}$  is equal to the number of positive singular values.

*Proof.* We assume  $m \geq n$ . The case  $m < n$  is similar. If  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^* = [\mathbf{u}_1, \dots, \mathbf{u}_m] \Sigma [\mathbf{v}_1, \dots, \mathbf{v}_n]^*$  is a singular value factorization of  $\mathbf{A}$  then  $\mathbf{A}^*\mathbf{A} = (\mathbf{U}\Sigma\mathbf{V}^*)^* (\mathbf{U}\Sigma\mathbf{V}^*) = \mathbf{V} \Sigma^T \mathbf{U}^* \mathbf{U} \Sigma \mathbf{V}^* = \mathbf{V} \Sigma^T \Sigma \mathbf{V}^*$  and part 1 follows. Part 2 is similar. Since these are spectral decompositions part 3 and 4 follow. Part 5 follows from part 5 of Theorem 6.5. □

**Theorem 6.7 (Existence of SVD)**

Suppose for  $m, n, r \in \mathbb{N}$  that  $\mathbf{A} \in \mathbb{C}^{m \times n}$  has rank  $r$ , and that  $(\lambda_j, \mathbf{v}_j)$  are orthonormal eigenpairs for  $\mathbf{A}^*\mathbf{A}$  with  $\lambda_1 \geq \dots \geq \lambda_r > 0 = \lambda_{r+1} = \dots = \lambda_n$ . Define

1.  $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_n] \in \mathbb{C}^{n \times n}$ ,
2.  $\Sigma \in \mathbb{R}^{m \times n}$  is a diagonal matrix with diagonal elements  $\sigma_j := \sqrt{\lambda_j}$  for  $j = 1, \dots, \min(m, n)$ ,

3.  $\mathbf{U} := [\mathbf{u}_1, \dots, \mathbf{u}_m] \in \mathbb{C}^{m \times m}$ , where  $\mathbf{u}_j = \sigma_j^{-1} \mathbf{A}\mathbf{v}_j$  for  $j = 1, \dots, r$  and  $\mathbf{u}_{r+1}, \dots, \mathbf{u}_m$  is an extension of  $\mathbf{u}_1, \dots, \mathbf{u}_r$  to an orthonormal basis  $\mathbf{u}_1, \dots, \mathbf{u}_m$  for  $\mathbb{C}^m$ .

Then  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$  is an ordered singular value decomposition of  $\mathbf{A}$ .

*Proof.* Let  $\mathbf{\Sigma}, \mathbf{U}, \mathbf{V}$  be as in the theorem. The vectors  $\mathbf{u}_1, \dots, \mathbf{u}_r$  are orthonormal since  $\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r$  are orthogonal and  $\sigma_j = \|\mathbf{A}\mathbf{v}_j\|_2 > 0$ ,  $j = 1, \dots, r$  by (6.5). But then  $\mathbf{U}$  and  $\mathbf{V}$  are unitary and  $\mathbf{\Sigma}$  is a nonnegative diagonal matrix. Moreover,

$$\mathbf{U}\mathbf{\Sigma} = \mathbf{U}[\sigma_1\mathbf{e}_1, \dots, \sigma_r\mathbf{e}_r, 0, \dots, 0] = [\sigma_1\mathbf{u}_1, \dots, \sigma_r\mathbf{u}_r, 0, \dots, 0] = [\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r].$$

Thus  $\mathbf{U}\mathbf{\Sigma} = \mathbf{A}\mathbf{V}$  and since  $\mathbf{V}$  is unitary we find  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^* = \mathbf{A}\mathbf{V}\mathbf{V}^* = \mathbf{A}$  and we have an ordered SVD of  $\mathbf{A}$ .  $\square$

### 6.1.3 The singular value factorization

Suppose  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$  is an ordered singular value decomposition of  $\mathbf{A}$  of rank  $r$ . The matrix  $\mathbf{\Sigma}$  can be partitioned in the form

$$\mathbf{\Sigma} = \begin{bmatrix} \mathbf{\Sigma}_1 & \mathbf{0}_{r,n-r} \\ \mathbf{0}_{m-r,r} & \mathbf{0}_{m-r,n-r} \end{bmatrix} \in \mathbb{R}^{m \times n}, \text{ where } \mathbf{\Sigma}_1 := \text{diag}(\sigma_1, \dots, \sigma_r), \quad (6.6)$$

Thus  $\mathbf{\Sigma}_1$  contains the  $r$  ordered positive singular values on the diagonal. Here, for  $k, l \geq 0$  the symbol  $\mathbf{0}_{k,l} = []$  denotes the empty matrix if  $k = 0$  or  $l = 0$ , and the zero matrix with  $k$  rows and  $l$  columns otherwise.

Using the block partitions

$$\begin{aligned} \mathbf{U} &= [\mathbf{U}_1, \mathbf{U}_2] \in \mathbb{C}^{m \times m}, & \mathbf{U}_1 &:= [\mathbf{u}_1, \dots, \mathbf{u}_r], & \mathbf{U}_2 &:= [\mathbf{u}_{r+1}, \dots, \mathbf{u}_m], \\ \mathbf{V} &= [\mathbf{V}_1, \mathbf{V}_2] \in \mathbb{C}^{n \times n}, & \mathbf{V}_1 &:= [\mathbf{v}_1, \dots, \mathbf{v}_r], & \mathbf{V}_2 &:= [\mathbf{v}_{r+1}, \dots, \mathbf{v}_n], \end{aligned} \quad (6.7)$$

we obtain by block multiplication

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^* = \mathbf{U}_1\mathbf{\Sigma}_1\mathbf{V}_1^*. \quad (6.8)$$

As an example:

$$\begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 0 \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \end{bmatrix}.$$

#### Definition 6.8 (SVF)

Let  $m, n, r \in \mathbb{N}$  with  $1 \leq r \leq \min(m, n)$ . A **singular value factorization (SVF)** is a factorization of  $\mathbf{A} \in \mathbb{C}^{m \times n}$  of the form  $\mathbf{A} = \mathbf{U}_1\mathbf{\Sigma}_1\mathbf{V}_1^*$ , where  $\mathbf{U}_1 \in \mathbb{C}^{m \times r}$  and  $\mathbf{V}_1 \in \mathbb{C}^{n \times r}$  have orthonormal columns, and  $\mathbf{\Sigma}_1 \in \mathbb{R}^{r \times r}$  is a diagonal matrix with positive diagonal elements. We say that the SVF is **ordered** if the diagonal elements of  $\mathbf{\Sigma}_1$  are ordered.

An SVD and an SVF of a matrix  $\mathbf{A}$  are closely related.

1. Let  $\mathbf{A}$  have rank  $r$  and let  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$  be an ordered SVD of  $\mathbf{A}$ . Then  $\mathbf{A} = \mathbf{U}_1\mathbf{\Sigma}_1\mathbf{V}_1^*$  is an ordered SVF of  $\mathbf{A}$ . Moreover,  $\mathbf{U}_1, \mathbf{V}_1$  contain the first  $r$  columns of  $\mathbf{U}, \mathbf{V}$  and  $\mathbf{\Sigma}_1$  is a diagonal matrix with the positive singular values on the diagonal.
2. Conversely, suppose  $\mathbf{A} = \mathbf{U}_1\mathbf{\Sigma}_1\mathbf{V}_1^*$  is a singular value factorization of  $\mathbf{A}$  with  $\mathbf{\Sigma}_1 \in \mathbb{R}^{r \times r}$ . Extend  $\mathbf{U}_1$  and  $\mathbf{V}_1$  in any way to unitary matrices  $\mathbf{U} \in \mathbb{C}^{m \times m}$  and  $\mathbf{V} \in \mathbb{C}^{n \times n}$ , and let  $\mathbf{\Sigma}$  be given by (6.6). Then  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$  is an SVD of  $\mathbf{A}$ . Moreover,  $r$  is uniquely given as the rank of  $\mathbf{A}$ .
3. If  $\mathbf{A} = [\mathbf{u}_1, \dots, \mathbf{u}_r] \text{diag}(\sigma_1, \dots, \sigma_r) [\mathbf{v}_1, \dots, \mathbf{v}_r]^*$  is a singular value factorization of  $\mathbf{A}$  then

$$\mathbf{A} = \sum_{j=1}^r \sigma_j \mathbf{u}_j \mathbf{v}_j^*. \quad (6.9)$$

This is known as the **outer product form** of the SVF.

4. We note that a nonsingular square matrix has full rank and only positive singular values. Thus the SVD and SVF are the same for a nonsingular matrix.

### Theorem 6.9 (Singular values of a normal matrix)

*The singular values of a symmetric positive semidefinite matrix are its eigenvalues. The singular values of a normal matrix are the absolute values of its eigenvalues.*

**Proof.** If  $\mathbf{A}$  is normal then by Theorem 5.21,  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^*$ , where  $\mathbf{U}^*\mathbf{U} = \mathbf{I}$  and  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$  contains the eigenvalues of  $\mathbf{A}$ . Now  $\mathbf{A}^*\mathbf{A} = \mathbf{U}\mathbf{D}^*\mathbf{D}\mathbf{U}^*$ , and  $\mathbf{D}^*\mathbf{D} = \text{diag}(|\lambda_1|^2, \dots, |\lambda_n|^2)$  and by Theorem 6.6  $\sigma_j = \sqrt{|\lambda_j|^2} = |\lambda_j|$  for  $j = 1, \dots, n$ . If  $\mathbf{A}$  is symmetric positive semidefinite then the eigenvalues are nonnegative.  $\square$

## 6.1.4 Examples

We use Theorem 6.7 to derive some singular value factorizations and decompositions.

### Example 6.10 (Nonsingular matrix)

*Derive the SVF and SVD of the matrix in (6.4). Discussion: Eigenpairs of  $\mathbf{B} := \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 97 & 96 \\ 96 & 153 \end{bmatrix} / 25$  are given by*

$$\mathbf{B} \begin{bmatrix} 3 \\ 4 \end{bmatrix} = 9 \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 4 \\ -3 \end{bmatrix} = \begin{bmatrix} 4 \\ -3 \end{bmatrix}.$$

Taking square roots and normalizing we find  $\sigma_1 = 3$ ,  $\sigma_2 = 1$ ,  $\mathbf{v}_1 = \begin{bmatrix} 3 \\ 4 \end{bmatrix}/5$ , and  $\mathbf{v}_2 = \begin{bmatrix} 4 \\ -3 \end{bmatrix}/5$ . Thus  $\mathbf{u}_1 := \mathbf{A}\mathbf{v}_1/\sigma_1 = \begin{bmatrix} 3 \\ 4 \end{bmatrix}/5$ ,  $\mathbf{u}_2 := \mathbf{A}\mathbf{v}_2/\sigma_2 = \begin{bmatrix} -4 \\ 3 \end{bmatrix}/5$ , and (6.4) follows. Since  $\mathbf{A}$  is nonsingular this is both an SVF and an SVD of  $\mathbf{A}$ .

**Example 6.11 (Full column rank)**

Find the SVF and SVD of

$$\mathbf{A} = \frac{1}{15} \begin{bmatrix} 14 & 2 \\ 4 & 22 \\ 16 & 13 \end{bmatrix} \in \mathbb{R}^{3,2}.$$

Discussion: Eigenpairs of

$$\mathbf{B} := \mathbf{A}^T \mathbf{A} = \frac{1}{25} \begin{bmatrix} 52 & 36 \\ 36 & 73 \end{bmatrix}$$

are found from

$$\mathbf{B} \begin{bmatrix} 3 \\ 4 \end{bmatrix} = 4 \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 4 \\ -3 \end{bmatrix} = 1 \begin{bmatrix} 4 \\ -3 \end{bmatrix}.$$

Thus  $\sigma_1 = 2$ ,  $\sigma_2 = 1$ , and  $\mathbf{V} = \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix}$ . Now  $\mathbf{u}_1 = \mathbf{A}\mathbf{u}/\sigma_1 = [1, 2, 2]^T/3$ ,  $\mathbf{u}_2 = \mathbf{A}\mathbf{v}_2/\sigma_2 = [2, -2, 1]^T/3$  giving the singular value factorization

$$\mathbf{A} = \frac{1}{3} \begin{bmatrix} 1 & 2 \\ 2 & -2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix}.$$

For an SVD we also need  $\mathbf{u}_3$  which should be orthogonal to  $\mathbf{u}_1$  and  $\mathbf{u}_2$ .  $\mathbf{u}_3 = [2, 1, -2]^T$  is such a vector and normalizing  $\mathbf{u}_3$  we obtain the singular value decomposition

$$\mathbf{A} = \frac{1}{3} \begin{bmatrix} 1 & 2 & 2 \\ 2 & -2 & 1 \\ 2 & 1 & -2 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix}. \quad (6.10)$$

**Example 6.12 (Full row rank)**

Find the SVF and SVD of

$$\mathbf{A}_1 := \frac{1}{15} \begin{bmatrix} 14 & 4 & 16 \\ 2 & 22 & 13 \end{bmatrix} \in \mathbb{R}^{2 \times 3}.$$

Discussion: Since  $\mathbf{A}_1 = \mathbf{A}^T$ , where  $\mathbf{A}$  is the matrix in Example 6.11 we can find an SVF and SVD of  $\mathbf{A}_1$  by simply transposing the corresponding factorization of

*A. Thus*

$$\begin{aligned} \mathbf{A}_1 &= (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T = \mathbf{V}\mathbf{\Sigma}^T\mathbf{U}^T = \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \frac{1}{3} \begin{bmatrix} 1 & 2 & 2 \\ 2 & -2 & 1 \\ 2 & 1 & -2 \end{bmatrix} \\ &= (\mathbf{U}_1\mathbf{\Sigma}_1\mathbf{V}_1^T)^T = \mathbf{V}_1\mathbf{\Sigma}_1^T\mathbf{U}_1^T = \frac{1}{5} \begin{bmatrix} 3 & 4 \\ 4 & -3 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \frac{1}{3} \begin{bmatrix} 1 & 2 & 2 \\ 2 & -2 & 1 \end{bmatrix}. \end{aligned} \quad (6.11)$$

**Example 6.13** ( $r < n < m$ )

Find the SVD of

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix}.$$

*Discussion: Eigenpairs of*

$$\mathbf{B} := \mathbf{A}^T\mathbf{A} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}$$

are derived from

$$\mathbf{B} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 4 \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \mathbf{B} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = 0 \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

and we find  $\sigma_1 = 2$ ,  $\sigma_2 = 0$ , Thus  $r = 1$ ,  $m = 3$ ,  $n = 2$  and

$$\mathbf{\Sigma} = \begin{bmatrix} \mathbf{\Sigma}_1 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{\Sigma}_1 = [2], \quad \mathbf{V} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

We find  $\mathbf{u}_1 = \mathbf{A}\mathbf{v}_1/\sigma_1 = \mathbf{s}_1/\sqrt{2}$ , where  $\mathbf{s}_1 = [1, 1, 0]^T$ , and the SVF of  $\mathbf{A}$  is given by

$$\mathbf{A} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} [2] \frac{1}{\sqrt{2}} [1 \quad 1].$$

To find an SVD we need to extend  $\mathbf{u}_1$  to an orthonormal basis for  $\mathbb{R}^3$ . We first extend  $\mathbf{s}_1$  to a basis  $\{\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3\}$  for  $\mathbb{R}^3$ , apply the Gram-Schmidt orthogonalization process to  $\{\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3\}$ , and then normalize. Choosing the basis

$$\mathbf{s}_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{s}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{s}_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$$

we find from (26)

$$\mathbf{w}_1 = \mathbf{s}_1, \quad \mathbf{w}_2 = \mathbf{s}_2 - \frac{\mathbf{s}_2^T\mathbf{w}_1}{\mathbf{w}_1^T\mathbf{w}_1}\mathbf{w}_1 = \begin{bmatrix} -1/2 \\ 1/2 \\ 0 \end{bmatrix}, \quad \mathbf{w}_3 = \mathbf{s}_3 - \frac{\mathbf{s}_3^T\mathbf{w}_1}{\mathbf{w}_1^T\mathbf{w}_1}\mathbf{w}_1 - \frac{\mathbf{s}_3^T\mathbf{w}_2}{\mathbf{w}_2^T\mathbf{w}_2}\mathbf{w}_2 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

Normalizing the  $\mathbf{w}_i$ 's we obtain  $\mathbf{u}_1 = \mathbf{w}_1/\|\mathbf{w}_1\|_2 = [1/\sqrt{2}, 1/\sqrt{2}, 0]^T$ ,  $\mathbf{u}_2 = \mathbf{w}_2/\|\mathbf{w}_2\|_2 = [-1/\sqrt{2}, 1/\sqrt{2}, 0]^T$ , and  $\mathbf{u}_3 = \mathbf{s}_3/\|\mathbf{s}_3\|_2 = [0, 0, 1]^T$ . Therefore,  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ , where

$$\mathbf{U} := \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} & 0 \\ 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 0 & 0 & 1 \end{bmatrix} \in \mathbb{R}^{3,3}, \quad \mathbf{\Sigma} := \begin{bmatrix} 2 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{3,2}, \quad \mathbf{V} := \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \in \mathbb{R}^{2,2}.$$

The method we used to find the singular value decomposition in the previous examples and exercises can be suitable for hand calculation with small matrices, but it is not appropriate as a basis for a general purpose numerical method. In particular, the Gram-Schmidt orthogonalization process is not numerically stable, and forming  $\mathbf{A}^*\mathbf{A}$  can lead to extra errors in the computation. Standard computer implementations of the singular value decomposition ([29]) first reduces  $\mathbf{A}$  to bidiagonal form and then use an adapted version of the QR algorithm where the matrix  $\mathbf{A}^*\mathbf{A}$  is not formed. The QR algorithm is discussed in Chapter 13.

### Exercise 6.14 (SVD examples)

Find the singular value decomposition of the following matrices

(a)  $\mathbf{A} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$ .

(b)  $\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 2 & 2 \\ 2 & 2 \end{bmatrix}$ .

### Exercise 6.15 (More SVD examples)

Find the singular value decomposition of the following matrices

(a)  $\mathbf{A} = \mathbf{e}_1$  the first unit vector in  $\mathbb{R}^m$ .

(b)  $\mathbf{A} = \mathbf{e}_n^T$  the last unit vector in  $\mathbb{R}^n$ .

(c)  $\mathbf{A} = \begin{bmatrix} -1 & 0 \\ 0 & 3 \end{bmatrix}$ .

## 6.2 SVD and the Four Fundamental Subspaces

The singular vectors form orthonormal bases for the four fundamental subspaces  $\text{span}(\mathbf{A})$ ,  $\text{ker}(\mathbf{A})$ ,  $\text{span}(\mathbf{A}^*)$ , and  $\text{ker}(\mathbf{A}^*)$ .

### Theorem 6.16 (Singular vectors and orthonormal bases)

For positive integers  $m, n$  let  $\mathbf{A} \in \mathbb{C}^{m \times n}$  have rank  $r$  and a singular value decomposition  $\mathbf{A} = [\mathbf{u}_1, \dots, \mathbf{u}_m]\mathbf{\Sigma}[\mathbf{v}_1, \dots, \mathbf{v}_n]^* = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ . Then the singular vectors satisfy

$$\begin{aligned} \mathbf{A}\mathbf{v}_i &= \sigma_i\mathbf{u}_i, \quad i = 1, \dots, r, & \mathbf{A}\mathbf{v}_i &= 0, \quad i = r + 1, \dots, n, \\ \mathbf{A}^*\mathbf{u}_i &= \sigma_i\mathbf{v}_i, \quad i = 1, \dots, r, & \mathbf{A}^*\mathbf{u}_i &= 0, \quad i = r + 1, \dots, m. \end{aligned} \quad (6.12)$$

Moreover,

1.  $\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$  is an orthonormal basis for  $\text{span}(\mathbf{A})$ ,
  2.  $\{\mathbf{u}_{r+1}, \dots, \mathbf{u}_m\}$  is an orthonormal basis for  $\ker(\mathbf{A}^*)$ ,
  3.  $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$  is an orthonormal basis for  $\text{span}(\mathbf{A}^*)$ ,
  4.  $\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$  is an orthonormal basis for  $\ker(\mathbf{A})$ .
- (6.13)

**Proof.** If  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$  then  $\mathbf{A}\mathbf{V} = \mathbf{U}\mathbf{\Sigma}$ , or in terms of the block partition (6.7)  $\mathbf{A}[\mathbf{V}_1, \mathbf{V}_2] = [\mathbf{U}_1, \mathbf{U}_2] \begin{bmatrix} \mathbf{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ . But then  $\mathbf{A}\mathbf{V}_1 = \mathbf{U}_1\mathbf{\Sigma}_1$ ,  $\mathbf{A}\mathbf{V}_2 = \mathbf{0}$ , and this implies the first part of (6.12). Taking conjugate transpose of  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$  gives  $\mathbf{A}^* = \mathbf{V}\mathbf{\Sigma}^T\mathbf{U}^*$  or  $\mathbf{A}^*\mathbf{U} = \mathbf{V}\mathbf{\Sigma}^T$ . Using the block partition as before we obtain the last part of (6.12).

It follows from Theorem 6.5 that  $\{\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r\}$  is an orthogonal basis for  $\text{span}(\mathbf{A})$  and  $\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_m\}$  is an orthonormal basis for  $\ker(\mathbf{A})$ . Applying this theorem to  $\mathbf{A}\mathbf{A}^*$  it also follows that  $\{\mathbf{A}^*\mathbf{u}_1, \dots, \mathbf{A}^*\mathbf{u}_r\}$  is an orthogonal basis for  $\text{span}(\mathbf{A}^*)$  and  $\{\mathbf{u}_{r+1}, \dots, \mathbf{u}_m\}$  is an orthonormal basis for  $\ker(\mathbf{A}^*)$ . By (6.12)  $\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$  is an orthonormal basis for  $\text{span}(\mathbf{A})$  and  $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$  is an orthonormal basis for  $\text{span}(\mathbf{A}^*)$ .  $\square$

### Exercise 6.17 (Counting dimensions of fundamental subspaces)

Suppose  $\mathbf{A} \in \mathbb{C}^{m \times n}$ . Show using SVD that

1.  $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^*)$ .
2.  $\text{rank}(\mathbf{A}) + \text{null}(\mathbf{A}) = n$ ,
3.  $\text{rank}(\mathbf{A}) + \text{null}(\mathbf{A}^*) = m$ ,

where  $\text{null}(\mathbf{A})$  is defined as the dimension of  $\ker(\mathbf{A})$ .

### Exercise 6.18 (Rank and nullity relations)

Use Theorem 6.5 to show that for any  $\mathbf{A} \in \mathbb{C}^{m \times n}$

1.  $\text{rank } \mathbf{A} = \text{rank}(\mathbf{A}^*\mathbf{A}) = \text{rank}(\mathbf{A}\mathbf{A}^*)$ ,
2.  $\text{null}(\mathbf{A}^*\mathbf{A}) = \text{null } \mathbf{A}$ , and  $\text{null}(\mathbf{A}\mathbf{A}^*) = \text{null}(\mathbf{A}^*)$ .

### Exercise 6.19 (Orthonormal bases example)

Let  $\mathbf{A}$  and  $\mathbf{B}$  be as in Example 6.11. Give orthonormal bases for  $\text{span}(\mathbf{B})$  and  $\ker(\mathbf{B})$ .

### Exercise 6.20 (Some spanning sets)

Show for any  $\mathbf{A} \in \mathbb{C}^{m \times n}$  that  $\text{span}(\mathbf{A}^*\mathbf{A}) = \text{span}(\mathbf{V}_1) = \text{span}(\mathbf{A}^*)$



**Exercise 6.21 (Singular values and eigenpair of composite matrix)**

Let  $\mathbf{A} \in \mathbb{C}^{m \times n}$  with  $m \geq n$  have singular values  $\sigma_1, \dots, \sigma_n$ , left singular vectors  $\mathbf{u}_1, \dots, \mathbf{u}_m \in \mathbb{C}^m$ , and right singular vectors  $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{C}^n$ . Show that the matrix

$$\mathbf{C} := \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^* & \mathbf{0} \end{bmatrix}$$

has the  $n + m$  eigenpairs

$$\{(\sigma_1, \mathbf{p}_1), \dots, (\sigma_n, \mathbf{p}_n), (-\sigma_1, \mathbf{q}_1), \dots, (-\sigma_n, \mathbf{q}_n), (0, \mathbf{r}_{n+1}), \dots, (0, \mathbf{r}_m)\},$$

where

$$\mathbf{p}_i = \begin{bmatrix} \mathbf{u}_i \\ \mathbf{v}_i \end{bmatrix}, \quad \mathbf{q}_i = \begin{bmatrix} \mathbf{u}_i \\ -\mathbf{v}_i \end{bmatrix}, \quad \mathbf{r}_j = \begin{bmatrix} \mathbf{u}_j \\ \mathbf{0} \end{bmatrix}, \quad \text{for } i = 1, \dots, n \text{ and } j = n + 1, \dots, m.$$

### 6.3 A Geometric Interpretation

The singular value decomposition and factorization give insight into the geometry of a linear transformation. Consider the linear transformation  $\mathbf{T} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  given by  $\mathbf{Tz} := \mathbf{Az}$  where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Assume that  $\text{rank}(\mathbf{A}) = n$ . The function  $\mathbf{T}$  maps the unit sphere  $\mathcal{S} := \{\mathbf{z} \in \mathbb{R}^n : \|\mathbf{z}\|_2 = 1\}$  onto an ellipsoid  $\mathcal{E} := \mathbf{AS} = \{\mathbf{Az} : \mathbf{z} \in \mathcal{S}\}$  in  $\mathbb{R}^m$ .

**Theorem 6.22 (SVF ellipse)**

Suppose  $\mathbf{A} \in \mathbb{R}^{m,n}$  has rank  $r = n$ , and let  $\mathbf{A} = \mathbf{U}_1 \Sigma_1 \mathbf{V}_1^T$  be a singular value factorization of  $\mathbf{A}$ . Then

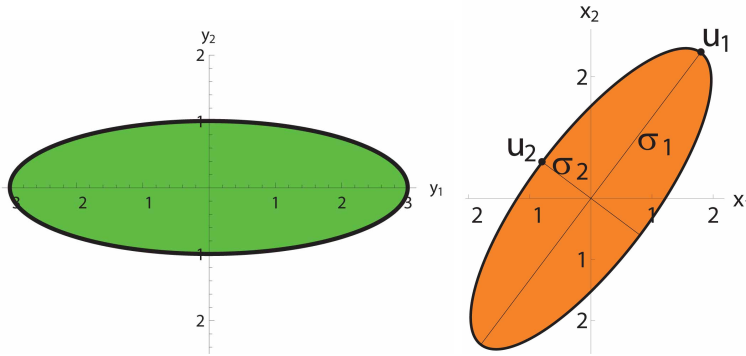
$$\mathcal{E} = \mathbf{U}_1 \tilde{\mathcal{E}} \text{ where } \tilde{\mathcal{E}} := \{\mathbf{y} = [y_1, \dots, y_n]^T \in \mathbb{R}^n : \frac{y_1^2}{\sigma_1^2} + \dots + \frac{y_n^2}{\sigma_n^2} = 1\}.$$

**Proof.** Suppose  $\mathbf{z} \in \mathcal{S}$ . Now  $\mathbf{Az} = \mathbf{U}_1 \Sigma_1 \mathbf{V}_1^T \mathbf{z} = \mathbf{U}_1 \mathbf{y}$ , where  $\mathbf{y} := \Sigma_1 \mathbf{V}_1^T \mathbf{z}$ . Since  $\text{rank}(\mathbf{A}) = n$  it follows that  $\mathbf{V}_1 = \mathbf{V}$  is square so that  $\mathbf{V}_1 \mathbf{V}_1^T = \mathbf{I}$ . But then  $\mathbf{V}_1 \Sigma_1^{-1} \mathbf{y} = \mathbf{z}$  and we obtain

$$1 = \|\mathbf{z}\|_2^2 = \|\mathbf{V}_1 \Sigma_1^{-1} \mathbf{y}\|_2^2 = \|\Sigma_1^{-1} \mathbf{y}\|_2^2 = \frac{y_1^2}{\sigma_1^2} + \dots + \frac{y_n^2}{\sigma_n^2}.$$

This implies that  $\mathbf{y} \in \tilde{\mathcal{E}}$ . Finally,  $\mathbf{x} = \mathbf{Az} = \mathbf{U}_1 \Sigma_1 \mathbf{V}_1^T \mathbf{z} = \mathbf{U}_1 \mathbf{y}$ , where  $\mathbf{y} \in \tilde{\mathcal{E}}$  implies that  $\mathcal{E} = \mathbf{U}_1 \tilde{\mathcal{E}}$ .  $\square$

The equation  $1 = \frac{y_1^2}{\sigma_1^2} + \dots + \frac{y_n^2}{\sigma_n^2}$  describes an ellipsoid in  $\mathbb{R}^n$  with semi-axes of length  $\sigma_j$  along the unit vectors  $\mathbf{e}_j$  for  $j = 1, \dots, n$ . Since the orthonormal transformation  $\mathbf{U}_1 \mathbf{y} \rightarrow \mathbf{x}$  preserves length, the image  $\mathcal{E} = \mathbf{AS}$  is a rotated ellipsoid



**Figure 6.1.** The ellipse  $y_1^2/9 + y_2^2 = 1$  (left) and the rotated ellipse  $\mathbf{A}\mathcal{S}$  (right).

with semiaxes along the left singular vectors  $\mathbf{u}_j = \mathbf{U}\mathbf{e}_j$ , of length  $\sigma_j$ ,  $j = 1, \dots, n$ . Since  $\mathbf{A}\mathbf{v}_j = \sigma_j\mathbf{u}_j$ , for  $j = 1, \dots, n$  the right singular vectors defines points in  $\mathcal{S}$  that are mapped onto the semiaxes of  $\mathcal{E}$ .

**Example 6.23 (Ellipse)**

Consider the transformation  $\mathbf{A} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  given by the matrix

$$\mathbf{A} := \frac{1}{25} \begin{bmatrix} 11 & 48 \\ 48 & 39 \end{bmatrix}$$

in Example 6.10. Recall that  $\sigma_1 = 3$ ,  $\sigma_2 = 1$ ,  $\mathbf{u}_1 = [3, 4]^T/5$  and  $\mathbf{u}_2 = [-4, 3]^T/5$ . The ellipses  $y_1^2/\sigma_1^2 + y_2^2/\sigma_2^2 = 1$  and  $\mathcal{E} = \mathbf{A}\mathcal{S} = \mathbf{U}_1\tilde{\mathcal{E}}$  are shown in Figure 6.1. Since  $\mathbf{y} = \mathbf{U}_1^T\mathbf{x} = [3/5x_1 + 4/5x_2, -4/5x_1 + 3/5x_2]^T$ , the equation for the ellipse on the right is

$$\frac{(\frac{3}{5}x_1 + \frac{4}{5}x_2)^2}{9} + \frac{(-\frac{4}{5}x_1 + \frac{3}{5}x_2)^2}{1} = 1,$$

## 6.4 Determining the Rank of a Matrix Numerically

In many elementary linear algebra courses a version of Gaussian elimination, called Gauss-Jordan elimination, is used to determine the rank of a matrix. To carry this out by hand for a large matrix can be a Herculean task and using a computer and floating point arithmetic the result will not be reliable. Entries, which in the final result should have been zero, will have nonzero values because of round-off errors. As an alternative we can use the singular value decomposition to determine rank. Although success is not at all guaranteed, the result will be more reliable than if Gauss-Jordan elimination is used.

By Theorem 6.7 the rank of a matrix is equal to the number of nonzero singular values and if we have computed the singular values, then all we have to do is to count the nonzero ones. The problem however is the same as for Gaussian elimination. Due to round-off errors none of the computed singular values are likely to be zero.

### 6.4.1 The Frobenius norm

This commonly occurring matrix norm will be used here in a discussion of how many of the computed singular values can possibly be considered to be zero. The **Frobenius norm**, of a matrix  $\mathbf{A} \in \mathbb{C}^{m \times n}$  is defined by

$$\|\mathbf{A}\|_F := \left( \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{1/2}. \quad (6.14)$$

There is a relation between the Frobenius norm of a matrix and its singular values. First we derive some elementary properties of this norm. A systematic study of matrix norms is given in the next chapter.

#### Lemma 6.24 (Frobenius norm properties)

For any  $m, n \in \mathbb{N}$  and any matrix  $\mathbf{A} \in \mathbb{C}^{m \times n}$

1.  $\|\mathbf{A}^*\|_F = \|\mathbf{A}\|_F$ ,
2.  $\|\mathbf{A}\|_F^2 = \sum_{j=1}^n \|\mathbf{a}_{:j}\|_2^2$ ,
3.  $\|\mathbf{U}\mathbf{A}\|_F = \|\mathbf{A}\mathbf{V}\|_F = \|\mathbf{A}\|_F$  for any unitary matrices  $\mathbf{U} \in \mathbb{C}^{m \times m}$  and  $\mathbf{V} \in \mathbb{C}^{n \times n}$ ,
4.  $\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F$  for any  $\mathbf{B} \in \mathbb{C}^{n,k}$ ,  $k \in \mathbb{N}$ ,
5.  $\|\mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{A}\|_F \|\mathbf{x}\|_2$ , for all  $\mathbf{x} \in \mathbb{C}^n$ .

*Proof.*

1.  $\|\mathbf{A}^*\|_F^2 = \sum_{j=1}^n \sum_{i=1}^m |\bar{a}_{ij}|^2 = \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 = \|\mathbf{A}\|_F^2$ .
2. This follows since the Frobenius norm is the Euclidian norm of a vector,  $\|\mathbf{A}\|_F := \|\text{vec}(\mathbf{A})\|_2$ , where  $\text{vec}(\mathbf{A}) \in \mathbb{C}^{mn}$  is the vector obtained by stacking the columns of  $\mathbf{A}$  on top of each other.
3. Recall that if  $\mathbf{U}^*\mathbf{U} = \mathbf{I}$  then  $\|\mathbf{U}\mathbf{x}\|_2 = \|\mathbf{x}\|_2$  for all  $\mathbf{x} \in \mathbb{C}^n$ . Applying this to each column  $\mathbf{a}_{:j}$  of  $\mathbf{A}$  we find  $\|\mathbf{U}\mathbf{A}\|_F^2 \stackrel{2.}{=} \sum_{j=1}^n \|\mathbf{U}\mathbf{a}_{:j}\|_2^2 = \sum_{j=1}^n \|\mathbf{a}_{:j}\|_2^2 \stackrel{2.}{=} \|\mathbf{A}\|_F^2$ . Similarly, since  $\mathbf{V}\mathbf{V}^* = \mathbf{I}$  we find  $\|\mathbf{A}\mathbf{V}\|_F \stackrel{1.}{=} \|\mathbf{V}^*\mathbf{A}\|_F = \|\mathbf{A}\|_F \stackrel{1.}{=} \|\mathbf{A}\|_F$ .

4. Using the Cauchy-Schwarz inequality and 2. we obtain

$$\|\mathbf{AB}\|_F^2 = \sum_{i=1}^m \sum_{j=1}^k |\mathbf{a}_i^T \mathbf{b}_{:j}|^2 \leq \sum_{i=1}^m \sum_{j=1}^k \|\mathbf{a}_i\|_2^2 \|\mathbf{b}_{:j}\|_2^2 = \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2.$$

5. Since  $\|\mathbf{v}\|_F = \|\mathbf{v}\|_2$  for a vector this follows by taking  $k = 1$  and  $\mathbf{B} = \mathbf{x}$  in 4.

□

### Theorem 6.25 (Frobenius norm and singular values)

We have  $\|\mathbf{A}\|_F = \sqrt{\sigma_1^2 + \cdots + \sigma_n^2}$ , where  $\sigma_1, \dots, \sigma_n$  are the singular values of  $\mathbf{A}$ .

*Proof.* Using Lemma 6.24 we find  $\|\mathbf{A}\|_F \stackrel{3.}{=} \|\mathbf{U}^* \mathbf{A} \mathbf{V}\|_F = \|\mathbf{\Sigma}\|_F = \sqrt{\sigma_1^2 + \cdots + \sigma_n^2}$ .  
□

## 6.4.2 Low rank approximation

Suppose  $m \geq n \geq 1$  and  $\mathbf{A} \in \mathbb{C}^{m \times n}$  has an ordered singular value decomposition  $\mathbf{A} = \mathbf{U} \begin{bmatrix} \mathbf{D} \\ \mathbf{0} \end{bmatrix} \mathbf{V}^*$ , where  $\mathbf{D} = \text{diag}(\sigma_1, \dots, \sigma_n)$ . We choose  $\epsilon > 0$  and let  $1 \leq r \leq n$  be the smallest integer such that  $\sigma_{r+1}^2 + \cdots + \sigma_n^2 < \epsilon^2$ . Define  $\mathbf{A}' := \mathbf{U} \begin{bmatrix} \mathbf{D}' \\ \mathbf{0} \end{bmatrix} \mathbf{V}^*$ , where  $\mathbf{D}' := \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0) \in \mathbb{R}^{n \times n}$ . By Lemma 6.24

$$\|\mathbf{A} - \mathbf{A}'\|_F = \|\mathbf{U} \begin{bmatrix} \mathbf{D} - \mathbf{D}' \\ \mathbf{0} \end{bmatrix} \mathbf{V}^*\|_F = \|\begin{bmatrix} \mathbf{D} - \mathbf{D}' \\ \mathbf{0} \end{bmatrix}\|_F = \sqrt{\sigma_{r+1}^2 + \cdots + \sigma_n^2} < \epsilon.$$

Thus, if  $\epsilon$  is small then  $\mathbf{A}$  is near a matrix  $\mathbf{A}'$  of rank  $r$ . This can be used to determine rank numerically. We choose an  $r$  such that  $\sqrt{\sigma_{r+1}^2 + \cdots + \sigma_n^2}$  is “small”. Then we postulate that  $\text{rank}(\mathbf{A}) = r$  since  $\mathbf{A}$  is close to a matrix of rank  $r$ .

The following theorem shows that of all  $m \times n$  matrices of rank  $r$ ,  $\mathbf{A}'$  is closest to  $\mathbf{A}$  measured in the Frobenius norm.

### Theorem 6.26 (Best low rank approximation)

Suppose  $\mathbf{A} \in \mathbb{R}^{m \times n}$  has singular values  $\sigma_1 \geq \cdots \geq \sigma_n \geq 0$ . For any  $r \leq \text{rank}(\mathbf{A})$  we have

$$\|\mathbf{A} - \mathbf{A}'\|_F = \min_{\substack{\mathbf{B} \in \mathbb{R}^{m \times n} \\ \text{rank}(\mathbf{B})=r}} \|\mathbf{A} - \mathbf{B}\|_F = \sqrt{\sigma_{r+1}^2 + \cdots + \sigma_n^2}.$$

For the proof of this theorem we refer to p. 322 of [29].

**Exercise 6.27 (Rank example)**

Consider the singular value decomposition

$$\mathbf{A} := \begin{bmatrix} 0 & 3 & 3 \\ 4 & 1 & -1 \\ 4 & 1 & -1 \\ 0 & 3 & 3 \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 6 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} \\ \frac{1}{\sqrt{3}} & -\frac{1}{\sqrt{3}} & -\frac{1}{\sqrt{3}} \\ \frac{1}{\sqrt{3}} & -\frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} \end{bmatrix}$$

- (a) Give orthonormal bases for  $\text{span}(\mathbf{A})$ ,  $\text{span}(\mathbf{A}^T)$ ,  $\ker(\mathbf{A})$ ,  $\ker(\mathbf{A}^T)$  and  $\text{span}(\mathbf{A})^\perp$ .
- (b) Explain why for all matrices  $\mathbf{B} \in \mathbb{R}^{4,3}$  of rank one we have  $\|\mathbf{A} - \mathbf{B}\|_F \geq 6$ .
- (c) Give a matrix  $\mathbf{A}_1$  of rank one such that  $\|\mathbf{A} - \mathbf{A}_1\|_F = 6$ .

**Exercise 6.28 (Another rank example)**

Let  $\mathbf{A}$  be the  $n \times n$  matrix that for  $n = 4$  takes the form

$$\mathbf{A} = \begin{bmatrix} 1 & -1 & -1 & -1 \\ 0 & 1 & -1 & -1 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

Thus  $\mathbf{A}$  is upper triangular with diagonal elements one and all elements above the diagonal equal to  $-1$ . Let  $\mathbf{B}$  be the matrix obtained from  $\mathbf{A}$  by changing the  $(n, 1)$  element from zero to  $-2^{2-n}$ .

- (a) Show that  $\mathbf{B}\mathbf{x} = \mathbf{0}$ , where  $\mathbf{x} := [2^{n-2}, 2^{n-3}, \dots, 2^0, 1]^T$ . Conclude that  $\mathbf{B}$  is singular,  $\det(\mathbf{A}) = 1$ , and  $\|\mathbf{A} - \mathbf{B}\|_F = 2^{2-n}$ . Thus even if  $\det(\mathbf{A})$  is not small the Frobenius norm of  $\mathbf{A} - \mathbf{B}$  is small for large  $n$ , and the matrix  $\mathbf{A}$  is very close to being singular for large  $n$ .
- (b) Use Theorem 6.26 to show that the smallest singular value  $\sigma_n$  of  $\mathbf{A}$  is bounded above by  $2^{2-n}$ .

## 6.5 The Minmax Theorem for Singular Values and the Hoffman-Wielandt Theorem

We have a minmax and maxmin characterization for singular values.

**Theorem 6.29 (The Courant-Fischer theorem for singular values)**

Suppose  $\mathbf{A} \in \mathbb{C}^{m \times n}$  has singular values  $\sigma_1, \sigma_2, \dots, \sigma_n$  ordered so that  $\sigma_1 \geq \dots \geq \sigma_n$ . Then for  $k = 1, \dots, n$

$$\sigma_k = \min_{\dim(S)=n-k+1} \max_{\substack{\mathbf{x} \in S \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2} = \max_{\dim(S)=k} \min_{\substack{\mathbf{x} \in S \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2}. \tag{6.15}$$

*Proof.* Since

$$\frac{\|\mathbf{Ax}\|_2^2}{\|\mathbf{x}\|_2^2} = \frac{(\mathbf{Ax})^*(\mathbf{Ax})}{\mathbf{x}^*\mathbf{x}} = \frac{\mathbf{x}^*(\mathbf{A}^*\mathbf{A})\mathbf{x}}{\mathbf{x}^*\mathbf{x}}$$

is the Rayleigh quotient  $R_{\mathbf{A}^*\mathbf{A}}(\mathbf{x})$  of  $\mathbf{A}^*\mathbf{A}$ , and since the singular values of  $\mathbf{A}$  are the nonnegative square roots of the eigenvalues of  $\mathbf{A}^*\mathbf{A}$ , the results follow from the Courant-Fischer Theorem for eigenvalues, see Theorem 5.29.  $\square$

By taking  $k = 1$  and  $k = n$  in (6.15) we obtain for any  $\mathbf{A} \in \mathbb{C}^{m \times n}$

$$\sigma_1 = \max_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2}, \quad \sigma_n = \min_{\substack{\mathbf{x} \in \mathbb{C}^n \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2}. \quad (6.16)$$

This follows since the only subspace of  $\mathbb{C}^n$  of dimension  $n$  is  $\mathbb{C}^n$  itself.

The Hoffman-Wielandt Theorem, see Theorem 5.32, for eigenvalues of Hermitian matrices can be written

$$\sum_{j=1}^n |\mu_j - \lambda_j|^2 \leq \|\mathbf{A} - \mathbf{B}\|_F^2 := \sum_{i=1}^n \sum_{j=1}^n |a_{ij} - b_{ij}|^2, \quad (6.17)$$

where  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$  are both Hermitian matrices with eigenvalues  $\lambda_1 \geq \dots \geq \lambda_n$  and  $\mu_1 \geq \dots \geq \mu_n$ , respectively.

For singular values we have a similar result, see also Section 11.6.

**Theorem 6.30 (Hoffman-Wielandt theorem for singular values)**

For any  $m, n \in \mathbb{N}$  and  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{m \times n}$  we have

$$\sum_{j=1}^n |\beta_j - \alpha_j|^2 \leq \|\mathbf{A} - \mathbf{B}\|_F^2. \quad (6.18)$$

where  $\alpha_1 \geq \dots \geq \alpha_n$  and  $\beta_1 \geq \dots \geq \beta_n$  are the singular values of  $\mathbf{A}$  and  $\mathbf{B}$ , respectively.

## 6.6 Proof of the Hoffman-Wielandt Theorem for Singular Values

We apply the Hoffman-Wielandt Theorem for eigenvalues to the Hermitian matrices

$$\mathbf{C} := \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^* & \mathbf{0} \end{bmatrix} \text{ and } \mathbf{D} := \begin{bmatrix} \mathbf{0} & \mathbf{B} \\ \mathbf{B}^* & \mathbf{0} \end{bmatrix} \in \mathbb{C}^{m+n \times m+n}.$$

If  $\mathbf{C}$  and  $\mathbf{D}$  has eigenvalues  $\lambda_1 \geq \dots \geq \lambda_{m+n}$  and  $\mu_1 \geq \dots \geq \mu_{m+n}$ , respectively then

$$\sum_{j=1}^{m+n} |\lambda_j - \mu_j|^2 \leq \|\mathbf{C} - \mathbf{D}\|_F^2. \quad (6.19)$$

Suppose  $\mathbf{A}$  has rank  $r$  and singular value decomposition  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ . We use (6.12) and determine the eigenpairs of  $\mathbf{C}$  as follows.

$$\begin{aligned} \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^* & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i \\ \mathbf{v}_i \end{bmatrix} &= \begin{bmatrix} \mathbf{A}\mathbf{v}_i \\ \mathbf{A}^*\mathbf{u}_i \end{bmatrix} = \begin{bmatrix} \alpha_i\mathbf{u}_i \\ \alpha_i\mathbf{v}_i \end{bmatrix} = \alpha_i \begin{bmatrix} \mathbf{u}_i \\ \mathbf{v}_i \end{bmatrix}, \quad i = 1, \dots, r, \\ \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^* & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i \\ -\mathbf{v}_i \end{bmatrix} &= \begin{bmatrix} -\mathbf{A}\mathbf{v}_i \\ \mathbf{A}^*\mathbf{u}_i \end{bmatrix} = \begin{bmatrix} -\alpha_i\mathbf{u}_i \\ \alpha_i\mathbf{v}_i \end{bmatrix} = -\alpha_i \begin{bmatrix} \mathbf{u}_i \\ -\mathbf{v}_i \end{bmatrix}, \quad i = 1, \dots, r, \\ \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^* & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u}_i \\ \mathbf{0} \end{bmatrix} &= \begin{bmatrix} \mathbf{0} \\ \mathbf{A}^*\mathbf{u}_i \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix} = 0 \begin{bmatrix} \mathbf{u}_i \\ \mathbf{0} \end{bmatrix}, \quad i = r + 1, \dots, m, \\ \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{A}^* & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{0} \\ \mathbf{v}_i \end{bmatrix} &= \begin{bmatrix} \mathbf{A}\mathbf{v}_i \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix} = 0 \begin{bmatrix} \mathbf{0} \\ \mathbf{v}_i \end{bmatrix}, \quad i = r + 1, \dots, n. \end{aligned}$$

Thus  $\mathbf{C}$  has the  $2r$  eigenvalues  $\alpha_1, -\alpha_1, \dots, \alpha_r, -\alpha_r$  and  $m + n - 2r$  additional zero eigenvalues. Similarly, if  $\mathbf{B}$  has rank  $s$  then  $\mathbf{D}$  has the  $2s$  eigenvalues  $\beta_1, -\beta_1, \dots, \beta_s, -\beta_s$  and  $m + n - 2s$  additional zero eigenvalues. Let

$$t := \max(r, s).$$

Then

$$\begin{aligned} \lambda_1 \geq \dots \geq \lambda_{m+n} &= \alpha_1 \geq \dots \geq \alpha_t \geq 0 = \dots = 0 \geq -\alpha_t \geq \dots \geq -\alpha_1, \\ \mu_1 \geq \dots \geq \mu_{m+n} &= \beta_1 \geq \dots \geq \beta_t \geq 0 = \dots = 0 \geq -\beta_t \geq \dots \geq -\beta_1. \end{aligned}$$

We find  $\sum_{j=1}^{m+n} |\lambda_j - \mu_j|^2 = 2 \sum_{i=1}^t |\alpha_i - \beta_i|^2$  and

$$\|\mathbf{C} - \mathbf{D}\|_F^2 = \left\| \begin{bmatrix} \mathbf{0} & \mathbf{A} - \mathbf{B} \\ \mathbf{A}^* - \mathbf{B}^* & \mathbf{0} \end{bmatrix} \right\|_F^2 = \|\mathbf{B} - \mathbf{A}\|_F^2 + \|(\mathbf{B} - \mathbf{A})^*\|_F^2 = 2\|\mathbf{B} - \mathbf{A}\|_F^2.$$

But then (6.19) implies  $\sum_{i=1}^t |\alpha_i - \beta_i|^2 \leq \|\mathbf{B} - \mathbf{A}\|_F^2$ . Since  $t \leq n$  and  $\alpha_i = \beta_i = 0$  for  $i = t + 1, \dots, n$  we obtain (6.18).

## 6.7 Review Questions

**6.7.1** Consider an SVD and an SVF of a matrix  $\mathbf{A}$ .

- What are the singular values of  $\mathbf{A}$ ?
- how is the SVD defined?
- how can we find an SVF if we know an SVD?
- how can we find an SVD if we know an SVF?
- what are the relations between the singular vectors?
- which singular vectors form bases for  $\text{span}(\mathbf{A})$  and  $\ker(\mathbf{A}^*)$ ?

- 6.7.2** How are the Frobenius norm and singular values related?
- 6.7.3** State the Courant-Fischer theorem for singular values.
- 6.7.4** State the Hoffman-Wieland theorem for singular values.



## Chapter 7

# Matrix Norms

To measure the size of a matrix we can use a matrix norm. In this chapter we initiate a systematic study of matrix norms.

### 7.1 Matrix Norms

For simplicity we consider only matrix norms on the vector space  $(\mathbb{C}^{m \times n}, \mathbb{C})$ . All results also holds for  $(\mathbb{R}^{m \times n}, \mathbb{R})$ .

#### Definition 7.1 (Matrix norms)

Suppose  $m, n$  are positive integers. A function  $\|\cdot\|: \mathbb{C}^{m \times n} \rightarrow \mathbb{R}$  is called a **matrix norm** on  $\mathbb{C}^{m \times n}$  if for all  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{m \times n}$  and all  $c \in \mathbb{C}$

1.  $\|\mathbf{A}\| \geq 0$  with equality if and only if  $\mathbf{A} = \mathbf{0}$ . (positivity)
2.  $\|c\mathbf{A}\| = |c| \|\mathbf{A}\|$ . (homogeneity)
3.  $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$ . (subadditivity)

A matrix norm is simply a vector norm on the finite dimensional vector space  $(\mathbb{C}^{m \times n}, \mathbb{C})$  of  $m \times n$  matrices. Adapting Theorem 0.19 to this situation gives

#### Theorem 7.2 (Matrix norm equivalence)

All matrix norms are equivalent. Thus, if  $\|\cdot\|$  and  $\|\cdot\|'$  are two matrix norms on  $\mathbb{C}^{m \times n}$  then there are positive constants  $\mu$  and  $M$  such that

$$\mu \|\mathbf{A}\| \leq \|\mathbf{A}\|' \leq M \|\mathbf{A}\|$$

holds for all  $\mathbf{A} \in \mathbb{C}^{m \times n}$ . Moreover, a matrix norm is a continuous function.

Any vector norm  $\|\cdot\|_V$  on  $\mathbb{C}^{mn}$  defines a matrix norm on  $\mathbb{C}^{m \times n}$  given by  $\|\mathbf{A}\| := \|\text{vec}(\mathbf{A})\|_V$ , where  $\text{vec}(\mathbf{A}) \in \mathbb{C}^{mn}$  is the vector obtained by stacking the columns of  $\mathbf{A}$  on top of each other. In particular, to the  $p$  vector norms for  $p = 1, 2, \infty$ , we have the corresponding **sum norm**, **Frobenius norm**, and **max norm** defined by

$$\|\mathbf{A}\|_S := \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|, \quad \|\mathbf{A}\|_F := \left( \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{1/2}, \quad \|\mathbf{A}\|_M := \max_{i,j} |a_{ij}|. \quad (7.1)$$



Ferdinand Georg Frobenius, (1849-1917).

Of these norms the Frobenius norm is the most useful. Some of its properties were derived in Lemma 6.24 and Theorem 6.25.

### 7.1.1 Consistent and subordinate matrix norms

Since matrices can be multiplied it is useful to have an analogue of subadditivity for matrix multiplication. For square matrices the product  $\mathbf{AB}$  is defined in a fixed space  $\mathbb{C}^{n \times n}$ , while in the rectangular case matrix multiplication combines matrices in different spaces. The following definition captures this distinction.

#### Definition 7.3 (Consistent matrix norms)

A matrix norm is called **consistent on**  $\mathbb{C}^{n \times n}$  if

$$4. \quad \|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\| \quad (\text{submultiplicativity})$$

holds for all  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$ . A matrix norm is **consistent** if it is defined on  $\mathbb{C}^{m \times n}$  for all  $m, n \in \mathbb{N}$ , and 4. holds for all matrices  $\mathbf{A}, \mathbf{B}$  for which the product  $\mathbf{AB}$  is defined.

Clearly the three norms in (7.1) are defined for all  $m, n \in \mathbb{N}$ . From Lemma 6.24 it follows that the Frobenius norm is consistent.

**Exercise 7.4 (Consistency of sum norm?)**

Show that the sum norm is consistent.

**Exercise 7.5 (Consistency of max norm?)**

Show that the max norm is not consistent by considering  $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ .

**Exercise 7.6 (Consistency of modified max norm)**

(a) Show that the norm

$$\|\mathbf{A}\| := \sqrt{mn}\|\mathbf{A}\|_M, \quad \mathbf{A} \in \mathbb{C}^{m \times n}$$

is a consistent matrix norm.

(b) Show that the constant  $\sqrt{mn}$  can be replaced by  $m$  and by  $n$ .

For a consistent matrix norm on  $\mathbb{C}^{n \times n}$  we have the inequality

$$\|\mathbf{A}^k\| \leq \|\mathbf{A}\|^k \text{ for } k \in \mathbb{N}. \quad (7.2)$$

When working with norms one often has to bound the vector norm of a matrix times a vector by the norm of the matrix times the norm of the vector. This leads to the following definition.

**Definition 7.7 (Subordinate matrix norms)**

Suppose  $m, n \in \mathbb{N}$  are given, let  $\|\cdot\|$  on  $\mathbb{C}^m$  and  $\|\cdot\|_\beta$  on  $\mathbb{C}^n$  be vector norms, and let  $\|\cdot\|$  be a matrix norm on  $\mathbb{C}^{m \times n}$ . We say that the matrix norm  $\|\cdot\|$  is **subordinate** to the vector norms  $\|\cdot\|$  and  $\|\cdot\|_\beta$  if  $\|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|_\beta$  for all  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and all  $\mathbf{x} \in \mathbb{C}^n$ . If  $\|\cdot\| = \|\cdot\|_\beta$  then we say that  $\|\cdot\|$  is subordinate to  $\|\cdot\|_\beta$ .

By Lemma 6.24 we have  $\|\mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{A}\|_F \|\mathbf{x}\|_2$ , for all  $\mathbf{x} \in \mathbb{C}^n$ . Thus the Frobenius norm is subordinate to the Euclidian vector norm.

**Exercise 7.8 (What is the sum norm subordinate to?)**

Show that the sum norm is subordinate to the  $l_1$ -norm.

**Exercise 7.9 (What is the max norm subordinate to?)**

- (a) Show that the max norm is subordinate to the  $\infty$  and 1 norm, i. e.,  $\|\mathbf{A}\mathbf{x}\|_\infty \leq \|\mathbf{A}\|_M \|\mathbf{x}\|_1$  holds for all  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and all  $\mathbf{x} \in \mathbb{C}^n$ .
- (b) Show that if  $\|\mathbf{A}\|_M = |a_{kl}|$ , then  $\|\mathbf{A}\mathbf{e}_l\|_\infty = \|\mathbf{A}\|_M \|\mathbf{e}_l\|_1$ .
- (c) Show that  $\|\mathbf{A}\|_M = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|_\infty}{\|\mathbf{x}\|_1}$ .

## 7.1.2 Operator norms

Corresponding to vector norms on  $\mathbb{C}^n$  and  $\mathbb{C}^m$  there is an induced matrix norm on  $\mathbb{C}^{m \times n}$  which we call the **operator norm**. It is possible to consider one vector norm on  $\mathbb{C}^m$  and another vector norm on  $\mathbb{C}^n$ , but we treat only the case of one vector norm defined on  $\mathbb{C}^n$  for all  $n \in \mathbb{N}$ <sup>9</sup>.

### Definition 7.10 (Operator norm)

Let  $\|\cdot\|$  be a vector norm defined on  $\mathbb{C}^n$  for all  $n \in \mathbb{N}$ . For given  $m, n \in \mathbb{N}$  and  $\mathbf{A} \in \mathbb{C}^{m \times n}$  we define

$$\|\mathbf{A}\| := \max_{\mathbf{x} \neq 0} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|}. \quad (7.3)$$

We call this the **operator norm** corresponding to the vector norm  $\|\cdot\|$ .

With a risk of confusion we use the same symbol for the operator norm and the corresponding vector norm. Before we show that the operator norm is a matrix norm we make some observations.

1. It is enough to take the max over subsets of  $\mathbb{C}^n$ . For example

$$\|\mathbf{A}\| = \max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|. \quad (7.4)$$

The set

$$\mathcal{S} := \{x \in \mathbb{C}^n : \|x\| = 1\} \quad (7.5)$$

is the unit sphere in  $\mathbb{C}^n$  with respect to the vector norm  $\|\cdot\|$ . It is enough to take the max over this unit sphere since

$$\max_{\mathbf{x} \neq 0} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\mathbf{x} \neq 0} \left\| \mathbf{A} \left( \frac{\mathbf{x}}{\|\mathbf{x}\|} \right) \right\| = \max_{\|\mathbf{y}\|=1} \|\mathbf{A}\mathbf{y}\|.$$

2. The operator norm is subordinate to the corresponding vector norm. Thus,

$$\|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{x}\| \text{ for all } \mathbf{A} \in \mathbb{C}^{m \times n} \text{ and } \mathbf{x} \in \mathbb{C}^n. \quad (7.6)$$

3. We can use max instead of sup in (7.3). This follows by the following compactness argument. The unit sphere  $\mathcal{S}$  given by (7.5) is bounded. It is also finite dimensional and closed, and hence compact. Moreover, since the vector norm  $\|\cdot\| : \mathcal{S} \rightarrow \mathbb{R}$  is a continuous function, it follows that the function  $f : \mathcal{S} \rightarrow \mathbb{R}$  given by  $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x}\|$  is continuous. But then  $f$  attains its max and min and we have

$$\|\mathbf{A}\| = \|\mathbf{A}\mathbf{x}^*\| \text{ for some } \mathbf{x}^* \in \mathcal{S}. \quad (7.7)$$

---

<sup>9</sup>In the case of one vector norm  $\|\cdot\|$  on  $\mathbb{C}^m$  and another vector norm  $\|\cdot\|_\beta$  on  $\mathbb{C}^n$  we would define  $\|\mathbf{A}\| := \max_{\mathbf{x} \neq 0} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|_\beta}$ .

**Lemma 7.11 (The operator norm is a matrix norm)**

For any vector norm the operator norm given by (7.3) is a consistent matrix norm. Moreover,  $\|\mathbf{I}\| = 1$ .

*Proof.* We use (7.4). In 2. and 3. below we take the max over the unit sphere  $\mathcal{S}$  given by (7.5).

1. Nonnegativity is obvious. If  $\|\mathbf{A}\| = 0$  then  $\|\mathbf{A}\mathbf{y}\| = 0$  for each  $\mathbf{y} \in \mathbb{C}^n$ . In particular, each column  $\mathbf{A}\mathbf{e}_j$  in  $\mathbf{A}$  is zero. Hence  $\mathbf{A} = \mathbf{0}$ .
2.  $\|c\mathbf{A}\| = \max_{\mathbf{x}} \|c\mathbf{A}\mathbf{x}\| = \max_{\mathbf{x}} |c| \|\mathbf{A}\mathbf{x}\| = |c| \|\mathbf{A}\|$ .
3.  $\|\mathbf{A} + \mathbf{B}\| = \max_{\mathbf{x}} \|(\mathbf{A} + \mathbf{B})\mathbf{x}\| \leq \max_{\mathbf{x}} \|\mathbf{A}\mathbf{x}\| + \max_{\mathbf{x}} \|\mathbf{B}\mathbf{x}\| = \|\mathbf{A}\| + \|\mathbf{B}\|$ .
4.  $\|\mathbf{AB}\| = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{AB}\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\mathbf{B}\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{AB}\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\mathbf{B}\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{AB}\mathbf{x}\|}{\|\mathbf{B}\mathbf{x}\|} \frac{\|\mathbf{B}\mathbf{x}\|}{\|\mathbf{x}\|}$   
 $\leq \max_{\mathbf{y} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{y}\|}{\|\mathbf{y}\|} \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{B}\mathbf{x}\|}{\|\mathbf{x}\|} = \|\mathbf{A}\| \|\mathbf{B}\|$ .

That  $\|\mathbf{I}\| = 1$  for any operator norm follows immediately from the definition.  $\square$

Since  $\|\mathbf{I}\|_F = \sqrt{n}$ , we see that the Frobenius norm is not an operator norm for  $n > 1$ .

**7.1.3 The operator  $p$ -norms**

Recall that the  $p$  or  $\ell_p$  vector norms (10) are given by

$$\|\mathbf{x}\|_p := \left( \sum_{j=1}^n |x_j|^p \right)^{1/p}, \quad p \geq 1, \quad \|\mathbf{x}\|_\infty := \max_{1 \leq j \leq n} |x_j|.$$

The operator norms  $\|\cdot\|_p$  defined from these  $p$ -vector norms are used quite frequently for  $p = 1, 2, \infty$ . We define for any  $1 \leq p \leq \infty$

$$\|\mathbf{A}\|_p := \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|_p}{\|\mathbf{x}\|_p} = \max_{\|\mathbf{y}\|_p=1} \|\mathbf{A}\mathbf{y}\|_p. \quad (7.8)$$

For  $p = 1, 2, \infty$  we have explicit expressions for these norms.

**Theorem 7.12 (on two norms)**

For  $\mathbf{A} \in \mathbb{C}^{m \times n}$  we have

$$\begin{aligned} \|\mathbf{A}\|_1 &:= \max_{1 \leq j \leq n} \|\mathbf{A}\mathbf{e}_j\|_1 = \max_{1 \leq j \leq n} \sum_{k=1}^m |a_{k,j}|, && (\text{max column sum}) \\ \|\mathbf{A}\|_2 &:= \sigma_1, && (\text{largest singular value of } \mathbf{A}) \\ \|\mathbf{A}\|_\infty &= \max_{1 \leq k \leq m} \|\mathbf{e}_k^T \mathbf{A}\|_1 = \max_{1 \leq k \leq m} \sum_{j=1}^n |a_{k,j}|. && (\text{max row sum}) \end{aligned} \quad (7.9)$$

The **two-norm**  $\|\mathbf{A}\|_2$  is also called the **spectral norm** of  $\mathbf{A}$ .

**Proof.** The result for  $p = 2$  follows from the minmax theorem for singular values. Indeed, by (6.16) we have  $\sigma_1 = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{x}\|_2}$ . For  $p = 1, \infty$  we proceed as follows:

(a) We derive a constant  $K_p$  such that  $\|\mathbf{Ax}\|_p \leq K_p$  for any  $\mathbf{x} \in \mathbb{C}^n$  with  $\|\mathbf{x}\|_p = 1$ .

(b) We give an extremal vector  $\mathbf{y}^* \in \mathbb{C}^n$  with  $\|\mathbf{y}^*\|_p = 1$  so that  $\|\mathbf{Ay}^*\|_p = K_p$ .

It then follows from (7.8) that  $\|\mathbf{A}\|_p = \|\mathbf{Ay}^*\|_p = K_p$ .

**1-norm:** Define  $K_1$ ,  $c$  and  $\mathbf{y}^*$  by  $K_1 := \|\mathbf{Ae}_c\|_1 = \max_{1 \leq j \leq n} \|\mathbf{Ae}_j\|_1$  and  $\mathbf{y}^* := \mathbf{e}_c$ , a unit vector. Then  $\|\mathbf{y}^*\|_1 = 1$  and we obtain

(a)

$$\|\mathbf{Ax}\|_1 = \sum_{k=1}^m \left| \sum_{j=1}^n a_{kj} x_j \right| \leq \sum_{k=1}^m \sum_{j=1}^n |a_{kj}| |x_j| = \sum_{j=1}^n \left( \sum_{k=1}^m |a_{kj}| \right) |x_j| \leq K_1.$$

(b)  $\|\mathbf{Ay}^*\|_1 = K_1$ .

**$\infty$ -norm:** Define  $K_\infty$ ,  $r$  and  $\mathbf{y}^*$  by  $K_\infty := \|\mathbf{e}_r^T \mathbf{A}\|_1 = \max_{1 \leq k \leq m} \|\mathbf{e}_k^T \mathbf{A}\|_1$  and  $\mathbf{y}^* := [e^{-i\theta_1}, \dots, e^{-i\theta_n}]^T$ , where  $a_{rj} = |a_{rj}| e^{i\theta_j}$  for  $j = 1, \dots, n$ .

(a)  $\|\mathbf{Ax}\|_\infty = \max_{1 \leq k \leq m} \left| \sum_{j=1}^n a_{kj} x_j \right| \leq \max_{1 \leq k \leq m} \sum_{j=1}^n |a_{kj}| |x_j| \leq K_\infty$ .

(b)  $\|\mathbf{Ay}^*\|_\infty = \max_{1 \leq k \leq m} \left| \sum_{j=1}^n a_{kj} e^{-i\theta_j} \right| = K_\infty$ .

The last equality is correct because  $\left| \sum_{j=1}^n a_{kj} e^{-i\theta_j} \right| \leq \sum_{j=1}^n |a_{kj}| \leq K_\infty$  with equality for  $k = r$ .

□

### Example 7.13 (Compare onetwoinfnorms)

The largest singular value of the matrix  $\mathbf{A} := \frac{1}{15} \begin{bmatrix} 14 & 4 & 16 \\ 2 & 22 & 13 \end{bmatrix}$ , is  $\sigma_1 = 2$  (cf. Example 6.12). We find

$$\|\mathbf{A}\|_1 = \frac{29}{15}, \quad \|\mathbf{A}\|_2 = 2, \quad \|\mathbf{A}\|_\infty = \frac{37}{15}, \quad \|\mathbf{A}\|_F = \sqrt{5}.$$

We observe that the values of these norms do not differ by much.

In some cases the spectral norm is equal to an eigenvalue of the matrix.

**Theorem 7.14 (Spectral norm)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  has singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$  and eigenvalues  $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$ . Then

$$\|\mathbf{A}\|_2 = \sigma_1 \text{ and } \|\mathbf{A}^{-1}\|_2 = \frac{1}{\sigma_n}, \quad (7.10)$$

$$\|\mathbf{A}\|_2 = \lambda_1 \text{ and } \|\mathbf{A}^{-1}\|_2 = \frac{1}{\lambda_n}, \quad \text{if } \mathbf{A} \text{ is symmetric positive definite,} \quad (7.11)$$

$$\|\mathbf{A}\|_2 = |\lambda_1| \text{ and } \|\mathbf{A}^{-1}\|_2 = \frac{1}{|\lambda_n|}, \quad \text{if } \mathbf{A} \text{ is normal.} \quad (7.12)$$

For the norms of  $\mathbf{A}^{-1}$  we assume of course that  $\mathbf{A}$  is nonsingular.

**Proof.** Since  $1/\sigma_n$  is the largest singular value of  $\mathbf{A}^{-1}$ , (7.10) follows. By Theorem 6.9 the singular values of a symmetric positive definite matrix (normal matrix) are equal to the eigenvalues (absolute value of the eigenvalues). This implies (7.11) and (7.12).  $\square$

The following result is sometimes useful.

**Theorem 7.15 (Spectral norm bound)**

For any  $\mathbf{A} \in \mathbb{C}^{m \times n}$  we have  $\|\mathbf{A}\|_2^2 \leq \|\mathbf{A}\|_1 \|\mathbf{A}\|_\infty$ .

**Proof.** Let  $(\sigma^2, \mathbf{v})$  be an eigenpair for  $\mathbf{A}^* \mathbf{A}$  corresponding to the largest singular value  $\sigma$  of  $\mathbf{A}$ . Then

$$\|\mathbf{A}\|_2^2 \|\mathbf{v}\|_1 = \sigma^2 \|\mathbf{v}\|_1 = \|\sigma^2 \mathbf{v}\|_1 = \|\mathbf{A}^* \mathbf{A} \mathbf{v}\|_1 \leq \|\mathbf{A}^*\|_1 \|\mathbf{A}\|_1 \|\mathbf{v}\|_1.$$

Observing that  $\|\mathbf{A}^*\|_1 = \|\mathbf{A}\|_\infty$  by Theorem 7.12 and canceling  $\|\mathbf{v}\|_1$  proves the result.  $\square$

**Exercise 7.16 (Spectral norm)**

Let  $m, n \in \mathbb{N}$  and  $\mathbf{A} \in \mathbb{C}^{m \times n}$ . Show that

$$\|\mathbf{A}\|_2 = \max_{\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 = 1} |\mathbf{y}^* \mathbf{A} \mathbf{x}|.$$

**Exercise 7.17 (Spectral norm of the inverse)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is nonsingular. Use (7.10) and (6.16) to show that

$$\|\mathbf{A}^{-1}\|_2 = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{x}\|_2}{\|\mathbf{A} \mathbf{x}\|_2}.$$

**Exercise 7.18 ( $p$ -norm example)**

Let

$$\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}.$$

Compute  $\|\mathbf{A}\|_p$  and  $\|\mathbf{A}^{-1}\|_p$  for  $p = 1, 2, \infty$ .

**7.1.4 Unitary invariant matrix norms****Definition 7.19 (Unitary invariant norm)**

A matrix norm  $\|\cdot\|$  on  $\mathbb{C}^{m \times n}$  is called **unitary invariant** if  $\|\mathbf{UAV}\| = \|\mathbf{A}\|$  for any  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and any unitary matrices  $\mathbf{U} \in \mathbb{C}^{m \times m}$  and  $\mathbf{V} \in \mathbb{C}^{n \times n}$ .

When an unitary invariant matrix norm is used, the size of a perturbation is not increased by a unitary transformation. Thus if  $\mathbf{U}$  and  $\mathbf{V}$  are unitary then  $\mathbf{U}(\mathbf{A} + \mathbf{E})\mathbf{V} = \mathbf{UAV} + \mathbf{F}$ , where  $\|\mathbf{F}\| = \|\mathbf{E}\|$ .

It follows from Lemma 6.24 that the Frobenius norm is unitary invariant. We show here that this also holds for the spectral norm. It can be shown that the spectral norm is the only unitary invariant operator norm, see [13] p. 308.

**Theorem 7.20 (Unitary invariant norms)**

The Frobenius norm and the spectral norm are unitary invariant. Moreover  $\|\mathbf{A}^*\|_F = \|\mathbf{A}\|_F$  and  $\|\mathbf{A}^*\|_2 = \|\mathbf{A}\|_2$ .

**Proof.** The results for the Frobenius norm follow from Lemma 6.24. Suppose  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and let  $\mathbf{U} \in \mathbb{C}^{m \times m}$  and  $\mathbf{V} \in \mathbb{C}^{n \times n}$  be unitary. Since the 2-vector norm is unitary invariant we obtain

$$\|\mathbf{UA}\|_2 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{UAx}\|_2 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{Ax}\|_2 = \|\mathbf{A}\|_2.$$

Now  $\mathbf{A}$  and  $\mathbf{A}^*$  have the same nonzero singular values, and it follows from Theorem 7.12 that  $\|\mathbf{A}^*\|_2 = \|\mathbf{A}\|_2$ . Moreover  $\mathbf{V}^*$  is unitary. Using these facts we find

$$\|\mathbf{AV}\|_2 = \|(\mathbf{AV})^*\|_2 = \|\mathbf{V}^*\mathbf{A}^*\|_2 = \|\mathbf{A}^*\|_2 = \|\mathbf{A}\|_2.$$

□

**Exercise 7.21 (Unitary invariance of the spectral norm)**

Show that  $\|\mathbf{VA}\|_2 = \|\mathbf{A}\|_2$  holds even for a rectangular  $\mathbf{V}$  as long as  $\mathbf{V}^*\mathbf{V} = \mathbf{I}$ .

**Exercise 7.22 ( $\|\mathbf{AU}\|_2$  rectangular  $\mathbf{A}$ )**

Find  $\mathbf{A} \in \mathbb{R}^{2 \times 2}$  and  $\mathbf{U} \in \mathbb{R}^{2 \times 1}$  with  $\mathbf{U}^T\mathbf{U} = \mathbf{I}$  such that  $\|\mathbf{AU}\|_2 < \|\mathbf{A}\|_2$ . Thus, in general,  $\|\mathbf{AU}\|_2 = \|\mathbf{A}\|_2$  does not hold for a rectangular  $\mathbf{U}$  even if  $\mathbf{U}^*\mathbf{U} = \mathbf{I}$ .



**Exercise 7.23 ( $p$ -norm of diagonal matrix)**

Show that  $\|\mathbf{A}\|_p = \rho(\mathbf{A}) := \max |\lambda_i|$  (the largest eigenvalue of  $\mathbf{A}$ ),  $1 \leq p \leq \infty$ , when  $\mathbf{A}$  is a diagonal matrix.

**Exercise 7.24 (spectral norm of a column vector)**

A vector  $\mathbf{a} \in \mathbb{C}^m$  can also be considered as a matrix  $\mathbf{A} \in \mathbb{C}^{m,1}$ .

(a) Show that the spectral matrix norm (2-norm) of  $\mathbf{A}$  equals the Euclidean vector norm of  $\mathbf{a}$ .

(b) Show that  $\|\mathbf{A}\|_p = \|\mathbf{a}\|_p$  for  $1 \leq p \leq \infty$ .

**7.1.5 Absolute and monotone norms**

A vector norm on  $\mathbb{C}^n$  is an **absolute norm** if  $\|\mathbf{x}\| = \|\mathbf{x}\|$  for all  $\mathbf{x} \in \mathbb{C}^n$ . Here  $|\mathbf{x}| := [|x_1|, \dots, |x_n|]^T$ , the absolute values of the components of  $\mathbf{x}$ . Clearly the vector  $p$  norms are absolute norms. We state without proof (see Theorem 5.5.10 of [13]) that a vector norm on  $\mathbb{C}^n$  is an absolute norm if and only if it is a **monotone norm**, i. e.,

$$|x_i| \leq |y_i|, i = 1, \dots, n \implies \|\mathbf{x}\| \leq \|\mathbf{y}\|, \text{ for all } \mathbf{x}, \mathbf{y} \in \mathbb{C}^n.$$

Absolute and monotone matrix norms are defined as for vector norms.

**Exercise 7.25 (Norm of absolute value matrix)**

If  $\mathbf{A} \in \mathbb{C}^{m \times n}$  has elements  $a_{ij}$ , let  $|\mathbf{A}| \in \mathbb{R}^{m \times n}$  be the matrix with elements  $|a_{ij}|$ .

(a) Compute  $|\mathbf{A}|$  if  $\mathbf{A} = \begin{bmatrix} 1+i & -2 \\ 1 & 1-i \end{bmatrix}$ ,  $i = \sqrt{-1}$ .

(b) Show that for any  $\mathbf{A} \in \mathbb{C}^{m \times n}$   $\|\mathbf{A}\|_F = \|\mathbf{A}\|_F$ ,  $\|\mathbf{A}\|_p = \|\mathbf{A}\|_p$  for  $p = 1, \infty$ .

(c) Show that for any  $\mathbf{A} \in \mathbb{C}^{m \times n}$   $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_2$ .

(d) Find a real symmetric  $2 \times 2$  matrix  $\mathbf{A}$  such that  $\|\mathbf{A}\|_2 < \|\mathbf{A}\|_2$ .

The study of matrix norms will be continued in Chapter 8.

**7.2 The Condition Number with Respect to Inversion**

Consider the system of two linear equations

$$\begin{aligned} x_1 + x_2 &= 20 \\ x_1 + (1 - 10^{-16})x_2 &= 20 - 10^{-15} \end{aligned}$$

whose exact solution is  $x_1 = x_2 = 10$ . If we replace the second equation by

$$x_1 + (1 + 10^{-16})x_2 = 20 - 10^{-15},$$

the exact solution changes to  $x_1 = 30$ ,  $x_2 = -10$ . Here a small change in one of the coefficients, from  $1 - 10^{-16}$  to  $1 + 10^{-16}$ , changed the exact solution by a large amount.

A mathematical problem in which the solution is very sensitive to changes in the data is called **ill-conditioned**. Such problems can be difficult to solve on a computer.

In this section we consider what effect a small change (perturbation) in the data  $\mathbf{A}, \mathbf{b}$  has on the solution  $\mathbf{x}$  of a linear system  $\mathbf{Ax} = \mathbf{b}$ . Suppose  $\mathbf{y}$  solves  $(\mathbf{A} + \mathbf{E})\mathbf{y} = \mathbf{b} + \mathbf{e}$  where  $\mathbf{E}$  is a (small)  $n \times n$  matrix and  $\mathbf{e}$  a (small) vector. How large can  $\mathbf{y} - \mathbf{x}$  be? To measure this we use vector and matrix norms. In this section  $\|\cdot\|$  will denote a vector norm on  $\mathbb{C}^n$  and also a matrix norm on  $\mathbb{C}^{n \times n}$  which for any  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$  and any  $\mathbf{x} \in \mathbb{C}^n$  satisfy

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\| \text{ and } \|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|.$$

This holds if the matrix norm is the operator norm corresponding to the given vector norm, but is also satisfied for the Frobenius matrix norm and the Euclidian vector norm. This follows from Lemma 6.24.

Suppose  $\mathbf{x}$  and  $\mathbf{y}$  are vectors in  $\mathbb{C}^n$  that we want to compare. The difference  $\|\mathbf{y} - \mathbf{x}\|$  measures the **absolute error** in  $\mathbf{y}$  as an approximation to  $\mathbf{x}$ , while  $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{x}\|$  and  $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{y}\|$  are measures for the **relative error**.

We consider first a perturbation in the right-hand side  $\mathbf{b}$ .

**Theorem 7.26 (Perturbation in the right-hand side)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is nonsingular,  $\mathbf{b}, \mathbf{e} \in \mathbb{C}^n$ ,  $\mathbf{b} \neq \mathbf{0}$  and  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{Ay} = \mathbf{b} + \mathbf{e}$ . Then

$$\frac{1}{K(\mathbf{A})} \frac{\|\mathbf{e}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(\mathbf{A}) \frac{\|\mathbf{e}\|}{\|\mathbf{b}\|}, \quad K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|. \quad (7.13)$$

*Proof.* Subtracting  $\mathbf{Ax} = \mathbf{b}$  from  $\mathbf{Ay} = \mathbf{b} + \mathbf{e}$  we have  $\mathbf{A}(\mathbf{y} - \mathbf{x}) = \mathbf{e}$  or  $\mathbf{y} - \mathbf{x} = \mathbf{A}^{-1}\mathbf{e}$ . Combining  $\|\mathbf{y} - \mathbf{x}\| = \|\mathbf{A}^{-1}\mathbf{e}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{e}\|$  and  $\|\mathbf{b}\| = \|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$  we obtain the upper bound in (7.13). Combining  $\|\mathbf{e}\| \leq \|\mathbf{A}\| \|\mathbf{y} - \mathbf{x}\|$  and  $\|\mathbf{x}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{b}\|$  we obtain the lower bound.  $\square$

Consider (7.13).  $\|\mathbf{e}\|/\|\mathbf{b}\|$  is a measure of the size of the perturbation  $\mathbf{e}$  relative to the size of  $\mathbf{b}$ . The upper bound says that  $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{x}\|$  in the worst case can be  $K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$  times as large as  $\|\mathbf{e}\|/\|\mathbf{b}\|$ .  $K(\mathbf{A})$  is called the **condition number with respect to inversion of a matrix**, or just the condition number, if it is clear from the context that we are talking about solving

linear systems or inverting a matrix. The condition number depends on the matrix  $\mathbf{A}$  and on the norm used. If  $K(\mathbf{A})$  is large,  $\mathbf{A}$  is called **ill-conditioned** (with respect to inversion). If  $K(\mathbf{A})$  is small,  $\mathbf{A}$  is called **well-conditioned** (with respect to inversion). We always have  $K(\mathbf{A}) \geq 1$ . For since  $\|\mathbf{x}\| = \|\mathbf{I}\mathbf{x}\| \leq \|\mathbf{I}\|\|\mathbf{x}\|$  for any  $\mathbf{x}$  we have  $\|\mathbf{I}\| \geq 1$  and therefore  $\|\mathbf{A}\|\|\mathbf{A}^{-1}\| \geq \|\mathbf{A}\mathbf{A}^{-1}\| = \|\mathbf{I}\| \geq 1$ .

Since all matrix norms are equivalent, the dependence of  $K(\mathbf{A})$  on the norm chosen is less important than the dependence on  $\mathbf{A}$ . Sometimes one chooses the spectral norm when discussing properties of the condition number, and the  $\ell_1$ ,  $\ell_\infty$ , or Frobenius norm when one wishes to compute it or estimate it.

The following explicit expressions for the 2-norm condition number follow from Theorem 7.14.

**Theorem 7.27 (Spectral condition number)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is nonsingular with singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0$  and eigenvalues  $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n| > 0$ . Then  $K_2(\mathbf{A}) := \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2 = \sigma_1/\sigma_n$ . Moreover,

$$K_2(\mathbf{A}) = \begin{cases} \lambda_1/\lambda_n, & \text{if } \mathbf{A} \text{ is symmetric positive definite,} \\ |\lambda_1|/|\lambda_n|, & \text{if } \mathbf{A} \text{ is normal.} \end{cases} \quad (7.14)$$

It follows that  $\mathbf{A}$  is ill-conditioned with respect to inversion if and only if  $\sigma_1/\sigma_n$  is large, or  $\lambda_1/\lambda_n$  is large when  $\mathbf{A}$  is symmetric positive definite.

Suppose we have computed an approximate solution  $\mathbf{y}$  to  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . The vector  $\mathbf{r}(\mathbf{y}) := \mathbf{A}\mathbf{y} - \mathbf{b}$  is called the **residual vector**, or just the residual. We can bound  $\mathbf{x} - \mathbf{y}$  in term of  $\mathbf{r}$ .

**Theorem 7.28 (Perturbation and residual)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$ ,  $\mathbf{b} \in \mathbb{C}^n$ ,  $\mathbf{A}$  is nonsingular and  $\mathbf{b} \neq \mathbf{0}$ . Let  $\mathbf{r}(\mathbf{y}) = \mathbf{A}\mathbf{y} - \mathbf{b}$  for any  $\mathbf{y} \in \mathbb{C}^n$ . If  $\mathbf{A}\mathbf{x} = \mathbf{b}$  then

$$\frac{1}{K(\mathbf{A})} \frac{\|\mathbf{r}(\mathbf{y})\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(\mathbf{A}) \frac{\|\mathbf{r}(\mathbf{y})\|}{\|\mathbf{b}\|}. \quad (7.15)$$

*Proof.* We simply take  $\mathbf{e} = \mathbf{r}(\mathbf{y})$  in Theorem 7.26.  $\square$

If  $\mathbf{A}$  is well-conditioned, (7.15) says that  $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{x}\| \approx \|\mathbf{r}(\mathbf{y})\|/\|\mathbf{b}\|$ . In other words, the accuracy in  $\mathbf{y}$  is about the same order of magnitude as the residual as long as  $\|\mathbf{b}\| \approx 1$ . If  $\mathbf{A}$  is ill-conditioned, anything can happen. We can for example have an accurate solution even if the residual is large.

Consider next the effect of a perturbation in the coefficient matrix. Suppose  $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n \times n}$  with  $\mathbf{A}$  nonsingular. We like to compare the solution  $\mathbf{x}$  and  $\mathbf{y}$  of the systems  $\mathbf{A}\mathbf{x} = \mathbf{b}$  and  $(\mathbf{A} + \mathbf{E})\mathbf{y} = \mathbf{b}$ . We expect  $\mathbf{A} + \mathbf{E}$  to be nonsingular if the elements of  $\mathbf{E}$  are sufficiently small and we need to address this question. Consider first the case where  $\mathbf{A} = \mathbf{I}$ .

**Theorem 7.29 (Nonsingularity of perturbation of identity)**

Suppose  $\mathbf{B} \in \mathbb{C}^{n \times n}$  and  $\|\mathbf{B}\| < 1$  for some consistent matrix norm on  $\mathbb{C}^{n \times n}$ . Then  $\mathbf{I} - \mathbf{B}$  is nonsingular and

$$\frac{1}{1 + \|\mathbf{B}\|} \leq \|(\mathbf{I} - \mathbf{B})^{-1}\| \leq \frac{1}{1 - \|\mathbf{B}\|}. \quad (7.16)$$

*Proof.* Suppose  $\mathbf{I} - \mathbf{B}$  is singular. Then  $(\mathbf{I} - \mathbf{B})\mathbf{x} = \mathbf{0}$  for some nonzero  $\mathbf{x} \in \mathbb{C}^n$ , and  $\mathbf{x} = \mathbf{B}\mathbf{x}$  so that  $\|\mathbf{x}\| = \|\mathbf{B}\mathbf{x}\| \leq \|\mathbf{B}\|\|\mathbf{x}\|$ . But then  $\|\mathbf{B}\| \geq 1$ . It follows that  $\mathbf{I} - \mathbf{B}$  is nonsingular if  $\|\mathbf{B}\| < 1$ . Next, since

$$\begin{aligned} \|\mathbf{I}\| &= \|(\mathbf{I} - \mathbf{B})(\mathbf{I} - \mathbf{B})^{-1}\| \leq \|\mathbf{I} - \mathbf{B}\| \|(\mathbf{I} - \mathbf{B})^{-1}\| \\ &\leq (\|\mathbf{I}\| + \|\mathbf{B}\|) \|(\mathbf{I} - \mathbf{B})^{-1}\|, \end{aligned}$$

and since  $\|\mathbf{I}\| \geq 1$ , we obtain the lower bound in (7.16):

$$\frac{1}{1 + \|\mathbf{B}\|} \leq \frac{\|\mathbf{I}\|}{\|\mathbf{I}\| + \|\mathbf{B}\|} \leq \|(\mathbf{I} - \mathbf{B})^{-1}\|. \quad (7.17)$$

Taking norms and using the inverse triangle inequality in

$$\mathbf{I} = (\mathbf{I} - \mathbf{B})(\mathbf{I} - \mathbf{B})^{-1} = (\mathbf{I} - \mathbf{B})^{-1} - \mathbf{B}(\mathbf{I} - \mathbf{B})^{-1}$$

implies

$$\|\mathbf{I}\| \geq \|(\mathbf{I} - \mathbf{B})^{-1}\| - \|\mathbf{B}(\mathbf{I} - \mathbf{B})^{-1}\| \geq (1 - \|\mathbf{B}\|) \|(\mathbf{I} - \mathbf{B})^{-1}\|.$$

If the matrix norm is an operator norm then  $\|\mathbf{I}\| = 1$  and the upper bound follows. We show in Section 8.4 that the upper bound also holds for the Frobenius norm, and more generally for any consistent matrix norm on  $\mathbb{C}^{n \times n}$ .  $\square$

**Theorem 7.30 (Nonsingularity of perturbation)**

Suppose  $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n \times n}$ ,  $\mathbf{b} \in \mathbb{C}^n$  with  $\mathbf{A}$  invertible and  $\mathbf{b} \neq \mathbf{0}$ . If  $r := \|\mathbf{A}^{-1}\mathbf{E}\| < 1$  for some matrix norm consistent on  $\mathbb{C}^{n \times n}$  then  $\mathbf{A} + \mathbf{E}$  is nonsingular. If  $\mathbf{A}\mathbf{x} = \mathbf{b}$  and  $(\mathbf{A} + \mathbf{E})\mathbf{y} = \mathbf{b}$  then

$$\frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{y}\|} \leq \|\mathbf{A}^{-1}\mathbf{E}\| \leq K(\mathbf{A}) \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}, \quad (7.18)$$

$$\frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq 2K(\mathbf{A}) \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}. \quad (7.19)$$

In (7.19) we have assumed that  $r \leq 1/2$ .

**Proof.** Since  $r < 1$  Theorem 7.29 implies that the matrix  $\mathbf{I} - \mathbf{B} := \mathbf{I} + \mathbf{A}^{-1}\mathbf{E}$  is nonsingular and then  $\mathbf{A} + \mathbf{E} = \mathbf{A}(\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})$  is nonsingular. Subtracting  $(\mathbf{A} + \mathbf{E})\mathbf{y} = \mathbf{b}$  from  $\mathbf{A}\mathbf{x} = \mathbf{b}$  gives  $\mathbf{A}(\mathbf{x} - \mathbf{y}) = \mathbf{E}\mathbf{y}$  or  $\mathbf{x} - \mathbf{y} = \mathbf{A}^{-1}\mathbf{E}\mathbf{y}$ . Taking norms and dividing by  $\|\mathbf{y}\|$  proves (7.18). Solving  $\mathbf{x} - \mathbf{y} = \mathbf{A}^{-1}\mathbf{E}\mathbf{y}$  for  $\mathbf{y}$  we obtain  $\mathbf{y} = (\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1}\mathbf{x}$ . By (7.16)

$$\|\mathbf{y}\| \leq \|(\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1}\|\|\mathbf{x}\| \leq \frac{\|\mathbf{x}\|}{1 - \|\mathbf{A}^{-1}\mathbf{E}\|} \leq 2\|\mathbf{x}\|.$$

But then (7.19) follows from (7.18).  $\square$

In Theorem 7.30 we gave bounds for the relative error in  $\mathbf{x}$  as an approximation to  $\mathbf{y}$  and the relative error in  $\mathbf{y}$  as an approximation to  $\mathbf{x}$ .  $\|\mathbf{E}\|/\|\mathbf{A}\|$  is a measure for the size of the perturbation  $\mathbf{E}$  in  $\mathbf{A}$  relative to the size of  $\mathbf{A}$ . The condition number again plays a crucial role.  $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{y}\|$  can be as large as  $K(\mathbf{A})$  times  $\|\mathbf{E}\|/\|\mathbf{A}\|$ . It can be shown that the upper bound can be attained for any  $\mathbf{A}$  and any  $\mathbf{b}$ . In deriving the upper bound we used the inequality  $\|\mathbf{A}^{-1}\mathbf{E}\mathbf{y}\| \leq \|\mathbf{A}^{-1}\|\|\mathbf{E}\|\|\mathbf{y}\|$ . For a more or less random perturbation  $\mathbf{E}$  this is not a severe overestimate for  $\|\mathbf{A}^{-1}\mathbf{E}\mathbf{y}\|$ . In the situation where  $\mathbf{E}$  is due to round-off errors (7.18) can give a fairly realistic estimate for  $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{y}\|$ .

We end this section with a perturbation result for the inverse matrix. Again the condition number plays an important role.

**Theorem 7.31 (Perturbation of inverse matrix)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is nonsingular and let  $\|\cdot\|$  be a consistent matrix norm on  $\mathbb{C}^{n \times n}$ . If  $\mathbf{E} \in \mathbb{C}^{n \times n}$  is so small that  $r := \|\mathbf{A}^{-1}\mathbf{E}\| < 1$  then  $\mathbf{A} + \mathbf{E}$  is nonsingular and

$$\|(\mathbf{A} + \mathbf{E})^{-1}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - r}. \quad (7.20)$$

If  $r < 1/2$  then

$$\frac{\|(\mathbf{A} + \mathbf{E})^{-1} - \mathbf{A}^{-1}\|}{\|\mathbf{A}^{-1}\|} \leq 2K(\mathbf{A}) \frac{\|\mathbf{E}\|}{\|\mathbf{A}\|}. \quad (7.21)$$

**Proof.** We showed in Theorem 7.30 that  $\mathbf{A} + \mathbf{E}$  is nonsingular and since  $(\mathbf{A} + \mathbf{E})^{-1} = (\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1}\mathbf{A}^{-1}$  we obtain

$$\|(\mathbf{A} + \mathbf{E})^{-1}\| \leq \|(\mathbf{I} + \mathbf{A}^{-1}\mathbf{E})^{-1}\|\|\mathbf{A}^{-1}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - \|\mathbf{A}^{-1}\mathbf{E}\|}$$

and (7.20) follows. Since

$$(\mathbf{A} + \mathbf{E})^{-1} - \mathbf{A}^{-1} = (\mathbf{I} - \mathbf{A}^{-1}(\mathbf{A} + \mathbf{E}))(\mathbf{A} + \mathbf{E})^{-1} = -\mathbf{A}^{-1}\mathbf{E}(\mathbf{A} + \mathbf{E})^{-1}$$

we obtain by (7.20)

$$\|(\mathbf{A} + \mathbf{E})^{-1} - \mathbf{A}^{-1}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{E}\| \|(\mathbf{A} + \mathbf{E})^{-1}\| \leq K(\mathbf{A}) \frac{\|\mathbf{E}\| \|\mathbf{A}^{-1}\|}{\|\mathbf{A}\| (1-r)}.$$

Dividing by  $\|\mathbf{A}^{-1}\|$  and setting  $r = 1/2$  proves (7.21).  $\square$

**Exercise 7.32 (Sharpness of perturbation bounds)**

The upper and lower bounds for  $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{x}\|$  given by (7.13) can be attained for any matrix  $\mathbf{A}$ , but only for special choices of  $\mathbf{b}$ . Suppose  $\mathbf{y}_{\mathbf{A}}$  and  $\mathbf{y}_{\mathbf{A}^{-1}}$  are vectors with  $\|\mathbf{y}_{\mathbf{A}}\| = \|\mathbf{y}_{\mathbf{A}^{-1}}\| = 1$  and  $\|\mathbf{A}\| = \|\mathbf{A}\mathbf{y}_{\mathbf{A}}\|$  and  $\|\mathbf{A}^{-1}\| = \|\mathbf{A}^{-1}\mathbf{y}_{\mathbf{A}^{-1}}\|$ .

(a) Show that the upper bound in (7.13) is attained if  $\mathbf{b} = \mathbf{A}\mathbf{y}_{\mathbf{A}}$  and  $\mathbf{e} = \mathbf{y}_{\mathbf{A}^{-1}}$ .

(b) Show that the lower bound is attained if  $\mathbf{b} = \mathbf{y}_{\mathbf{A}^{-1}}$  and  $\mathbf{e} = \mathbf{A}\mathbf{y}_{\mathbf{A}}$ .

**Exercise 7.33 (Condition number of 2. derivative matrix)**

In this exercise we will show that for  $m \geq 1$

$$\frac{4}{\pi^2}(m+1)^2 - 2/3 < \text{cond}_p(\mathbf{T}) \leq \frac{1}{2}(m+1)^2, \quad p = 1, 2, \infty, \quad (7.22)$$

where  $\mathbf{T} := \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m \times m}$  and  $\text{cond}_p(\mathbf{T}) := \|\mathbf{T}\|_p \|\mathbf{T}^{-1}\|_p$  is the  $p$ -norm condition number of  $\mathbf{T}$ . The  $p$  matrix norm is given by (7.8). You will need the explicit inverse of  $\mathbf{T}$  given by (1.10) and the eigenvalues given in Lemma 3.8. As usual we define  $h := 1/(m+1)$ .

a) Show that for  $m \geq 3$

$$\text{cond}_1(\mathbf{T}) = \text{cond}_\infty(\mathbf{T}) = \frac{1}{2} \begin{cases} h^{-2}, & m \text{ odd}, \\ h^{-2} - 1, & m \text{ even}. \end{cases} \quad (7.23)$$

and that  $\text{cond}_1(\mathbf{T}) = \text{cond}_\infty(\mathbf{T}) = 3$  for  $m = 2$ .

b) Show that for  $p = 2$  and  $m \geq 1$  we have

$$\text{cond}_2(\mathbf{T}) = \cot^2\left(\frac{\pi h}{2}\right) = 1/\tan^2\left(\frac{\pi h}{2}\right).$$

c) Show the bounds

$$\frac{4}{\pi^2}h^{-2} - \frac{2}{3} < \text{cond}_2(\mathbf{T}) < \frac{4}{\pi^2}h^{-2}. \quad (7.24)$$

*Hint:* For the upper bound use the inequality  $\tan x > x$  valid for  $0 < x < \pi/2$ . For the lower bound we use the inequality  $\cot^2 x > \frac{1}{x^2} - \frac{2}{3}$  for  $x > 0$ . This can be derived for  $0 < x < \pi$  by first showing that the second derivative of  $\cot^2 x$  is positive and then use Taylor's theorem.

d) Show (7.22).

## 7.3 Proof that the $p$ -Norms are Norms

We want to show

### Theorem 7.34 (The $p$ vector norms are norms)

Let for  $1 \leq p \leq \infty$  and  $\mathbf{x} \in \mathbb{C}^n$

$$\|\mathbf{x}\|_p := \left( \sum_{j=1}^n |x_j|^p \right)^{1/p}, \quad \|\mathbf{x}\|_\infty := \max_{1 \leq j \leq n} |x_j|.$$

Then for all  $1 \leq p \leq \infty$ ,  $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$  and all  $a \in \mathbb{C}$

1.  $\|\mathbf{x}\|_p \geq 0$  with equality if and only if  $\mathbf{x} = \mathbf{0}$ . (positivity)
2.  $\|a\mathbf{x}\|_p = |a| \|\mathbf{x}\|_p$ . (homogeneity)
3.  $\|\mathbf{x} + \mathbf{y}\|_p \leq \|\mathbf{x}\|_p + \|\mathbf{y}\|_p$ . (subadditivity)

Positivity and homogeneity follows immediately. To show the subadditivity we need some elementary properties of convex functions.

### Definition 7.35 (Convex function)

Let  $I \subset \mathbb{R}$  be an interval. A function  $f : I \rightarrow \mathbb{R}$  is convex if

$$f((1 - \lambda)x_1 + \lambda x_2) \leq (1 - \lambda)f(x_1) + \lambda f(x_2) \quad (7.25)$$

for all  $x_1, x_2 \in I$  with  $x_1 < x_2$  and all  $\lambda \in [0, 1]$ . The sum  $\sum_{j=1}^n \lambda_j x_j$  is called a **convex combination** of  $x_1, \dots, x_n$  if  $\lambda_j \geq 0$  for  $j = 1, \dots, n$  and  $\sum_{j=1}^n \lambda_j = 1$ .

The convexity condition is illustrated in Figure 7.1.

### Lemma 7.36 (A sufficient condition for convexity)

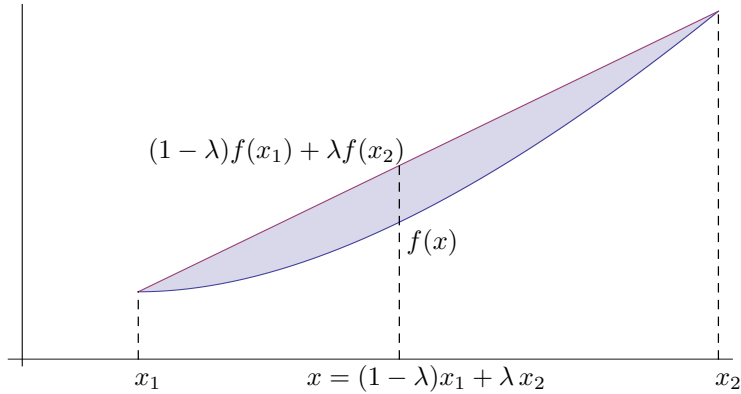
If  $f \in C^2[a, b]$  and  $f''(x) \geq 0$  for  $x \in [a, b]$  then  $f$  is convex.

**Proof.** We recall the formula for linear interpolation with remainder, (cf a book on numerical methods) For any  $a \leq x_1 \leq x \leq x_2 \leq b$  there is a  $c \in [x_1, x_2]$  such that

$$\begin{aligned} f(x) &= \frac{x_2 - x}{x_2 - x_1} f(x_1) + \frac{x - x_1}{x_2 - x_1} f(x_2) + (x - x_1)(x - x_2) f''(c)/2 \\ &= (1 - \lambda) f(x_1) + \lambda f(x_2) + (x_2 - x_1)^2 \lambda(\lambda - 1) f''(c)/2, \quad \lambda := \frac{x - x_1}{x_2 - x_1}. \end{aligned}$$

Since  $\lambda \in [0, 1]$  the remainder term is not positive. Moreover,

$$x = \frac{x_2 - x}{x_2 - x_1} x_1 + \frac{x - x_1}{x_2 - x_1} x_2 = (1 - \lambda)x_1 + \lambda x_2$$



**Figure 7.1.** A convex function.

so that (7.25) holds, and  $f$  is convex.  $\square$

The following inequality is elementary, but can be used to prove many non-trivial inequalities.

**Theorem 7.37 (Jensen's inequality)**

Suppose  $I \subseteq \mathbb{R}$  is an interval and  $f : I \rightarrow \mathbb{R}$  is convex. Then for all  $n \in \mathbb{N}$ , all  $\lambda_1, \dots, \lambda_n$  with  $\lambda_j \geq 0$  for  $j = 1, \dots, n$  and  $\sum_{j=1}^n \lambda_j = 1$ , and all  $z_1, \dots, z_n \in I$

$$f\left(\sum_{j=1}^n \lambda_j z_j\right) \leq \sum_{j=1}^n \lambda_j f(z_j).$$

**Proof.** We use induction on  $n$ . The result is trivial for  $n = 1$ . Let  $n \geq 2$ , assume the inequality holds for  $n - 1$ , and let  $\lambda_j, z_j$  for  $j = 1, \dots, n$  be given as in the theorem. Since  $n \geq 2$  we have  $\lambda_i < 1$  for at least one  $i$  so assume without loss of generality that  $\lambda_1 < 1$ , and define  $u := \sum_{j=2}^n \frac{\lambda_j}{1-\lambda_1} z_j$ . Since  $\sum_{j=2}^n \lambda_j = 1 - \lambda_1$  this is a convex combination of  $n - 1$  terms and the induction hypothesis implies that  $f(u) \leq \sum_{j=2}^n \frac{\lambda_j}{1-\lambda_1} f(z_j)$ . But then by the convexity of  $f$

$$f\left(\sum_{j=1}^n \lambda_j z_j\right) = f(\lambda_1 z_1 + (1 - \lambda_1)u) \leq \lambda_1 f(z_1) + (1 - \lambda_1)f(u) \leq \sum_{j=1}^n \lambda_j f(z_j)$$

and the inequality holds for  $n$ .  $\square$



**Corollary 7.38 (Weighted geometric/arithmetic mean inequality)**

Suppose  $\sum_{j=1}^n \lambda_j a_j$  is a convex combination of nonnegative numbers  $a_1, \dots, a_n$ . Then

$$a_1^{\lambda_1} a_2^{\lambda_2} \cdots a_n^{\lambda_n} \leq \sum_{j=1}^n \lambda_j a_j, \quad (7.26)$$

where  $0^0 := 0$ .

**Proof.** The result is trivial if one or more of the  $a_j$ 's are zero so assume  $a_j > 0$  for all  $j$ . Consider the function  $f : (0, \infty)$  given by  $f(x) = -\log x$ . Since  $f''(x) = 1/x^2 > 0$  for  $x \in (0, \infty)$ , this function is convex. By Jensen's inequality

$$-\log \left( \sum_{j=1}^n \lambda_j a_j \right) \leq -\sum_{j=1}^n \lambda_j \log(a_j) = -\log(a_1^{\lambda_1} \cdots a_n^{\lambda_n})$$

or  $\log(a_1^{\lambda_1} \cdots a_n^{\lambda_n}) \leq \log(\sum_{j=1}^n \lambda_j a_j)$ . The inequality follows since  $\exp(\log x) = x$  for  $x > 0$  and the exponential function is monotone increasing.  $\square$

Taking  $\lambda_j = \frac{1}{n}$  for all  $j$  in (7.26) we obtain the classical **geometric/arithmetic mean inequality**

$$(a_1 a_2 \cdots a_n)^{\frac{1}{n}} \leq \frac{1}{n} \sum_{j=1}^n a_j. \quad (7.27)$$

**Corollary 7.39 (Hölder's inequality)**

For  $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$  and  $1 \leq p \leq \infty$

$$\sum_{j=1}^n |x_j y_j| \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q, \text{ where } \frac{1}{p} + \frac{1}{q} = 1.$$

**Proof.** We leave the proof for  $p = 1$  and  $p = \infty$  as an exercise so assume  $1 < p < \infty$ . For any  $a, b \geq 0$  the weighted arithmetic/geometric mean inequality implies that

$$a^{\frac{1}{p}} b^{\frac{1}{q}} \leq \frac{1}{p} a + \frac{1}{q} b, \text{ where } \frac{1}{p} + \frac{1}{q} = 1. \quad (7.28)$$

If  $\mathbf{x} = \mathbf{0}$  or  $\mathbf{y} = \mathbf{0}$  there is nothing to prove so assume that both  $\mathbf{x}$  and  $\mathbf{y}$  are nonzero. Using 7.28 on each term we obtain

$$\frac{1}{\|\mathbf{x}\|_p \|\mathbf{y}\|_q} \sum_{j=1}^n |x_j y_j| = \sum_{j=1}^n \left( \frac{|x_j|^p}{\|\mathbf{x}\|_p^p} \right)^{\frac{1}{p}} \left( \frac{|y_j|^q}{\|\mathbf{y}\|_q^q} \right)^{\frac{1}{q}} \leq \sum_{j=1}^n \left( \frac{1}{p} \frac{|x_j|^p}{\|\mathbf{x}\|_p^p} + \frac{1}{q} \frac{|y_j|^q}{\|\mathbf{y}\|_q^q} \right) = 1$$

and the proof of the inequality is complete.  $\square$

**Corollary 7.40 (Minkowski's inequality)**

For  $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$  and  $1 \leq p \leq \infty$

$$\|\mathbf{x} + \mathbf{y}\|_p \leq \|\mathbf{x}\|_p + \|\mathbf{y}\|_p.$$

**Proof.** We leave the proof for  $p = 1$  and  $p = \infty$  as an exercise so assume  $1 < p < \infty$ . We write

$$\|\mathbf{x} + \mathbf{y}\|_p^p = \sum_{j=1}^n |x_j + y_j|^p \leq \sum_{j=1}^n |x_j| |x_j + y_j|^{p-1} + \sum_{j=1}^n |y_j| |x_j + y_j|^{p-1}.$$

We apply Hölder's inequality with exponent  $p$  and  $q$  to each sum. In view of the relation  $(p-1)q = p$  the result is

$$\|\mathbf{x} + \mathbf{y}\|_p^p \leq \|\mathbf{x}\|_p \|\mathbf{x} + \mathbf{y}\|_p^{p/q} + \|\mathbf{y}\|_p \|\mathbf{x} + \mathbf{y}\|_p^{p/q} = (\|\mathbf{x}\|_p + \|\mathbf{y}\|_p) \|\mathbf{x} + \mathbf{y}\|_p^{p-1},$$

and canceling the common factor, the inequality follows.  $\square$

It is possible to characterize the  $p$ -norms that are derived from an inner product. We start with the following identity.

**Theorem 7.41 (Parallelogram identity)**

For all  $\mathbf{x}, \mathbf{y}$  in a real or complex inner product space

$$\|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2 = 2\|\mathbf{x}\|^2 + 2\|\mathbf{y}\|^2. \quad (7.29)$$

**Proof.** We set  $a = \pm 1$  in (22) and add the two equations.  $\square$

**Theorem 7.42 (When is a norm an inner product norm?)**

To a given norm on a real or complex vector space  $\mathcal{V}$  there exists an inner product on  $\mathcal{V}$  such that  $\langle \mathbf{x}, \mathbf{x} \rangle = \|\mathbf{x}\|^2$  if and only if the parallelogram identity (7.29) holds for all  $\mathbf{x}, \mathbf{y} \in \mathcal{V}$ .

**Proof.** If  $\langle \mathbf{x}, \mathbf{x} \rangle = \|\mathbf{x}\|^2$  then Theorem 7.41 shows that the parallelogram identity holds. For the converse we prove the real case and leave the complex case as an exercise. Suppose (7.29) holds for all  $\mathbf{x}, \mathbf{y}$  in the real vector space  $\mathcal{V}$ . We show that

$$\langle \mathbf{x}, \mathbf{y} \rangle := \frac{1}{4} (\|\mathbf{x} + \mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2), \quad \mathbf{x}, \mathbf{y} \in \mathcal{V} \quad (7.30)$$

defines an inner product on  $\mathcal{V}$ . Clearly 1. and 2. in Definition 0.20 hold. The hard part is to show 3. We need to show that

$$\langle \mathbf{x}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z} \rangle = \langle \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle, \quad \mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{V}, \quad (7.31)$$

$$\langle a\mathbf{x}, \mathbf{y} \rangle = a\langle \mathbf{x}, \mathbf{y} \rangle, \quad a \in \mathbb{R}, \quad \mathbf{x}, \mathbf{y} \in \mathcal{V}. \quad (7.32)$$

Now

$$\begin{aligned} & 4\langle \mathbf{x}, \mathbf{z} \rangle + 4\langle \mathbf{y}, \mathbf{z} \rangle \stackrel{(7.30)}{=} \|\mathbf{x} + \mathbf{z}\|^2 - \|\mathbf{x} - \mathbf{z}\|^2 + \|\mathbf{y} + \mathbf{z}\|^2 - \|\mathbf{y} - \mathbf{z}\|^2 \\ &= \left\| \left( \mathbf{z} + \frac{\mathbf{x} + \mathbf{y}}{2} \right) + \frac{\mathbf{x} - \mathbf{y}}{2} \right\|^2 - \left\| \left( \mathbf{z} - \frac{\mathbf{x} + \mathbf{y}}{2} \right) + \frac{\mathbf{y} - \mathbf{x}}{2} \right\|^2 \\ &+ \left\| \left( \mathbf{z} + \frac{\mathbf{x} + \mathbf{y}}{2} \right) - \frac{\mathbf{x} - \mathbf{y}}{2} \right\|^2 - \left\| \left( \mathbf{z} - \frac{\mathbf{x} + \mathbf{y}}{2} \right) - \frac{\mathbf{y} - \mathbf{x}}{2} \right\|^2 \\ &\stackrel{(7.29)}{=} 2\left\| \mathbf{z} + \frac{\mathbf{x} + \mathbf{y}}{2} \right\|^2 + 2\left\| \frac{\mathbf{x} - \mathbf{y}}{2} \right\|^2 - 2\left\| \mathbf{z} - \frac{\mathbf{x} + \mathbf{y}}{2} \right\|^2 - 2\left\| \frac{\mathbf{y} - \mathbf{x}}{2} \right\|^2 \\ &\stackrel{(7.30)}{=} 8\left\langle \frac{\mathbf{x} + \mathbf{y}}{2}, \mathbf{z} \right\rangle, \end{aligned}$$

or

$$\langle \mathbf{x}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z} \rangle = 2\left\langle \frac{\mathbf{x} + \mathbf{y}}{2}, \mathbf{z} \right\rangle, \quad \mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{V}.$$

In particular, since  $\mathbf{y} = \mathbf{0}$  implies  $\langle \mathbf{y}, \mathbf{z} \rangle = 0$  we obtain  $\langle \mathbf{x}, \mathbf{z} \rangle = 2\langle \frac{\mathbf{x}}{2}, \mathbf{z} \rangle$  for all  $\mathbf{x}, \mathbf{z} \in \mathcal{V}$ . This means that  $2\langle \frac{\mathbf{x} + \mathbf{y}}{2}, \mathbf{z} \rangle = \langle \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle$  for all  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{V}$  and (7.31) follows.

We first show (7.32) when  $a = n$  is a positive integer. By induction

$$\langle n\mathbf{x}, \mathbf{y} \rangle = \langle (n-1)\mathbf{x} + \mathbf{x}, \mathbf{y} \rangle \stackrel{(7.31)}{=} \langle (n-1)\mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{y} \rangle = n\langle \mathbf{x}, \mathbf{y} \rangle. \quad (7.33)$$

If  $m, n \in \mathbb{N}$  then

$$m^2 \left\langle \frac{n}{m}\mathbf{x}, \mathbf{y} \right\rangle \stackrel{(7.33)}{=} m \langle n\mathbf{x}, \mathbf{y} \rangle \stackrel{(7.33)}{=} mn \langle \mathbf{x}, \mathbf{y} \rangle,$$

implying that (7.32) holds for positive rational numbers

$$\left\langle \frac{n}{m}\mathbf{x}, \mathbf{y} \right\rangle = \frac{n}{m} \langle \mathbf{x}, \mathbf{y} \rangle.$$

Now if  $a > 0$  there is a sequence  $\{a_n\}$  of positive rational numbers converging to  $a$ . For each  $n$

$$a_n \langle \mathbf{x}, \mathbf{y} \rangle = \langle a_n \mathbf{x}, \mathbf{y} \rangle \stackrel{(7.30)}{=} \frac{1}{4} (\|a_n \mathbf{x} + \mathbf{y}\|^2 - \|a_n \mathbf{x} - \mathbf{y}\|^2).$$

Taking limits and using continuity of norms we obtain  $a\langle \mathbf{x}, \mathbf{y} \rangle = \langle a\mathbf{x}, \mathbf{y} \rangle$ . This also holds for  $a = 0$ . Finally, if  $a < 0$  then  $(-a) > 0$  and from what we just showed

$$(-a)\langle \mathbf{x}, \mathbf{y} \rangle = \langle (-a)\mathbf{x}, \mathbf{y} \rangle \stackrel{(7.30)}{=} \frac{1}{4} (\| -a\mathbf{x} + \mathbf{y} \|^2 - \| -a\mathbf{x} - \mathbf{y} \|^2) = -\langle a\mathbf{x}, \mathbf{y} \rangle,$$

so (7.32) also holds for negative  $a$ .  $\square$

**Corollary 7.43 (Are the  $p$ -norms inner product norms?)**

For the  $p$  vector norms on  $\mathcal{V} = \mathbb{R}^n$  or  $\mathcal{V} = \mathbb{C}^n$ ,  $1 \leq p \leq \infty$ ,  $n \geq 2$ , there is an inner product on  $\mathcal{V}$  such that  $\langle \mathbf{x}, \mathbf{x} \rangle = \|\mathbf{x}\|_p^2$  for all  $\mathbf{x} \in \mathcal{V}$  if and only if  $p = 2$ .

**Proof.** For  $p = 2$  the  $p$ -norm is the Euclidian norm which corresponds to the standard inner product. If  $p \neq 2$  then the parallelogram identity (7.29) does not hold for say  $\mathbf{x} := \mathbf{e}_1$  and  $\mathbf{y} := \mathbf{e}_2$ .  $\square$

**Exercise 7.44 (When is a complex norm an inner product norm?)**

Given a vector norm in a complex vector space  $\mathcal{V}$ , and suppose (7.29) holds for all  $\mathbf{x}, \mathbf{y}$ . Show that

$$\langle \mathbf{x}, \mathbf{y} \rangle := \frac{1}{4} (\|\mathbf{x} + \mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2 + i\|\mathbf{x} + i\mathbf{y}\|^2 - i\|\mathbf{x} - i\mathbf{y}\|^2), \quad (7.34)$$

defines an inner product on  $\mathcal{V}$ , where  $i = \sqrt{-1}$ . The identity (7.34) is called the **polarization identity**.<sup>10</sup>

**Exercise 7.45 ( $p$  norm for  $p = 1$  and  $p = \infty$ )**

Show that  $\|\cdot\|_p$  is a vector norm in  $\mathbb{R}^n$  for  $p = 1, p = \infty$ .

**Exercise 7.46 (The  $p$ - norm unit sphere)**

The set

$$S_p = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_p = 1\}$$

is called the unit sphere in  $\mathbb{R}^n$  with respect to  $p$ . Draw  $S_p$  for  $p = 1, 2, \infty$  for  $n = 2$ .

**Exercise 7.47 (Sharpness of  $p$ -norm inequality)**

For  $p \geq 1$ , and any  $\mathbf{x} \in \mathbb{C}^n$  we have  $\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_p \leq n^{1/p} \|\mathbf{x}\|_\infty$  (cf. (14)).

Produce a vector  $\mathbf{x}_l$  such that  $\|\mathbf{x}_l\|_\infty = \|\mathbf{x}_l\|_p$  and another vector  $\mathbf{x}_u$  such that  $\|\mathbf{x}_u\|_p = n^{1/p} \|\mathbf{x}_u\|_\infty$ . Thus, these inequalities are sharp.

**Exercise 7.48 ( $p$ -norm inequalities for arbitrary  $p$ )**

If  $1 \leq q \leq p \leq \infty$  then

$$\|\mathbf{x}\|_p \leq \|\mathbf{x}\|_q \leq n^{1/q-1/p} \|\mathbf{x}\|_p, \quad \mathbf{x} \in \mathbb{C}^n.$$

*Hint:* For the rightmost inequality use Jensen's inequality Cf. Theorem 7.37 with  $f(z) = z^{p/q}$  and  $z_i = |x_i|^q$ . For the left inequality consider first  $y_i = x_i / \|\mathbf{x}\|_\infty$ ,  $i = 1, 2, \dots, n$ .

<sup>10</sup>Hint: We have  $\langle \mathbf{x}, \mathbf{y} \rangle = s(\mathbf{x}, \mathbf{y}) + is(\mathbf{x}, i\mathbf{y})$ , where  $s(\mathbf{x}, \mathbf{y}) := \frac{1}{4} (\|\mathbf{x} + \mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2)$ .

## 7.4 Review Questions

7.4.1]

- What is a consistent matrix norm?
- what is a subordinate matrix norm?
- is an operator norm consistent?
- why is the Frobenius norm not an operator norm?
- what is the spectral norm of a matrix?
- how do we compute  $\|\mathbf{A}\|_\infty$ ?
- what is the spectral condition number of a symmetric positive definite matrix?

**7.4.2** Why is  $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F$  for any matrix  $\mathbf{A}$ ?

**7.4.3** What is the spectral norm of the inverse of a normal matrix?



## **Part III**

# **Iterative Methods for Large Linear Systems**





## Chapter 8

# The Classical Iterative Methods

Gaussian elimination and Cholesky factorization are **direct methods**. In absence of rounding errors they find the exact solution using a finite number of arithmetic operations. In an **iterative method** we start with an approximation  $\mathbf{x}_0$  to the exact solution  $\mathbf{x}$  and then compute a sequence  $\{\mathbf{x}_k\}$  such that hopefully  $\mathbf{x}_k \rightarrow \mathbf{x}$ . Iterative methods are mainly used for large sparse systems, i. e., where many of the elements in the coefficient matrix are zero. The main advantages of iterative methods are reduced storage requirements and ease of implementation. In an iterative method the main work in each iteration is a matrix times vector multiplication, an operation which often does not need storing the matrix, not even in sparse form.

In this chapter we consider the classical iterative methods of Richardson, Jacobi, Gauss-Seidel and an accelerated version of Gauss-Seidel's method called successive overrelaxation (SOR). David Young developed in his thesis a beautiful theory describing the convergence rate of SOR, see [35].

We give the main points of this theory specialized to the discrete Poisson matrix. With a careful choice of an acceleration parameter the amount of work using SOR on the discrete Poisson problem is the same as for the fast Poisson solver without FFT (cf. Algorithm 4.1). Moreover, SOR is not restricted to constant coefficient methods on a rectangle. However, to obtain fast convergence using SOR it is necessary to have a good estimate for an acceleration parameter.

For convergence we study convergence of powers of matrices.

## 8.1 Classical Iterative Methods; Component Form

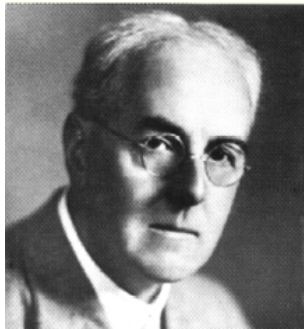
We start with an example showing how a linear system can be solved using an iterative method.

### Example 8.1 (Iterative methods on a special $2 \times 2$ matrix)

Solving for the diagonal elements the linear system  $\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} y \\ z \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$  can be written in component form as  $y = (z + 1)/2$  and  $z = (y + 1)/2$ . Starting with  $y_0, z_0$  we generate two sequences  $\{y_k\}$  and  $\{z_k\}$  using the difference equations  $y_{k+1} = (z_k + 1)/2$  and  $z_{k+1} = (y_k + 1)/2$ . This is known as Jacobi's method. If  $y_0 = z_0 = 0$  then we find  $y_1 = z_1 = 1/2$  and in general  $y_k = z_k = 1 - 2^{-k}$  for  $k = 0, 1, 2, 3, \dots$ . The iteration converges to the exact solution  $[1, 1]^T$ , and the error is halved in each iteration.

We can improve the convergence rate by using the most current approximation in each iteration. This leads to Gauss-Seidel's method:  $y_{k+1} = (z_k + 1)/2$  and  $z_{k+1} = (y_{k+1} + 1)/2$ . If  $y_0 = z_0 = 0$  then we find  $y_1 = 1/2$ ,  $z_1 = 3/4$ ,  $y_2 = 7/8$ ,  $z_2 = 15/16$ , and in general  $y_k = 1 - 2 \cdot 4^{-k}$  and  $z_k = 1 - 4^{-k}$  for  $k = 1, 2, 3, \dots$ . The error is now reduced by a factor 4 in each iteration.

Consider the general case. Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is nonsingular and  $\mathbf{b} \in \mathbb{C}^n$ . Suppose we know an approximation  $\mathbf{x}_k = [\mathbf{x}_k(1), \dots, \mathbf{x}_k(n)]^T$  to the exact solution  $\mathbf{x}$  of  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .



Lewis Fry Richardson, 1881-1953 (left), Carl Gustav Jacob Jacobi, 1804-1851 (right).



Philipp Ludwig von Seidel, 1821-1896 (left), David M. Young Jr., 1923-2008 (right)

In **Richardson's method (R method)** we pick a positive parameter  $\alpha$  and compute a new approximation by adding a multiple of the residual vector  $\mathbf{r}_k := \mathbf{b} - \mathbf{A}\mathbf{x}_k$ :

$$\mathbf{x}_{k+1}(i) = \mathbf{x}_k(i) + \alpha \mathbf{r}_k(i), \quad \mathbf{r}_k(i) := b_i - \sum_{j=1}^n a_{ij} \mathbf{x}_k(j), \quad \text{for } i = 1, 2, \dots, n. \quad (8.1)$$

Richardson considered the simplest case  $\alpha = 1$ . The parameter  $\alpha$  is added to get faster convergence.

For the other methods we need to assume that  $\mathbf{A}$  has nonzero diagonal elements. Solving the  $i$ th equation of  $\mathbf{A}\mathbf{x} = \mathbf{b}$  for  $\mathbf{x}(i)$ , we obtain a **fixed-point form** of  $\mathbf{A}\mathbf{x} = \mathbf{b}$

$$\mathbf{x}(i) = \left( - \sum_{j=1}^{i-1} a_{ij} \mathbf{x}(j) - \sum_{j=i+1}^n a_{ij} \mathbf{x}(j) + b_i \right) / a_{ii}, \quad i = 1, 2, \dots, n. \quad (8.2)$$

1. In **Jacobi's method (J method)** we substitute  $\mathbf{x}_k$  into the right hand side of (8.2) and compute a new approximation by

$$\mathbf{x}_{k+1}(i) = \left( - \sum_{j=1}^{i-1} a_{ij} \mathbf{x}_k(j) - \sum_{j=i+1}^n a_{ij} \mathbf{x}_k(j) + b_i \right) / a_{ii}, \quad \text{for } i = 1, 2, \dots, n. \quad (8.3)$$

2. **Gauss-Seidel's method (GS method)** is a modification of Jacobi's method, where we use the new  $\mathbf{x}_{k+1}(i)$  immediately after it has been computed.

$$\mathbf{x}_{k+1}(i) = \left( - \sum_{j=1}^{i-1} a_{ij} \mathbf{x}_{k+1}(j) - \sum_{j=i+1}^n a_{ij} \mathbf{x}_k(j) + b_i \right) / a_{ii}, \quad \text{for } i = 1, 2, \dots, n. \quad (8.4)$$

3. The **Successive overrelaxation method (SOR method)** is obtained by introducing an acceleration parameter  $0 < \omega < 2$  in the GS method. We write  $\mathbf{x}(i) = \omega\mathbf{x}(i) + (1 - \omega)\mathbf{x}(i)$  and this leads to the method

$$\mathbf{x}_{k+1}(i) = \omega \left( - \sum_{j=1}^{i-1} a_{ij} \mathbf{x}_{k+1}(j) - \sum_{j=i+1}^n a_{ij} \mathbf{x}_k(j) + b_i \right) / a_{ii} + (1 - \omega) \mathbf{x}_k(i). \quad (8.5)$$

The SOR method reduces to the Gauss-Seidel method for  $\omega = 1$ . Denoting the right hand side of (8.4) by  $\mathbf{x}_{k+1}^{gs}$  we can write (8.5) as  $\mathbf{x}_{k+1} = \omega \mathbf{x}_{k+1}^{gs} + (1 - \omega) \mathbf{x}_k$ , and we see that  $\mathbf{x}_{k+1}$  is located on the straight line passing through the two points  $\mathbf{x}_{k+1}^{gs}$  and  $\mathbf{x}_k$ . The restriction  $0 < \omega < 2$  is necessary for convergence (cf. Theorem 8.14). Normally, the best results are obtained for the relaxation parameter  $\omega$  in the range  $1 \leq \omega < 2$  and then  $\mathbf{x}_{k+1}$  is computed by linear extrapolation, i. e., it is not located between  $\mathbf{x}_{k+1}^{gs}$  and  $\mathbf{x}_k$ .

4. We mention also briefly the symmetric successive overrelaxation method **SSOR**. One iteration in SSOR consists of two SOR sweeps. A forward SOR sweep (8.5), computing an approximation denoted  $\mathbf{x}_{k+1/2}$  instead of  $\mathbf{x}_{k+1}$ , is followed by a back SOR sweep computing

$$\mathbf{x}_{k+1}(i) = \omega \left( - \sum_{j=1}^{i-1} a_{ij} \mathbf{x}_{k+1/2}(j) - \sum_{j=i+1}^n a_{ij} \mathbf{x}_{k+1}(j) + b_i \right) / a_{ii} + (1 - \omega) \mathbf{x}_{k+1/2}(i) \quad (8.6)$$

in the order  $i = n, n - 1, \dots, 1$ . The method is slower and more complicated than the SOR method. Its main use is as a symmetric preconditioner. For if  $\mathbf{A}$  is symmetric then SSOR combines the two SOR steps in such a way that the resulting iteration matrix is similar to a symmetric matrix. We will not discuss this method any further here and refer to Section 9.6 for an alternative example of a preconditioner.

We will refer to the R, J, GS and SOR methods as the **classical (iteration) methods**. The R method will be discussed later, see Section 8.3.1.

### 8.1.1 The discrete Poisson system

Consider the classical methods applied to the discrete Poisson matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  given by (3.7). Let  $n = m^2$  and set  $h = 1/(m + 1)$ . In component form the linear system  $\mathbf{Ax} = \mathbf{b}$  can be written (cf. (3.4))

$$4\mathbf{v}(i, j) - \mathbf{v}(i-1, j) - \mathbf{v}(i+1, j) - \mathbf{v}(i, j-1) - \mathbf{v}(i, j+1) = h^2 f_{i,j}, \quad i, j = 1, \dots, m,$$

with homogenous boundary conditions also given in (3.4). Solving for  $\mathbf{v}(i, j)$  we obtain the **fixed point form**

$$\mathbf{v}(i, j) = (\mathbf{v}(i-1, j) + \mathbf{v}(i+1, j) + \mathbf{v}(i, j-1) + \mathbf{v}(i, j+1) + \mathbf{e}_{i,j})/4, \quad \mathbf{e}_{i,j} := f_{i,j}/(m+1)^2. \quad (8.7)$$

The R, J, GS, and SOR methods takes the form

$$\begin{aligned} R : \mathbf{v}_{k+1}(i, j) &= \mathbf{v}_k(i, j) + \alpha(\mathbf{v}_k(i-1, j) + \mathbf{v}_k(i, j-1) + \mathbf{v}_k(i+1, j) \\ &\quad + \mathbf{v}_k(i, j+1) - 4\mathbf{v}_k(i, j) + \mathbf{e}(i, j)) \\ J : \mathbf{v}_{k+1}(i, j) &= (\mathbf{v}_k(i-1, j) + \mathbf{v}_k(i, j-1) + \mathbf{v}_k(i+1, j) + \mathbf{v}_k(i, j+1) \\ &\quad + \mathbf{e}(i, j))/4 \\ GS : \mathbf{v}_{k+1}(i, j) &= (\mathbf{v}_{k+1}(i-1, j) + \mathbf{v}_{k+1}(i, j-1) + \mathbf{v}_k(i+1, j) + \mathbf{v}_k(i, j+1) \\ &\quad + \mathbf{e}(i, j))/4 \\ SOR : \mathbf{v}_{k+1}(i, j) &= \omega(\mathbf{v}_{k+1}(i-1, j) + \mathbf{v}_{k+1}(i, j-1) + \mathbf{v}_k(i+1, j) + \mathbf{v}_k(i, j+1) \\ &\quad + \mathbf{e}(i, j))/4 + (1 - \omega)\mathbf{v}_k(i, j). \end{aligned} \quad (8.8)$$

We note that

- For  $\alpha = 1/4$  the R and J methods in (8.8) are identical.
- For a general system the R and J methods are identical if  $\mathbf{A}$  is constant and nonzero on the diagonal and  $\alpha$  is chosen as the inverse of this constant. See Exercise 8.2.
- For the discrete Poisson problem the choice  $\alpha = 1/4$  in the J method is optimal in a way to become clear in Section 8.3.1.
- For GS and SOR we have used the **natural ordering**, i. e.,  $(i_1, j_1) < (i_2, j_2)$  if and only if  $j_1 \leq j_2$  and  $i_1 < i_2$  if  $j_1 = j_2$ . For the J method any ordering can be used.

### Exercise 8.2 (Richardson and Jacobi)

Show that if  $a_{ii} = d \neq 0$  for all  $i$  then Richardson's method with  $\alpha := 1/d$  is the same as Jacobi's method.

In Algorithm 8.3 we give a Matlab program to test the convergence of Jacobi's method on the discrete Poisson problem. We carry out Jacobi iterations on the linear system (8.7) with  $\mathbf{F} = (f_{ij}) \in \mathbb{R}^{m \times m}$ , starting with  $\mathbf{V}_0 = \mathbf{0} \in \mathbb{R}^{(m+2) \times (m+2)}$ . The output is the number of iterations  $k$ , to obtain  $\|\mathbf{V}^{(k)} - \mathbf{U}\|_M := \max_{i,j} |v_{ij} - u_{ij}| < tol$ . Here  $[u_{ij}] \in \mathbb{R}^{(m+2) \times (m+2)}$  is the "exact" solution of (8.7) computed using the fast Poisson solver in Algorithm 4.1. We set  $k = K + 1$

|     | $k_{100}$ | $k_{2500}$ | $k_{10\ 000}$ | $k_{40\ 000}$ | $k_{160\ 000}$ |
|-----|-----------|------------|---------------|---------------|----------------|
| J   | 385       | 8386       |               |               |                |
| GS  | 194       | 4194       |               |               |                |
| SOR | 35        | 164        | 324           | 645           | 1286           |

**Table 8.1.** *The number of iterations  $k_n$  to solve the discrete Poisson problem with  $n$  unknowns using the methods of Jacobi, Gauss-Seidel, and SOR (see text) with a tolerance  $10^{-8}$ .*

if convergence is not obtained in  $K$  iterations. In Table 8.1 we show the output  $k = k_n$  from this algorithm using  $\mathbf{F} = \mathbf{ones}(m, m)$  for  $m = 10, 50$ ,  $K = 10^4$ , and  $tol = 10^{-8}$ . We also show the number of iterations for Gauss-Seidel and SOR with a value of  $\omega$  known as the optimal acceleration parameter  $\omega^* := 2/(1 + \sin(\frac{\pi}{m+1}))$ . We will derive this value later.

### Algorithm 8.3 (Jacobi)

```

1 function k=jdp(F,K,tol)
2 m=length(F); U=fastpoisson(F); V=zeros(m+2,m+2); E=F/(m+1)^2;
3 for k=1:K
4     V(2:m+1,2:m+1)=(V(1:m,2:m+1)+V(3:m+2,2:m+1)...
5         +V(2:m+1,1:m)+V(2:m+1,3:m+2)+E)/4;
6     if max(max(abs(V-U)))<tol, return
7     end
8 end
9 k=K+1;

```

For the GS and SOR methods we have used Algorithm 8.4. This is the analog of Algorithm 8.3 using SOR instead of J to solve the discrete Poisson problem.  $w$  is an acceleration parameter with  $0 < w < 2$ . For  $w = 1$  we obtain Gauss-Seidel's method.

**Algorithm 8.4 (SOR)**

```

1 function k=sordp(F,K,w,tol)
2 m=length(F); U=fastpoisson(F); V=zeros(m+2,m+2); E=F/(m+1)^2;
3 for k=1:K
4     for j=2:m+1
5         for i=2:m+1
6             V(i,j)=w*(V(i-1,j)+V(i+1,j)+V(i,j-1)...
7                 +V(i,j+1)+E(i-1,j-1))/4+(1-w)*V(i,j);
8         end
9     end
10    if max(max(abs(V-U)))<tol, return
11    end
12 end
13 k=K+1;

```

We make several remarks about these programs and the results in Table 8.1.

1. The rate (speed) of convergence is quite different for the four methods. The R, J and GS methods converge, but rather slowly. The R method with  $\alpha = 1/4$  and the J method needs about twice as many iterations as the GS method. The improvement using the SOR method with optimal  $\omega$  is spectacular.
2. We show in Section 8.3.4 that the number of iterations  $k_n$  for a size  $n$  problem is  $k_n = O(n)$  for the J and GS method and  $k_n = O(\sqrt{n})$  for SOR with optimal  $\omega$ . The choice of *tol* will only influence the constants multiplying  $n$  or  $\sqrt{n}$ .
3. From (8.8) it follows that each iteration requires  $O(n)$  arithmetic operations. Thus the number of arithmetic operations to achieve a given tolerance is  $O(k_n \times n)$ . Therefore the number of arithmetic operations for the J and GS method is  $O(n^2)$ , while it is only  $O(n^{3/2})$  for the SOR method with optimal  $\omega$ . Asymptotically, for J and GS this is the same as using banded Cholesky, while SOR competes with the fast method (without FFT).
4. We do not need to store the coefficient matrix so the storage requirements for these methods on the discrete Poisson problem is  $O(n)$ , asymptotically the same as for the fast methods.
5. Jacobi's method has the advantage that it can be easily parallelized.

## 8.2 Classical Iterative Methods; Matrix Form

To study convergence we need matrix formulations of the classical methods.

### 8.2.1 Fixed-point form

In general we can construct an iterative method by choosing a nonsingular matrix  $M$  and write  $Ax = b$  in the equivalent form  $M^{-1}Ax = M^{-1}b$ . This system can be written  $x = x - M^{-1}Ax + M^{-1}b$ , and we obtain  $Ax = b$  in a **fixed-point form**

$$x = Gx + c, \quad G = I - M^{-1}A, \quad c = M^{-1}b. \quad (8.9)$$

For a general  $G \in \mathbb{C}^{n \times n}$  and  $c \in \mathbb{C}^n$  a solution of  $x = Gx + c$  is called a **fixed-point**. The fixed-point is unique if  $I - G$  is nonsingular.

The corresponding iterative method is given by

$$x_{k+1} := Gx_k + c. \quad (8.10)$$

This is known as a **fixed-point iteration**. Starting with  $x_0$  this defines a sequence  $\{x_k\}$  of vectors in  $\mathbb{C}^n$ . If  $\lim_{k \rightarrow \infty} x_k = x$  for some  $x \in \mathbb{C}^n$  then  $x$  is a fixed point since

$$x = \lim_{k \rightarrow \infty} x_{k+1} = \lim_{k \rightarrow \infty} (Gx_k + c) = G \lim_{k \rightarrow \infty} x_k + c = Gx + c.$$

### 8.2.2 The preconditioning and splitting matrix

Different choices of  $M$  in (8.9) lead to different iterative methods. The matrix  $M$  can be interpreted in two ways. It is a **preconditioning matrix** since a good choice of  $M$  can lead to a preconditioned system  $M^{-1}Ax = M^{-1}b$  with smaller condition number. It is also known as a **splitting matrix**, since if we split  $A$  in the form  $A = M + (A - M)$  then  $Ax = b$  can be written  $Mx = (M - A)x + b$ , and this leads to the iterative method

$$Mx_{k+1} = (M - A)x_k + b \quad (8.11)$$

which is equivalent to (8.9).

### 8.2.3 The splitting matrices for the classical methods

We now derive  $M$  for the classical methods. For J, GS and SOR it is convenient to write  $A$  as a sum of three matrices,  $A = D - A_L - A_R$ , where  $-A_L$ ,  $D$ , and  $-A_R$  are the lower, diagonal, and upper part of  $A$ , respectively. Thus  $D := \text{diag}(a_{11}, \dots, a_{nn})$ ,

$$A_L := \begin{bmatrix} 0 & & & & \\ -a_{2,1} & 0 & & & \\ \vdots & \ddots & \ddots & & \\ -a_{n,1} & \cdots & -a_{n,n-1} & 0 & \end{bmatrix}, \quad A_R := \begin{bmatrix} 0 & -a_{1,2} & \cdots & -a_{1,n} \\ & \ddots & \ddots & \vdots \\ & & 0 & -a_{n-1,n} \\ & & & 0 \end{bmatrix}. \quad (8.12)$$



**Theorem 8.5 (Splitting matrices for R, J, and SOR)**

The splitting and iteration matrices for the R, J and SOR methods are given by

$$\begin{aligned} M_R &= \alpha^{-1}I, \quad M_J = D, \quad M_\omega = \omega^{-1}D - A_L, \\ G_R &= I - \alpha A, \quad G_J = I - D^{-1}A, \quad G_\omega = I - (\omega^{-1}D - A_L)^{-1}A. \end{aligned} \quad (8.13)$$

We obtain the matrices  $M_1$  and  $G_1$  for the GS method by letting  $\omega = 1$  in  $M_\omega$  and  $G_\omega$

**Proof.** To find  $M$  we write the methods in the form (8.11). The formulas for  $G$  then follows immediately from (8.9).

The matrix form of the R method is  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha(\mathbf{b} - \mathbf{A}\mathbf{x}_k)$  or  $\alpha^{-1}\mathbf{x}_{k+1} = \alpha^{-1}(\mathbf{I} - \mathbf{A})\mathbf{x}_k + \mathbf{b}$ , and the formulas for  $M_R$  and  $G_R$  follows. The equation  $\mathbf{A}\mathbf{x} = \mathbf{b}$  can be written  $D\mathbf{x} - A_L\mathbf{x} - A_R\mathbf{x} = \mathbf{b}$  or  $D\mathbf{x} = A_L\mathbf{x} + A_R\mathbf{x} + \mathbf{b}$ . This leads to

$$\begin{aligned} J: \quad D\mathbf{x}_{k+1} &= A_L\mathbf{x}_k + A_R\mathbf{x}_k + \mathbf{b}, \quad \text{or} \\ M_J\mathbf{x}_{k+1} &= (A_L + A_R)\mathbf{x}_k + \mathbf{b}, \\ \text{SOR: } D\mathbf{x}_{k+1} &= \omega(A_L\mathbf{x}_{k+1} + A_R\mathbf{x}_k + \mathbf{b}) + (1 - \omega)D\mathbf{x}_k, \quad \text{or} \\ M_\omega\mathbf{x}_{k+1} &= (A_R + (\omega^{-1} - 1)D)\mathbf{x}_k + \mathbf{b}. \end{aligned} \quad (8.14)$$

□

**Example 8.6 (Splitting matrices)**

For the system

$$\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

we find

$$A_L = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}, \quad D = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad A_R = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix},$$

and

$$M_J = D = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad M_\omega = \omega^{-1}D - A_L = \begin{bmatrix} 2\omega^{-1} & 0 \\ -1 & 2\omega^{-1} \end{bmatrix}.$$

The iteration matrix  $G_\omega = I - M_\omega^{-1}A$  is given by

$$G_\omega = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} \omega/2 & 0 \\ \omega^2/4 & \omega/2 \end{bmatrix} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} 1 - \omega & \omega/2 \\ \omega(1 - \omega)/2 & 1 - \omega + \omega^2/4 \end{bmatrix}. \quad (8.15)$$

For the J and GS method we have

$$G_J = I - D^{-1}A = \begin{bmatrix} 0 & 1/2 \\ 1/2 & 0 \end{bmatrix}, \quad G_1 = \begin{bmatrix} 0 & 1/2 \\ 0 & 1/4 \end{bmatrix}. \quad (8.16)$$

We could have derived these matrices directly from the component form of the iteration. For example, for the GS method we have the component form

$$\mathbf{x}_{k+1}(1) = \frac{1}{2}\mathbf{x}_k(2) + \frac{1}{2}, \quad \mathbf{x}_{k+1}(2) = \frac{1}{2}\mathbf{x}_{k+1}(1) + \frac{1}{2}.$$

Substituting the value of  $\mathbf{x}_{k+1}(1)$  from the first equation into the second equation we find

$$\mathbf{x}_{k+1}(2) = \frac{1}{2}\left(\frac{1}{2}\mathbf{x}_k(2) + \frac{1}{2}\right) + \frac{1}{2} = \frac{1}{4}\mathbf{x}_k(2) + \frac{3}{4}.$$

Thus

$$\mathbf{x}_{k+1} = \begin{bmatrix} \mathbf{x}_{k+1}(1) \\ \mathbf{x}_{k+1}(2) \end{bmatrix} = \begin{bmatrix} 0 & 1/2 \\ 0 & 1/4 \end{bmatrix} \begin{bmatrix} \mathbf{x}_k(1) \\ \mathbf{x}_k(2) \end{bmatrix} + \begin{bmatrix} 1/2 \\ 3/4 \end{bmatrix} = \mathbf{G}_1 \mathbf{x}_k + \mathbf{c}.$$

### 8.3 Convergence

**Definition 8.7 (Convergence of  $\mathbf{x}_{k+1} := \mathbf{G}\mathbf{x}_k + \mathbf{c}$ )**

We say that the iterative method  $\mathbf{x}_{k+1} := \mathbf{G}\mathbf{x}_k + \mathbf{c}$  converges if the sequence  $\{\mathbf{x}_k\}$  converges for **any** starting vector  $\mathbf{x}_0$ .

We have the following necessary and sufficient condition for convergence:

**Theorem 8.8 (Convergence of an iterative method)**

The iterative method  $\mathbf{x}_{k+1} := \mathbf{G}\mathbf{x}_k + \mathbf{c}$  converges if and only if  $\lim_{k \rightarrow \infty} \mathbf{G}^k = \mathbf{0}$ .

*Proof.* We subtract  $\mathbf{x} = \mathbf{G}\mathbf{x} + \mathbf{c}$  from  $\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}$ . The vector  $\mathbf{c}$  cancels and we obtain  $\mathbf{x}_{k+1} - \mathbf{x} = \mathbf{G}(\mathbf{x}_k - \mathbf{x})$ . By induction on  $k$

$$\mathbf{x}_k - \mathbf{x} = \mathbf{G}^k(\mathbf{x}_0 - \mathbf{x}), \quad k = 0, 1, 2, \dots \quad (8.17)$$

Clearly  $\mathbf{x}_k - \mathbf{x} \rightarrow \mathbf{0}$  if  $\mathbf{G}^k \rightarrow \mathbf{0}$ . The converse follows by choosing  $\mathbf{x}_0 - \mathbf{x} = \mathbf{e}_j$ , the  $j$ th unit vector for  $j = 1, \dots, n$ .  $\square$

**Theorem 8.9 (Sufficient condition for convergence)**

If  $\|\mathbf{G}\| < 1$  for some consistent matrix norm on  $\mathbb{C}^{n \times n}$ , then the iteration  $\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}$  converges.

*Proof.* We have

$$\|\mathbf{x}_k - \mathbf{x}\| = \|\mathbf{G}^k(\mathbf{x}_0 - \mathbf{x})\| \leq \|\mathbf{G}^k\| \|\mathbf{x}_0 - \mathbf{x}\| \leq \|\mathbf{G}\|^k \|\mathbf{x}_0 - \mathbf{x}\| \rightarrow \mathbf{0}, \quad k \rightarrow \infty.$$

$\square$

A necessary and sufficient condition for convergence involves the eigenvalues of  $\mathbf{G}$ . We define the **spectral radius** of a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  as the maximum absolute value of its eigenvalues.

$$\rho(\mathbf{A}) := \max_{\lambda \in \sigma(\mathbf{A})} |\lambda|. \quad (8.18)$$

**Theorem 8.10 (When does an iterative method converge?)**

Suppose  $\mathbf{G} \in \mathbb{C}^{n \times n}$  and  $\mathbf{c} \in \mathbb{C}^n$ . The iteration  $\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}$  converges if and only if  $\rho(\mathbf{G}) < 1$ .

We will prove this theorem using Theorem 8.27 in Section 8.4.

### 8.3.1 Convergence of Richardson's method.

Recall that Richardson's method can be written in the form

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{r}_k, \quad \mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k. \quad (8.19)$$

We will assume that  $\alpha$  is real. If all eigenvalues of  $\mathbf{A}$  have positive real parts then the R method converges provided  $\alpha$  is positive and sufficiently small. We show this result for positive eigenvalues and leave the more general case to Exercise 8.13.

**Theorem 8.11 (Convergence of Richardson's method)**

If  $\mathbf{A}$  has positive eigenvalues then the R method converges if and only if  $0 < \alpha < 2/\rho(\mathbf{A})$ . Moreover,

$$\min_{\alpha} \mathbf{G}(\alpha) = \mathbf{G}(\alpha^*), \quad \alpha^* := \frac{2}{\lambda_{max} + \lambda_{min}}, \quad \text{and } \rho(\mathbf{G}(\alpha^*)) = \frac{\lambda_{max} - \lambda_{min}}{\lambda_{max} + \lambda_{min}} = \frac{\kappa - 1}{\kappa + 1}, \quad (8.20)$$

where  $\lambda_{max}$  and  $\lambda_{min}$  are the largest and smallest eigenvalue of  $\mathbf{A}$  and  $\kappa := \lambda_{max}/\lambda_{min}$ .

**Proof.** The eigenvalues of  $\mathbf{G}(\alpha) = \mathbf{I} - \alpha\mathbf{A}$  are  $\mu_j(\alpha) = 1 - \alpha\lambda_j$ ,  $j = 1, \dots, n$ . We have  $\max_j \mu_j < 1$  if and only if  $\alpha > 0$  and  $\min_j \mu_j = 1 - \alpha\rho(\mathbf{A}) > -1$  if and only if  $\alpha < 2/\rho(\mathbf{A})$ . The method converges if and only if  $\rho(\mathbf{G}(\alpha)) < 1$  which from what we have shown is equivalent to  $0 < \alpha < 2/\rho(\mathbf{A})$ . Since  $1 - \alpha^*\lambda_{min} = \alpha^*\lambda_{max} - 1$  we have we have

$$\rho(\mathbf{G}(\alpha^*)) = 1 - \alpha^*\lambda_{min} = \frac{\lambda_{max} - \lambda_{min}}{\lambda_{max} + \lambda_{min}} = \frac{\kappa - 1}{\kappa + 1}.$$

Now  $\mathbf{G}(\alpha) \geq \mathbf{G}(\alpha^*)$ , for if  $0 < \alpha < \alpha^*$  then from what we showed  $\rho(\mathbf{G}(\alpha)) \geq 1 - \alpha\lambda_{min} > 1 - \alpha^*\lambda_{min} = \rho(\mathbf{G}(\alpha^*))$ , and if  $\alpha^* < \alpha < 2/\rho(\mathbf{A})$  then  $-\rho(\mathbf{G}(\alpha)) \leq 1 - \alpha\lambda_{max} < 1 - \alpha^*\lambda_{max} = -\rho(\mathbf{G}(\alpha^*))$ , and again  $\rho(\mathbf{G}(\alpha)) > \rho(\mathbf{G}(\alpha^*))$ .  $\square$

**Corollary 8.12 (Rate of convergence for the R method)**

Suppose  $\mathbf{A}$  is symmetric positive definite with largest and smallest eigenvalue  $\lambda_{max}$  and  $\lambda_{min}$ , respectively. Richardson's method with acceleration parameter  $\alpha^* := \frac{2}{\lambda_{max} + \lambda_{min}}$  converges. More precisely

$$\|\mathbf{x}_k - \mathbf{x}\|_2 \leq \left(\frac{\kappa - 1}{\kappa + 1}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2, \quad k = 0, 1, 2, \dots \quad (8.21)$$

where  $\kappa := \lambda_{max}/\lambda_{min}$  is the spectral condition number of  $\mathbf{A}$ .

**Proof.** The spectral norm  $\|\cdot\|_2$  is consistent and therefore  $\|\mathbf{x}_k - \mathbf{x}\|_2 \leq \|\mathbf{G}(\alpha^*)\|_2^k \|\mathbf{x}_0 - \mathbf{x}\|_2$ . But for a symmetric positive definite matrix the spectral norm is equal to the spectral radius and the result follows from (8.20).  $\square$

**Exercise 8.13 (Convergence of the R-method when eigenvalues have positive real parts)**

Suppose all eigenvalues  $\lambda_j$  of  $\mathbf{A}$  have positive real parts  $u_j$  for  $j = 1, \dots, n$  and that  $\alpha$  is real. Show that the R method converges if and only if  $0 < \alpha < \min_j (2u_j/|\lambda_j|^2)$ .

**8.3.2 Convergence of SOR**

The condition  $\omega \in (0, 2)$  is necessary for convergence of the SOR method.

**Theorem 8.14 (Necessary condition for convergence of SOR)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is nonsingular with nonzero diagonal elements. If the SOR method applied to  $\mathbf{A}$  converges then  $\omega \in (0, 2)$ .

**Proof.** We have (cf. (8.14))  $\mathbf{D}\mathbf{x}_{k+1} = \omega(\mathbf{A}_L\mathbf{x}_{k+1} + \mathbf{A}_R\mathbf{x}_k + \mathbf{b}) + (1 - \omega)\mathbf{D}\mathbf{x}_k$  or  $\mathbf{x}_{k+1} = \omega(\mathbf{L}\mathbf{x}_{k+1} + \mathbf{R}\mathbf{x}_k + \mathbf{D}^{-1}\mathbf{b}) + (1 - \omega)\mathbf{x}_k$ , where  $\mathbf{L} := \mathbf{D}^{-1}\mathbf{A}_L$  and  $\mathbf{R} := \mathbf{D}^{-1}\mathbf{A}_R$ . Thus  $(\mathbf{I} - \omega\mathbf{L})\mathbf{x}_{k+1} = (\omega\mathbf{R} + (1 - \omega)\mathbf{I})\mathbf{x}_k + \mathbf{D}^{-1}\mathbf{b}$  so the following form of the iteration matrix is obtained

$$\mathbf{G}_\omega = (\mathbf{I} - \omega\mathbf{L})^{-1}(\omega\mathbf{R} + (1 - \omega)\mathbf{I}). \quad (8.22)$$

We next compute the determinant of  $\mathbf{G}_\omega$ . Since  $\mathbf{I} - \omega\mathbf{L}$  is lower triangular with ones on the diagonal, the same holds for the inverse by Lemma 1.22, and therefore the determinant of this matrix is equal to one. The matrix  $\omega\mathbf{R} + (1 - \omega)\mathbf{I}$  is upper triangular with  $1 - \omega$  on the diagonal and therefore its determinant equals  $(1 - \omega)^n$ . It follows that  $\det(\mathbf{G}_\omega) = (1 - \omega)^n$ . Since the determinant of a matrix equals the product of its eigenvalues we must have  $|\lambda| \geq |1 - \omega|$  for at least one eigenvalue  $\lambda$  of  $\mathbf{G}_\omega$  and we conclude that  $\rho(\mathbf{G}_\omega) \geq |1 - \omega|$ . But then  $\rho(\mathbf{G}_\omega) \geq 1$  if  $\omega$  is not in the interval  $(0, 2)$  and by Theorem 8.10 SOR diverges.  $\square$

The SOR method always converges for a symmetric positive definite matrix.

**Theorem 8.15 (SOR on positive definite matrix)**

*SOR converges for a symmetric positive definite matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  if and only if  $0 < \omega < 2$ . In particular, Gauss-Seidel's method converges for a symmetric positive definite matrix.*

**Proof.** By Theorem 8.14 convergence implies  $0 < \omega < 2$ . Suppose now  $0 < \omega < 2$  and let  $(\lambda, \mathbf{x})$  be an eigenpair for  $\mathbf{G}_\omega$ . Note that  $\lambda$  and  $\mathbf{x}$  can be complex. We need to show that  $|\lambda| < 1$ . The following identity will be shown below:

$$\omega^{-1}(2 - \omega)|1 - \lambda|^2 \mathbf{x}^* \mathbf{D} \mathbf{x} = (1 - |\lambda|^2) \mathbf{x}^* \mathbf{A} \mathbf{x}, \quad (8.23)$$

where  $\mathbf{D} := \text{diag}(a_{11}, \dots, a_{nn})$ . Now  $\mathbf{x}^* \mathbf{A} \mathbf{x}$  and  $\mathbf{x}^* \mathbf{D} \mathbf{x}$  are positive for all nonzero  $\mathbf{x} \in \mathbb{C}^n$  since a positive definite matrix has positive diagonal elements  $a_{ii} = \mathbf{e}_i^T \mathbf{A} \mathbf{e}_i > 0$ . It follows that the left hand side of (8.23) is nonnegative and then the right hand side must be nonnegative as well. This implies  $|\lambda| \leq 1$ . It remains to show that we cannot have  $\lambda = 1$ . By (8.13) the eigenpair equation  $\mathbf{G}_\omega \mathbf{x} = \lambda \mathbf{x}$  can be written  $\mathbf{x} - (\omega^{-1} \mathbf{D} - \mathbf{A}_L)^{-1} \mathbf{A} \mathbf{x} = \lambda \mathbf{x}$  or

$$\mathbf{A} \mathbf{x} = (\omega^{-1} \mathbf{D} - \mathbf{A}_L) \mathbf{y}, \quad \mathbf{y} := (1 - \lambda) \mathbf{x}. \quad (8.24)$$

Now  $\mathbf{A} \mathbf{x} \neq \mathbf{0}$  implies that  $\lambda \neq 1$ .

To prove equation (8.23) consider the matrix  $\mathbf{E} := \omega^{-1} \mathbf{D} + \mathbf{A}_R - \mathbf{D}$ . Since  $\mathbf{A}_R - \mathbf{D} = -\mathbf{A}_L - \mathbf{A}$  we find  $\mathbf{E} \mathbf{y} = (\omega^{-1} \mathbf{D} - \mathbf{A}_L - \mathbf{A}) \mathbf{y} \stackrel{(8.24)}{=} \mathbf{A} \mathbf{x} - \mathbf{A} \mathbf{y} = \lambda \mathbf{A} \mathbf{x}$ . Observe that  $(\omega^{-1} \mathbf{D} - \mathbf{A}_L)^* = \omega^{-1} \mathbf{D} - \mathbf{A}_R$  so that by (8.24)

$$\begin{aligned} (\mathbf{A} \mathbf{x})^* \mathbf{y} + \mathbf{y}^* (\lambda \mathbf{A} \mathbf{x}) &= \mathbf{y}^* (\omega^{-1} \mathbf{D} - \mathbf{A}_R) \mathbf{y} + \mathbf{y}^* \mathbf{E} \mathbf{y} = \mathbf{y}^* (2\omega^{-1} - 1) \mathbf{D} \mathbf{y} \\ &= \omega^{-1} (2 - \omega) |1 - \lambda|^2 \mathbf{x}^* \mathbf{D} \mathbf{x}. \end{aligned}$$

Since  $(\mathbf{A} \mathbf{x})^* = \mathbf{x}^* \mathbf{A}$ ,  $\mathbf{y} := (1 - \lambda) \mathbf{x}$  and  $\mathbf{y}^* = (1 - \bar{\lambda}) \mathbf{x}^*$  this also equals

$$(\mathbf{A} \mathbf{x})^* \mathbf{y} + \mathbf{y}^* (\lambda \mathbf{A} \mathbf{x}) = (1 - \lambda) \mathbf{x}^* \mathbf{A} \mathbf{x} + \lambda (1 - \bar{\lambda}) \mathbf{x}^* \mathbf{A} \mathbf{x} = (1 - |\lambda|^2) \mathbf{x}^* \mathbf{A} \mathbf{x},$$

and (8.23) follows.  $\square$

**Exercise 8.16 (Example: GS converges, J diverges)**

Show (by finding its eigenvalues) that the matrix

$$\begin{bmatrix} 1 & a & a \\ a & 1 & a \\ a & a & 1 \end{bmatrix}$$

is symmetric positive definite for  $-1/2 < a < 1$ . Thus, GS converges for these values of  $a$ . Show that the J method does not converge for  $1/2 < a < 1$ .

**Exercise 8.17 (Divergence example for J and GS)**

Show that both Jacobi's method and Gauss-Seidel's method diverge for  $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ .

**Exercise 8.18 (Strictly diagonally dominance; The J method)**

Show that the J method converges if  $|a_{ii}| > \sum_{j \neq i} |a_{ij}|$  for  $i = 1, \dots, n$ .

**Exercise 8.19 (Strictly diagonally dominance; The GS method)**

Consider the GS method. Suppose  $r := \max_i r_i < 1$ , where  $r_i = \sum_{j \neq i} \frac{|a_{ij}|}{|a_{ii}|}$ . Show using induction on  $i$  that  $|\epsilon_{k+1}(j)| \leq r \|\epsilon_k\|_\infty$  for  $j = 1, \dots, i$ . Conclude that Gauss-Seidel's method is convergent when  $\mathbf{A}$  is strictly diagonally dominant.

### 8.3.3 Convergence of the classical methods for the discrete Poisson matrix

We know the eigenvalues of the discrete Poisson matrix  $\mathbf{A}$  given by (3.7) and we can use this to estimate the number of iterations necessary to achieve a given accuracy for the various methods.

Recall that by (3.20) the eigenvalues  $\lambda_{j,k}$  of  $\mathbf{A}$  are

$$\lambda_{j,k} = 4 - 2 \cos(j\pi h) - 2 \cos(k\pi h), \quad j, k = 1, \dots, m, h = 1/(m+1).$$

It follows that the largest and smallest eigenvalue of  $\mathbf{A}$ , and the spectral condition number  $\kappa$  of  $\mathbf{A}$ , are given by

$$\lambda_{max} = 8 \cos^2 w, \quad \lambda_{min} = 8 \sin^2 w, \quad \kappa := \frac{\cos^2 w}{\sin^2 w}, \quad w := \frac{\pi}{2(m+1)}. \quad (8.25)$$

Consider first the J method. The matrix  $\mathbf{G}_J = \mathbf{I} - \mathbf{D}^{-1}\mathbf{A} = \mathbf{I} - \mathbf{A}/4$  has eigenvalues

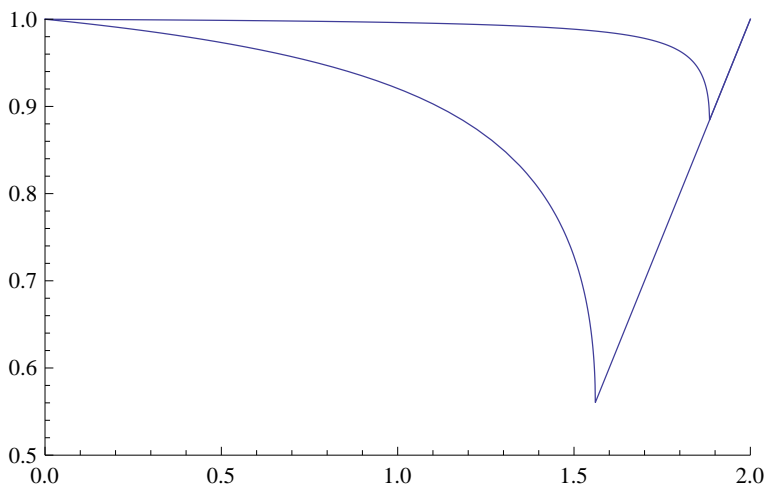
$$\mu_{j,k} = 1 - \frac{1}{4}\lambda_{j,k} = \frac{1}{2} \cos(j\pi h) + \frac{1}{2} \cos(k\pi h), \quad j, k = 1, \dots, m. \quad (8.26)$$

It follows that  $\rho(\mathbf{G}_J) = \cos(\pi h) < 1$ . Since  $\mathbf{G}_J$  is symmetric it is normal, and the spectral norm is equal to the spectral radius (cf. Theorem 7.14). We obtain

$$\|\mathbf{x}_k - \mathbf{x}\|_2 \leq \|\mathbf{G}_J\|_2^k \|\mathbf{x}_0 - \mathbf{x}\|_2 = \cos^k(\pi h) \|\mathbf{x}_0 - \mathbf{x}\|_2, \quad k = 0, 1, 2, \dots \quad (8.27)$$

The R method given by  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{r}_k$  with  $\alpha = 2/(\lambda_{max} + \lambda_{min}) = 1/4$  is the same as the J-method so (8.27) holds in this case as well. This also follows from Corollary 8.12 with  $\kappa$  given by (8.25).

For the SOR method it is possible to explicitly determine  $\rho(\mathbf{G}_\omega)$  for any  $\omega \in (0, 2)$ . The following result will be shown in Section 8.3.2.



**Figure 8.1.**  $\rho(\mathbf{G}_\omega)$  with  $\omega \in [0, 2]$  for  $n = 100$ , (lower curve) and  $n = 2500$  (upper curve).

**Theorem 8.20 (The spectral radius of SOR matrix)**

Consider the SOR iteration (8.8), with the natural ordering. The spectral radius of  $\mathbf{G}_\omega$  is

$$\rho(\mathbf{G}_\omega) = \begin{cases} \frac{1}{4} \left( \omega\beta + \sqrt{(\omega\beta)^2 - 4(\omega - 1)} \right)^2, & \text{for } 0 < \omega \leq \omega^*, \\ \omega - 1, & \text{for } \omega^* < \omega < 2, \end{cases} \quad (8.28)$$

where  $\beta := \rho(\mathbf{G}_J) = \cos(\pi h)$  and

$$\omega^* := \frac{2}{1 + \sqrt{1 - \beta^2}} > 1. \quad (8.29)$$

Moreover,

$$\rho(\mathbf{G}_\omega) > \rho(\mathbf{G}_{\omega^*}) \text{ for } \omega \in (0, 2) \setminus \{\omega^*\}. \quad (8.30)$$

A plot of  $\rho(\mathbf{G}_\omega)$  as a function of  $\omega \in (0, 2)$  is shown in Figure 8.1 for  $n = 100$  (lower curve) and  $n = 2500$  (upper curve). As  $\omega$  increases the spectral radius of  $\mathbf{G}_\omega$  decreases monotonically to the minimum  $\omega^*$ . Then it increases linearly to the value one for  $\omega = 2$ . We call  $\omega^*$  the **optimal relaxation parameter**.

For the discrete Poisson problem we have  $\beta = \cos(\pi h)$  and it follows from (8.28),(8.29) that

$$\omega^* = \frac{2}{1 + \sin(\pi h)}, \quad \rho(\mathbf{G}_{\omega^*}) = \omega^* - 1 = \frac{1 - \sin(\pi h)}{1 + \sin(\pi h)}, \quad h = \frac{1}{m + 1}. \quad (8.31)$$

|     | n=100    | n=2500   | $k_{100}$ | $k_{2500}$ |
|-----|----------|----------|-----------|------------|
| J   | 0.959493 | 0.998103 | 446       | 9703       |
| GS  | 0.920627 | 0.99621  | 223       | 4852       |
| SOR | 0.56039  | 0.88402  | 32        | 150        |

**Table 8.2.** Spectral radii for  $\mathbf{G}_J$ ,  $\mathbf{G}_1$ ,  $\mathbf{G}_{\omega^*}$  and the smallest integer  $k_n$  such that  $\rho(\mathbf{G})^{k_n} \leq 10^{-8}$ .

Letting  $\omega = 1$  in (8.28) we find  $\rho(\mathbf{G}_1) = \beta^2 = \rho(\mathbf{G}_J)^2 = \cos^2(\pi h)$  for the GS method. Thus, for the discrete Poisson problem the J method needs twice as many iterations as the GS method for a given accuracy.

The values of  $\rho(\mathbf{G}_J)$ ,  $\rho(\mathbf{G}_1)$ , and  $\rho(\mathbf{G}_{\omega^*}) = \omega^* - 1$  are shown in Table 8.2 for  $n = 100$  and  $n = 2500$ . We also show the smallest integer  $k_n$  such that  $\rho(\mathbf{G})^{k_n} \leq 10^{-8}$ . This is an estimate for the number of iteration needed to obtain an accuracy of  $10^{-8}$ . These values are comparable to the exact values given in Table 8.1.

### 8.3.4 Number of iterations

Consider next the **rate of convergence** of the iteration  $\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}$ . We like to know how fast the iterative method converges. Suppose  $\|\cdot\|$  is a matrix norm that is subordinate to a vector norm also denoted by  $\|\cdot\|$ . Recall that  $\mathbf{x}_k - \mathbf{x} = \mathbf{G}^k(\mathbf{x}_0 - \mathbf{x})$ . For  $k$  sufficiently large

$$\|\mathbf{x}_k - \mathbf{x}\| \leq \|\mathbf{G}^k\| \|\mathbf{x}_0 - \mathbf{x}\| \approx \rho(\mathbf{G})^k \|\mathbf{x}_0 - \mathbf{x}\|.$$

For the last formula we apply Theorem 8.30 which says that  $\lim_{k \rightarrow \infty} \|\mathbf{G}^k\|^{1/k} = \rho(\mathbf{G})$ . For Jacobi's method and the spectral norm we have  $\|\mathbf{G}_J^k\|_2 = \rho(\mathbf{G}_J)^k$  (cf. (8.27)).

For fast convergence we should use a  $\mathbf{G}$  with small spectral radius.

**Lemma 8.22 (Number of iterations)**

Suppose  $\rho(\mathbf{G}) = 1 - \eta$  for some  $0 < \eta < 1$ ,  $\|\cdot\|$  a consistent matrix norm on  $\mathbb{C}^{n \times n}$ , and let  $s \in \mathbb{N}$ . Then

$$\tilde{k} := \frac{s \log(10)}{\eta} \tag{8.32}$$

is an estimate for the smallest number of iterations  $k$  so that  $\rho(\mathbf{G})^k \leq 10^{-s}$ .

**Proof.** The estimate  $\tilde{k}$  is an approximate solution of the equation  $\rho(\mathbf{G})^k = 10^{-s}$ . Thus, since  $-\log(1 - \eta) \approx \eta$  when  $\eta$  is small

$$k = -\frac{s \log(10)}{\log(1 - \eta)} \approx \frac{s \log(10)}{\eta} = \tilde{k}.$$



□

The following estimates are obtained. They agree with those we found numerically in Section 8.1.1.

- R and J:  $\rho(\mathbf{G}_J) = \cos(\pi h) = 1 - \eta$ ,  $\eta = 1 - \cos(\pi h) = \frac{1}{2}\pi^2 h^2 + O(h^4) = \frac{\pi^2}{2}/n + O(n^{-2})$ . Thus

$$\tilde{k}_n = \frac{2\log(10)s}{\pi^2}n + O(n^{-1}) = O(n).$$

- GS:  $\rho(\mathbf{G}_1) = \cos^2(\pi h) = 1 - \eta$ ,  $\eta = 1 - \cos^2(\pi h) = \sin^2 \pi h = \pi^2 h^2 + O(h^4) = \pi^2/n + O(n^{-2})$ . Thus

$$\tilde{k}_n = \frac{\log(10)s}{\pi^2}n + O(n^{-1}) = O(n).$$

- SOR:  $\rho(\mathbf{G}_{\omega^*}) = \frac{1 - \sin(\pi h)}{1 + \sin(\pi h)} = 1 - 2\pi h + O(h^2)$ . Thus,

$$\tilde{k}_n = \frac{\log(10)s}{2\pi}\sqrt{n} + O(n^{-1/2}) = O(\sqrt{n}).$$

### Exercise 8.23 (Convergence example for fix point iteration)

Consider for  $a \in \mathbb{C}$

$$\mathbf{x} := \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 & a \\ a & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 1 - a \\ 1 - a \end{bmatrix} =: \mathbf{G}\mathbf{x} + \mathbf{c}.$$

Starting with  $\mathbf{x}_0 = \mathbf{0}$  show by induction

$$\mathbf{x}_k(1) = \mathbf{x}_k(2) = 1 - a^k, \quad k \geq 0,$$

and conclude that the iteration converges to the fixed-point  $\mathbf{x} = [1, 1]^T$  for  $|a| < 1$  and diverges for  $|a| > 1$ . Show that  $\rho(\mathbf{G}) = 1 - \eta$  with  $\eta = 1 - |a|$ . Compute the estimate (8.32) for the rate of convergence for  $a = 0.9$  and  $s = 16$  and compare with the true number of iterations determined from  $|a|^k \leq 10^{-16}$ .

### Exercise 8.24 (Estimate in Lemma 8.22 can be exact)

Consider the iteration in Example 8.6. Show that  $\rho(\mathbf{G}_J) = 1/2$ . Then show that  $\mathbf{x}_k(1) = \mathbf{x}_k(2) = 1 - 2^{-k}$  for  $k \geq 0$ . Thus the estimate in Lemma 8.22 is exact in this case.

We note that

1. The convergence depends on the behavior of the powers  $\mathbf{G}^k$  as  $k$  increases. The matrix  $\mathbf{M}$  should be chosen so that all elements in  $\mathbf{G}^k$  converge quickly to zero and such that the linear system (8.11) is easy to solve for  $\mathbf{x}_{k+1}$ . These are conflicting demands.  $\mathbf{M}$  should be an approximation to  $\mathbf{A}$  to obtain a  $\mathbf{G}$  with small elements, but then (8.11) might not be easy to solve for  $\mathbf{x}_{k+1}$ .
2. The convergence  $\lim_{k \rightarrow \infty} \|\mathbf{G}^k\|^{1/k} = \rho(\mathbf{G})$  can be quite slow (cf. Exercise 8.25).

**Exercise 8.25 (Slow spectral radius convergence)**

The convergence  $\lim_{k \rightarrow \infty} \|\mathbf{A}^k\|^{1/k} = \rho(\mathbf{A})$  can be quite slow. Consider

$$\mathbf{A} := \begin{bmatrix} \lambda & a & 0 & \cdots & 0 & 0 \\ 0 & \lambda & a & \cdots & 0 & 0 \\ 0 & 0 & \lambda & \cdots & 0 & 0 \\ \vdots & & & & \vdots & \\ 0 & 0 & 0 & \cdots & \lambda & a \\ 0 & 0 & 0 & \cdots & 0 & \lambda \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

If  $|\lambda| = \rho(\mathbf{A}) < 1$  then  $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0}$  for any  $a \in \mathbb{R}$ . We show below that the  $(1, n)$  element of  $\mathbf{A}^k$  is given by  $f(k) := \binom{k}{n-1} a^{n-1} \lambda^{k-n+1}$  for  $k \geq n-1$ .

- (a) Pick an  $n$ , e.g.  $n = 5$ , and make a plot of  $f(k)$  for  $\lambda = 0.9$ ,  $a = 10$ , and  $n-1 \leq k \leq 200$ . Your program should also compute  $\max_k f(k)$ . Use your program to determine how large  $k$  must be before  $f(k) < 10^{-8}$ .
- (b) We can determine the elements of  $\mathbf{A}^k$  explicitly for any  $k$ . Let  $\mathbf{E} := (\mathbf{A} - \lambda \mathbf{I})/a$ . Show by induction that  $\mathbf{E}^k = \begin{bmatrix} \mathbf{0} & \mathbf{I}_{n-k} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$  for  $1 \leq k \leq n-1$  and that  $\mathbf{E}^n = \mathbf{0}$ .
- (c) We have  $\mathbf{A}^k = (a\mathbf{E} + \lambda\mathbf{I})^k = \sum_{j=0}^{\min\{k, n-1\}} \binom{k}{j} a^j \lambda^{k-j} \mathbf{E}^j$  and conclude that the  $(1, n)$  element is given by  $f(k)$  for  $k \geq n-1$ .

### 8.3.5 Stopping the iteration

In Algorithms 8.3 and 8.4 we had access to the exact solution and could stop the iteration when the error was sufficiently small in the infinity norm. The decision when to stop is obviously more complicated when the exact solution is not known. One possibility is to choose a vector norm, keep track of  $\|\mathbf{x}_{k+1} - \mathbf{x}_k\|$ , and stop when this number is sufficiently small. The following result indicates that  $\|\mathbf{x}_k - \mathbf{x}\|$  can be quite large if  $\|\mathbf{G}\|$  is close to one.

**Lemma 8.26 (Be careful when stopping)**

Suppose  $\|\mathbf{G}\| < 1$  for some consistent matrix norm on  $\mathbb{C}^{n \times n}$  which is subordinate

to a vector norm also denoted by  $\|\cdot\|$ . If  $\mathbf{x}_k = \mathbf{G}\mathbf{x}_{k-1} + \mathbf{c}$  and  $\mathbf{x} = \mathbf{G}\mathbf{x} + \mathbf{c}$ . Then

$$\|\mathbf{x}_k - \mathbf{x}_{k-1}\| \geq \frac{1 - \|\mathbf{G}\|}{\|\mathbf{G}\|} \|\mathbf{x}_k - \mathbf{x}\|, \quad k \geq 1. \quad (8.33)$$

*Proof.* We find

$$\begin{aligned} \|\mathbf{x}_k - \mathbf{x}\| &= \|\mathbf{G}(\mathbf{x}_{k-1} - \mathbf{x})\| \leq \|\mathbf{G}\| \|\mathbf{x}_{k-1} - \mathbf{x}\| \\ &= \|\mathbf{G}\| \|\mathbf{x}_{k-1} - \mathbf{x}_k + \mathbf{x}_k - \mathbf{x}\| \leq \|\mathbf{G}\| (\|\mathbf{x}_{k-1} - \mathbf{x}_k\| + \|\mathbf{x}_k - \mathbf{x}\|). \end{aligned}$$

Thus  $(1 - \|\mathbf{G}\|)\|\mathbf{x}_k - \mathbf{x}\| \leq \|\mathbf{G}\| \|\mathbf{x}_{k-1} - \mathbf{x}_k\|$  which implies (8.33).  $\square$

Another possibility is to stop when the residual vector  $\mathbf{r}_k := \mathbf{b} - \mathbf{A}\mathbf{x}_k$  is sufficiently small in some norm. To use the residual vector for stopping it is convenient to write the iterative method (8.10) in an alternative form. If  $\mathbf{M}$  is the splitting matrix of the method then by (8.11) we have  $\mathbf{M}\mathbf{x}_{k+1} = \mathbf{M}\mathbf{x}_k - \mathbf{A}\mathbf{x}_k + \mathbf{b}$ . This leads to

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{M}^{-1}\mathbf{r}_k, \quad \mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k. \quad (8.34)$$

Testing on  $\mathbf{r}_k$  works fine if  $\mathbf{A}$  is well conditioned, but Theorem 7.28 shows that the relative error in the solution can be much larger than the relative error in  $\mathbf{r}_k$  if  $\mathbf{A}$  is ill-conditioned.

## 8.4 Powers of a matrix

Let  $\mathbf{A} \in \mathbb{C}^{n \times n}$  be a square matrix. In this section we consider the special matrix sequence  $\{\mathbf{A}^k\}$  of powers of  $\mathbf{A}$ . We want to know when this sequence converges to the zero matrix. Such a sequence occurs in iterative methods (cf. (8.17)), in Markov processes in statistics, in the converge of geometric series of matrices (Neumann series cf. Section 8.4.2) and in many other applications.

### 8.4.1 The spectral radius

In this section we show the following theorem.

**Theorem 8.27 (When is  $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0}$ ?)**

For any  $\mathbf{A} \in \mathbb{C}^{n \times n}$  we have

$$\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0} \iff \rho(\mathbf{A}) < 1,$$

where  $\rho(\mathbf{A})$  is the spectral radius of  $\mathbf{A}$  given by (8.18).

Clearly  $\rho(\mathbf{A}) < 1$  is a necessary condition for  $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0}$ . For if  $(\lambda, \mathbf{x})$  is an eigenpair of  $\mathbf{A}$  with  $|\lambda| \geq 1$  and  $\|\mathbf{x}\|_2 = 1$  then  $\mathbf{A}^k \mathbf{x} = \lambda^k \mathbf{x}$ , and this implies  $\|\mathbf{A}^k\|_2 \geq \|\mathbf{A}^k \mathbf{x}\|_2 = \|\lambda^k \mathbf{x}\|_2 = |\lambda|^k$ , and it follows that  $\mathbf{A}^k$  does not tend to zero.

The sufficiency condition is harder to show. We construct a consistent matrix norm on  $\mathbb{C}^{n \times n}$  such that  $\|\mathbf{A}\| < 1$  and then use Theorems 8.8 and 8.9.

We start with

**Theorem 8.28 (Any consistent norm majorizes the spectral radius)**

For any matrix norm  $\|\cdot\|$  that is consistent on  $\mathbb{C}^{n \times n}$  and any  $\mathbf{A} \in \mathbb{C}^{n \times n}$  we have  $\rho(\mathbf{A}) \leq \|\mathbf{A}\|$ .

*Proof.* Let  $(\lambda, \mathbf{x})$  be an eigenpair for  $\mathbf{A}$  and define  $\mathbf{X} := [\mathbf{x}, \dots, \mathbf{x}] \in \mathbb{C}^{n \times n}$ . Then  $\lambda\mathbf{X} = \mathbf{A}\mathbf{X}$ , which implies  $|\lambda| \|\mathbf{X}\| = \|\lambda\mathbf{X}\| = \|\mathbf{A}\mathbf{X}\| \leq \|\mathbf{A}\| \|\mathbf{X}\|$ . Since  $\|\mathbf{X}\| \neq 0$  we obtain  $|\lambda| \leq \|\mathbf{A}\|$ .  $\square$

The next theorem shows that if  $\rho(\mathbf{A}) < 1$  then  $\|\mathbf{A}\| < 1$  for some consistent matrix norm on  $\mathbb{C}^{n \times n}$ , thus completing the proof of Theorem 8.27.

**Theorem 8.29 (The spectral radius can be approximated by a norm)**

Let  $\mathbf{A} \in \mathbb{C}^{n \times n}$  and  $\epsilon > 0$  be given. There is a consistent matrix norm  $\|\cdot\|$  on  $\mathbb{C}^{n \times n}$  such that  $\rho(\mathbf{A}) \leq \|\mathbf{A}\| \leq \rho(\mathbf{A}) + \epsilon$ .

*Proof.* Let  $\mathbf{A}$  have eigenvalues  $\lambda_1, \dots, \lambda_n$ . By the Schur Triangulation Theorem 5.13 there is a unitary matrix  $\mathbf{U}$  and an upper triangular matrix  $\mathbf{R} = [r_{ij}]$  such that  $\mathbf{U}^* \mathbf{A} \mathbf{U} = \mathbf{R}$ . For  $t > 0$  we define  $\mathbf{D}_t := \text{diag}(t, t^2, \dots, t^n) \in \mathbb{R}^{n \times n}$ , and note that the  $(i, j)$  element in  $\mathbf{D}_t \mathbf{R} \mathbf{D}_t^{-1}$  is given by  $t^{i-j} r_{ij}$  for all  $i, j$ . For  $n = 3$

$$\mathbf{D}_t \mathbf{R} \mathbf{D}_t^{-1} = \begin{bmatrix} \lambda_1 & t^{-1} r_{12} & t^{-2} r_{13} \\ 0 & \lambda_2 & t^{-1} r_{23} \\ 0 & 0 & \lambda_3 \end{bmatrix}.$$

For each  $\mathbf{B} \in \mathbb{C}^{n \times n}$  and  $t > 0$  we use the one norm to define the matrix norm  $\|\mathbf{B}\|_t := \|\mathbf{D}_t \mathbf{U}^* \mathbf{B} \mathbf{U} \mathbf{D}_t^{-1}\|_1$ . We leave it as an exercise to show that  $\|\cdot\|_t$  is a consistent matrix norm on  $\mathbb{C}^{n \times n}$ . We define  $\|\mathbf{B}\| := \|\mathbf{B}\|_t$ , where  $t$  is chosen so large that the sum of the absolute values of all off-diagonal elements in  $\mathbf{D}_t \mathbf{R} \mathbf{D}_t^{-1}$  is less than  $\epsilon$ . Then

$$\begin{aligned} \|\mathbf{A}\| &= \|\mathbf{D}_t \mathbf{U}^* \mathbf{A} \mathbf{U} \mathbf{D}_t^{-1}\|_1 = \|\mathbf{D}_t \mathbf{R} \mathbf{D}_t^{-1}\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |(\mathbf{D}_t \mathbf{R} \mathbf{D}_t^{-1})_{ij}| \\ &\leq \max_{1 \leq j \leq n} (|\lambda_j| + \epsilon) = \rho(\mathbf{A}) + \epsilon. \end{aligned}$$

$\square$

A consistent matrix norm of a matrix can be much larger than the spectral radius. However the following result holds.

**Theorem 8.30 (Spectral radius convergence)**

For any consistent matrix norm  $\|\cdot\|$  on  $\mathbb{C}^{n \times n}$  and any  $\mathbf{A} \in \mathbb{C}^{n \times n}$  we have

$$\lim_{k \rightarrow \infty} \|\mathbf{A}^k\|^{1/k} = \rho(\mathbf{A}). \quad (8.35)$$

**Proof.** By Theorems 5.1 and 8.28 we obtain  $\rho(\mathbf{A})^k = \rho(\mathbf{A}^k) \leq \|\mathbf{A}^k\|$  for any  $k \in \mathbb{N}$  so that  $\rho(\mathbf{A}) \leq \|\mathbf{A}^k\|^{1/k}$ . Let  $\epsilon > 0$  and consider the matrix  $\mathbf{B} := (\rho(\mathbf{A}) + \epsilon)^{-1} \mathbf{A}$ . Then  $\rho(\mathbf{B}) = \rho(\mathbf{A}) / (\rho(\mathbf{A}) + \epsilon) < 1$  and  $\|\mathbf{B}^k\| \rightarrow 0$  by Theorem 8.27 as  $k \rightarrow \infty$ . Choose  $N \in \mathbb{N}$  such that  $\|\mathbf{B}^k\| < 1$  for all  $k \geq N$ . Then for  $k \geq N$

$$\|\mathbf{A}^k\| = \|(\rho(\mathbf{A}) + \epsilon)^k \mathbf{B}^k\| = (\rho(\mathbf{A}) + \epsilon)^k \|\mathbf{B}^k\| < (\rho(\mathbf{A}) + \epsilon)^k.$$

We have shown that  $\rho(\mathbf{A}) \leq \|\mathbf{A}^k\|^{1/k} \leq \rho(\mathbf{A}) + \epsilon$  for  $k \geq N$ . Since  $\epsilon$  is arbitrary the result follows.  $\square$

**Exercise 8.31 (A special norm)**

Show that  $\|\mathbf{B}\|_t := \|\mathbf{D}_t \mathbf{U}^* \mathbf{B} \mathbf{U} \mathbf{D}_t^{-1}\|_1$  defined in the proof of Theorem 8.29 is a consistent matrix norm on  $\mathbb{C}^{n \times n}$ .

**8.4.2 Neumann series**

Carl Neumann., 1832–1925. He studied potential theory. The Neumann boundary conditions are named after him.

Let  $\mathbf{B}$  be a square matrix. In this section we consider the **Neumann series**  $\sum_{k=0}^{\infty} \mathbf{B}^k$  which is a matrix analogue of a geometric series of numbers.

Consider an infinite series  $\sum_{k=0}^{\infty} \mathbf{A}_k$  of matrices in  $\mathbb{C}^{n \times n}$ . We say that the series converges if the sequence of partial sums  $\{\mathbf{S}_m\}$  given by  $\mathbf{S}_m = \sum_{k=0}^m \mathbf{A}_k$  converges. The series converges if and only if  $\{\mathbf{S}_m\}$  is a Cauchy sequence, i.e. to each  $\epsilon > 0$  there exists an integer  $N$  so that  $\|\mathbf{S}_l - \mathbf{S}_m\| < \epsilon$  for all  $l > m \geq N$ .

**Theorem 8.32 (Neumann series)**

Suppose  $\mathbf{B} \in \mathbb{C}^{n \times n}$ . Then

1. The series  $\sum_{k=0}^{\infty} \mathbf{B}^k$  converges if and only if  $\rho(\mathbf{B}) < 1$ .

2. If  $\rho(\mathbf{B}) < 1$  then  $(\mathbf{I} - \mathbf{B})$  is nonsingular and  $(\mathbf{I} - \mathbf{B})^{-1} = \sum_{k=0}^{\infty} \mathbf{B}^k$ .
3. If  $\|\mathbf{B}\| < 1$  for some consistent matrix norm  $\|\cdot\|$  on  $\mathbb{C}^{n \times n}$  then

$$\|(\mathbf{I} - \mathbf{B})^{-1}\| \leq \frac{1}{1 - \|\mathbf{B}\|}. \quad (8.36)$$

**Proof.**

1. Suppose  $\rho(\mathbf{B}) < 1$ . We show that  $\mathbf{S}_m := \sum_{k=0}^m \mathbf{B}^k$  is a Cauchy sequence and hence convergent. Let  $\epsilon > 0$ . By Theorem 8.29 there is a consistent matrix norm  $\|\cdot\|$  on  $\mathbb{C}^{n \times n}$  such that  $\|\mathbf{B}\| < 1$ . Then for  $l > m$

$$\|\mathbf{S}_l - \mathbf{S}_m\| = \left\| \sum_{k=m+1}^l \mathbf{B}^k \right\| \leq \sum_{k=m+1}^l \|\mathbf{B}\|^k \leq \|\mathbf{B}\|^{m+1} \sum_{k=0}^{\infty} \|\mathbf{B}\|^k = \frac{\|\mathbf{B}\|^{m+1}}{1 - \|\mathbf{B}\|}.$$

But then  $\{\mathbf{S}_m\}$  is a Cauchy sequence provided  $N$  is such that  $\frac{\|\mathbf{B}\|^{N+1}}{1 - \|\mathbf{B}\|} < \epsilon$ .

Conversely, suppose  $(\lambda, \mathbf{x})$  is an eigenpair for  $\mathbf{B}$  with  $|\lambda| \geq 1$ . We find  $\mathbf{S}_m \mathbf{x} = \sum_{k=0}^m \mathbf{B}^k \mathbf{x} = (\sum_{k=0}^m \lambda^k) \mathbf{x}$ . Since  $\lambda^k$  does not tend to zero the series  $\sum_{k=0}^{\infty} \lambda^k$  is not convergent and therefore  $\{\mathbf{S}_m \mathbf{x}\}$  and hence  $\{\mathbf{S}_m\}$  does not converge.

2. We have

$$\left( \sum_{k=0}^m \mathbf{B}^k \right) (\mathbf{I} - \mathbf{B}) = \mathbf{I} + \mathbf{B} + \cdots + \mathbf{B}^m - (\mathbf{B} + \cdots + \mathbf{B}^{m+1}) = \mathbf{I} - \mathbf{B}^{m+1}. \quad (8.37)$$

Since  $\rho(\mathbf{B}) < 1$  we conclude that  $\mathbf{B}^{m+1} \rightarrow 0$  and hence taking limits in (8.37) we obtain  $(\sum_{k=0}^{\infty} \mathbf{B}^k)(\mathbf{I} - \mathbf{B}) = \mathbf{I}$  which completes the proof of 2.

3. By 2:  $\|(\mathbf{I} - \mathbf{B})^{-1}\| = \|\sum_{k=0}^{\infty} \mathbf{B}^k\| \leq \sum_{k=0}^{\infty} \|\mathbf{B}\|^k = \frac{1}{1 - \|\mathbf{B}\|}$ .

□

### Exercise 8.33 (When is $\mathbf{A} + \mathbf{E}$ nonsingular?)

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is nonsingular and  $\mathbf{E} \in \mathbb{C}^{n \times n}$ . Show that  $\mathbf{A} + \mathbf{E}$  is nonsingular if and only if  $\rho(\mathbf{A}^{-1} \mathbf{E}) < 1$ .

## 8.5 The Optimal SOR Parameter $\omega$

The following analysis is only carried out for the discrete Poisson matrix. It also holds for the averaging matrix given by (3.9). A more general theory is presented

in [35]. We will compare the eigenpair equations for  $\mathbf{G}_J$  and  $\mathbf{G}_\omega$ . It is convenient to write these equations using the matrix formulation  $\mathbf{T}\mathbf{V} + \mathbf{V}\mathbf{T} = h^2\mathbf{F}$ . If  $\mathbf{G}_J\mathbf{v} = \mu\mathbf{v}$  is an eigenpair of  $\mathbf{G}_J$  then

$$\frac{1}{4}(v_{i-1,j} + v_{i,j-1} + v_{i+1,j} + v_{i,j+1}) = \mu v_{i,j}, \quad i, j = 1, \dots, m, \quad (8.38)$$

where  $\mathbf{V} := \text{vec}(\mathbf{v}) \in \mathbb{R}^{m \times m}$  and  $v_{i,j} = 0$  if  $i \in \{0, m+1\}$  or  $j \in \{0, m+1\}$ .

Suppose  $(\lambda, \mathbf{w})$  is an eigenpair for  $\mathbf{G}_\omega$ . By (8.22)  $(\mathbf{I} - \omega\mathbf{L})^{-1}(\omega\mathbf{R} + (1 - \omega)\mathbf{I})\mathbf{w} = \lambda\mathbf{w}$  or

$$(\omega\mathbf{R} + \lambda\omega\mathbf{L})\mathbf{w} = (\lambda + \omega - 1)\mathbf{w}, \quad (8.39)$$

where  $l_{i,i-m} = r_{i,i+m} = 1/4$  for all  $i$ , and all other elements in  $\mathbf{L}$  and  $\mathbf{R}$  are equal to zero. Let  $\mathbf{w} = \text{vec}(\mathbf{W})$ , where  $\mathbf{W} \in \mathbb{C}^{m \times m}$ . Then (8.39) can be written

$$\frac{\omega}{4}(\lambda w_{i-1,j} + \lambda w_{i,j-1} + w_{i+1,j} + w_{i,j+1}) = (\lambda + \omega - 1)w_{i,j}, \quad (8.40)$$

where  $w_{i,j} = 0$  if  $i \in \{0, m+1\}$  or  $j \in \{0, m+1\}$ .

#### Theorem 8.34 (The optimal $\omega$ )

Consider the SOR method applied to the discrete Poisson matrix (3.9), where we use the natural ordering. Moreover, assume  $\omega \in (0, 2)$ .

1. If  $\lambda \neq 0$  is an eigenvalue of  $\mathbf{G}_\omega$  then

$$\mu := \frac{\lambda + \omega - 1}{\omega\lambda^{1/2}} \quad (8.41)$$

is an eigenvalue of  $\mathbf{G}_J$ .

2. If  $\mu$  is an eigenvalue of  $\mathbf{G}_J$  and  $\lambda$  satisfies the equation

$$\mu\omega\lambda^{1/2} = \lambda + \omega - 1 \quad (8.42)$$

then  $\lambda$  is an eigenvalue of  $\mathbf{G}_\omega$ .

**Proof.** Suppose  $(\lambda, \mathbf{w})$  is an eigenpair for  $\mathbf{G}_\omega$ . We claim that  $(\mu, \mathbf{v})$  is an eigenpair for  $\mathbf{G}_J$ , where  $\mu$  is given by (8.41) and  $\mathbf{v} = \overline{(\mathbf{V})}$  with  $v_{i,j} := \lambda^{-(i+j)/2}w_{i,j}$ . Indeed, replacing  $w_{i,j}$  by  $\lambda^{(i+j)/2}v_{i,j}$  in (8.40) and cancelling the common factor  $\lambda^{(i+j)/2}$  we obtain

$$\frac{\omega}{4}(v_{i-1,j} + v_{i,j-1} + v_{i+1,j} + v_{i,j+1}) = \lambda^{-1/2}(\lambda + \omega - 1)v_{i,j}.$$

But then

$$\mathbf{G}_J\mathbf{v} = (\mathbf{L} + \mathbf{R})\mathbf{v} = \frac{\lambda + \omega - 1}{\omega\lambda^{1/2}}\mathbf{v} = \mu\mathbf{v}.$$

For the converse let  $(\mu, \mathbf{v})$  be an eigenpair for  $\mathbf{G}_J$  and let  $\lambda$  be a solution of (8.42). We define as before  $\mathbf{v} =: \text{vec}(\mathbf{V})$ ,  $\mathbf{W} = \text{vec}(\mathbf{W})$  with  $w_{i,j} := \lambda^{(i+j)/2} v_{i,j}$ . Inserting this in (8.38) and canceling  $\lambda^{-(i+j)/2}$  we obtain

$$\frac{1}{4}(\lambda^{1/2}w_{i-1,j} + \lambda^{1/2}w_{i,j-1} + \lambda^{-1/2}w_{i+1,j} + \lambda^{-1/2}w_{i,j+1}) = \mu w_{i,j}.$$

Multiplying by  $\omega\lambda^{1/2}$  we obtain

$$\frac{\omega}{4}(\lambda w_{i-1,j} + \lambda w_{i,j-1} + w_{i+1,j} + w_{i,j+1}) = \omega\mu\lambda^{1/2}w_{i,j},$$

Thus, if  $\omega\mu^{1/2}\lambda^{1/2} = \lambda + \omega - 1$  then by (8.40)  $(\lambda, \mathbf{w})$  is an eigenpair for  $\mathbf{G}_\omega$ .  $\square$

### Proof of Theorem 8.20

Combining statement 1 and 2 in Theorem 8.34 we see that  $\rho(\mathbf{G}_\omega) = |\lambda(\mu)|$ , where  $\lambda(\mu)$  is an eigenvalue of  $\mathbf{G}_\omega$  satisfying (8.42) for some eigenvalue  $\mu$  of  $\mathbf{G}_J$ . The eigenvalues of  $\mathbf{G}_J$  are  $\frac{1}{2}\cos(j\pi h) + \frac{1}{2}\cos(k\pi h)$ ,  $j, k = 1, \dots, m$ , so  $\mu$  is real and both  $\mu$  and  $-\mu$  are eigenvalues. Thus, to compute  $\rho(\mathbf{G}_\omega)$  it is enough to consider (8.42) for a positive eigenvalue  $\mu$  of  $\mathbf{G}_J$ . Solving (8.42) for  $\lambda = \lambda(\mu)$  gives

$$\lambda(\mu) := \frac{1}{4}\left(\omega\mu \pm \sqrt{(\omega\mu)^2 - 4(\omega - 1)}\right)^2. \quad (8.43)$$

Both roots  $\lambda(\mu)$  are eigenvalues of  $\mathbf{G}_\omega$ . The discriminant

$$d(\omega) := (\omega\mu)^2 - 4(\omega - 1).$$

is strictly decreasing on  $(0, 2)$  since

$$d'(\omega) = 2(\omega\mu^2 - 2) < 2(\omega - 2) < 0.$$

Moreover  $d(0) = 4 > 0$  and  $d(2) = 4\mu^2 - 4 < 0$ . As a function of  $\omega$ ,  $\lambda(\mu)$  changes from real to complex when  $d(\omega) = 0$ . The root in  $(0, 2)$  is

$$\omega = \tilde{\omega}(\mu) := 2\frac{1 - \sqrt{1 - \mu^2}}{\mu^2} = \frac{2}{1 + \sqrt{1 - \mu^2}}. \quad (8.44)$$

In the complex case we find

$$|\lambda(\mu)| = \frac{1}{4}\left((\omega\mu)^2 + 4(\omega - 1) - (\omega\mu)^2\right) = \omega - 1, \quad \tilde{\omega}(\mu) < \omega < 2.$$

In the real case both roots of (8.43) are positive and the larger one is

$$\lambda(\mu) = \frac{1}{4}\left(\omega\mu + \sqrt{(\omega\mu)^2 - 4(\omega - 1)}\right)^2, \quad 0 < \omega \leq \tilde{\omega}(\mu). \quad (8.45)$$



Both  $\lambda(\mu)$  and  $\tilde{\omega}(\mu)$  are strictly increasing as functions of  $\mu$ . It follows that  $|\lambda(\mu)|$  is maximized for  $\mu = \rho(\mathbf{G}_J) =: \beta$  and for this value of  $\mu$  we obtain (8.28) for  $0 < \omega \leq \tilde{\omega}(\beta) = \omega^*$ .

Evidently  $\rho(\mathbf{G}_\omega) = \omega - 1$  is strictly increasing in  $\omega^* < \omega < 2$ . Equation (8.30) will follow if we can show that  $\rho(\mathbf{G}_\omega)$  is strictly decreasing in  $0 < \omega < \omega^*$ . By differentiation

$$\frac{d}{d\omega} \left( \omega\beta + \sqrt{(\omega\beta)^2 - 4(\omega - 1)} \right) = \frac{\beta\sqrt{(\omega\beta)^2 - 4(\omega - 1)} + \omega\beta^2 - 2}{\sqrt{(\omega\beta)^2 - 4(\omega - 1)}}.$$

Since  $\beta^2(\omega^2\beta^2 - 4\omega + 4) < (2 - \omega\beta^2)^2$  the numerator is negative and the strict decrease of  $\rho(\mathbf{G}_\omega)$  in  $0 < \omega < \omega^*$  follows.

## 8.6 Review Questions

**8.6.1** Consider a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  with nonzero diagonal elements.

- Define the J and GS method in component form,
- Do they always converge?
- Give a necessary and sufficient condition that  $\mathbf{A}^n \rightarrow \mathbf{0}$ .
- Is there a matrix norm  $\| \cdot \|$  consistent on  $\mathbb{C}^{n \times n}$  such that  $\| \mathbf{A} \| < \rho(\mathbf{A})$ ?

**8.6.2** What is a Neumann series? when does it converge?

**8.6.3** How do we define convergence of a fixed point iteration  $\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}$ ?  
When does it converge?

**8.6.4** Define Richardson's method.



## Chapter 9

# The Conjugate Gradient Method

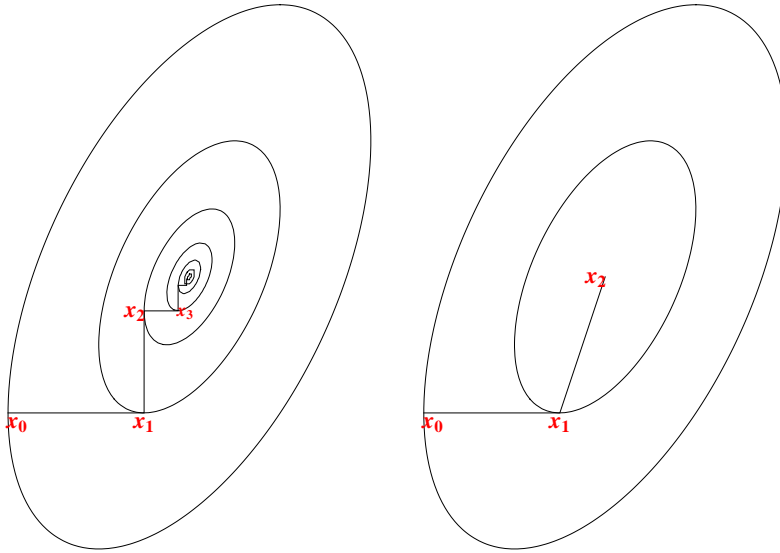


Magnus Rudolph Hestenes, 1906-1991 (left), Eduard L. Stiefel, 1909-1978 (right).

The **conjugate gradient method** was published by Hestenes and Stiefel in 1952, [11] as a direct method for solving linear systems. Today its main use is as an iterative method for solving large sparse linear systems. On a test problem we show that it performs as well as the SOR method with optimal acceleration parameter, and we do not have to estimate any such parameter. However the conjugate gradient method is restricted to symmetric positive definite systems. We also consider the mathematical formulation of the **preconditioned conjugate gradient method**. It is used to speed up convergence of the conjugate gradient method and we study this on a partial differential equation example.

The conjugate gradient method can also be used for minimization and is related to a method known as **steepest descent**. This method and the conjugate gradient method are both minimization methods and iterative methods for solving linear equations.

Throughout this chapter  $\mathbf{A} \in \mathbb{R}^{n \times n}$  will be a symmetric positive definite matrix. Thus,  $\mathbf{A}^T = \mathbf{A}$  and  $\mathbf{y}^T \mathbf{A} \mathbf{y} > 0$  for all nonzero  $\mathbf{y} \in \mathbb{R}^n$ . We recall that  $\mathbf{A}$  has positive eigenvalues and that the spectral (2-norm) condition number of  $\mathbf{A}$  is given by  $\kappa := \frac{\lambda_{max}}{\lambda_{min}}$ , where  $\lambda_{max}$  and  $\lambda_{min}$  are the largest and smallest eigenvalue



**Figure 9.1.** Level curves for  $Q(x, y)$  given by (9.2). Also shown is a steepest descent iteration (left) and a conjugate gradient iteration (right) to find the minimum of  $Q$ . (cf. Examples 9.3, 9.6)

of  $\mathbf{A}$ .

## 9.1 Quadratic Minimization and Steepest Descent

We start by discussing some aspect of quadratic minimization and its relation to solving linear systems.

Consider for  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\mathbf{b} \in \mathbb{R}^n$  the quadratic function  $Q : \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$Q(\mathbf{y}) := \frac{1}{2} \mathbf{y}^T \mathbf{A} \mathbf{y} - \mathbf{b}^T \mathbf{y}. \quad (9.1)$$

As an example, some level curves of

$$Q(x, y) := \frac{1}{2} \begin{bmatrix} x & y \end{bmatrix} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = x^2 - xy + y^2 \quad (9.2)$$

are shown in Figure 9.1. The level curves are ellipses and the graph of  $Q$  is a paraboloid (cf. Exercise 9.1).

**Exercise 9.1 (Paraboloid)**

Let  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^T$  be the spectral decomposition of  $\mathbf{A}$ , i. e.,  $\mathbf{U}$  is orthonormal and  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$  is diagonal. Define new variables  $\mathbf{v} = [v_1, \dots, v_n]^T := \mathbf{U}^T \mathbf{y}$ , and set  $\mathbf{c} := \mathbf{U}^T \mathbf{b} = [c_1, \dots, c_n]^T$ . Show that

$$Q(\mathbf{y}) = \frac{1}{2} \sum_{j=1}^n \lambda_j v_j^2 - \sum_{j=1}^n c_j v_j.$$

Minimizing a quadratic function is equivalent to solving a linear system.

**Lemma 9.2 (Quadratic function)**

A vector  $\mathbf{x} \in \mathbb{R}^n$  minimizes  $Q$  given by (9.1) if and only if  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Moreover, the residual  $\mathbf{r}(\mathbf{y}) := \mathbf{b} - \mathbf{A}\mathbf{y}$  at any  $\mathbf{y} \in \mathbb{R}^n$  is equal to the negative gradient, i. e.,

$$\mathbf{r}(\mathbf{y}) = -\nabla Q(\mathbf{y}), \text{ where } \nabla := \left[ \frac{\partial}{\partial y_1}, \dots, \frac{\partial}{\partial y_n} \right]^T.$$

*Proof.* Expanding  $Q(\mathbf{y} + \varepsilon \mathbf{h}) := \frac{1}{2}(\mathbf{y} + \varepsilon \mathbf{h})^T \mathbf{A}(\mathbf{y} + \varepsilon \mathbf{h}) - \mathbf{b}^T(\mathbf{y} + \varepsilon \mathbf{h})$  we find for any  $\mathbf{y}, \mathbf{h} \in \mathbb{R}^n$  and  $\varepsilon \in \mathbb{R}$

$$Q(\mathbf{y} + \varepsilon \mathbf{h}) = Q(\mathbf{y}) - \varepsilon \mathbf{h}^T \mathbf{r}(\mathbf{y}) + \frac{1}{2} \varepsilon^2 \mathbf{h}^T \mathbf{A} \mathbf{h}, \text{ where } \mathbf{r}(\mathbf{y}) := \mathbf{b} - \mathbf{A}\mathbf{y}. \quad (9.3)$$

If  $\mathbf{y} = \mathbf{x}$ ,  $\varepsilon = 1$ , and  $\mathbf{A}\mathbf{x} = \mathbf{b}$  then (9.3) simplifies to  $Q(\mathbf{x} + \mathbf{h}) = Q(\mathbf{x}) + \frac{1}{2} \mathbf{h}^T \mathbf{A} \mathbf{h}$ , and since  $\mathbf{A}$  is symmetric positive definite  $Q(\mathbf{x} + \mathbf{h}) > Q(\mathbf{x})$  for all nonzero  $\mathbf{h} \in \mathbb{R}^n$ . It follows that  $\mathbf{x}$  is the unique minimum of  $Q$ . Conversely, if  $\mathbf{A}\mathbf{x} \neq \mathbf{b}$  and  $\mathbf{h} := \mathbf{r}(\mathbf{x})$ , then by (9.3),  $Q(\mathbf{x} + \varepsilon \mathbf{h}) - Q(\mathbf{x}) = -\varepsilon(\mathbf{h}^T \mathbf{r}(\mathbf{x}) - \frac{1}{2} \varepsilon \mathbf{h}^T \mathbf{A} \mathbf{h}) < 0$  for  $\varepsilon > 0$  sufficiently small. Thus  $\mathbf{x}$  does not minimize  $Q$ . By (9.3) for  $\mathbf{y} \in \mathbb{R}^n$

$$\begin{aligned} \frac{\partial}{\partial y_i} Q(\mathbf{y}) &:= \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} (Q(\mathbf{y} + \varepsilon \mathbf{e}_i) - Q(\mathbf{y})) \\ &= \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} \left( -\varepsilon \mathbf{e}_i^T \mathbf{r}(\mathbf{y}) + \frac{1}{2} \varepsilon^2 \mathbf{e}_i^T \mathbf{A} \mathbf{e}_i \right) = -\mathbf{e}_i^T \mathbf{r}(\mathbf{y}), \quad i = 1, \dots, n, \end{aligned}$$

showing that  $\mathbf{r}(\mathbf{y}) = -\nabla Q(\mathbf{y})$ .  $\square$

A general class of minimization algorithms for  $Q$  and solution algorithms for a linear system is given as follows:

1. Choose  $\mathbf{x}_0 \in \mathbb{R}^n$ .
2. For  $k = 0, 1, 2, \dots$

Choose a “search direction”  $\mathbf{p}_k$ ,  
 Choose a “step length”  $\alpha_k$ ,  
 Compute  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$ .

(9.4)

We would like to generate a sequence  $\{\mathbf{x}_k\}$  that converges quickly to the minimum  $\mathbf{x}$  of  $Q$ .

For a fixed direction  $\mathbf{p}_k$  we say that  $\alpha_k$  is **optimal** if  $Q(\mathbf{x}_{k+1})$  is as small as possible, i.e.

$$Q(\mathbf{x}_{k+1}) = Q(\mathbf{x}_k + \alpha_k \mathbf{p}_k) = \min_{\alpha \in \mathbb{R}} Q(\mathbf{x}_k + \alpha \mathbf{p}_k).$$

By (9.3) we have  $Q(\mathbf{x}_k + \alpha \mathbf{p}_k) = Q(\mathbf{x}_k) - \alpha \mathbf{p}_k^T \mathbf{r}_k + \frac{1}{2} \alpha^2 \mathbf{p}_k^T \mathbf{A} \mathbf{p}_k$ , where  $\mathbf{r}_k := \mathbf{b} - \mathbf{A} \mathbf{x}_k$ . Since  $\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k \geq 0$  we find a minimum  $\alpha_k$  by solving  $\frac{\partial}{\partial \alpha} Q(\mathbf{x}_k + \alpha \mathbf{p}_k) = 0$ . It follows that the optimal  $\alpha_k$  is uniquely given by

$$\alpha_k := \frac{\mathbf{p}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}. \quad (9.5)$$

In the method of **Steepest descent**, also known as the **Gradient method** we choose  $\mathbf{p}_k = \mathbf{r}_k$  the negative gradient, and the optimal  $\alpha_k$ . Starting from  $\mathbf{x}_0$  we compute for  $k = 0, 1, 2, \dots$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \left( \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k} \right) \mathbf{r}_k. \quad (9.6)$$

This is similar to Richardson's method (8.19), but in that method we used a constant step length. Computationally, a step in the steepest descent iteration can be organized as follows

$$\begin{array}{l} \mathbf{t}_k = \mathbf{A} \mathbf{r}_k, \\ \alpha_k = (\mathbf{r}_k^T \mathbf{r}_k) / (\mathbf{r}_k^T \mathbf{t}_k), \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{r}_k, \\ \mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{t}_k. \end{array} \quad (9.7)$$

Here, and in general, the following updating of the residual is used:

$$\mathbf{r}_{k+1} = \mathbf{b} - \mathbf{A} \mathbf{x}_{k+1} = \mathbf{b} - \mathbf{A}(\mathbf{x}_k + \alpha_k \mathbf{p}_k) = \mathbf{r}_k - \alpha_k \mathbf{t}_k, \quad \mathbf{t}_k := \mathbf{A} \mathbf{p}_k. \quad (9.8)$$

### Example 9.3 (Steepest descent iteration)

Suppose  $Q(x, y)$  is given by (9.2). Starting with  $\mathbf{x}_0 = [-1, -1/2]^T$  and  $\mathbf{r}_0 = -\mathbf{A} \mathbf{x}_0 = [3/2, 0]^T$  we find

$$\begin{array}{l} \mathbf{t}_0 = 3 \begin{bmatrix} -1 \\ -1/2 \end{bmatrix}, \quad \alpha_0 = \frac{1}{2}, \quad \mathbf{x}_1 = -4^{-1} \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad \mathbf{r}_1 = 3 * 4^{-1} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\ \mathbf{t}_1 = 3 * 4^{-1} \begin{bmatrix} -1 \\ 2 \end{bmatrix}, \quad \alpha_1 = \frac{1}{2}, \quad \mathbf{x}_2 = -4^{-1} \begin{bmatrix} 1 \\ 1/2 \end{bmatrix}, \quad \mathbf{r}_2 = 3 * 4^{-1} \begin{bmatrix} 1/2 \\ 0 \end{bmatrix}, \end{array}$$

and in general for  $k \geq 1$

$$\begin{aligned} \mathbf{t}_{2k-2} &= 3 * 4^{1-k} \begin{bmatrix} 1 \\ -1/2 \end{bmatrix}, & \mathbf{x}_{2k-1} &= -4^{-k} \begin{bmatrix} 1 \\ 2 \end{bmatrix}, & \mathbf{r}_{2k-1} &= 3 * 4^{-k} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\ \mathbf{t}_{2k-1} &= 3 * 4^{-k} \begin{bmatrix} -1 \\ 2 \end{bmatrix}, & \mathbf{x}_{2k} &= -4^{-k} \begin{bmatrix} 1 \\ 1/2 \end{bmatrix}, & \mathbf{r}_{2k} &= 3 * 4^{-k} \begin{bmatrix} 1/2 \\ 0 \end{bmatrix}. \end{aligned}$$

Since  $\alpha_k = 1/2$  is constant for all  $k$  the methods of Richardson, Jacobi and steepest descent are the same on this simple problem. See the left part of Figure 9.1. The rate of convergence is determined from  $\|\mathbf{x}_{j+1}\|_2/\|\mathbf{x}_j\| = \|\mathbf{r}_{j+1}\|_2/\|\mathbf{r}_j\| = 1/2$  for all  $j$ .

#### Exercise 9.4 (Steepest descent iteration)

Verify the numbers in Example 9.3.

## 9.2 The Conjugate Gradient Method

In the steepest descent method the choice  $\mathbf{p}_k = \mathbf{r}_k$  implies that the last two gradients are orthogonal. Indeed, by (9.8),  $\mathbf{r}_{k+1}^T \mathbf{r}_k = (\mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{r}_k)^T \mathbf{p}_k = 0$  since  $\alpha_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}$  and  $\mathbf{A}$  is symmetric. In the conjugate gradient method all gradients are orthogonal<sup>11</sup>. We achieve this by using **A-orthogonal search directions** i. e.,  $\mathbf{p}_i^T \mathbf{A} \mathbf{p}_j = 0$  for all  $i \neq j$ .

### 9.2.1 Derivation of the method

As in the steepest descent method we choose a starting vector  $\mathbf{x}_0 \in \mathbb{R}^n$ . If  $\mathbf{r}_0 = \mathbf{b} - \mathbf{A} \mathbf{x}_0 = \mathbf{0}$  then  $\mathbf{x}_0$  is the exact solution and we are finished, otherwise we initially make a steepest descent step. It follows that  $\mathbf{r}_1^T \mathbf{r}_0 = 0$  and  $\mathbf{p}_0 := \mathbf{r}_0$ .

For the general case we define for  $j \geq 0$

$$\mathbf{p}_j := \mathbf{r}_j - \sum_{i=0}^{j-1} \left( \frac{\mathbf{r}_j^T \mathbf{A} \mathbf{p}_i}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i} \right) \mathbf{p}_i, \quad (9.9)$$

$$\mathbf{x}_{j+1} := \mathbf{x}_j + \alpha_j \mathbf{p}_j \quad \alpha_j := \frac{\mathbf{r}_j^T \mathbf{r}_j}{\mathbf{p}_j^T \mathbf{A} \mathbf{p}_j}, \quad (9.10)$$

$$\mathbf{r}_{j+1} = \mathbf{r}_j - \alpha_j \mathbf{A} \mathbf{p}_j. \quad (9.11)$$

We note that

1.  $\mathbf{p}_j$  is computed by the Gram-Schmidt orthogonalization process applied to the residuals  $\mathbf{r}_0, \dots, \mathbf{r}_j$  using the  $\mathbf{A}$ -inner product. The search directions are therefore  $\mathbf{A}$ -orthogonal and nonzero as long as the residuals are linearly independent.

<sup>11</sup>It is this property that has given the method its name.

2. Equation (9.11) follows from (9.8).
3. It can be shown that the step length  $\alpha_j$  is optimal for all  $j$  (cf. Exercise 9.9)).

**Lemma 9.5 (The residuals are orthogonal)**

Suppose that for some  $k \geq 0$  that  $\mathbf{x}_j$  is well defined,  $\mathbf{r}_j \neq 0$ , and  $\mathbf{r}_i^T \mathbf{r}_j = 0$  for  $i, j = 0, 1, \dots, k$ ,  $i \neq j$ . Then  $\mathbf{x}_{k+1}$  is well defined and  $\mathbf{r}_{k+1}^T \mathbf{r}_j = 0$  for  $j = 0, 1, \dots, k$ .

**Proof.** Since the residuals  $\mathbf{r}_j$  are orthogonal and nonzero for  $j \leq k$ , they are linearly independent, and it follows from the Gram-Schmidt Theorem 0.29 that  $\mathbf{p}_k$  is nonzero and  $\mathbf{p}_k^T \mathbf{A} \mathbf{p}_i = 0$  for  $i < k$ . But then  $\mathbf{x}_{k+1}$  and  $\mathbf{r}_{k+1}$  are well defined. Now

$$\begin{aligned} \mathbf{r}_{k+1}^T \mathbf{r}_j &\stackrel{(9.11)}{=} (\mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k)^T \mathbf{r}_j \\ &\stackrel{(9.9)}{=} \mathbf{r}_k^T \mathbf{r}_j - \alpha_k \mathbf{p}_k^T \mathbf{A} (\mathbf{p}_j + \sum_{i=0}^{j-1} (\frac{\mathbf{r}_j^T \mathbf{A} \mathbf{p}_i}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i}) \mathbf{p}_i) \\ \mathbf{p}_k^T \mathbf{A} \mathbf{p}_i &\stackrel{=}{=} 0 \quad \mathbf{r}_k^T \mathbf{r}_j - \alpha_k \mathbf{p}_k^T \mathbf{A} \mathbf{p}_j = 0, \quad j = 0, 1, \dots, k. \end{aligned}$$

That the final expression is equal to zero follows by orthogonality and  $\mathbf{A}$ -orthogonality for  $j < k$  and by the definition of  $\alpha_k$  for  $j = k$ . This completes the proof.  $\square$

The expression (9.9) for  $\mathbf{p}_k$  can be greatly simplified. All terms except the last one vanish, since by orthogonality of the residuals

$$\mathbf{r}_j^T \mathbf{A} \mathbf{p}_i \stackrel{(9.11)}{=} \mathbf{r}_j^T \left( \frac{\mathbf{r}_i - \mathbf{r}_{i+1}}{\alpha_i} \right) = 0, \quad i = 0, 1, \dots, j-2.$$

For the last term with  $k = j-1$

$$\beta_k := -\frac{\mathbf{r}_{k+1}^T \mathbf{A} \mathbf{p}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k} \stackrel{(9.11)}{=} \frac{\mathbf{r}_{k+1}^T (\mathbf{r}_{k+1} - \mathbf{r}_k)}{\alpha_k \mathbf{p}_k^T \mathbf{A} \mathbf{p}_k} \stackrel{(9.10)}{=} \frac{\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{r}_k^T \mathbf{r}_k}. \quad (9.12)$$

To summarize, in the **conjugate gradient method** we start with  $\mathbf{x}_0$ ,  $\mathbf{p}_0 = \mathbf{r}_0 = \mathbf{b} - \mathbf{A} \mathbf{x}_0$  and then generate a sequence of vectors  $\{\mathbf{x}_k\}$  as follows:

For  $k = 0, 1, 2, \dots$

$$\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k, \quad \alpha_k := \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}, \quad (9.13)$$

$$\mathbf{r}_{k+1} := \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k, \quad (9.14)$$

$$\mathbf{p}_{k+1} := \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k, \quad \beta_k := \frac{\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{r}_k^T \mathbf{r}_k}. \quad (9.15)$$



The residuals and search directions are orthogonal and  $\mathbf{A}$ -orthogonal, respectively.

The conjugate gradient method is also a direct method. Since  $\dim \mathbb{R}^n = n$  the  $n + 1$  residuals  $\mathbf{r}_0, \dots, \mathbf{r}_n$  cannot all be nonzero and for orthogonal residuals we find the exact solution in at most  $n$  iterations.

For computation we organize the iterations as follows for  $k = 0, 1, 2, \dots$

$$\begin{array}{l}
 \mathbf{t}_k = \mathbf{A}\mathbf{p}_k, \\
 \alpha_k = (\mathbf{r}_k^T \mathbf{r}_k) / (\mathbf{p}_k^T \mathbf{t}_k), \\
 \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k, \\
 \mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{t}_k, \\
 \beta_k = (\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}) / (\mathbf{r}_k^T \mathbf{r}_k), \\
 \mathbf{p}_{k+1} := \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k.
 \end{array} \tag{9.16}$$

### Example 9.6 (Conjugate gradient iteration)

Consider (9.16) applied to the positive definite linear system  $\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ . Starting as in Example 9.3 with  $\mathbf{x}_0 = \begin{bmatrix} -1 \\ -1/2 \end{bmatrix}$  we find  $\mathbf{p}_0 = \mathbf{r}_0 = \begin{bmatrix} 3/2 \\ 0 \end{bmatrix}$  and then

$$\begin{array}{l}
 \mathbf{t}_0 = \begin{bmatrix} -3 \\ -3/2 \end{bmatrix}, \quad \alpha_0 = 1/2, \quad \mathbf{x}_1 = \begin{bmatrix} -1/4 \\ -1/2 \end{bmatrix}, \quad \mathbf{r}_1 = \begin{bmatrix} 0 \\ 3/4 \end{bmatrix}, \quad \beta_0 = 1/4, \quad \mathbf{p}_1 = \begin{bmatrix} 3/8 \\ 3/4 \end{bmatrix}, \\
 \mathbf{t}_1 = \begin{bmatrix} 0 \\ 9/8 \end{bmatrix}, \quad \alpha_1 = 2/3, \quad \mathbf{x}_2 = \mathbf{0}, \quad \mathbf{r}_2 = \mathbf{0}.
 \end{array}$$

Thus  $\mathbf{x}_2$  is the exact solution as illustrated in the right part of Figure 9.1.

### Exercise 9.7 (Conjugate gradient iteration, II)

Do one iteration with the conjugate gradient method when  $\mathbf{x}_0 = \mathbf{0}$ . (Answer:  $\mathbf{x}_1 = \left( \frac{\mathbf{b}^T \mathbf{b}}{\mathbf{b}^T \mathbf{A} \mathbf{b}} \right) \mathbf{b}$ .)

### Exercise 9.8 (Conjugate gradient iteration, III)

Do two conjugate gradient iterations for the system

$$\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 3 \end{bmatrix}$$

starting with  $\mathbf{x}_0 = \mathbf{0}$ .

### Exercise 9.9 (The cg step length is optimal)

Show that the step length  $\alpha_k$  in the conjugate gradient method is optimal<sup>12</sup>.

<sup>12</sup>Hint: use induction on  $k$  to show that  $\mathbf{p}_k = \mathbf{r}_k + \sum_{j=0}^{k-1} a_{k,j} \mathbf{r}_j$  for some constants  $a_{k,j}$ .

**Exercise 9.10 (Starting value in cg)**

Show that the conjugate gradient method (9.16) for  $\mathbf{Ax} = \mathbf{b}$  starting with  $\mathbf{x}_0$  is the same as applying the method to the system  $\mathbf{Ay} = \mathbf{r}_0 := \mathbf{b} - \mathbf{Ax}_0$  starting with  $\mathbf{y}_0 = \mathbf{0}$ .<sup>13</sup>

**9.2.2 The conjugate gradient algorithm**

In this section we give numerical examples and discuss implementation.

The formulas in (9.16) form a basis for an algorithm.

**Algorithm 9.11 (Conjugate gradient iteration)**

The symmetric positive definite linear system  $\mathbf{Ax} = \mathbf{b}$  is solved by the conjugate gradient method.  $\mathbf{x}$  is a starting vector for the iteration. The iteration is stopped when  $\|\mathbf{r}_k\|_2/\|\mathbf{b}\|_2 \leq \text{tol}$  or  $k > \text{itmax}$ .  $K$  is the number of iterations used.

```

1 function [x,K]=cg(A,b,x,tol,itmax)
2 r=b-A*x; p=r; rho0=b'*b; rho=r'*r;
3 for k=0:itmax
4     if sqrt(rho/rho0)<= tol
5         K=k; return
6     end
7     t=A*p; a=rho/(p'*t);
8     x=x+a*p; r=r-a*t;
9     rhos=rho; rho=r'*r;
10    p=r+(rho/rhos)*p;
11 end
12 K=itmax+1;

```

The work involved in each iteration is

1. one matrix times vector ( $\mathbf{t} = \mathbf{Ap}$ ),
2. two inner products ( $(\mathbf{p}^T \mathbf{t}$  and  $\mathbf{r}^T \mathbf{r}$ ),
3. three vector-plus-scalar-times-vector ( $\mathbf{x} = \mathbf{x} + \mathbf{ap}$ ,  $\mathbf{r} = \mathbf{r} - \mathbf{at}$  and  $\mathbf{p} = \mathbf{r} + (\text{rho}/\text{rhos})\mathbf{p}$ ),

The dominating part is the computation of  $\mathbf{t} = \mathbf{Ap}$ .

**9.2.3 Numerical example**

We test the conjugate gradient method on two examples. For a similar test for the steepest descent method see Exercise 9.17. Consider the matrix given by the

<sup>13</sup>Hint: The conjugate gradient method for  $\mathbf{Ay} = \mathbf{r}_0$  can be written  $\mathbf{y}_{k+1} := \mathbf{y}_k + \gamma_k \mathbf{q}_k$ ,  $\gamma_k := \frac{\mathbf{s}_k^T \mathbf{s}_k}{\mathbf{q}_k^T \mathbf{A} \mathbf{q}_k}$ ,  $\mathbf{s}_{k+1} := \mathbf{s}_k - \gamma_k \mathbf{A} \mathbf{q}_k$ ,  $\mathbf{q}_{k+1} := \mathbf{s}_{k+1} + \delta_k \mathbf{q}_k$ ,  $\delta_k := \frac{\mathbf{s}_{k+1}^T \mathbf{s}_{k+1}}{\mathbf{s}_k^T \mathbf{s}_k}$ . Show that  $\mathbf{y}_k = \mathbf{x}_k - \mathbf{x}_0$ ,  $\mathbf{s}_k = \mathbf{r}_k$ , and  $\mathbf{q}_k = \mathbf{p}_k$ , for  $k = 0, 1, 2, \dots$

|     |       |        |        |           |           |
|-----|-------|--------|--------|-----------|-----------|
| $n$ | 2 500 | 10 000 | 40 000 | 1 000 000 | 4 000 000 |
| $K$ | 19    | 18     | 18     | 16        | 15        |

**Table 9.12.** *The number of iterations  $K$  for the averaging problem on a  $\sqrt{n} \times \sqrt{n}$  grid for various  $n$*

Kronecker sum  $\mathbf{T}_2 := \mathbf{T}_1 \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{T}_1$  where  $\mathbf{T}_1 = \text{tridiag}_m(a, d, a)$ . We recall that this matrix is symmetric positive definite if  $d > 0$  and  $d \geq 2|a|$ . We set  $h = 1/(m+1)$  and  $\mathbf{f} = [1, \dots, 1]^T \in \mathbb{R}^n$ .

We consider two problems.

1.  $a = 1/9$ ,  $d = 5/18$ , the Averaging matrix.
2.  $a = -1$ ,  $d = 2$ , the Poisson matrix.

### 9.2.4 Implementation issues

Note that for our test problems  $\mathbf{T}_2$  only has  $O(5n)$  nonzero elements. Therefore, taking advantage of the sparseness of  $\mathbf{T}_2$  we can compute  $\mathbf{t}$  in Algorithm 9.11 in  $O(n)$  arithmetic operations. With such an implementation the total number of arithmetic operations in one iteration is  $O(n)$ . We also note that it is not necessary to store the matrix  $\mathbf{T}_2$ .

To use the Conjugate Gradient Algorithm on the test matrix for large  $n$  it is advantageous to use a matrix equation formulation. We define matrices  $\mathbf{V}, \mathbf{R}, \mathbf{P}, \mathbf{B}, \mathbf{T} \in \mathbb{R}^{m \times m}$  by  $\mathbf{x} = \text{vec}(\mathbf{V})$ ,  $\mathbf{r} = \text{vec}(\mathbf{R})$ ,  $\mathbf{p} = \text{vec}(\mathbf{P})$ ,  $\mathbf{t} = \text{vec}(\mathbf{T})$ , and  $h^2 \mathbf{f} = \text{vec}(\mathbf{B})$ . Then  $\mathbf{T}_2 \mathbf{x} = h^2 \mathbf{f} \iff \mathbf{T}_1 \mathbf{V} + \mathbf{V} \mathbf{T}_1 = \mathbf{B}$ , and  $\mathbf{t} = \mathbf{T}_2 \mathbf{p} \iff \mathbf{T} = \mathbf{T}_1 \mathbf{P} + \mathbf{P} \mathbf{T}_1$ .

This leads to the following algorithm for testing the conjugate gradient algorithm on the matrix

$$\mathbf{A} = \text{tridiag}_m(a, d, a) \otimes \mathbf{I}_m + \mathbf{I}_m \otimes \text{tridiag}_m(a, d, a) \in \mathbb{R}^{(m^2) \times (m^2)}.$$

|              |       |        |        |         |
|--------------|-------|--------|--------|---------|
| $n$          | 2 500 | 10 000 | 40 000 | 160 000 |
| $K$          | 94    | 188    | 370    | 735     |
| $K/\sqrt{n}$ | 1.88  | 1.88   | 1.85   | 1.84    |

**Table 9.14.** *The number of iterations  $K$  for the Poisson problem on a  $\sqrt{n} \times \sqrt{n}$  grid for various  $n$*

### Algorithm 9.13 (Testing conjugate gradient)

```

1 function [V,K]=cgtest(m,a,d,tol,itmax)
2 R=ones(m)/(m+1)^2; rho=sum(sum(R.*R)); rho0=rho; P=R;
3 V=zeros(m,m); T1=sparse(tridiagonal(a,d,a,m));
4 for k=1:itmax
5     if sqrt(rho/rho0)<= tol
6         K=k; return
7     end
8     T=T1*P+P*T1;
9     a=rho/sum(sum(P.*T)); V=V+a*P; R=R-a*T;
10    rhos=rho; rho=sum(sum(R.*R)); P=R+(rho/rhos)*P;
11 end
12 K=itmax+1;

```

For both the averaging- and Poisson matrix we use  $tol = 10^{-8}$ .

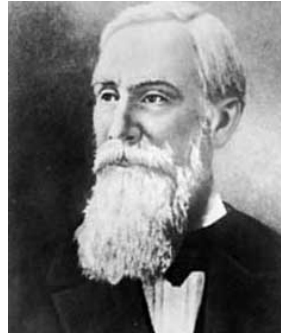
For the averaging matrix we obtain the values in Table 9.12.

The convergence is quite rapid. It appears that the number of iterations can be bounded independently of  $n$ , and therefore we solve the problem in  $O(n)$  operations. This is the best we can do for a problem with  $n$  unknowns.

Consider next the Poisson problem. In Table 9.14 we list  $K$ , the required number of iterations, and  $K/\sqrt{n}$ .

The results show that  $K$  is much smaller than  $n$  and appears to be proportional to  $\sqrt{n}$ . This is the same speed as for SOR and we don't have to estimate any acceleration parameter.

## 9.3 Convergence



Leonid Vitaliyevich Kantorovich, 1912-1986 (left), Aleksey Nikolaevich Krylov, 1863-1945 (center), Pafnuty Lvovich Chebyshev, 1821-1894 (right)

### 9.3.1 The $\mathbf{A}$ -norm

The convergence analysis for both steepest descent and conjugate gradients is in terms of a special inner product. We define the  $\mathbf{A}$ -inner product and the corresponding  $\mathbf{A}$ -norm by

$$\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^T \mathbf{A} \mathbf{y}, \quad \|\mathbf{y}\|_{\mathbf{A}} := \sqrt{\mathbf{y}^T \mathbf{A} \mathbf{y}}, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad (9.17)$$

#### Exercise 9.15 (The $\mathbf{A}$ -inner product)

Show that if  $\mathbf{A}$  is symmetric positive definite then the  $\mathbf{A}$ -inner product is indeed an inner product.

### 9.3.2 The Main Theorem

The following theorem gives upper bounds for the  $\mathbf{A}$ -norm of the error in both methods.

#### Theorem 9.16 (Error bound for steepest descent and conjugate gradients)

Suppose  $\mathbf{A}$  is symmetric positive definite. For the  $\mathbf{A}$ -norms of the errors in the steepest descent method (9.6) the following upper bounds hold

$$\frac{\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}}{\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}} \leq \left( \frac{\kappa - 1}{\kappa + 1} \right)^k < e^{-\frac{2}{\kappa}k}, \quad \kappa > 0, \quad (9.18)$$

while for the conjugate gradient method we have

$$\frac{\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}}{\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}} \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k < 2e^{-\frac{2}{\sqrt{\kappa}}k}, \quad \kappa \geq 0. \quad (9.19)$$

Here  $\kappa = \text{cond}_2(\mathbf{A}) := \lambda_{\max}/\lambda_{\min}$  is the spectral condition number of  $\mathbf{A}$ , and  $\lambda_{\max}$  and  $\lambda_{\min}$  are the largest and smallest eigenvalue of  $\mathbf{A}$ , respectively.

Theorem 9.16 implies

1. Since  $\frac{\kappa-1}{\kappa+1} < 1$  the steepest descent method always converges for a symmetric positive definite matrix. The convergence can be slow when  $\frac{\kappa-1}{\kappa+1}$  is close to one, and this happens even for a moderately ill-conditioned  $\mathbf{A}$ .
2. The rate of convergence for the conjugate gradient method appears to be determined by the square root of the spectral condition number. This is much better than the estimate for the steepest descent method. Especially for problems with large condition numbers.
3. The proofs of the estimates in (9.18) and (9.19) are quite different. This is in spite of their similar appearance.

### 9.3.3 The number of iterations for the model problems

Consider the test matrix

$$\mathbf{T}_2 := \text{tridiag}_m(a, d, a) \otimes \mathbf{I}_m + \mathbf{I}_m \otimes \text{tridiag}_m(a, d, a) \in \mathbb{R}^{(m^2) \times (m^2)}.$$

The eigenvalues were given in (3.20) as

$$\lambda_{j,k} = 2d + 2a \cos(j\pi h) + 2a \cos(k\pi h), \quad j, k = 1, \dots, m. \quad (9.20)$$

For the averaging problem given by  $d = 5/18$ ,  $a = 1/9$ , the largest and smallest eigenvalue of  $\mathbf{T}_2$  are given by  $\lambda_{\max} = \frac{5}{9} + \frac{4}{9} \cos(\pi h)$  and  $\lambda_{\min} = \frac{5}{9} - \frac{4}{9} \cos(\pi h)$ . Thus

$$\kappa_A = \frac{5 + 4 \cos(\pi h)}{5 - 4 \cos(\pi h)} \leq 9,$$

and the condition number is bounded independently of  $n$ . It follows from (9.19) that the number of iterations can be bounded independently of the size  $n$  of the problem, and this is in agreement with what we observed in Table 9.12.

For the Poisson problem we have by (8.25) the condition number

$$\kappa_P = \frac{\lambda_{\max}}{\lambda_{\min}} = \frac{\cos^2(\pi h/2)}{\sin^2(\pi h/2)} \quad \text{and} \quad \sqrt{\kappa_P} = \frac{\cos(\pi h/2)}{\sin(\pi h/2)} \approx \frac{2}{\pi h} \approx \frac{2}{\pi} \sqrt{n}.$$

Thus, (see also Exercise 7.33) we solve the discrete Poisson problem in  $O(n^{3/2})$  arithmetic operations using the conjugate gradient method. This is the same as for the SOR method and for the fast method without the FFT. In comparison the Cholesky Algorithm requires  $O(n^2)$  arithmetic operations both for the averaging and the Poisson problem.

**Exercise 9.17 (Program code for testing steepest descent)**

Write a function  $K=sdtest(m,a,d,tol,itmax)$  to test the Steepest descent method on the matrix  $\mathbf{T}_2$ . Make the analogues of Table 9.12 and Table 9.14. For Table 9.14 it is enough to test for say  $n = 100, 400, 1600, 2500$ , and tabulate  $K/n$  instead of  $K/\sqrt{n}$  in the last row. Conclude that the upper bound (9.18) is realistic. Compare also with the number of iterations for the J and GS method in Table 8.1.

**Exercise 9.18 (Using cg to solve normal equations)**

Consider solving the linear system  $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$  by using the conjugate gradient method. Here  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{b} \in \mathbb{R}^m$  and  $\mathbf{A}^T \mathbf{A}$  is positive definite<sup>14</sup>. Explain why only the following modifications in Algorithm 9.11 are necessary

1.  $r = \mathbf{A}^T(\mathbf{b} - \mathbf{A}^* \mathbf{x}); p = r;$
2.  $a = \text{rho} / (t^* t);$
3.  $r = r - a^* \mathbf{A}^* t;$

Note that the condition number of the normal equations is  $\text{cond}_2(\mathbf{A})^2$ , the square of the condition number of  $\mathbf{A}$ .

## 9.4 Proof of the Convergence Estimates

### 9.4.1 Convergence proof for steepest descent

For the proof of (9.18) the following inequality will be used.

**Theorem 9.19 (Kantorovich inequality)**

For any symmetric positive definite matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$

$$1 \leq \frac{(\mathbf{y}^T \mathbf{A} \mathbf{y})(\mathbf{y}^T \mathbf{A}^{-1} \mathbf{y})}{(\mathbf{y}^T \mathbf{y})^2} \leq \frac{(M + m)^2}{4Mm} \quad \mathbf{y} \neq \mathbf{0}, \mathbf{y} \in \mathbb{R}^n, \quad (9.21)$$

where  $M := \lambda_{\max}$  and  $m := \lambda_{\min}$  are the largest and smallest eigenvalue of  $\mathbf{A}$ , respectively.

**Proof.** For  $j = 1, \dots, n$  let  $(\lambda_j, \mathbf{u}_j)$  be orthonormal eigenpairs of  $\mathbf{A}$  and  $\mathbf{y} \in \mathbb{R}^n$ . By Theorem 5.1  $(\lambda_j^{-1}, \mathbf{u}_j)$  are eigenpairs for  $\mathbf{A}^{-1}$ . Let  $\mathbf{y} = \sum_{j=1}^n c_j \mathbf{u}_j$  be the corresponding eigenvector expansion of  $\mathbf{y}$ . By orthonormality, (cf. (5.6))

$$a := \frac{\mathbf{y}^T \mathbf{A} \mathbf{y}}{\mathbf{y}^T \mathbf{y}} = \sum_{i=1}^n t_i \lambda_i, \quad b := \frac{\mathbf{y}^T \mathbf{A}^{-1} \mathbf{y}}{\mathbf{y}^T \mathbf{y}} = \sum_{i=1}^n \frac{t_i}{\lambda_i}, \quad (9.22)$$

<sup>14</sup>This system known as the **normal equations** appears in linear least squares problems and will be considered in this context in Chapter 11.

where

$$t_i = \frac{c_i^2}{\sum_{j=1}^n c_j^2} \geq 0, \quad i = 1, \dots, n \text{ and } \sum_{i=1}^n t_i = 1. \quad (9.23)$$

Thus  $a$  and  $b$  are **convex combinations** of the eigenvalues of  $\mathbf{A}$  and  $\mathbf{A}^{-1}$ , respectively. Let  $c$  be a positive constant to be chosen later. By the geometric/arithmetical mean inequality (7.27) and (9.22)

$$\sqrt{ab} = \sqrt{(ac)(b/c)} \leq (ac + b/c)/2 = \frac{1}{2} \sum_{i=1}^n t_i (\lambda_i c + 1/(\lambda_i c)) = \frac{1}{2} \sum_{i=1}^n t_i f(\lambda_i c),$$

where  $f : [mc, Mc] \rightarrow \mathbb{R}$  is given by  $f(x) := x + 1/x$ . By (9.23)

$$\sqrt{ab} \leq \frac{1}{2} \max_{mc \leq x \leq Mc} f(x).$$

Since  $f \in C^2$  and  $f''$  is positive it follows from Lemma 7.36 that  $f$  is a convex function. But a convex function takes its maximum at one of the endpoints of the range (cf. Exercise 9.20) and we obtain

$$\sqrt{ab} \leq \frac{1}{2} \max\{f(mc), f(Mc)\}. \quad (9.24)$$

Choosing  $c := 1/\sqrt{mM}$  we find  $f(mc) = f(Mc) = \sqrt{\frac{M}{m}} + \sqrt{\frac{m}{M}} = \frac{M+m}{\sqrt{mM}}$ . By (9.24) we obtain

$$\frac{(\mathbf{y}^T \mathbf{A} \mathbf{y})(\mathbf{y}^T \mathbf{A}^{-1} \mathbf{y})}{(\mathbf{y}^T \mathbf{y})^2} = ab \leq \frac{(M+m)^2}{4Mm},$$

the upper bound in (9.21). For the lower bound we use the Cauchy-Schwarz inequality as follows

$$1 = \left( \sum_{i=1}^n t_i \right)^2 = \left( \sum_{i=1}^n (t_i \lambda_i)^{1/2} (t_i / \lambda_i)^{1/2} \right)^2 \leq \left( \sum_{i=1}^n t_i \lambda_i \right) \left( \sum_{i=1}^n t_i / \lambda_i \right) = ab.$$

□

### Exercise 9.20 (Maximum of a convex function)

Show that if  $f : [a, b] \rightarrow \mathbb{R}$  is convex then  $\max_{a \leq x \leq b} f(x) \leq \max\{f(a), f(b)\}$ .

### Proof of (9.18)

Let  $\boldsymbol{\epsilon}_j := \mathbf{x} - \mathbf{x}_j$ ,  $j = 0, 1, \dots$ , where  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . It is enough to show that

$$\frac{\|\boldsymbol{\epsilon}_{k+1}\|_{\mathbf{A}}^2}{\|\boldsymbol{\epsilon}_k\|_{\mathbf{A}}^2} \leq \left( \frac{\kappa - 1}{\kappa + 1} \right)^2, \quad k = 0, 1, 2, \dots, \quad (9.25)$$



for then  $\|\epsilon_k\|_{\mathbf{A}} \leq \left(\frac{\kappa-1}{\kappa+1}\right) \|\epsilon_{k-1}\| \leq \dots \leq \left(\frac{\kappa-1}{\kappa+1}\right)^k \|\epsilon_0\|$ . It follows from (9.6) that

$$\epsilon_{k+1} = \epsilon_k - \alpha_k \mathbf{r}_k, \quad \alpha_k := \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k}.$$

We find

$$\begin{aligned} \|\epsilon_k\|_{\mathbf{A}}^2 &= \epsilon_k^T \mathbf{A} \epsilon_k = \mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k, \\ \|\epsilon_{k+1}\|_{\mathbf{A}}^2 &= (\epsilon_k - \alpha_k \mathbf{r}_k)^T \mathbf{A} (\epsilon_k - \alpha_k \mathbf{r}_k) \\ &= \epsilon_k^T \mathbf{A} \epsilon_k - 2\alpha_k \mathbf{r}_k^T \mathbf{A} \epsilon_k + \alpha_k^2 \mathbf{r}_k^T \mathbf{A} \mathbf{r}_k = \|\epsilon_k\|_{\mathbf{A}}^2 - \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k}. \end{aligned}$$

Combining these and using Kantorovich inequality

$$\frac{\|\epsilon_{k+1}\|_{\mathbf{A}}^2}{\|\epsilon_k\|_{\mathbf{A}}^2} = 1 - \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{(\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k)(\mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k)} \leq 1 - \frac{4\lambda_{\min}\lambda_{\max}}{(\lambda_{\min} + \lambda_{\max})^2} = \left(\frac{\kappa-1}{\kappa+1}\right)^2$$

and (9.25) is proved.

The inequality

$$\frac{x-1}{x+1} < e^{-2/x} \quad \text{for } x > 1 \quad (9.26)$$

follows from the familiar series expansion of the exponential function. Indeed, with  $y = 1/x$ , using  $2^k/k! = 2$ ,  $k = 1, 2$ , and  $2^k/k! < 2$  for  $k > 2$ , we find

$$e^{2/x} = e^{2y} = \sum_{k=0}^{\infty} \frac{(2y)^k}{k!} < 1 + 2 \sum_{k=1}^{\infty} y^k = \frac{1+y}{1-y} = \frac{x+1}{x-1}$$

and (9.26) follows.  $\square$

## 9.4.2 Krylov spaces and the best approximation property

For the convergence analysis of the conjugate gradient method certain subspaces of  $\mathbb{R}^n$  called **Krylov spaces** play a central role. In fact the iterates in the conjugate gradient method are best approximation of the solution from these subspaces using the  $\mathbf{A}$ -norm to measure the error.

The Krylov spaces are defined by  $\mathbb{W}_0 = \{\mathbf{0}\}$  and

$$\mathbb{W}_k = \text{span}(\mathbf{r}_0, \mathbf{A}\mathbf{r}_0, \mathbf{A}^2\mathbf{r}_0, \dots, \mathbf{A}^{k-1}\mathbf{r}_0), \quad k = 1, 2, 3, \dots$$

They are nested subspaces

$$\mathbb{W}_0 \subset \mathbb{W}_1 \subset \mathbb{W}_2 \subset \dots \subset \mathbb{W}_n \subset \mathbb{R}^n$$

with  $\dim(\mathbb{W}_k) \leq k$  for all  $k \geq 0$ . Moreover, if  $\mathbf{v} \in \mathbb{W}_k$  then  $\mathbf{A}\mathbf{v} \in \mathbb{W}_{k+1}$ .

**Lemma 9.21 (Krylov space)**

For the iterates in the conjugate gradient method we have

$$\mathbf{x}_k - \mathbf{x}_0 \in \mathbb{W}_k, \quad \mathbf{r}_k, \mathbf{p}_k \in \mathbb{W}_{k+1}, \quad k = 0, 1, \dots, \quad (9.27)$$

and

$$\mathbf{r}_k^T \mathbf{w} = \mathbf{p}_k^T \mathbf{A} \mathbf{w} = 0, \quad \mathbf{w} \in \mathbb{W}_k. \quad (9.28)$$

**Proof.** (9.27) clearly holds for  $k = 0$  since  $\mathbf{p}_0 = \mathbf{r}_0$ . Suppose it holds for some  $k \geq 0$ . Then  $\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k \in \mathbb{W}_{k+2}$  and by  $\mathbf{p}_{k+1} = \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k \in \mathbb{W}_{k+2}$  and  $\mathbf{x}_{k+1} - \mathbf{x}_0 \stackrel{(9.10)}{=} \mathbf{x}_k - \mathbf{x}_0 + \alpha_k \mathbf{p}_k \in \mathbb{W}_{k+1}$ . Thus (9.27) follows by induction. Since any  $\mathbf{w} \in \mathbb{W}_k$  is a linear combination of  $\{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{k-1}\}$  and also  $\{\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{k-1}\}$ , (9.28) follows.  $\square$

**Theorem 9.22 (Best approximation property)**

Suppose  $\mathbf{A} \mathbf{x} = \mathbf{b}$ , where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is symmetric positive definite and  $\{\mathbf{x}_k\}$  is generated by the conjugate gradient method (cf. (9.13)). Then

$$\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}} = \min_{\mathbf{w} \in \mathbb{W}_k} \|\mathbf{x} - \mathbf{x}_0 - \mathbf{w}\|_{\mathbf{A}}. \quad (9.29)$$

**Proof.** Fix  $k$ , let  $\mathbf{w} \in \mathbb{W}_k$  and  $\mathbf{u} := \mathbf{x}_k - \mathbf{x}_0 - \mathbf{w}$ . By (9.27)  $\mathbf{u} \in \mathbb{W}_k$  and then (9.28) implies that  $\mathbf{r}_k^T \mathbf{u} = 0$ . Since  $(\mathbf{x} - \mathbf{x}_k)^T \mathbf{A} \mathbf{u} = \mathbf{r}_k^T \mathbf{u}$  we find

$$\begin{aligned} \|\mathbf{x} - \mathbf{x}_0 - \mathbf{w}\|_{\mathbf{A}}^2 &= (\mathbf{x} - \mathbf{x}_k + \mathbf{u})^T \mathbf{A} (\mathbf{x} - \mathbf{x}_k + \mathbf{u}) \\ &= (\mathbf{x} - \mathbf{x}_k)^T \mathbf{A} (\mathbf{x} - \mathbf{x}_k) + 2\mathbf{r}_k^T \mathbf{u} + \mathbf{u}^T \mathbf{A} \mathbf{u} \\ &= \|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}^2 + \|\mathbf{u}\|_{\mathbf{A}}^2 \geq \|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}^2. \end{aligned}$$

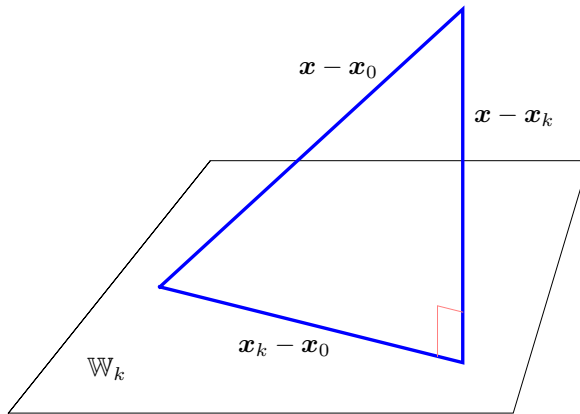
Taking square roots the result follows.  $\square$

If  $\mathbf{x}_0 = \mathbf{0}$  then (9.29) says that  $\mathbf{x}_k$  is the element in  $\mathbb{W}_k$  that is closest to the solution  $\mathbf{x}$  in the  $\mathbf{A}$ -norm. More generally, if  $\mathbf{x}_0 \neq \mathbf{0}$  then  $\mathbf{x} - \mathbf{x}_k = (\mathbf{x} - \mathbf{x}_0) - (\mathbf{x}_k - \mathbf{x}_0)$  and  $\mathbf{x}_k - \mathbf{x}_0$  is the element in  $\mathbb{W}_k$  that is closest to  $\mathbf{x} - \mathbf{x}_0$  in the  $\mathbf{A}$ -norm. This is the orthogonal projection of  $\mathbf{x} - \mathbf{x}_0$  into  $\mathbb{W}_k$ , see Figure 9.2.

Recall that to each polynomial  $p(t) := \sum_{j=0}^m a_j t^j$  there corresponds a matrix polynomial  $p(\mathbf{A}) := a_0 \mathbf{I} + a_1 \mathbf{A} + \dots + a_m \mathbf{A}^m$ . Moreover, if  $(\lambda_j, \mathbf{u}_j)$  are eigenpairs of  $\mathbf{A}$  then  $(p(\lambda_j), \mathbf{u}_j)$  are eigenpairs of  $p(\mathbf{A})$  for  $j = 1, \dots, n$ .

**Lemma 9.23 (Krylov space and polynomials)**

Suppose  $\mathbf{A} \mathbf{x} = \mathbf{b}$  where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is symmetric positive definite with orthonormal eigenpairs  $(\lambda_j, \mathbf{u}_j)$ ,  $j = 1, 2, \dots, n$ , and let  $\mathbf{r}_0 := \mathbf{b} - \mathbf{A} \mathbf{x}_0$  for some  $\mathbf{x}_0 \in \mathbb{R}^n$ .



**Figure 9.2.** The orthogonal projection of  $\mathbf{x} - \mathbf{x}_0$  into  $\mathbb{W}_k$ .

To each  $\mathbf{w} \in \mathbb{W}_k$  there corresponds a polynomial  $P(t) := \sum_{j=0}^{k-1} a_j t^{k-1}$  such that  $\mathbf{w} = P(\mathbf{A})\mathbf{r}_0$ . Moreover, if  $\mathbf{r}_0 = \sum_{j=1}^n \sigma_j \mathbf{u}_j$  then

$$\|\mathbf{x} - \mathbf{x}_0 - \mathbf{w}\|_{\mathbf{A}}^2 = \sum_{j=1}^n \frac{\sigma_j^2}{\lambda_j} Q(\lambda_j)^2, \quad Q(t) := 1 - tP(t). \quad (9.30)$$

**Proof.** If  $\mathbf{w} \in \mathbb{W}_k$  then  $\mathbf{w} = a_0 \mathbf{r}_0 + a_1 \mathbf{A} \mathbf{r}_0 + \cdots + a_{k-1} \mathbf{A}^{k-1} \mathbf{r}_0$  for some scalars  $a_0, \dots, a_{k-1}$ . But then  $\mathbf{w} = P(\mathbf{A})\mathbf{r}_0$ . We find

$$\mathbf{A}(\mathbf{x} - \mathbf{x}_0 - \mathbf{w}) = \mathbf{A}(\mathbf{x} - \mathbf{x}_0 - P(\mathbf{A})\mathbf{r}_0) = \mathbf{r}_0 - \mathbf{A}P(\mathbf{A})\mathbf{r}_0 = Q(\mathbf{A})\mathbf{r}_0,$$

and so  $\|\mathbf{x} - \mathbf{x}_0 - \mathbf{w}\|_{\mathbf{A}}^2 = \mathbf{c}^T \mathbf{A}^{-1} \mathbf{c}$ , where  $\mathbf{c} = Q(\mathbf{A})\mathbf{r}_0$ . Using the eigenvector expansion for  $\mathbf{r}_0$  we obtain

$$\mathbf{c} = \sum_{j=1}^n \sigma_j Q(\lambda_j) \mathbf{u}_j, \quad \mathbf{A}^{-1} \mathbf{c} = \sum_{i=1}^n \sigma_i \frac{Q(\lambda_i)}{\lambda_i} \mathbf{u}_i. \quad (9.31)$$

Now (9.30) follows by the orthonormality of the eigenvectors.  $\square$

We will use the following theorem to estimate the rate of convergence.

**Theorem 9.24 (cg and best polynomial approximation)**

Suppose  $[a, b]$  with  $0 < a < b$  is an interval containing all the eigenvalues of  $\mathbf{A}$ .

Then in the conjugate gradient method

$$\frac{\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}}{\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}} = \min_{\substack{Q \in \Pi_k \\ Q(0)=1}} \max_{a \leq x \leq b} |Q(x)|, \quad (9.32)$$

where  $\Pi_k$  denotes the class of univariate polynomials of degree  $\leq k$  with real coefficients.

**Proof.** With the notation in Lemma 9.23 we find  $\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}^2 = \mathbf{r}_0 \mathbf{A}^{-1} \mathbf{r}_0 = \sum_{j=1}^n \frac{\sigma_j^2}{\lambda_j}$ . Therefore, by the best approximation property and (9.30), for any  $\mathbf{w} \in \mathbb{W}_k$

$$\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}^2 \leq \|\mathbf{x} - \mathbf{x}_0 - \mathbf{w}\|_{\mathbf{A}}^2 \leq \max_{a \leq x \leq b} |Q(x)|^2 \sum_{j=1}^n \frac{\sigma_j^2}{\lambda_j} = \max_{a \leq x \leq b} |Q(x)|^2 \|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}^2,$$

where  $Q \in \Pi_k$  and  $Q(0) = 1$ . Minimizing over such polynomials  $Q$  and taking square roots the result follows.  $\square$

In the next section we use properties of the Chebyshev polynomials to show that

$$\frac{\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}}{\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}} \leq \min_{\substack{Q \in \Pi_k \\ Q(0)=1}} \max_{\lambda_{\min} \leq x \leq \lambda_{\max}} |Q(x)| = \frac{2}{\left(\frac{\sqrt{\kappa}+1}{\sqrt{\kappa}-1}\right)^k + \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^k}, \quad (9.33)$$

where  $\kappa = \lambda_{\max}/\lambda_{\min}$  is the spectral condition number of  $\mathbf{A}$ . Ignoring the second term in the denominator this implies the first inequality in (9.19). The second inequality follows from (9.26).

### Exercise 9.25 (Krylov space and cg iterations)

Consider the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  where

$$\mathbf{A} = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}, \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} 4 \\ 0 \\ 0 \end{bmatrix}.$$

a) Determine the vectors defining the Krylov spaces for  $k \leq 3$  taking as initial

$$\text{approximation } \mathbf{x} = \mathbf{0}. \text{ Answer: } [\mathbf{b}, \mathbf{A}\mathbf{b}, \mathbf{A}^2\mathbf{b}] = \begin{bmatrix} 4 & 8 & 20 \\ 0 & -4 & -16 \\ 0 & 0 & 4 \end{bmatrix}.$$

b) Carry out three CG-iterations on  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Answer:

$$[\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3] = \begin{bmatrix} 0 & 2 & 8/3 & 3 \\ 0 & 0 & 4/3 & 2 \\ 0 & 0 & 0 & 1 \end{bmatrix},$$

$$\begin{aligned}
 [\mathbf{r}_0, \mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3] &= \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 4/3 & 0 \end{bmatrix}, \\
 [\mathbf{A}\mathbf{p}_0, \mathbf{A}\mathbf{p}_1, \mathbf{A}\mathbf{p}_2] &= \begin{bmatrix} 8 & 0 & 0 \\ -4 & 3 & 0 \\ 0 & -2 & 16/9 \end{bmatrix}, \\
 [\mathbf{p}_0, \mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3] &= \begin{bmatrix} 4 & 1 & 4/9 & 0 \\ 0 & 2 & 8/9 & 0 \\ 0 & 0 & 12/9 & 0 \end{bmatrix},
 \end{aligned}$$

c) Verify that

- $\dim(\mathbb{W}_k) = k$  for  $k = 0, 1, 2, 3$ .
- $\mathbf{x}_3$  is the exact solution of  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .
- $\mathbf{r}_0, \dots, \mathbf{r}_{k-1}$  is an orthogonal basis for  $\mathbb{W}_k$  for  $k = 1, 2, 3$ .
- $\mathbf{p}_0, \dots, \mathbf{p}_{k-1}$  is an  $\mathbf{A}$ -orthogonal basis for  $\mathbb{W}_k$  for  $k = 1, 2, 3$ .
- $\{\|\mathbf{r}_k\|_2\}$  is monotonically decreasing.
- $\{\|\mathbf{x}_k - \mathbf{x}\|_2\}$  is monotonically decreasing.

### 9.4.3 Chebyshev polynomials

The proof of the estimate (9.33) for the error in the conjugate gradient method is based on an extremal property of the Chebyshev polynomials. Suppose  $a < b$ ,  $c \notin [a, b]$  and  $k \in \mathbb{N}$ . Consider the set  $\mathcal{S}_k$  of all polynomials  $Q$  of degree  $\leq k$  such that  $Q(c) = 1$ . For any continuous function  $f$  on  $[a, b]$  we define

$$\|f\|_\infty = \max_{a \leq x \leq b} |f(x)|.$$

We want to find a polynomial  $Q^* \in \mathcal{S}_k$  such that

$$\|Q^*\|_\infty = \min_{Q \in \mathcal{S}_k} \|Q\|_\infty.$$

We will show that  $Q^*$  is uniquely given as a suitably shifted and normalized version of the **Chebyshev polynomial**. The Chebyshev polynomial  $T_n$  of degree  $n$  can be defined recursively by

$$T_{n+1}(t) = 2tT_n(t) - T_{n-1}(t), \quad n \geq 1, \quad t \in \mathbb{R},$$

starting with  $T_0(t) = 1$  and  $T_1(t) = t$ . Thus  $T_2(t) = 2t^2 - 1$ ,  $T_3(t) = 4t^3 - 3t$  etc. In general  $T_n$  is a polynomial of degree  $n$ .

There are some convenient closed form expressions for  $T_n$ .

**Lemma 9.26 (Closed forms of Chebyshev polynomials)**

For  $n \geq 0$

1.  $T_n(t) = \cos(\arccos t)$  for  $t \in [-1, 1]$ ,
2.  $T_n(t) = \frac{1}{2}[(t + \sqrt{t^2 - 1})^n + (t + \sqrt{t^2 - 1})^{-n}]$  for  $|t| \geq 1$ .

**Proof.** 1. With  $P_n(t) = \cos(n \arccos t)$  we have  $P_n(t) = \cos n\phi$ , where  $t = \cos \phi$ . Therefore,

$$P_{n+1}(t) + P_{n-1}(t) = \cos(n+1)\phi + \cos(n-1)\phi = 2 \cos \phi \cos n\phi = 2tP_n(t),$$

and it follows that  $P_n$  satisfies the same recurrence relation as  $T_n$ . Since  $P_0 = T_0$  and  $P_1 = T_1$  we have  $P_n = T_n$  for all  $n \geq 0$ .

2. Fix  $t$  with  $|t| \geq 1$  and let  $x_n := T_n(t)$  for  $n \geq 0$ . The recurrence relation for the Chebyshev polynomials can then be written

$$x_{n+1} - 2tx_n + x_{n-1} = 0 \text{ for } n \geq 1, \text{ with } x_0 = 1, x_1 = t. \quad (9.34)$$

To solve this difference equation we insert  $x_n = z^n$  into (9.34) and obtain  $z^{n+1} - 2tz^n + z^{n-1} = 0$  or  $z^2 - 2tz + 1 = 0$ . The roots of this equation are

$$z_1 = t + \sqrt{t^2 - 1}, \quad z_2 = t - \sqrt{t^2 - 1} = (t + \sqrt{t^2 - 1})^{-1}.$$

Now  $z_1^n, z_2^n$  and more generally  $c_1 z_1^n + c_2 z_2^n$  are solutions of (9.34) for any constants  $c_1$  and  $c_2$ . We find these constants from the initial conditions  $x_0 = c_1 + c_2 = 1$  and  $x_1 = c_1 z_1 + c_2 z_2 = t$ . Since  $z_1 + z_2 = 2t$  the solution is  $c_1 = c_2 = \frac{1}{2}$ .  $\square$

We show that the unique solution to our minimization problem is

$$Q^*(x) = \frac{T_k(u(x))}{T_k(u(c))}, \quad u(x) = \frac{b + a - 2x}{b - a}. \quad (9.35)$$

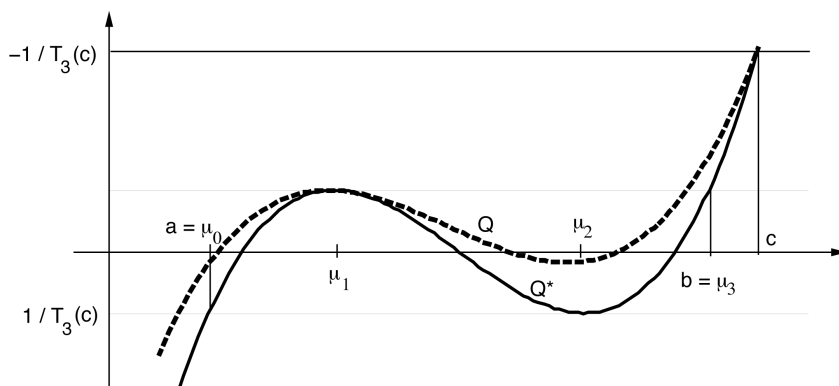
Clearly  $Q^* \in S_k$ .

**Theorem 9.27 (A minimal norm problem)**

Suppose  $a < b$ ,  $c \notin [a, b]$  and  $k \in \mathbb{N}$ . If  $Q \in S_k$  and  $Q \neq Q^*$  then  $\|Q\|_\infty > \|Q^*\|_\infty$ .

**Proof.** Recall that a nonzero polynomial  $p$  of degree  $k$  can have at most  $k$  zeros. If  $p(z) = p'(z) = 0$ , we say that  $p$  has a double zero at  $z$ . Counting such a zero as two zeros it is still true that a nonzero polynomial of degree  $k$  has at most  $k$  zeros.

$|Q^*|$  takes on its maximum  $1/|T_k(u(c))|$  at the  $k+1$  points  $\mu_0, \dots, \mu_k$  in  $[a, b]$  such that  $u(\mu_i) = \cos(i\pi/k)$  for  $i = 0, 1, \dots, k$ . Suppose  $Q \in S_k$  and that  $\|Q\| \leq \|Q^*\|$ . We have to show that  $Q \equiv Q^*$ . Let  $f \equiv Q - Q^*$ . We show that  $f$



**Figure 9.3.** This is an illustration of the proof of Theorem 9.27 for  $k = 3$ .  $f \equiv Q - Q^*$  has a double zero at  $\mu_1$  and one zero between  $\mu_2$  and  $\mu_3$ .

has at least  $k$  zeros in  $[a, b]$ . Since  $f$  is a polynomial of degree  $\leq k$  and  $f(c) = 0$ , this means that  $f \equiv 0$  or equivalently  $Q \equiv Q^*$ .

Consider  $I_j = [\mu_{j-1}, \mu_j]$  for a fixed  $j$ . Let

$$\sigma_j = f(\mu_{j-1})f(\mu_j).$$

We have  $\sigma_j \leq 0$ . For if say  $Q^*(\mu_j) > 0$  then

$$Q(\mu_j) \leq \|Q\|_\infty \leq \|Q^*\|_\infty = Q^*(\mu_j)$$

so that  $f(\mu_j) \leq 0$ . Moreover,

$$-Q(\mu_{j-1}) \leq \|Q\|_\infty \leq \|Q^*\|_\infty = -Q^*(\mu_{j-1}).$$

Thus  $f(\mu_{j-1}) \geq 0$  and it follows that  $\sigma_j \leq 0$ . Similarly,  $\sigma_j \leq 0$  if  $Q^*(\mu_j) < 0$ .

If  $\sigma_j < 0$ ,  $f$  must have a zero in  $I_j$  since it is continuous. Suppose  $\sigma_j = 0$ . Then  $f(\mu_{j-1}) = 0$  or  $f(\mu_j) = 0$ . If  $f(\mu_j) = 0$  then  $Q(\mu_j) = Q^*(\mu_j)$ . But then  $\mu_j$  is a maximum or minimum both for  $Q$  and  $Q^*$ . If  $\mu_j \in (a, b)$  then  $Q'(\mu_j) = Q^{*'}(\mu_j) = 0$ . Thus  $f(\mu_j) = f'(\mu_j) = 0$ , and  $f$  has a double zero at  $\mu_j$ . We can count this as one zero for  $I_j$  and one for  $I_{j+1}$ . If  $\mu_j = b$ , we still have a zero in  $I_j$ . Similarly, if  $f(\mu_{j-1}) = 0$ , a double zero of  $f$  at  $\mu_{j-1}$  appears if  $\mu_{j-1} \in (a, b)$ . We count this as one zero for  $I_{j-1}$  and one for  $I_j$ .

In this way we associate one zero of  $f$  for each of the  $k$  intervals  $I_j$ ,  $j = 1, 2, \dots, k$ . We conclude that  $f$  has at least  $k$  zeros in  $[a, b]$ .  $\square$

**Exercise 9.28 (Another explicit formula for the Chebyshev polynomial)**

Show that

$$T_n(t) = \cosh(n \operatorname{arccosh} t) \text{ for } |t| \geq 1,$$

where  $\operatorname{arccosh}$  is the inverse function of  $\cosh x := (e^x + e^{-x})/2$ .

Theorem 9.27 with  $a$ , and  $b$ , the smallest and largest eigenvalue of  $\mathbf{A}$ , and  $c = 0$  implies that the minimizing polynomial in (9.33) is given by

$$Q^*(x) = T_k \left( \frac{b+a-2x}{b-a} \right) / T_k \left( \frac{b+a}{b-a} \right). \quad (9.36)$$

By Lemma 9.26

$$\max_{a \leq x \leq b} \left| T_k \left( \frac{b+a-2x}{b-a} \right) \right| = \max_{-1 \leq t \leq 1} |T_k(t)| = 1. \quad (9.37)$$

Moreover with  $t = (b+a)/(b-a)$  we have

$$t + \sqrt{t^2 - 1} = \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1}, \quad \kappa = b/a.$$

Thus again by Lemma 9.26 we find

$$T_k \left( \frac{b+a}{b-a} \right) = T_k \left( \frac{\kappa+1}{\kappa-1} \right) = \frac{1}{2} \left[ \left( \frac{\sqrt{\kappa}+1}{\sqrt{\kappa}-1} \right)^k + \left( \frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1} \right)^k \right] \quad (9.38)$$

and (9.33) follows.

**9.4.4 Monotonicity of the error**

The error analysis for the conjugate gradient method is based on the  $\mathbf{A}$ -norm. We end this chapter by considering the Euclidian norm of the error, and show that it is strictly decreasing.

**Theorem 9.29 (The error in cg is strictly decreasing)**

Let in the conjugate gradient method  $m$  be the smallest integer such that  $\mathbf{r}_{m+1} = \mathbf{0}$ . For  $k \leq m$  we have  $\|\boldsymbol{\epsilon}_{k+1}\|_2 < \|\boldsymbol{\epsilon}_k\|_2$ . More precisely,

$$\|\boldsymbol{\epsilon}_k\|_2^2 - \|\boldsymbol{\epsilon}_{k+1}\|_2^2 = \frac{\|\mathbf{p}_k\|_2^2}{\|\mathbf{p}_k\|_{\mathbf{A}}^2} (\|\boldsymbol{\epsilon}_k\|_{\mathbf{A}}^2 + \|\boldsymbol{\epsilon}_{k+1}\|_{\mathbf{A}}^2)$$

where  $\boldsymbol{\epsilon}_j = \mathbf{x} - \mathbf{x}_j$  and  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .



**Proof.** For  $j \leq m$

$$\boldsymbol{\epsilon}_j = \mathbf{x}_{m+1} - \mathbf{x}_j = \mathbf{x}_m - \mathbf{x}_j + \alpha_m \mathbf{p}_m = \mathbf{x}_{m-1} - \mathbf{x}_j + \alpha_{m-1} \mathbf{p}_{m-1} + \alpha_m \mathbf{p}_m = \dots$$

so that

$$\boldsymbol{\epsilon}_j = \sum_{i=j}^m \alpha_i \mathbf{p}_i, \quad \alpha_i = \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i}. \quad (9.39)$$

By (9.39) and  $\mathbf{A}$ -orthogonality

$$\|\boldsymbol{\epsilon}_j\|_{\mathbf{A}}^2 = \boldsymbol{\epsilon}_j^T \mathbf{A} \boldsymbol{\epsilon}_j = \sum_{i=j}^m \alpha_i^2 \mathbf{p}_i^T \mathbf{A} \mathbf{p}_i = \sum_{i=j}^m \frac{(\mathbf{r}_i^T \mathbf{r}_i)^2}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i}. \quad (9.40)$$

By (9.15) and Lemma 9.21

$$\mathbf{p}_i^T \mathbf{p}_k = (\mathbf{r}_i + \beta_{i-1} \mathbf{p}_{i-1})^T \mathbf{p}_k = \beta_{i-1} \mathbf{p}_{i-1}^T \mathbf{p}_k = \dots = \beta_{i-1} \dots \beta_k (\mathbf{p}_k^T \mathbf{p}_k),$$

and since  $\beta_{i-1} \dots \beta_k = (\mathbf{r}_i^T \mathbf{r}_i) / (\mathbf{r}_k^T \mathbf{r}_k)$  we obtain

$$\mathbf{p}_i^T \mathbf{p}_k = \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{r}_k^T \mathbf{r}_k} \mathbf{p}_k^T \mathbf{p}_k, \quad i \geq k. \quad (9.41)$$

Since

$$\|\boldsymbol{\epsilon}_k\|_2^2 = \|\boldsymbol{\epsilon}_{k+1} + \mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 = \|\boldsymbol{\epsilon}_{k+1} + \alpha_k \mathbf{p}_k\|_2^2,$$

we obtain

$$\begin{aligned} \|\boldsymbol{\epsilon}_k\|_2^2 - \|\boldsymbol{\epsilon}_{k+1}\|_2^2 &= \alpha_k (2\mathbf{p}_k^T \boldsymbol{\epsilon}_{k+1} + \alpha_k \mathbf{p}_k^T \mathbf{p}_k) \\ &\stackrel{(9.39)}{=} \alpha_k \left( 2 \sum_{i=k+1}^m \alpha_i \mathbf{p}_i^T \mathbf{p}_k + \alpha_k \mathbf{p}_k^T \mathbf{p}_k \right) = \left( \sum_{i=k}^m + \sum_{i=k+1}^m \right) \alpha_k \alpha_i \mathbf{p}_i^T \mathbf{p}_k \\ &\stackrel{(9.41)}{=} \left( \sum_{i=k}^m + \sum_{i=k+1}^m \right) \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k} \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i} \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{r}_k^T \mathbf{r}_k} \mathbf{p}_k^T \mathbf{p}_k \\ &\stackrel{(9.40)}{=} \frac{\|\mathbf{p}_k\|_2^2}{\|\mathbf{p}_k\|_{\mathbf{A}}^2} (\|\boldsymbol{\epsilon}_k\|_{\mathbf{A}}^2 + \|\boldsymbol{\epsilon}_{k+1}\|_{\mathbf{A}}^2). \end{aligned}$$

and the Theorem is proved.  $\square$

## 9.5 Preconditioning

For problems  $\mathbf{A}\mathbf{x} = \mathbf{b}$  of size  $n$ , where both  $n$  and  $\text{cond}_2(\mathbf{A})$  are large, it is often possible to improve the performance of the conjugate gradient method by using a

technique known as **preconditioning**. Instead of  $\mathbf{Ax} = \mathbf{b}$  we consider an equivalent system  $\mathbf{BAx} = \mathbf{Bb}$ , where  $\mathbf{B}$  is nonsingular and  $\text{cond}_2(\mathbf{BA})$  is smaller than  $\text{cond}_2(\mathbf{A})$ . The matrix  $\mathbf{B}$  will in many cases be the inverse of another matrix,  $\mathbf{B} = \mathbf{M}^{-1}$ . We cannot use CG on  $\mathbf{BAx} = \mathbf{Bb}$  directly since  $\mathbf{BA}$  in general is not symmetric even if both  $\mathbf{A}$  and  $\mathbf{B}$  are. But if  $\mathbf{B}$  (and hence  $\mathbf{M}$ ) is symmetric positive definite then we can apply CG to a symmetrized system and then transform the recurrence formulas to an iterative method for the original system  $\mathbf{Ax} = \mathbf{b}$ . This iterative method is known as the **preconditioned conjugate gradient method**. We shall see that the convergence properties of this method is determined by the eigenvalues of  $\mathbf{BA}$ .

Suppose  $\mathbf{B}$  is symmetric positive definite. By Theorem 2.38 there is a nonsingular matrix  $\mathbf{C}$  such that  $\mathbf{B} = \mathbf{C}^T \mathbf{C}$ . ( $\mathbf{C}$  is only needed for the derivation and will not appear in the final formulas). Now

$$\mathbf{BAx} = \mathbf{Bb} \Leftrightarrow \mathbf{C}^T (\mathbf{CAC}^T) \mathbf{C}^{-T} \mathbf{x} = \mathbf{C}^T \mathbf{Cb} \Leftrightarrow (\mathbf{CAC}^T) \mathbf{y} = \mathbf{Cb}, \text{ \& } \mathbf{x} = \mathbf{C}^T \mathbf{y}.$$

We have 3 linear systems

$$\mathbf{Ax} = \mathbf{b} \tag{9.42}$$

$$\mathbf{BAx} = \mathbf{Bb} \tag{9.43}$$

$$(\mathbf{CAC}^T) \mathbf{y} = \mathbf{Cb}, \text{ \& } \mathbf{x} = \mathbf{C}^T \mathbf{y}. \tag{9.44}$$

Note that (9.42) and (9.44) are symmetric positive definite linear systems. In addition to being symmetric positive definite the matrix  $\mathbf{CAC}^T$  is similar to  $\mathbf{BA}$ . Indeed,

$$\mathbf{C}^T (\mathbf{CAC}^T) \mathbf{C}^{-T} = \mathbf{BA}.$$

Thus  $\mathbf{CAC}^T$  and  $\mathbf{BA}$  have the same eigenvalues. Therefore if we apply the conjugate gradient method to (9.44) then the rate of convergence will be determined by the eigenvalues of  $\mathbf{BA}$ .

We apply the conjugate gradient method to  $(\mathbf{CAC}^T) \mathbf{y} = \mathbf{Cb}$ . Denoting the search direction by  $\mathbf{q}_k$  and the residual by  $\mathbf{z}_k = \mathbf{Cb} - \mathbf{CAC}^T \mathbf{y}_k$  we obtain the following from (9.13), (9.14), and (9.15).

$$\begin{aligned} \mathbf{y}_{k+1} &= \mathbf{y}_k + \alpha_k \mathbf{q}_k, & \alpha_k &= \mathbf{z}_k^T \mathbf{z}_k / \mathbf{q}_k^T (\mathbf{CAC}^T) \mathbf{q}_k, \\ \mathbf{z}_{k+1} &= \mathbf{z}_k - \alpha_k (\mathbf{CAC}^T) \mathbf{q}_k, \\ \mathbf{q}_{k+1} &= \mathbf{z}_{k+1} + \beta_k \mathbf{q}_k, & \beta_k &= \mathbf{z}_{k+1}^T \mathbf{z}_{k+1} / \mathbf{z}_k^T \mathbf{z}_k. \end{aligned}$$

With

$$\mathbf{x}_k := \mathbf{C}^T \mathbf{y}_k, \quad \mathbf{p}_k := \mathbf{C}^T \mathbf{q}_k, \quad \mathbf{s}_k := \mathbf{C}^T \mathbf{z}_k, \quad \mathbf{r}_k := \mathbf{C}^{-1} \mathbf{z}_k \tag{9.45}$$

this can be transformed into

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k, \quad \alpha_k = \frac{\mathbf{s}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}, \quad (9.46)$$

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k, \quad (9.47)$$

$$\mathbf{s}_{k+1} = \mathbf{s}_k - \alpha_k \mathbf{B} \mathbf{A} \mathbf{p}_k, \quad (9.48)$$

$$\mathbf{p}_{k+1} = \mathbf{s}_{k+1} + \beta_k \mathbf{p}_k, \quad \beta_k = \frac{\mathbf{s}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{s}_k^T \mathbf{r}_k}. \quad (9.49)$$

Here  $\mathbf{x}_k$  will be an approximation to the solution  $\mathbf{x}$  of  $\mathbf{A}\mathbf{x} = \mathbf{b}$ ,  $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$  is the residual in the original system, and  $\mathbf{s}_k = \mathbf{B}\mathbf{b} - \mathbf{B}\mathbf{A}\mathbf{x}_k$  is the residual in the preconditioned system. This follows since by (9.45)

$$\mathbf{r}_k = \mathbf{C}^{-1} \mathbf{z}_k = \mathbf{b} - \mathbf{C}^{-1} \mathbf{C} \mathbf{A} \mathbf{C}^T \mathbf{y}_k = \mathbf{b} - \mathbf{A} \mathbf{x}_k$$

and  $\mathbf{s}_k = \mathbf{C}^T \mathbf{z}_k = \mathbf{C}^T \mathbf{C} \mathbf{r}_k = \mathbf{B} \mathbf{r}_k$ . We start with  $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$ ,  $\mathbf{p}_0 = \mathbf{s}_0 = \mathbf{B}\mathbf{r}_0$  and obtain the following preconditioned conjugate gradient algorithm for determining approximations  $\mathbf{x}_k$  to the solution of a symmetric positive definite system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .

**Algorithm 9.30 (Preconditioned conjugate gradient )**

The symmetric positive definite linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  is solved by the preconditioned conjugate gradient method on the system  $\mathbf{B}\mathbf{A}\mathbf{x} = \mathbf{B}\mathbf{b}$ , where  $\mathbf{B}$  is symmetric positive definite.  $\mathbf{x}$  is a starting vector for the iteration. The iteration is stopped when  $\|\mathbf{r}_k\|_2 / \|\mathbf{b}\|_2 \leq \text{tol}$  or  $k > \text{itmax}$ .  $K$  is the number of iterations used.

```

1 function [x,K]=pcg(A,B,b,x,tol,itmax)
2 r=b-A*x; p=B*r; s=p; rho=s'*r; rho0=b'*b;
3 for k=0:itmax
4     if sqrt(rho/rho0)<= tol
5         K=k; return
6     end
7     t=A*p; a=rho/(p'*t);
8     x=x+a*p; r=r-a*t;
9     w=B*t; s=s-a*w;
10    rhos=rho; rho=s'*r;
11    p=r+(rho/rhos)*p;
12 end
13 K=itmax+1;
```

Appart from the calculation of  $\rho$  this algorithm is quite similar to Algorithm 9.11. The main additional work is contained in  $w = B * t$ . We'll discuss this further in connection with an example. There the inverse of  $\mathbf{B}$  is known and we have to solve a linear system to find  $\mathbf{w}$ .

We have the following convergence result for this algorithm.

**Theorem 9.31 (Error bound preconditioned cg)**

Suppose we apply a symmetric positive definite preconditioner  $\mathbf{B}$  to the symmetric positive definite system  $\mathbf{Ax} = \mathbf{b}$ . Then the quantities  $\mathbf{x}_k$  computed in Algorithm 9.30 satisfy the following bound:

$$\frac{\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}}{\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}} \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k, \quad \text{for } k \geq 0,$$

where  $\kappa = \lambda_{\max}/\lambda_{\min}$  is the ratio of the largest and smallest eigenvalue of  $\mathbf{BA}$ .

**Proof.** Since Algorithm 9.30 is equivalent to solving (9.44) by the conjugate gradient method Theorem 9.16 implies that

$$\frac{\|\mathbf{y} - \mathbf{y}_k\|_{\mathbf{CAC}^T}}{\|\mathbf{y} - \mathbf{y}_0\|_{\mathbf{CAC}^T}} \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k, \quad \text{for } k \geq 0,$$

where  $\mathbf{y}_k$  is the conjugate gradient approximation to the solution  $\mathbf{y}$  of (9.44) and  $\kappa$  is the ratio of the largest and smallest eigenvalue of  $\mathbf{CAC}^T$ . Since  $\mathbf{BA}$  and  $\mathbf{CAC}^T$  are similar this is the same as the  $\kappa$  in the theorem. By (9.45) we have

$$\begin{aligned} \|\mathbf{y} - \mathbf{y}_k\|_{\mathbf{CAC}^T}^2 &= (\mathbf{y} - \mathbf{y}_k)^T (\mathbf{CAC}^T) (\mathbf{y} - \mathbf{y}_k) \\ &= (\mathbf{C}^T (\mathbf{y} - \mathbf{y}_k))^T \mathbf{A} (\mathbf{C}^T (\mathbf{y} - \mathbf{y}_k)) = \|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{A}}^2 \end{aligned}$$

and the proof is complete.  $\square$

We conclude that  $\mathbf{B}$  should satisfy the following requirements for a problem of size  $n$ :

1. The eigenvalues of  $\mathbf{BA}$  should be located in a narrow interval. Preferably one should be able to bound the length of the interval independently of  $n$ .
2. The evaluation of  $\mathbf{Bx}$  for a given vector  $\mathbf{x}$  should not be expensive in storage and arithmetic operations, ideally  $O(n)$  for both.

## 9.6 Preconditioning Example

### 9.6.1 A variable coefficient problem

Consider the problem

$$\begin{aligned} -\frac{\partial}{\partial x} \left( c(x, y) \frac{\partial u}{\partial x} \right) - \frac{\partial}{\partial y} \left( c(x, y) \frac{\partial u}{\partial y} \right) &= f(x, y) & (x, y) \in \Omega = (0, 1)^2 \\ u(x, y) &= 0 & (x, y) \in \partial\Omega. \end{aligned} \quad (9.50)$$

Here  $\Omega$  is the open unit square while  $\partial\Omega$  is the boundary of  $\Omega$ . The functions  $f$  and  $c$  are given and we seek a function  $u = u(x, y)$  such that (9.50) holds. We

assume that  $c$  and  $f$  are defined and continuous on  $\Omega$  and that  $c(x, y) > 0$  for all  $(x, y) \in \Omega$ . The problem (9.50) reduces to the Poisson problem in the special case where  $c(x, y) = 1$  for  $(x, y) \in \Omega$ .

As for the Poisson problem we solve (9.50) numerically on a grid of points

$$\{(jh, kh) : j, k = 0, 1, \dots, m+1\}, \quad \text{where } h = 1/(m+1),$$

and where  $m$  is a positive integer. Let  $(x, y)$  be one of the interior grid points. For univariate functions  $f, g$  we use the central difference approximations

$$\begin{aligned} \frac{\partial}{\partial t} \left( f(t) \frac{\partial}{\partial t} g(t) \right) &\approx \left( f(t + \frac{h}{2}) \frac{\partial}{\partial t} g(t + h/2) - f(t - \frac{h}{2}) \frac{\partial}{\partial t} g(t - \frac{h}{2}) \right) / h \\ &\approx \left( f(t + \frac{h}{2}) (g(t + h) - g(t)) - f(t - \frac{h}{2}) (g(t) - g(t - h)) \right) / h^2 \end{aligned}$$

to obtain

$$\frac{\partial}{\partial x} \left( c \frac{\partial u}{\partial x} \right)_{j,k} \approx \frac{c_{j+\frac{1}{2},k} (v_{j+1,k} - v_{j,k}) - c_{j-\frac{1}{2},k} (v_{j,k} - v_{j-1,k})}{h^2}$$

and

$$\frac{\partial}{\partial y} \left( c \frac{\partial u}{\partial y} \right)_{j,k} \approx \frac{c_{j,k+\frac{1}{2}} (v_{j,k+1} - v_{j,k}) - c_{j,k-\frac{1}{2}} (v_{j,k} - v_{j,k-1})}{h^2},$$

where  $c_{p,q} = c(ph, qh)$  and  $v_{j,k} \approx u(jh, kh)$ . With these approximations the discrete analog of (9.50) turns out to be

$$\begin{aligned} -(\mathbf{P}_h v)_{j,k} &= h^2 f_{j,k} & j, k = 1, \dots, m \\ v_{j,k} &= 0 & j = 0, m+1 \text{ all } k \text{ or } k = 0, m+1 \text{ all } j, \end{aligned} \quad (9.51)$$

where

$$\begin{aligned} -(\mathbf{P}_h v)_{j,k} &= (c_{j,k-\frac{1}{2}} + c_{j-\frac{1}{2},k} + c_{j+\frac{1}{2},k} + c_{j,k+\frac{1}{2}}) v_{j,k} \\ &\quad - c_{j,k-\frac{1}{2}} v_{j,k-1} - c_{j-\frac{1}{2},k} v_{j-1,k} - c_{j+\frac{1}{2},k} v_{j+1,k} - c_{j,k+\frac{1}{2}} v_{j,k+1} \end{aligned} \quad (9.52)$$

and  $f_{j,k} = f(jh, kh)$ .

As before we let  $\mathbf{V} = (v_{j,k}) \in \mathbb{R}^{m \times m}$  and  $\mathbf{F} = (f_{j,k}) \in \mathbb{R}^{m \times m}$ . The corresponding linear system can be written  $\mathbf{A}\mathbf{x} = \mathbf{b}$  where  $\mathbf{x} = \text{vec}(\mathbf{V})$ ,  $\mathbf{b} = h^2 \text{vec}(\mathbf{F})$ , and the  $n$ -by- $n$  coefficient matrix  $\mathbf{A}$  is given by

$$\begin{aligned} a_{i,i} &= c_{j_i, k_i - \frac{1}{2}} + c_{j_i - \frac{1}{2}, k_i} + c_{j_i + \frac{1}{2}, k_i} + c_{j_i, k_i + \frac{1}{2}}, & i = 1, 2, \dots, n \\ a_{i+1,i} &= a_{i,i+1} = -c_{j_i + \frac{1}{2}, k_i}, & i \bmod m \neq 0 \\ a_{i+m,i} &= a_{i,i+m} = -c_{j_i, k_i + \frac{1}{2}}, & i = 1, 2, \dots, n-m \\ a_{i,j} &= 0 & \text{otherwise,} \end{aligned} \quad (9.53)$$

where  $(j_i, k_i)$  with  $1 \leq j_i, k_i \leq m$  is determined uniquely from the equation  $i = j_i + (k_i - 1)m$  for  $i = 1, \dots, n$ . If  $c(x, y) = 1$  for all  $(x, y) \in \Omega$  we recover the Poisson matrix.

In general we cannot write  $\mathbf{A}$  as a Kronecker sum. But we can show that  $\mathbf{A}$  is symmetric and it is positive definite as long as the function  $c$  is positive on  $\Omega$ .

**Theorem 9.32 (Positive definite matrix)**

If  $c(x, y) > 0$  for  $(x, y) \in \Omega$  then the matrix  $\mathbf{A}$  given by (9.53) is symmetric positive definite.

**Proof.**

To each  $x \in \mathbb{R}^n$  there corresponds a matrix  $\mathbf{V} \in \mathbb{R}^{m \times m}$  such that  $x = \text{vec}(\mathbf{V})$ . We claim that

$$x^T \mathbf{A} x = \sum_{j=1}^m \sum_{k=0}^m c_{j, k+\frac{1}{2}} (v_{j, k+1} - v_{j, k})^2 + \sum_{k=1}^m \sum_{j=0}^m c_{j+\frac{1}{2}, k} (v_{j+1, k} - v_{j, k})^2, \quad (9.54)$$

where  $v_{0, k} = v_{m+1, k} = v_{j, 0} = v_{j, m+1} = 0$  for  $j, k = 0, 1, \dots, m+1$ . Since  $c_{j+\frac{1}{2}, k}$  and  $c_{j, k+\frac{1}{2}}$  correspond to values of  $c$  in  $\Omega$  for the values of  $j, k$  in the sums it follows that they are positive and from (9.54) we see that  $x^T \mathbf{A} x \geq 0$  for all  $x \in \mathbb{R}^n$ . Moreover if  $x^T \mathbf{A} x = 0$  then all quadratic factors are zero and  $v_{j, k+1} = v_{j, k}$  for  $k = 0, 1, \dots, m$  and  $j = 1, \dots, m$ . Now  $v_{j, 0} = v_{j, m+1} = 0$  implies that  $\mathbf{V} = \mathbf{0}$  and hence  $x = 0$ . Thus  $\mathbf{A}$  is symmetric positive definite.

It remains to prove (9.54). From the connection between (9.52) and (9.53) we have

$$\begin{aligned} x^T \mathbf{A} x &= \sum_{j=1}^m \sum_{k=1}^m -(\mathbf{P}_h v)_{j, k} v_{j, k} \\ &= \sum_{j=1}^m \sum_{k=1}^m \left( c_{j, k-\frac{1}{2}} v_{j, k}^2 + c_{j-\frac{1}{2}, k} v_{j, k}^2 + c_{j+\frac{1}{2}, k} v_{j, k}^2 + c_{j, k+\frac{1}{2}} v_{j, k}^2 \right. \\ &\quad \left. - c_{j, k-\frac{1}{2}} v_{j, k-1} v_{j, k} - c_{j, k+\frac{1}{2}} v_{j, k} v_{j, k+1} \right. \\ &\quad \left. - c_{j-\frac{1}{2}, k} v_{j-1, k} v_{j, k} - c_{j+\frac{1}{2}, k} v_{j, k} v_{j+1, k} \right). \end{aligned}$$

Using the homogenous boundary conditions we obtain

$$\begin{aligned}
\sum_{j=1}^m \sum_{k=1}^m c_{j,k-\frac{1}{2}} v_{j,k}^2 &= \sum_{j=1}^m \sum_{k=0}^m c_{j,k+\frac{1}{2}} v_{j,k+1}^2, \\
\sum_{j=1}^m \sum_{k=1}^m c_{j,k-\frac{1}{2}} v_{j,k-1} v_{j,k} &= \sum_{j=1}^m \sum_{k=0}^m c_{j,k+\frac{1}{2}} v_{j,k+1} v_{j,k}, \\
\sum_{j=1}^m \sum_{k=1}^m c_{j-\frac{1}{2},k} v_{j,k}^2 &= \sum_{k=1}^m \sum_{j=0}^m c_{j+\frac{1}{2},k} v_{j+1,k}^2, \\
\sum_{j=1}^m \sum_{k=1}^m c_{j-\frac{1}{2},k} v_{j-,k} v_{j,k} &= \sum_{k=1}^m \sum_{j=0}^m c_{j+\frac{1}{2},k} v_{j+1,k} v_{j,k}.
\end{aligned}$$

It follows that

$$\begin{aligned}
\mathbf{x}^T \mathbf{A} \mathbf{x} &= \sum_{j=1}^m \sum_{k=0}^m c_{j,k+\frac{1}{2}} (v_{j,k}^2 + v_{j,k+1}^2 - 2v_{j,k} v_{j,k+1}) \\
&\quad + \sum_{k=1}^m \sum_{j=0}^m c_{j+\frac{1}{2},k} (v_{j,k}^2 + v_{j+1,k}^2 - 2v_{j,k} v_{j+1,k})
\end{aligned}$$

and (9.54) follows.  $\square$

### 9.6.2 Applying preconditioning

Consider solving  $\mathbf{A} \mathbf{x} = \mathbf{b}$ , where  $\mathbf{A}$  is given by (9.53) and  $\mathbf{b} \in \mathbb{R}^n$ . Since  $\mathbf{A}$  is positive definite it is nonsingular and the system has a unique solution  $\mathbf{x} \in \mathbb{R}^n$ . Moreover we can use either Cholesky factorization or the block tridiagonal solver to find  $\mathbf{x}$ . Since the bandwidth of  $\mathbf{A}$  is  $m = \sqrt{n}$  both of these methods require  $O(n^2)$  arithmetic operations for large  $n$ .

If we choose  $c(x, y) \equiv 1$  in (9.50), we get the Poisson problem. With this in mind, we may think of the coefficient matrix  $\mathbf{A}_p$  arising from the discretization of the Poisson problem as an approximation to the matrix (9.53). This suggests using  $\mathbf{B} = \mathbf{A}_p^{-1}$ , the inverse of the discrete Poisson matrix as a preconditioner for the system (9.51).

Consider Algorithm 9.30. With this preconditioner the calculation  $\mathbf{w} = \mathbf{B} \mathbf{t}$  takes the form  $\mathbf{A}_p \mathbf{w}_k = \mathbf{t}_k$ .

In Section 4.2 we developed a Simple fast Poisson Solver, Cf. Algorithm 4.1. This method can be utilized to solve  $\mathbf{A}_p \mathbf{w} = \mathbf{t}$ .

Consider the specific problem where

$$c(x, y) = e^{-x+y} \text{ and } f(x, y) = 1.$$

|              |      |       |       |       |       |
|--------------|------|-------|-------|-------|-------|
| $n$          | 2500 | 10000 | 22500 | 40000 | 62500 |
| $K$          | 222  | 472   | 728   | 986   | 1246  |
| $K/\sqrt{n}$ | 4.44 | 4.72  | 4.85  | 4.93  | 4.98  |
| $K_{pre}$    | 22   | 23    | 23    | 23    | 23    |

**Table 9.33.** *The number of iterations  $K$  (no preconditioning) and  $K_{pre}$  (with preconditioning) for the problem (9.50) using the discrete Poisson problem as a preconditioner.*

We have used Algorithm 9.11 (conjugate gradient without preconditioning), and Algorithm 9.30 (conjugate gradient with preconditioning) to solve the problem (9.50). We used  $\mathbf{x}_0 = 0$  and  $\epsilon = 10^{-8}$ . The results are shown in Table 9.33.

Without preconditioning the number of iterations still seems to be more or less proportional to  $\sqrt{n}$  although the convergence is slower than for the constant coefficient problem. Using preconditioning speeds up the convergence considerably. The number of iterations appears to be bounded independently of  $n$ .

Using a preconditioner increases the work in each iteration. For the present example the number of arithmetic operations in each iteration changes from  $O(n)$  without preconditioning to  $O(n^{3/2})$  or  $O(n \log_2 n)$  with preconditioning. This is not a large increase and both the number of iterations and the computing time is reduced significantly.

Let us finally show that the number  $\kappa = \lambda_{max}/\lambda_{min}$  which determines the rate of convergence for the preconditioned conjugate gradient method applied to (9.50) can be bounded independently of  $n$ .

**Theorem 9.34 (Eigvalues of preconditioned matrix)**

*Suppose  $0 < c_0 \leq c(x, y) \leq c_1$  for all  $(x, y) \in [0, 1]^2$ . For the eigenvalues of the matrix  $\mathbf{B}\mathbf{A} = \mathbf{A}_p^{-1}\mathbf{A}$  just described we have*

$$\kappa = \frac{\lambda_{max}}{\lambda_{min}} \leq \frac{c_1}{c_0}.$$

**Proof.**

Suppose  $\mathbf{A}_p^{-1}\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  for some  $\mathbf{x} \in \mathbb{R}^n \setminus \{0\}$ . Then  $\mathbf{A}\mathbf{x} = \lambda\mathbf{A}_p\mathbf{x}$ . Multiplying this by  $\mathbf{x}^T$  and solving for  $\lambda$  we find

$$\lambda = \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{A}_p \mathbf{x}}.$$

We computed  $\mathbf{x}^T \mathbf{A} \mathbf{x}$  in (9.54) and we obtain  $\mathbf{x}^T \mathbf{A}_p \mathbf{x}$  by setting all the  $c$ 's there



equal to one

$$\mathbf{x}^T \mathbf{A}_p \mathbf{x} = \sum_{i=1}^m \sum_{j=0}^m (v_{i,j+1} - v_{i,j})^2 + \sum_{j=1}^m \sum_{i=0}^m (v_{i+1,j} - v_{i,j})^2.$$

Thus  $\mathbf{x}^T \mathbf{A}_p \mathbf{x} > 0$  and bounding all the  $c$ 's in (9.54) from below by  $c_0$  and above by  $c_1$  we find

$$c_0(\mathbf{x}^T \mathbf{A}_p \mathbf{x}) \leq \mathbf{x}^T \mathbf{A} \mathbf{x} \leq c_1(\mathbf{x}^T \mathbf{A}_p \mathbf{x})$$

which implies that  $c_0 \leq \lambda \leq c_1$  for all eigenvalues  $\lambda$  of  $\mathbf{B}\mathbf{A} = \mathbf{A}_p^{-1}\mathbf{A}$ .  $\square$

Using  $c(x, y) = e^{-x+y}$  as above, we find  $c_0 = e^{-2}$  and  $c_1 = 1$ . Thus  $\kappa \leq e^2 \approx 7.4$ , a quite acceptable matrix condition number which explains the convergence results from our numerical experiment.

## 9.7 Review Questions

**9.7.1** Does the steepest descent and conjugate gradient method always converge?

**9.7.2** What kind of orthogonalities occur in the conjugate gradient method?

**9.7.3** What is a Krylow space?

**9.7.4** What is a convex function?

**9.7.5** How do SOR and conjugate gradient compare?



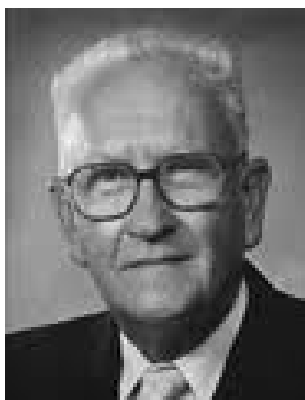
## **Part IV**

# **Orthonormal Transformations and Least Squares**



## Chapter 10

# Orthonormal and Unitary Transformations



Alston Scott Householder, 1904-1993 (left), James Hardy Wilkinson, 1919-1986 (right). Householder and Wilkinson are two of the founders of modern numerical analysis and scientific computing.

Gauss transformations, (cf. Theorem 2.60, the PLU theorem) are used in Gaussian elimination to reduce a matrix to triangular form. These are not the only kind of transformations that can be used for such a task. In this chapter we study how transformations by orthonormal and unitary matrices can be used to reduce a square matrix to upper triangular form and more generally a rectangular matrix to upper triangular (also called upper trapezoidal) form. This leads to a decomposition of the matrix known as a **QR decomposition** and a reduced form which we refer to as a **QR factorization**. The QR decomposition and factorization will be used in later chapters to solve least squares- and eigenvalue problems.

It cannot be repeated too often that orthonormal transformations have the advantage that they preserve the Euclidean norm of a vector, and the spectral norm and Frobenius norm of a matrix, see Lemma 6.24 and Theorem 7.20. This means that when an orthonormal transformation is applied to an inaccurate vector or matrix then the error will not grow. Thus in general an orthonormal transfor-

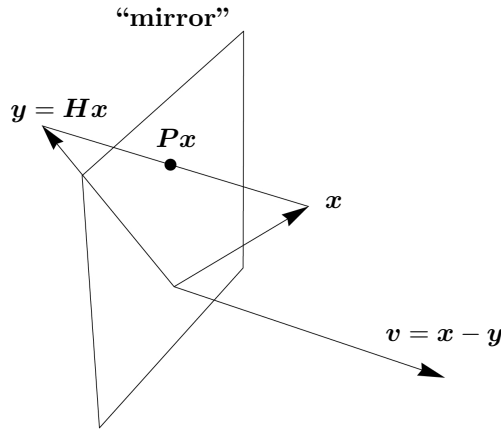


Figure 10.1. The Householder transformation in Exercise 10.2

mation is numerically stable.

## 10.1 The Householder Transformation

**Definition 10.1 (Householder transformation)**

A matrix  $\mathbf{H} \in \mathbb{C}^{n \times n}$  of the form

$$\mathbf{H} := \mathbf{I} - \mathbf{u}\mathbf{u}^*, \text{ where } \mathbf{u} \in \mathbb{C}^n \text{ and } \mathbf{u}^*\mathbf{u} = 2$$

is called a **Householder transformation**. The name **elementary reflector** is also used.

In the real case and for  $n = 2$  we find  $\mathbf{H} = \begin{bmatrix} 1-u_1^2 & -u_1u_2 \\ -u_2u_1 & 1-u_2^2 \end{bmatrix}$ . A Householder transformation is Hermitian and unitary. Indeed,  $\mathbf{H}^* = (\mathbf{I} - \mathbf{u}\mathbf{u}^*)^* = \mathbf{H}$  and

$$\mathbf{H}^*\mathbf{H} = \mathbf{H}^2 = (\mathbf{I} - \mathbf{u}\mathbf{u}^*)(\mathbf{I} - \mathbf{u}\mathbf{u}^*) = \mathbf{I} - 2\mathbf{u}\mathbf{u}^* + \mathbf{u}(\mathbf{u}^*\mathbf{u})\mathbf{u}^* = \mathbf{I}.$$

In the real case  $\mathbf{H}$  is symmetric and orthonormal.

There are several ways to represent a Householder transformation. Householder used  $\mathbf{I} - 2\mathbf{u}\mathbf{u}^*$ , where  $\mathbf{u}^*\mathbf{u} = 1$ . For any nonzero  $\mathbf{v} \in \mathbb{R}^n$  the matrix

$$\mathbf{H} := \mathbf{I} - 2\frac{\mathbf{v}\mathbf{v}^*}{\mathbf{v}^*\mathbf{v}} \tag{10.1}$$

is a Householder transformation. Indeed,  $\mathbf{H} = \mathbf{I} - \mathbf{u}\mathbf{u}^*$ , where  $\mathbf{u} := \sqrt{2}\frac{\mathbf{v}}{\|\mathbf{v}\|_2}$  has length  $\sqrt{2}$ . Moreover, if  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  with  $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2$  and  $\mathbf{v} := \mathbf{x} - \mathbf{y} \neq \mathbf{0}$ , then  $\mathbf{H}\mathbf{x} = \mathbf{y}$  (Cf. Exercise 10.2).

**Exercise 10.2 (Reflector)**

Suppose  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  with  $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2$  and  $\mathbf{v} := \mathbf{x} - \mathbf{y} \neq \mathbf{0}$ .

- (a) Show that  $\mathbf{H}\mathbf{x} := (\mathbf{I} - 2\frac{\mathbf{v}\mathbf{v}^T}{\mathbf{v}^T\mathbf{v}})\mathbf{x} = \mathbf{y}$ .<sup>15</sup>
- (b) Let  $\mathcal{M} := \{\mathbf{w} \in \mathbb{R}^n : \mathbf{w}^T\mathbf{v} = 0\}$  and  $\mathbf{P} := \mathbf{I} - \frac{\mathbf{v}\mathbf{v}^T}{\mathbf{v}^T\mathbf{v}}$ . Show that  $\mathbf{P}\mathbf{x} = (\mathbf{x} + \mathbf{y})/2 \in \mathcal{M}$ . Thus  $\mathbf{y}$  is the reflected image of  $\mathbf{x}$ , where  $\mathcal{M}$  is the "mirror". See Figure 10.1.
- (c) Determine the matrices  $\mathbf{H}, \mathbf{P}$  and the "mirror"  $\mathcal{M}$  when  $\mathbf{x} := [1, 0, 1]^T$  and  $\mathbf{y} := [-1, 0, 1]^T$ .

A main use of Householder transformations is to produce zeros in vectors.

**Theorem 10.3 (Zeros in vectors)**

Suppose  $\mathbf{x} \in \mathbb{C}^n$  is nonzero and define  $\rho \in \mathbb{C}$  and  $\mathbf{z}, \mathbf{u} \in \mathbb{C}^n$  by

$$\rho := \begin{cases} x_1/|x_1|, & \text{if } x_1 \neq 0, \\ 1, & \text{otherwise.} \end{cases}, \quad \mathbf{z} := \bar{\rho}\mathbf{x}/\|\mathbf{x}\|_2, \quad \mathbf{u} := \frac{\mathbf{z} + \mathbf{e}_1}{\sqrt{1 + z_1}}. \quad (10.2)$$

Then  $\mathbf{u}^*\mathbf{u} = 2$ ,  $\mathbf{x} = \rho\|\mathbf{x}\|_2\mathbf{z}$  and

$$\mathbf{H}\mathbf{x} := (\mathbf{I} - \mathbf{u}\mathbf{u}^*)\mathbf{x} = a\mathbf{e}_1, \quad a := -\rho\|\mathbf{x}\|_2. \quad (10.3)$$

**Proof.** Since  $|\rho| = 1$  we have  $\rho\|\mathbf{x}\|_2\mathbf{z} = |\rho|^2\mathbf{x} = \mathbf{x}$ . Moreover,  $\|\mathbf{z}\|_2 = 1$  and  $z_1 = |x_1|/\|\mathbf{x}\|_2$  is real so that  $\mathbf{u}^*\mathbf{u} = \frac{(\mathbf{z} + \mathbf{e}_1)^*(\mathbf{z} + \mathbf{e}_1)}{1 + z_1} = \frac{2 + 2z_1}{1 + z_1} = 2$ . Finally,

$$\begin{aligned} \mathbf{H}\mathbf{x} &= \mathbf{x} - (\mathbf{u}^*\mathbf{x})\mathbf{u} = \rho\|\mathbf{x}\|_2(\mathbf{z} - (\mathbf{u}^*\mathbf{z})\mathbf{u}) = \rho\|\mathbf{x}\|_2(\mathbf{z} - \frac{(\mathbf{z}^* + \mathbf{e}_1^*)\mathbf{z}}{1 + z_1}(\mathbf{z} + \mathbf{e}_1)) \\ &= \rho\|\mathbf{x}\|_2(\mathbf{z} - (\mathbf{z} + \mathbf{e}_1)) = -\rho\|\mathbf{x}\|_2\mathbf{e}_1 = a\mathbf{e}_1. \end{aligned}$$

□

The formulas in Theorem 10.3 are implemented in the following algorithm adapted from [26]. To any given  $\mathbf{x} \in \mathbb{C}^n$  a number  $a$  and a vector  $\mathbf{u}$  with  $\mathbf{u}^*\mathbf{u} = 2$  is computed so that  $(\mathbf{I} - \mathbf{u}\mathbf{u}^*)\mathbf{x} = a\mathbf{e}_1$ .

<sup>15</sup>Hint: Show first that  $\mathbf{v}^T\mathbf{v} = 2\mathbf{v}^T\mathbf{x}$

**Algorithm 10.4 (Generate a Householder transformation)**

```

1 function [u, a]=housegen(x)
2 a=norm(x);
3 if a==0
4     u=x; u(1)=sqrt(2); return;
5 end
6 if x(1)==0
7     r=1;
8 else
9     r=x(1)/abs(x(1));
10 end
11 u=conj(r)*x/a;
12 u(1)=u(1)+1;
13 u=u/sqrt(u(1));
14 a=-r*a;
15 end

```

Note that

- If  $\mathbf{x} = \mathbf{0}$  then any  $\mathbf{u}$  with  $\|\mathbf{u}\|_2 = \sqrt{2}$  can be used in the Householder transformation. In the algorithm we use  $\mathbf{u} = \sqrt{2}\mathbf{e}_1$  in this case.
- In Theorem 10.3 the first component of  $\mathbf{z}$  is  $z_1 = |x_1|/\|\mathbf{x}\|_2 \geq 0$ . Since  $\|\mathbf{z}\|_2 = 1$  we have  $1 \leq 1 + z_1 \leq 2$ . It follows that  $\mathbf{u}$  is well defined and we avoid cancellation error when computing  $1 + z_1$ .

**Exercise 10.5 (What does algorithm housegen do when  $\mathbf{x} = \mathbf{e}_1$ ?)**

Determine  $\mathbf{H}$  in Algorithm 10.4 when  $\mathbf{x} = \mathbf{e}_1$ .

Householder transformations can also be used to zero out only the lower part of a vector. Suppose  $\mathbf{x}^T := [\mathbf{y}, \mathbf{z}]^T$ , where  $\mathbf{y} \in \mathbb{C}^k$ ,  $\mathbf{z} \in \mathbb{C}^{n-k}$  for some  $1 \leq k < n$ . The command  $[\hat{\mathbf{u}}, \hat{a}] := \text{housegen}(\mathbf{z})$  defines a Householder transformation  $\hat{\mathbf{H}} = \mathbf{I} - \hat{\mathbf{u}}\hat{\mathbf{u}}^*$  so that  $\hat{\mathbf{H}}\mathbf{z} = \hat{a}\mathbf{e}_1$ . With  $\mathbf{u}^T := [\mathbf{0}, \hat{\mathbf{u}}]^T \in \mathbb{C}^n$  we see that  $\mathbf{u}^*\mathbf{u} = \hat{\mathbf{u}}^*\hat{\mathbf{u}} = 2$ , and

$$\mathbf{H}\mathbf{x} = \begin{bmatrix} \mathbf{y} \\ \hat{a}\mathbf{e}_1 \end{bmatrix}, \text{ where } \mathbf{H} := \mathbf{I} - \mathbf{u}\mathbf{u}^* = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} - \begin{bmatrix} \mathbf{0} \\ \hat{\mathbf{u}} \end{bmatrix} \begin{bmatrix} \mathbf{0} & \hat{\mathbf{u}}^* \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \hat{\mathbf{H}} \end{bmatrix},$$

defines a Householder transformation that produces zeros in the lower part of  $\mathbf{x}$ .

**Exercise 10.6 (Examples of Householder transformations)**

If  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  with  $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2$  and  $\mathbf{v} := \mathbf{x} - \mathbf{y} \neq \mathbf{0}$  then it follows from Exercise 10.2 that  $(\mathbf{I} - 2\frac{\mathbf{v}\mathbf{v}^T}{\mathbf{v}^T\mathbf{v}})\mathbf{x} = \mathbf{y}$ . Use this to construct a Householder transformation  $\mathbf{H}$  such that  $\mathbf{H}\mathbf{x} = \mathbf{y}$  in the following cases.



$$\text{a) } \mathbf{x} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}.$$

$$\text{b) } \mathbf{x} = \begin{bmatrix} 2 \\ 2 \\ 1 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 0 \\ 3 \\ 0 \end{bmatrix}.$$

**Exercise 10.7** ( $2 \times 2$  Householder transformation)

Show that a real  $2 \times 2$  Householder transformation can be written in the form

$$\mathbf{H} = \begin{bmatrix} -\cos \phi & \sin \phi \\ \sin \phi & \cos \phi \end{bmatrix}.$$

Find  $\mathbf{H}\mathbf{x}$  if  $\mathbf{x} = [\cos \phi, \sin \phi]^T$ .

## 10.2 Householder Triangulation

We say that a matrix  $\mathbf{R} \in \mathbb{C}^{m \times n}$  is **upper trapezoidal**, if  $r_{i,j} = 0$  for  $j < i$  and  $i = 1, 2, \dots, m$ . Upper trapezoidal matrices corresponding to  $m < n$ ,  $m = n$ , and  $m > n$  look as follows:

$$\begin{bmatrix} x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \end{bmatrix}, \quad \begin{bmatrix} x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & 0 & x \end{bmatrix}, \quad \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & 0 \end{bmatrix}.$$

In this section we consider a method for bringing a matrix to upper trapezoidal form using Householder transformations. We treat the cases  $m > n$  and  $m \leq n$  separately and consider first  $m > n$ . We describe how to find a sequence  $\mathbf{H}_1, \dots, \mathbf{H}_n$  of Householder transformations such that

$$\mathbf{A}_{n+1} := \mathbf{H}_n \mathbf{H}_{n-1} \cdots \mathbf{H}_1 \mathbf{A} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix} = \mathbf{R},$$

and where  $\mathbf{R}_1$  is upper triangular. We define

$$\mathbf{A}_1 := \mathbf{A}, \quad \mathbf{A}_{k+1} = \mathbf{H}_k \mathbf{A}_k, \quad k = 1, 2, \dots, n.$$

Suppose  $\mathbf{A}_k$  is upper trapezoidal in its first  $k-1$  columns (which is true for  $k=1$ )

$$\mathbf{A}_k = \left[ \begin{array}{ccc|ccc} a_{1,1}^1 & \cdots & a_{1,k-1}^1 & a_{1,k}^1 & \cdots & a_{1,j}^1 & \cdots & a_{1,n}^1 \\ & & \vdots & \vdots & & \vdots & & \vdots \\ & & \ddots & \vdots & & \vdots & & \vdots \\ & & & a_{k-1,k-1}^{k-1} & a_{k-1,k}^{k-1} & \cdots & a_{k-1,j}^{k-1} & \cdots & a_{k-1,n}^{k-1} \\ \hline & & & a_{k,k}^k & \cdots & a_{k,j}^k & \cdots & a_{k,n}^k \\ & & & \vdots & & \vdots & & \vdots \\ & & & a_{i,k}^k & \cdots & a_{i,j}^k & \cdots & a_{i,n}^k \\ & & & \vdots & & \vdots & & \vdots \\ & & & a_{m,k}^k & \cdots & a_{m,j}^k & \cdots & a_{m,n}^k \end{array} \right] \quad (10.4)$$

$$= \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \\ \mathbf{0} & \mathbf{D}_k \end{bmatrix}.$$

Let  $\hat{\mathbf{H}}_k := \mathbf{I} - \hat{\mathbf{u}}_k \hat{\mathbf{u}}_k^*$  be a Householder transformation that maps the first column  $[a_{k,k}^k, \dots, a_{m,k}^k]^T$  of  $\mathbf{D}_k$  to a multiple of  $\mathbf{e}_1$ ,  $\hat{\mathbf{H}}_k(\mathbf{D}_k \mathbf{e}_1) = a_k \mathbf{e}_1$ . Using Algorithm 10.4 we have  $[\hat{\mathbf{u}}_k, a_k] = \text{housegen}(\mathbf{D}_k \mathbf{e}_1)$ . Then  $\mathbf{H}_k := \begin{bmatrix} \mathbf{I}_{k-1} & \mathbf{0} \\ \mathbf{0} & \hat{\mathbf{H}}_k \end{bmatrix}$  is a Householder transformation and

$$\mathbf{A}_{k+1} := \mathbf{H}_k \mathbf{A}_k = \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \\ \mathbf{0} & \hat{\mathbf{H}}_k \mathbf{D}_k \end{bmatrix} = \begin{bmatrix} \mathbf{B}_{k+1} & \mathbf{C}_{k+1} \\ \mathbf{0} & \mathbf{D}_{k+1} \end{bmatrix},$$

where  $\mathbf{B}_{k+1} \in \mathbb{C}^{k \times k}$  is upper triangular and  $\mathbf{D}_{k+1} \in \mathbb{C}^{(m-k) \times (n-k)}$ . Thus  $\mathbf{A}_{k+1}$  is upper trapezoidal in its first  $k$  columns and the reduction has been carried one step further. At the end  $\mathbf{R} := \mathbf{A}_{n+1} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix}$ , where  $\mathbf{R}_1$  is upper triangular.

The process can also be applied to  $\mathbf{A} \in \mathbb{C}^{m \times n}$  if  $m \leq n$ . In this case  $m-1$  Householder transformations will suffice and  $\mathbf{H}_{m-1} \cdots \mathbf{H}_1 \mathbf{A}$  is upper trapezoidal.

In an algorithm we can store most of the vectors  $\hat{\mathbf{u}}_k = [u_{k,k}, \dots, u_{m,k}]^T$  and  $\mathbf{A}_k$  in  $\mathbf{A}$ . However, the elements  $u_{k,k}$  and  $a_k = r_{k,k}$  have to compete for the diagonal in  $\mathbf{A}$ . For  $m=4$  and  $n=3$  the two possibilities look as follows:

$$\mathbf{A} = \begin{bmatrix} u_{11} & r_{12} & r_{13} \\ u_{21} & u_{22} & r_{23} \\ u_{31} & u_{32} & u_{33} \\ u_{41} & u_{42} & u_{43} \end{bmatrix} \quad \text{or} \quad \mathbf{A} = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ u_{21} & r_{22} & r_{23} \\ u_{31} & u_{32} & r_{33} \\ u_{41} & u_{42} & u_{43} \end{bmatrix}.$$

Whatever alternative is chosen, if the looser is needed, it has to be stored in a separate vector. In the following algorithm we store  $a_k = r_{k,k}$  in  $\mathbf{A}$ . We also apply the Householder transformations to a second matrix  $\mathbf{B}$ . The algorithm can then be used to solve linear systems and least squares problems with one or more right hand sides, or to compute the product of the Householder transformations by choosing  $\mathbf{B} = \mathbf{I}$ .

**Algorithm 10.8 (Householder triangulation)**

Suppose  $\mathbf{A} \in \mathbb{C}^{m \times n}$ ,  $\mathbf{B} \in \mathbb{C}^{m \times r}$  and let  $s := \min(n, m - 1)$ . The algorithm uses `housegen` to compute Householder transformations  $\mathbf{H}_1, \dots, \mathbf{H}_s$  such that  $\mathbf{R} = \mathbf{H}_s \cdots \mathbf{H}_1 \mathbf{A}$  is upper trapezoidal and  $\mathbf{C} = \mathbf{H}_s \cdots \mathbf{H}_1 \mathbf{B}$ . If  $\mathbf{B}$  is the empty matrix then  $\mathbf{C}$  is the empty matrix with  $m$  rows and 0 columns.

```

1 function [R,C] = housetriang(A,B)
2 [m,n]=size(A); r=size(B,2); A=[A,B];
3 for k=1:min(n,m-1)
4     [v,A(k,k)]=housegen(A(k:m,k));
5     C=A(k:m,k+1:n+r); A(k:m,k+1:n+r)=C-v*(v'*C);
6 end
7 R=triu(A(:,1:n)); C=A(:,n+1:n+r);

```

Here  $v = \hat{u}_k$  and we have used  $\hat{H}_k \mathbf{C} = (\mathbf{I} - \mathbf{v}\mathbf{v}^*)\mathbf{C} = \mathbf{C} - \mathbf{v}(\mathbf{v}^*\mathbf{C})$  for the update. The Matlab command `triu` extracts the upper triangular part of  $\mathbf{A}$  putting zeros in rows  $n + 1, \dots, m$ .

**10.2.1 Solving linear systems using unitary transformations**

Consider now the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , where  $\mathbf{A}$  is square. Using Algorithm 10.8 we obtain an upper triangular system  $\mathbf{R}\mathbf{x} = \mathbf{c}$  that is nonsingular if  $\mathbf{A}$  is nonsingular. Thus, it can be solved by back substitution and we have a method for solving linear systems that is an alternative to Gaussian elimination. The two methods are similar since they both reduce  $\mathbf{A}$  to upper triangular form using certain transformations and they both work for nonsingular systems.

Which method is better? Here is a short discussion.

- Advantages with Householder:
  - Row interchanges are not necessary, but see [5].
  - Numerically stable.
- Advantages with Gauss
  - Half the number of arithmetic operations compared to Householder.
  - Row interchanges are often not necessary.
  - Usually stable (but no guarantee).

Linear systems can be constructed where Gaussian elimination will fail numerically even if row interchanges are used, see [34]. On the other hand the transformations used in Householder triangulation are unitary so the method is quite stable. So why is Gaussian elimination more popular than Householder triangulation? One reason is that the number of arithmetic operations in (10.5)

when  $m = n$  is  $4n^3/3 = 2G_n$ , which is twice the number for Gaussian elimination. We show this below. Numerical stability can be a problem with Gaussian elimination, but years and years of experience shows that it works well for most practical problems and pivoting is often not necessary. Also Gaussian elimination often wins for banded and sparse problems.

## 10.2.2 The number of arithmetic operations

The bulk of the work in Algorithm 10.8 is the computation of  $\mathbf{C} - \mathbf{v} * (\mathbf{v}^T * \mathbf{C})$  for each  $k$ . In the real case it can be determined from the following lemma.

### Lemma 10.9 (Updating a Householder transformation)

Suppose  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{u} \in \mathbb{R}^m$  and  $\mathbf{v} \in \mathbb{R}^n$ . The computation of  $\mathbf{A} - \mathbf{u}(\mathbf{u}^T \mathbf{A})$  and  $\mathbf{A} - (\mathbf{A}\mathbf{v})\mathbf{v}^T$  both cost approximately  $4mn$  arithmetic operations.

**Proof.** It costs  $2mn$  arithmetic operations to compute  $\mathbf{w}^T := \mathbf{u}^T \mathbf{A}$ ,  $mn$  arithmetic operations to compute  $\mathbf{W} = \mathbf{u}\mathbf{w}^T$  and  $mn$  arithmetic operations for the final subtraction  $\mathbf{A} - \mathbf{W}$ , a total of  $4mn$  arithmetic operations. Taking the transpose we obtain the same count for  $\mathbf{A} - (\mathbf{A}\mathbf{v})\mathbf{v}^T$ .  $\square$

Since in Algorithm 10.8,  $\mathbf{C} \in \mathbb{C}^{(m-k+1) \times (n+r-k)}$  and  $m \geq n$  the cost of computing the update  $\mathbf{C} - \mathbf{v} * (\mathbf{v}^T * \mathbf{C})$  is  $4(m-k)(n+r-k)$  arithmetic operations. This implies that the work in Algorithm 10.8 can be estimated as

$$\int_0^n 4(m-k)(n+r-k)dk = 2m(n+r)^2 - \frac{2}{3}(n+r)^3. \quad (10.5)$$

For  $m = n$  and  $r = 0$  this gives  $4n^3/3$  for the number of arithmetic operations to bring a matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  to upper triangular form using Householder transformations.

## 10.3 The QR Decomposition and QR Factorization

Gaussian elimination without row interchanges results in an LU factorization  $\mathbf{A} = \mathbf{L}\mathbf{U}$  of  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Consider Householder triangulation of  $\mathbf{A}$ . Applying Algorithm 10.8 gives  $\mathbf{R} = \mathbf{H}_{n-1} \cdots \mathbf{H}_1 \mathbf{A}$  implying the factorization  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ , where  $\mathbf{Q} = \mathbf{H}_1 \cdots \mathbf{H}_{n-1}$  is orthonormal and  $\mathbf{R}$  is upper triangular. This is known as a QR-factorization of  $\mathbf{A}$ .

### 10.3.1 Existence

For a rectangular matrix we define the following.

**Definition 10.10 (QR decomposition)**

Let  $\mathbf{A} \in \mathbb{C}^{m \times n}$  with  $m, n \in \mathbb{N}$ . We say that  $\mathbf{A} = \mathbf{QR}$  is a **QR decomposition** of  $\mathbf{A}$  if  $\mathbf{Q} \in \mathbb{C}^{m, m}$  is square and unitary and  $\mathbf{R}$  is upper trapezoidal. If  $m \geq n$  then  $\mathbf{R}$  takes the form

$$\mathbf{R} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0}_{m-n, n} \end{bmatrix}$$

where  $\mathbf{R}_1 \in \mathbb{C}^{n \times n}$  is upper triangular and  $\mathbf{0}_{m-n, n}$  is the zero matrix with  $m - n$  rows and  $n$  columns. For  $m \geq n$  we call  $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$  a **QR factorization** of  $\mathbf{A}$  if  $\mathbf{Q}_1 \in \mathbb{C}^{m \times n}$  has orthonormal columns and  $\mathbf{R}_1 \in \mathbb{C}^{n \times n}$  is upper triangular.

Suppose  $m \geq n$ . A QR factorization is obtained from a QR decomposition  $\mathbf{A} = \mathbf{QR}$  by simply using the first  $n$  columns of  $\mathbf{Q}$  and the first  $n$  rows of  $\mathbf{R}$ . Indeed, if we partition  $\mathbf{Q}$  as  $[\mathbf{Q}_1, \mathbf{Q}_2]$  and  $\mathbf{R} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix}$ , where  $\mathbf{Q}_1 \in \mathbb{R}^{m \times n}$  and  $\mathbf{R}_1 \in \mathbb{R}^{n \times n}$  then  $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$  is a QR factorization of  $\mathbf{A}$ . On the other hand a QR factorization  $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$  of  $\mathbf{A}$  can be turned into a QR decomposition by extending the set of columns  $\{\mathbf{q}_1, \dots, \mathbf{q}_n\}$  of  $\mathbf{Q}_1$  into an orthonormal basis  $\{\mathbf{q}_1, \dots, \mathbf{q}_n, \mathbf{q}_{n+1}, \dots, \mathbf{q}_m\}$  for  $\mathbb{R}^m$  and adding  $m - n$  rows of zeros to  $\mathbf{R}_1$ . We then obtain the QR decomposition  $\mathbf{A} = \mathbf{QR}$ , where  $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_m]$  and  $\mathbf{R} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix}$ .

**Example 10.11 (QR decomposition and factorization)**

An example of a QR decomposition is

$$\mathbf{A} = \begin{bmatrix} 1 & 3 & 1 \\ 1 & 3 & 7 \\ 1 & -1 & -4 \\ 1 & -1 & 2 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 & 1 & -1 & -1 \\ 1 & 1 & 1 & 1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{bmatrix} \times \begin{bmatrix} 2 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 6 \\ 0 & 0 & 0 \end{bmatrix} = \mathbf{QR},$$

while a QR factorization  $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$  is obtained by dropping the last column of  $\mathbf{Q}$  and the last row of  $\mathbf{R}$ , so that

$$\mathbf{A} = \frac{1}{2} \begin{bmatrix} 1 & 1 & -1 \\ 1 & 1 & 1 \\ 1 & -1 & -1 \\ 1 & -1 & 1 \end{bmatrix} \times \begin{bmatrix} 2 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 6 \end{bmatrix} = \mathbf{Q}_1 \mathbf{R}_1.$$

Consider existence and uniqueness.

**Theorem 10.12 (Existence of QR decomposition)**

Any matrix  $\mathbf{A} \in \mathbb{C}^{m \times n}$  with  $m, n \in \mathbb{N}$  has a QR decomposition.

**Proof.** The function `housegen(x)` returns a Householder transformation for any  $\mathbf{x} \in \mathbb{C}^n$ . Thus with  $\mathbf{B} = \mathbf{I}$  in Algorithm 10.8 we obtain a QR decomposition  $\mathbf{A} = \mathbf{QR}$ , where  $\mathbf{Q} = \mathbf{C}^* = \mathbf{H}_1 \cdots \mathbf{H}_s$ , is unitary. Thus a QR decomposition always exists.  $\square$

**Theorem 10.13 (Uniqueness of QR factorization)**

If  $m \geq n$  and  $\mathbf{A}$  is real then the QR factorization is unique if  $\mathbf{A}$  has linearly independent columns and  $\mathbf{R}$  has positive diagonal elements.

*Proof.* Let  $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$  be a QR factorization of  $\mathbf{A}$ . Note that  $\mathbf{A}^T \mathbf{A} = \mathbf{R}_1^T \mathbf{Q}_1^T \mathbf{Q}_1 \mathbf{R}_1 = \mathbf{R}_1^T \mathbf{R}_1$ . Since  $\mathbf{A}^T \mathbf{A}$  is symmetric positive definite the matrix  $\mathbf{R}_1$  is nonsingular, and if its diagonal elements are positive this is the Cholesky factorization of  $\mathbf{A}^T \mathbf{A}$ . Since the Cholesky factorization is unique it follows that  $\mathbf{R}_1$  is unique and since necessarily  $\mathbf{Q}_1 = \mathbf{A} \mathbf{R}_1^{-1}$ , it must also be unique.  $\square$

**Example 10.14 (QR decomposition and factorization)**

Consider finding the QR decomposition and factorization of the matrix  $\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$  using the method of the uniqueness proof of Theorem 10.12. We find  $\mathbf{B} := \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 5 & -4 \\ -4 & 5 \end{bmatrix}$ . The Cholesky factorization of  $\mathbf{B} = \mathbf{R}^T \mathbf{R}$  is given by  $\mathbf{R} = \frac{1}{\sqrt{5}} \begin{bmatrix} 5 & -4 \\ 0 & 3 \end{bmatrix}$ . Now  $\mathbf{R}^{-1} = \frac{1}{3\sqrt{5}} \begin{bmatrix} 3 & 4 \\ 0 & 5 \end{bmatrix}$  so  $\mathbf{Q} = \mathbf{A} \mathbf{R}^{-1} = \frac{1}{\sqrt{5}} \begin{bmatrix} 2 & 1 \\ -1 & 2 \end{bmatrix}$ . Since  $\mathbf{A}$  is square  $\mathbf{A} = \mathbf{Q} \mathbf{R}$  is both the QR decomposition and QR factorization of  $\mathbf{A}$ .

The QR factorization can be used to prove a classical determinant inequality.

**Theorem 10.15 (Hadamard's inequality)**

For any  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{C}^{n \times n}$  we have

$$|\det(\mathbf{A})| \leq \prod_{j=1}^n \|\mathbf{a}_j\|_2. \quad (10.6)$$

Equality holds if and only if  $\mathbf{A}$  has a zero column or the columns of  $\mathbf{A}$  are orthogonal.

*Proof.* Let  $\mathbf{A} = \mathbf{Q} \mathbf{R}$  be a QR factorization of  $\mathbf{A}$ . Since

$$1 = \det(\mathbf{I}) = \det(\mathbf{Q}^* \mathbf{Q}) = \det(\mathbf{Q}^*) \det(\mathbf{Q}) = \det(\mathbf{Q})^* \det(\mathbf{Q}) = |\det(\mathbf{Q})|^2$$

we have  $|\det(\mathbf{Q})| = 1$ . Let  $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_n]$ . Then  $(\mathbf{A}^* \mathbf{A})_{jj} = \|\mathbf{a}_j\|_2^2 = (\mathbf{R}^* \mathbf{R})_{jj} = \|\mathbf{r}_j\|_2^2$ , and

$$|\det(\mathbf{A})| = |\det(\mathbf{Q} \mathbf{R})| = |\det(\mathbf{R})| = \prod_{j=1}^n |r_{jj}| \leq \prod_{j=1}^n \|\mathbf{r}_j\|_2 = \prod_{j=1}^n \|\mathbf{a}_j\|_2.$$

The inequality is proved. We clearly have equality if  $\mathbf{A}$  has a zero column, for then both sides of (10.6) are zero. Suppose the columns are nonzero. We have equality if and only if  $r_{jj} = \|\mathbf{r}_j\|_2$  for  $j = 1, \dots, n$ . This happens if and only if  $\mathbf{R}$  is diagonal. But then  $\mathbf{A}^* \mathbf{A} = \mathbf{R}^* \mathbf{R}$  is diagonal, which means that the columns of  $\mathbf{A}$  are orthogonal.  $\square$

**Exercise 10.16 (QR decomposition)**

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 1 & 2 \\ 1 & 0 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{Q} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{bmatrix}, \quad \mathbf{R} = \begin{bmatrix} 2 & 2 \\ 0 & 2 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}.$$

Show that  $\mathbf{Q}$  is orthonormal and that  $\mathbf{QR}$  is a QR decomposition of  $\mathbf{A}$ . Find a QR factorization of  $\mathbf{A}$ .

**Exercise 10.17 (Householder triangulation)**

a) Let

$$\mathbf{A} := \begin{bmatrix} 1 & 0 & 1 \\ -2 & -1 & 0 \\ 2 & 2 & 1 \end{bmatrix}.$$

Find Householder transformations  $\mathbf{H}_1, \mathbf{H}_2 \in \mathbb{R}^{3 \times 3}$  such that  $\mathbf{H}_2 \mathbf{H}_1 \mathbf{A}$  is upper triangular.

b) Find the QR factorization of  $\mathbf{A}$ , when  $\mathbf{R}$  has positive diagonal elements.

**10.3.2 QR and Gram-Schmidt**

The Gram-Schmidt orthogonalization of the columns of  $\mathbf{A}$  can be used to find the QR factorization of  $\mathbf{A}$ .

**Theorem 10.18 (QR and Gram-Schmidt)**

Suppose  $\mathbf{A} \in \mathbb{R}^{m \times n}$  has rank  $n$  and let  $\mathbf{v}_1, \dots, \mathbf{v}_n$  be the result of applying Gram Schmidt to the columns  $\mathbf{a}_1, \dots, \mathbf{a}_n$  of  $\mathbf{A}$ , i. e.,

$$\mathbf{v}_1 = \mathbf{a}_1, \quad \mathbf{v}_j = \mathbf{a}_j - \sum_{i=1}^{j-1} \frac{\mathbf{a}_j^T \mathbf{v}_i}{\mathbf{v}_i^T \mathbf{v}_i} \mathbf{v}_i, \quad \text{for } j = 2, \dots, n. \quad (10.7)$$

Let

$$\mathbf{Q}_1 := [\mathbf{q}_1, \dots, \mathbf{q}_n], \quad \mathbf{q}_j := \frac{\mathbf{v}_j}{\|\mathbf{v}_j\|_2}, \quad j = 1, \dots, n \text{ and}$$

$$\mathbf{R}_1 := \begin{bmatrix} \|\mathbf{v}_1\|_2 & \mathbf{a}_2^T \mathbf{q}_1 & \mathbf{a}_3^T \mathbf{q}_1 & \cdots & \mathbf{a}_{n-1}^T \mathbf{q}_1 & \mathbf{a}_n^T \mathbf{q}_1 \\ 0 & \|\mathbf{v}_2\|_2 & \mathbf{a}_3^T \mathbf{q}_2 & \cdots & \mathbf{a}_{n-1}^T \mathbf{q}_2 & \mathbf{a}_n^T \mathbf{q}_2 \\ & 0 & \|\mathbf{v}_3\|_2 & \cdots & \mathbf{a}_{n-1}^T \mathbf{q}_3 & \mathbf{a}_n^T \mathbf{q}_3 \\ & & & \ddots & \vdots & \vdots \\ & & & & \ddots & \|\mathbf{v}_{n-1}\|_2 & \mathbf{a}_n^T \mathbf{q}_{n-1} \\ & & & & & 0 & \|\mathbf{v}_n\|_2 \end{bmatrix}. \quad (10.8)$$

Then  $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$  is the unique QR factorization of  $\mathbf{A}$ .

**Proof.** Let  $\mathbf{Q}_1$  and  $\mathbf{R}_1$  be given by (10.8). The matrix  $\mathbf{Q}_1$  is well defined and has orthonormal columns, since  $\{\mathbf{q}_1, \dots, \mathbf{q}_n\}$  is an orthonormal basis for  $\text{span}(\mathbf{A})$  by Theorem 0.29. By (10.7)

$$\mathbf{a}_j = \mathbf{v}_j + \sum_{i=1}^{j-1} \frac{\mathbf{a}_j^T \mathbf{v}_i}{\mathbf{v}_i^T \mathbf{v}_i} \mathbf{v}_i = r_{jj} \mathbf{q}_j + \sum_{i=1}^{j-1} \mathbf{q}_i r_{ij} = \mathbf{Q}_1 \mathbf{R}_1 \mathbf{e}_j, \quad j = 1, \dots, n.$$

Clearly  $\mathbf{R}_1$  has positive diagonal elements and the factorization is unique.  $\square$

### Example 10.19 (QR using Gram-Schmidt)

Consider finding the QR decomposition and factorization of the matrix  $\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} = [\mathbf{a}_1, \mathbf{a}_2]$  using Gram-Schmidt. Using (10.7) we find  $\mathbf{v}_1 = \mathbf{a}_1$  and  $\mathbf{v}_2 = \mathbf{a}_2 - \frac{\mathbf{a}_2^T \mathbf{v}_1}{\mathbf{v}_1^T \mathbf{v}_1} \mathbf{v}_1 = \frac{3}{5} \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ . Thus  $\mathbf{Q} = [\mathbf{q}_1, \mathbf{q}_2]$ , where  $\mathbf{q}_1 = \frac{1}{\sqrt{5}} \begin{bmatrix} 2 \\ -1 \end{bmatrix}$  and  $\mathbf{q}_2 = \frac{1}{\sqrt{5}} \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ . By (10.8) we find

$$\mathbf{R}_1 = \mathbf{R} = \begin{bmatrix} \|\mathbf{v}_1\|_2 & \mathbf{a}_2^T \mathbf{q}_1 \\ 0 & \|\mathbf{v}_2\|_2 \end{bmatrix} = \frac{1}{\sqrt{5}} \begin{bmatrix} 5 & -4 \\ 0 & 3 \end{bmatrix}$$

and this agrees with what we found in Example 10.14.

### Exercise 10.20 (QR using Gram-Schmidt, II)

Construct  $\mathbf{Q}_1$  and  $\mathbf{R}_1$  in Example 10.11 using Gram-Schmidt orthogonalization.

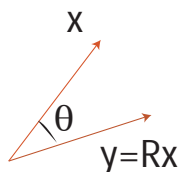
**Warning.** The Gram-Schmidt orthogonalization process should not be used to compute the QR factorization numerically. The columns of  $\mathbf{Q}_1$  computed in floating point arithmetic using Gram-Schmidt orthogonalization will often be far from orthogonal. There is a modified version of Gram-Schmidt which behaves better numerically, see [2]. Here we only considered Householder transformations (cf. Algorithm 10.8).

## 10.4 Givens Rotations

In some applications, the matrix we want to triangulate has a special structure. Suppose for example that  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is square and upper Hessenberg as illustrated by a **Wilkinson diagram** for  $n = 4$

$$\mathbf{A} = \begin{bmatrix} x & x & x & x \\ x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \end{bmatrix}.$$





**Figure 10.2.** A plane rotation.

Only one element in each column needs to be annihilated and a full Householder transformation will be inefficient. In this case we can use a simpler transformation.

**Definition 10.21 (Givens rotation, plane rotation)**

A **plane rotation** (also called a **Given's rotation**) is a matrix  $\mathbf{P} \in \mathbb{R}^{2,2}$  of the form

$$\mathbf{P} := \begin{bmatrix} c & s \\ -s & c \end{bmatrix}, \text{ where } c^2 + s^2 = 1.$$

A plane rotation is orthonormal and there is a unique angle  $\theta \in [0, 2\pi)$  such that  $c = \cos \theta$  and  $s = \sin \theta$ . Moreover, the identity matrix is a plane rotation corresponding to  $\theta = 0$ .

**Exercise 10.22 (Plane rotation)**

Show that if  $\mathbf{x} = \begin{bmatrix} r \cos \alpha \\ r \sin \alpha \end{bmatrix}$  then  $\mathbf{P}\mathbf{x} = \begin{bmatrix} r \cos(\alpha - \theta) \\ r \sin(\alpha - \theta) \end{bmatrix}$ . Thus  $\mathbf{P}$  rotates a vector  $\mathbf{x}$  in the plane an angle  $\theta$  clockwise. See Figure 10.2.

Suppose

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \neq \mathbf{0}, \quad c := \frac{x_1}{r}, \quad s := \frac{x_2}{r}, \quad r := \|\mathbf{x}\|_2.$$

Then

$$\mathbf{P}\mathbf{x} = \frac{1}{r} \begin{bmatrix} x_1 & x_2 \\ -x_2 & x_1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \frac{1}{r} \begin{bmatrix} x_1^2 + x_2^2 \\ 0 \end{bmatrix} = \begin{bmatrix} r \\ 0 \end{bmatrix},$$

and we have introduced a zero in  $\mathbf{x}$ . We can take  $\mathbf{P} = \mathbf{I}$  when  $\mathbf{x} = \mathbf{0}$ .

For an  $n$ -vector  $\mathbf{x} \in \mathbb{R}^n$  and  $1 \leq i < j \leq n$  we define a **rotation in the  $i, j$ -plane** as a matrix  $\mathbf{P}_{ij} = (p_{kl}) \in \mathbb{R}^{n \times n}$  by  $p_{kl} = \delta_{kl}$  except for positions  $ii, jj, ij, ji$ , which are given by

$$\begin{bmatrix} p_{ii} & p_{ij} \\ p_{ji} & p_{jj} \end{bmatrix} = \begin{bmatrix} c & s \\ -s & c \end{bmatrix}, \text{ where } c^2 + s^2 = 1.$$

Thus, for  $n = 4$ ,

$$\mathbf{P}_{1,2} = \begin{bmatrix} c & s & 0 & 0 \\ -s & c & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{P}_{13} = \begin{bmatrix} c & 0 & s & 0 \\ 0 & 1 & 0 & 0 \\ -s & 0 & c & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{P}_{23} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & s & c & 0 \\ 0 & -s & c & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$



Karl Adolf Hessenberg, 1904-1959 (left), James Wallace Givens, Jr., 1910-1993 (right)

Premultiplying a matrix by a rotation in the  $i, j$ -plane changes only rows  $i$  and  $j$  of the matrix, while postmultiplying the matrix by such a rotation only changes column  $i$  and  $j$ . In particular, if  $\mathbf{B} = \mathbf{P}_{ij}\mathbf{A}$  and  $\mathbf{C} = \mathbf{A}\mathbf{P}_{ij}$  then  $\mathbf{B}(k, :) = \mathbf{A}(k, :)$ ,  $\mathbf{C}(:, k) = \mathbf{A}(:, k)$  for all  $k \neq i, j$  and

$$\begin{bmatrix} \mathbf{B}(i, :) \\ \mathbf{B}(j, :) \end{bmatrix} = \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} \mathbf{A}(i, :) \\ \mathbf{A}(j, :) \end{bmatrix}, \quad [\mathbf{C}(:, i) \quad \mathbf{C}(:, j)] = [\mathbf{A}(:, i) \quad \mathbf{A}(:, j)] \begin{bmatrix} c & s \\ -s & c \end{bmatrix}. \quad (10.9)$$

An upper Hessenberg matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  can be transformed to upper triangular form using rotations  $\mathbf{P}_{i,i+1}$  for  $i = 1, \dots, n-1$ . For  $n = 4$  the process can be illustrated as follows.

$$\mathbf{A} = \begin{bmatrix} x & x & x & x \\ x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \end{bmatrix} \xrightarrow{\mathbf{P}_{12}} \begin{bmatrix} r_{11} & r_{12} & r_{13} & r_{14} \\ \mathbf{0} & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & x & x \end{bmatrix} \xrightarrow{\mathbf{P}_{23}} \begin{bmatrix} r_{11} & r_{12} & r_{13} & r_{14} \\ 0 & r_{22} & r_{23} & r_{24} \\ 0 & \mathbf{0} & x & x \\ 0 & 0 & x & x \end{bmatrix} \xrightarrow{\mathbf{P}_{34}} \begin{bmatrix} r_{11} & r_{12} & r_{13} & r_{14} \\ 0 & r_{22} & r_{23} & r_{24} \\ 0 & 0 & r_{33} & r_{34} \\ 0 & 0 & \mathbf{0} & r_{44} \end{bmatrix}.$$

For an algorithm see Exercise 10.23.

### Exercise 10.23 (Solving upper Hessenberg system using rotations)

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be upper Hessenberg and nonsingular, and let  $\mathbf{b} \in \mathbb{R}^n$ . The following algorithm solves the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  using rotations  $\mathbf{P}_{k,k+1}$  for  $k = 1, \dots, n-1$ . It uses the back solve algorithm 2.2. Determine the number of arithmetic operations of this algorithm.

**Algorithm 10.24 (Upper Hessenberg linear system)**

```

1 function x=rothesstri(A,b)
2 n=length(A); A=[A b];
3 for k=1:n-1
4     r=norm([A(k,k),A(k+1,k)]);
5     if r>0
6         c=A(k,k)/r; s=A(k+1,k)/r;
7         A([k k+1],k+1:n+1)=[c s;-s c]*A([k k+1],k+1:n+1);
8     end
9     A(k,k)=r; A(k+1,k)=0;
10 end
11 x=backsolve(A(:,1:n),A(:,n+1),n);

```

**10.5 Review Questions****10.5.1** What is a Householder transformation?**10.5.2** Why are they good for numerical work?**10.5.3** What are the main differences between solving a linear system by Gaussian elimination and Householder transformations?**10.5.4** What are the differences between a QR decomposition and a QR factorization?**10.5.5** Does any matrix have a QR decomposition?**10.5.6** What is a Givens transformation?



## Chapter 11

# Least Squares

Consider the linear system  $\mathbf{Ax} = \mathbf{b}$  of  $m$  equations in  $n$  unknowns. It is overdetermined, if  $m > n$ , square, if  $m = n$ , and underdetermined, if  $m < n$ . In either case the system can only be solved approximately if  $\mathbf{b} \notin \text{span}(\mathbf{A})$ . One way to solve  $\mathbf{Ax} = \mathbf{b}$  approximately is to select a vector norm  $\|\cdot\|$ , say a  $p$ -norm, and look for  $\mathbf{x} \in \mathbb{C}^n$  which minimizes  $\|\mathbf{Ax} - \mathbf{b}\|$ . The use of the one and  $\infty$  norm can be formulated as linear programming problems, while the Euclidian norm leads to a linear system. Only this norm is considered here.

### Definition 11.1 (Least squares problem)

Suppose  $m, n \in \mathbb{N}$ ,  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and  $\mathbf{b} \in \mathbb{C}^m$ . To find  $\mathbf{x} \in \mathbb{C}^n$  that minimizes  $E : \mathbb{C}^n \rightarrow \mathbb{R}$  given by

$$E(\mathbf{x}) := \|\mathbf{Ax} - \mathbf{b}\|_2^2,$$

is called the **least squares problem**. A minimizer  $\mathbf{x}$  is called a **least squares solution**.

Since the square root function is monotone, minimizing  $E(\mathbf{x})$  or  $\sqrt{E(\mathbf{x})}$  is equivalent.

One way to solve the least squares problem is to write  $E$  as a quadratic function and set partial derivatives equal to zero. If  $\mathbf{A}$  and  $\mathbf{b}$  have real components we find

$$E(\mathbf{x}) := (\mathbf{Ax} - \mathbf{b})^T(\mathbf{Ax} - \mathbf{b}) = \mathbf{x}^T \mathbf{Bx} - 2\mathbf{c}^T \mathbf{x} + \beta,$$

where

$$\mathbf{B} = \mathbf{A}^T \mathbf{A}, \quad \mathbf{c} = \mathbf{A}^T \mathbf{b}, \quad \beta = \mathbf{b}^T \mathbf{b}.$$

By Lemma 9.2 all minimums are solutions of the linear system  $\mathbf{Bx} = \mathbf{c}$  or

$$\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}$$

known as the **normal equations**. The coefficient matrix is symmetric positive semidefinite. If it is positive definite then the least square problem has a unique solution. By Corollary 2.36  $\mathbf{A}^T \mathbf{A}$  is positive definite if and only if  $\mathbf{A}$  has linearly independent columns. In particular,  $m \geq n$  is necessary for a unique solution.

## 11.1 Numerical Examples

### Example 11.2 (Average)

Consider the least squares problem defined by

$$\begin{array}{l} x_1 = 1 \\ x_1 = 1, \\ x_1 = 2 \end{array} \quad \mathbf{A} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \quad \mathbf{x} = [x_1], \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix},$$

We find

$$\|\mathbf{Ax} - \mathbf{b}\|_2^2 = (x_1 - 1)^2 + (x_1 - 1)^2 + (x_1 - 2)^2 = 3x_1^2 - 8x_1 + 6.$$

Setting the first derivative with respect to  $x_1$  equal to zero we obtain  $6x_1 - 8 = 0$  or  $x_1 = 4/3$ , the average of  $b_1, b_2, b_3$ . The second derivative is positive and  $x_1 = 4/3$  is a global minimum. The normal equation is  $3x_1 = 4$ .

### Example 11.3 (Input/output model)

Suppose we have a simple input/output model. To every input  $\mathbf{u} \in \mathbb{R}^n$  we obtain an output  $y \in \mathbb{R}$ . Assuming we have a linear relation

$$y = \mathbf{u}^T \mathbf{x} = \sum_{i=1}^n u_i x_i,$$

between  $\mathbf{u}$  and  $y$ , how can we determine  $\mathbf{x}$ ?

Performing  $m \geq n$  experiments we obtain a table of values

$$\begin{array}{c|c|c|c|c} \mathbf{u} & \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_m \\ \hline y & y_1 & y_2 & \cdots & y_m \end{array}.$$

We would like to find  $\mathbf{x}$  such that

$$\mathbf{Ax} = \begin{bmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \vdots \\ \mathbf{u}_m^T \end{bmatrix} \mathbf{x} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \mathbf{b}.$$

We can estimate  $\mathbf{x}$  by solving the least squares problem  $\min \|\mathbf{Ax} - \mathbf{b}\|_2^2$ .

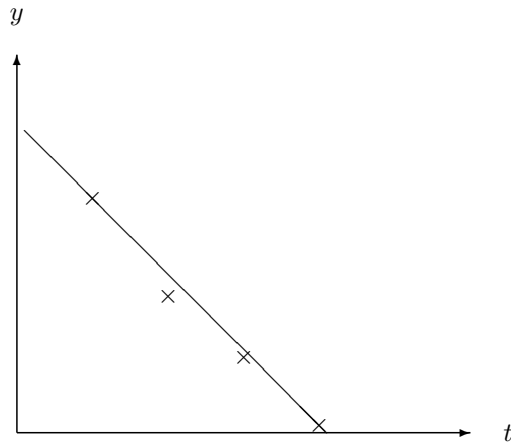
## 11.2 Curve Fitting

Given

- size:  $1 \leq n \leq m$ ,
- sites:  $\mathcal{S} := \{t_1, t_2, \dots, t_m\} \subset [a, b]$ ,
- $y$ -values:  $\mathbf{y} = [y_1, y_2, \dots, y_m]^T \in \mathbb{R}^m$ ,
- functions:  $\phi_j : [a, b] \rightarrow \mathbb{R}$ ,  $j = 1, \dots, n$ .

Find a function (curve fit)  $p : [a, b] \rightarrow \mathbb{R}$  given by  $p := \sum_{j=1}^n x_j \phi_j$  such that  $p(t_k) \approx y_k$  for  $k = 1, \dots, m$ .

An example is shown in Figure 11.1. Here  $\phi_1(t) = 1$  and  $\phi_2(t) = t$  and  $p$  is a straight line (linear regression).



**Figure 11.1.** *A least squares fit to data.*

The curve fitting problem can be defined from an overdetermined linear system:

$$\begin{bmatrix} p(t_1) \\ \vdots \\ p(t_m) \end{bmatrix} = \mathbf{A}\mathbf{x} := \begin{bmatrix} \phi_1(t_1) & \cdots & \phi_n(t_1) \\ \vdots & & \vdots \\ \phi_1(t_m) & \cdots & \phi_n(t_m) \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} =: \mathbf{b}. \quad (11.1)$$

Then we find  $\mathbf{x} \in \mathbb{R}^n$  as a solution of the corresponding least squares problem given by

$$E(\mathbf{x}) := \|\mathbf{Ax} - \mathbf{b}\|_2^2 = \sum_{k=1}^m \left( \sum_{j=1}^n x_j \phi_j(t_k) - y_k \right)^2. \quad (11.2)$$

Typical examples of functions  $\phi_j$  are polynomials, trigonometric functions, exponential functions, or splines.

In (11.2) one can also include **weights**  $w_k > 0$  for  $k = 1, \dots, m$  and minimize

$$E(\mathbf{x}) := \sum_{k=1}^m w_k \left( \sum_{j=1}^n x_j \phi_j(t_k) - y_k \right)^2.$$

If  $y_k$  is an accurate observation, we can choose a large weight  $w_k$ . This will force  $p(t_k) - y_k$  to be small. Similarly, a small  $w_k$  will allow  $p(t_k) - y_k$  to be large. If an estimate for the standard deviation  $\delta y_k$  in  $y_k$  is known for each  $k$ , we can choose  $w_k = 1/(\delta y_k)^2$ ,  $k = 1, 2, \dots, m$ . For simplicity we will assume in the following that  $w_k = 1$  for all  $k$ .

#### Lemma 11.4 (Curve fitting)

Let  $\mathbf{A}$  be given by (11.1). The matrix  $\mathbf{A}^T \mathbf{A}$  is symmetric positive definite if and only if  $\{\phi_1, \dots, \phi_n\}$  is linearly independent on  $\mathcal{S}$ , i. e.,

$$p(t_k) := \sum_{j=1}^n x_j \phi_j(t_k) = 0, \quad k = 1, \dots, m \Rightarrow x_1 = \dots = x_n = 0. \quad (11.3)$$

**Proof.**  $\mathbf{A}$  is positive definite if and only if  $\mathbf{A}$  has linearly independent columns. Since  $(\mathbf{Ax})_k = \sum_{j=1}^n x_j \phi_j(t_k)$ ,  $k = 1, \dots, m$  this is equivalent to (11.3).  $\square$

#### Example 11.5 (Straight line fit)

Consider  $m \geq n = 2$ ,  $\phi_1(t) = 1$ , and  $\phi_2(t) = t$ . The normal equations can be written

$$\begin{bmatrix} m & \sum t_k \\ \sum t_k & \sum t_k^2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \sum y_k \\ \sum t_k y_k \end{bmatrix}. \quad (11.4)$$

Here  $k$  ranges from 1 to  $m$  in the sums. Recall that a nonzero polynomial of degree at most  $n$  has at most  $n$  roots. Therefore, by the Lemma 11.4, this  $2 \times 2$  system is symmetric positive definite if and only if there are at least 2 distinct sites  $t_k$ . With the data

|     |     |     |     |     |
|-----|-----|-----|-----|-----|
| $t$ | 1.0 | 2.0 | 3.0 | 4.0 |
| $y$ | 3.1 | 1.8 | 1.0 | 0.1 |

the normal equations (11.4) become  $\begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 10.1 \end{bmatrix}$ . The data and the least squares polynomial  $p(t) = x_1 + x_2 t = 3.95 - 0.98t$  are shown in Figure 11.1.



**Example 11.6 (Ill conditioning and the Hilbert matrix)**

The normal equations can be extremely ill-conditioned. Consider the curve fitting problem using the polynomials  $\phi_j(t) := t^{j-1}$ , for  $j = 1, \dots, n$  and equidistant sites  $t_k = (k-1)/(m-1)$  for  $k = 1, \dots, m$ . The normal equations are  $\mathbf{B}_n \mathbf{x} = \mathbf{c}_n$ , where for  $n = 3$

$$\mathbf{B}_3 \mathbf{x} := \begin{bmatrix} m & \sum t_k & \sum t_k^2 \\ \sum t_k & \sum t_k^2 & \sum t_k^3 \\ \sum t_k^2 & \sum t_k^3 & \sum t_k^4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} \sum y_k \\ \sum t_k y_k \\ \sum t_k^2 y_k \end{bmatrix}.$$

$\mathbf{B}_n$  is symmetric positive definite if at least  $n$  of the  $t$ 's are distinct. However  $\mathbf{B}_n$  is extremely ill-conditioned even for moderate  $n$ . Indeed,  $\frac{1}{m} \mathbf{B}_n \approx \mathbf{H}_n$ , where  $\mathbf{H}_n \in \mathbb{R}^{n \times n}$  is the **Hilbert Matrix** with  $i, j$  element  $1/(i+j-1)$ . Thus for  $n = 3$

$$\mathbf{H}_3 = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \end{bmatrix}.$$

The elements of  $\frac{1}{m} \mathbf{B}_n$  are related to Riemann sums approximations to the elements of  $\mathbf{H}_n$ . In fact,

$$\frac{1}{m} b_{i,j} = \frac{1}{m} \sum_{k=1}^m t_k^{i+j-2} = \frac{1}{m} \sum_{k=1}^m \left( \frac{k-1}{m-1} \right)^{i+j-2} \approx \int_0^1 x^{i+j-2} dx = \frac{1}{i+j-1} = h_{i,j}.$$

The elements of  $\mathbf{H}_n^{-1}$  are determined in Exercise 0.51. We find  $K_1(\mathbf{H}_6) \approx 3 \cdot 10^7$ . It appears that  $\frac{1}{m} \mathbf{B}_n$  and hence  $\mathbf{B}_n$  is ill-conditioned for moderate  $n$  at least if  $m$  is large. The cure for this problem is to use a different basis for polynomials. Orthogonal polynomials are an excellent choice. Another possibility is to use the shifted power basis  $(t - \tilde{t})^{j-1}$ ,  $j = 1, \dots, n$ , for a suitable  $\tilde{t}$ , see Exercise 11.8.

**Exercise 11.7 (Straight line fit (linear regression))**

Suppose  $(t_i, y_i)_{i=1}^m$  are  $m$  points in the plane. We consider the over-determined systems

$$\begin{array}{ll} \text{(i)} & \begin{array}{l} x_1 = y_1 \\ x_1 = y_2 \\ \vdots \\ x_1 = y_m \end{array} & \text{(ii)} & \begin{array}{l} x_1 + t_1 x_2 = y_1 \\ x_1 + t_2 x_2 = y_2 \\ \vdots \\ x_1 + t_m x_2 = y_m \end{array} \end{array}$$

- a) Find the normal equations for (i) and the least squares solution.
- b) Find the normal equations for (ii) and give a geometric interpretation of the least squares solution.

**Exercise 11.8 (Straight line fit using shifted power form)**

Related to (ii) in Exercise 11.7 we have the overdetermined system

$$(iii) \quad x_1 + (t_i - \hat{t})x_2 = y_i, \quad i = 1, 2, \dots, m,$$

where  $\hat{t} = (t_1 + \dots + t_m)/m$ .

- a) Find the normal equations for (iii) and give a geometric interpretation of the least squares solution.
- b) Fit a straight line to the points  $(t_i, y_i)$ : (998.5, 1), (999.5, 1.9), (1000.5, 3.1) and (1001.5, 3.5) using a). Draw a sketch of the solution.

**Exercise 11.9 (Fitting a circle to points)**

In this problem we derive an algorithm to fit a circle  $(t - c_1)^2 + (y - c_2)^2 = r^2$  to  $m \geq 3$  given points  $(t_i, y_i)_{i=1}^m$  in the  $(t, y)$ -plane. We obtain the overdetermined system

$$(t_i - c_1)^2 + (y_i - c_2)^2 = r^2, \quad i = 1, \dots, m, \quad (11.5)$$

of  $m$  equations in the three unknowns  $c_1, c_2$  and  $r$ . This system is nonlinear, but it can be solved from the linear system

$$t_i x_1 + y_i x_2 + x_3 = t_i^2 + y_i^2, \quad i = 1, \dots, m, \quad (11.6)$$

and then setting  $c_1 = x_1/2$ ,  $c_2 = x_2/2$  and  $r^2 = c_1^2 + c_2^2 + x_3$ .

- a) Derive (11.6) from (11.5). Explain how we can find  $c_1, c_2, r$  once  $[x_1, x_2, x_3]$  is determined.
- b) Formulate (11.6) as a linear least squares problem for suitable  $\mathbf{A}$  and  $\mathbf{b}$ .
- c) Does the matrix  $\mathbf{A}$  in b) have linearly independent columns?
- d) Use (11.6) to find the circle passing through the three points (1, 4), (3, 2), (1, 0).

## 11.3 Least Squares and Singular Value Decomposition and Factorization

The singular value decomposition and factorization can be used to characterize all solution of the least squares problems. We first consider orthogonal projections.

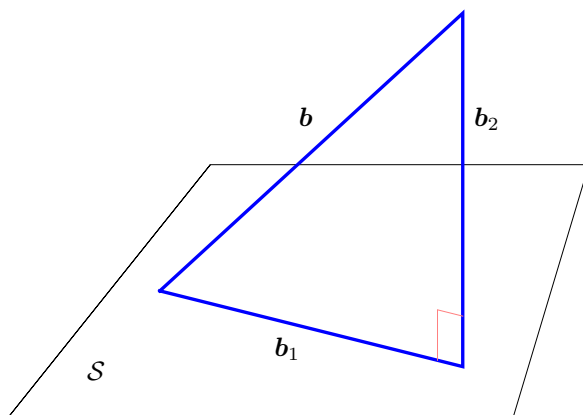


Figure 11.2. The orthogonal projection of  $\mathbf{b}$  into  $\mathcal{S}$ .

### 11.3.1 Sum of subspaces and orthogonal projections

Suppose  $\mathcal{S}$  and  $\mathcal{T}$  are subspaces of  $\mathbb{R}^n$  or  $\mathbb{C}^n$  endowed with an inner product  $\langle \cdot, \cdot \rangle$ . The following subsets are subspaces of  $\mathbb{R}^n$  or  $\mathbb{C}^n$ .

- The set  $\mathcal{S} + \mathcal{T} := \{\mathbf{s} + \mathbf{t} : \mathbf{s} \in \mathcal{S} \text{ and } \mathbf{t} \in \mathcal{T}\}$  is called the **sum** of  $\mathcal{S}$  and  $\mathcal{T}$ .
- If  $\mathcal{S} \cap \mathcal{T} = \{\mathbf{0}\}$  then  $\mathcal{S} + \mathcal{T}$  is called a **direct sum** and denoted  $\mathcal{S} \oplus \mathcal{T}$ .
- If  $\langle \mathbf{s}, \mathbf{t} \rangle = 0$  for all  $\mathbf{s} \in \mathcal{S}$  and  $\mathbf{t} \in \mathcal{T}$  then  $\mathcal{S} + \mathcal{T}$  is called an **orthogonal sum** and denoted  $\mathcal{S} \oplus^\perp \mathcal{T}$ .

Every  $\mathbf{x} \in \mathcal{S} \oplus \mathcal{T}$  can be decomposed uniquely in the form  $\mathbf{x} = \mathbf{s} + \mathbf{t}$ , where  $\mathbf{s} \in \mathcal{S}$  and  $\mathbf{t} \in \mathcal{T}$ . For if  $\mathbf{x} = \mathbf{s}_1 + \mathbf{t}_1 = \mathbf{s}_2 + \mathbf{t}_2$  for  $\mathbf{s}_1, \mathbf{s}_2 \in \mathcal{S}$  and  $\mathbf{t}_1, \mathbf{t}_2 \in \mathcal{T}$ , then  $\mathbf{s}_1 - \mathbf{s}_2 = \mathbf{t}_2 - \mathbf{t}_1$  and it follows that  $\mathbf{s}_1 - \mathbf{s}_2$  and  $\mathbf{t}_2 - \mathbf{t}_1$  belong to both  $\mathcal{S}$  and  $\mathcal{T}$  and hence to  $\mathcal{S} \cap \mathcal{T}$ . But then  $\mathbf{s}_1 - \mathbf{s}_2 = \mathbf{t}_2 - \mathbf{t}_1 = \mathbf{0}$  so  $\mathbf{s}_1 = \mathbf{s}_2$  and  $\mathbf{t}_2 = \mathbf{t}_1$ .

An orthogonal sum is a direct sum. For if  $\mathbf{b} \in \mathcal{S} \cap \mathcal{T}$  then  $\mathbf{b}$  is orthogonal to itself,  $\langle \mathbf{b}, \mathbf{b} \rangle = 0$ , which implies that  $\mathbf{b} = \mathbf{0}$ . Thus, every  $\mathbf{b} \in \mathcal{S} \oplus^\perp \mathcal{T}$  can be written uniquely as  $\mathbf{b} = \mathbf{b}_1 + \mathbf{b}_2$ , where  $\mathbf{b}_1 \in \mathcal{S}$  and  $\mathbf{b}_2 \in \mathcal{T}$ . The vectors  $\mathbf{b}_1$  and  $\mathbf{b}_2$  are called the **orthogonal projections** of  $\mathbf{b}$  into  $\mathcal{S}$  and  $\mathcal{T}$ . For any  $\mathbf{s} \in \mathcal{S}$  we have  $\langle \mathbf{b} - \mathbf{b}_1, \mathbf{s} \rangle = \langle \mathbf{b}_2, \mathbf{s} \rangle = 0$ , see Figure 11.2.

Consider a singular value decomposition of  $\mathbf{A}$  and the corresponding singular value factorization:

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^* = [\mathbf{U}_1, \mathbf{U}_2] \begin{bmatrix} \mathbf{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^* \\ \mathbf{V}_2^* \end{bmatrix} = \mathbf{U}_1\mathbf{\Sigma}_1\mathbf{V}_1^*, \quad \mathbf{\Sigma}_1 = \text{diag}(\sigma_1, \dots, \sigma_r),$$

where  $\mathbf{A}$  has rank  $r$  so that  $\sigma_1 \geq \dots \geq \sigma_r > 0$  and  $\mathbf{U}_1 = [\mathbf{u}_1, \dots, \mathbf{u}_r]$ ,  $\mathbf{V}_1 = [\mathbf{v}_1, \dots, \mathbf{v}_r]$ . We recall (cf. Theorem 6.16)

- the set of columns of  $\mathbf{U}_1$  is an orthonormal basis for the column space  $\text{span}(\mathbf{A})$ ,
- the set of columns of  $\mathbf{U}_2$  is an orthonormal basis for the null space  $\ker(\mathbf{A}^*)$ ,

**Theorem 11.10 (Orthogonal projection and least squares solution)**

Let  $\mathbf{A} \in \mathbb{C}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{C}^m$ ,  $\mathcal{S} := \text{span}(\mathbf{A})$  and  $\mathcal{T} := \ker(\mathbf{A}^*)$ . If  $\mathbf{A} = \mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^*$  is a singular value factorization of  $\mathbf{A}$  then:

1.  $\mathbb{C}^m = \mathcal{S} \oplus \mathcal{T}$  is an **orthogonal decomposition** of  $\mathbb{C}^m$  with respect to the usual inner product  $\langle \mathbf{s}, \mathbf{t} \rangle = \mathbf{t}^* \mathbf{s}$ .
2. The orthogonal projection  $\mathbf{b}_1$  of  $\mathbf{b}$  into  $\mathcal{S}$  is

$$\mathbf{b}_1 = \mathbf{U}_1 \mathbf{U}_1^* \mathbf{b} = \mathbf{A} \mathbf{A}^\dagger \mathbf{b}, \text{ where } \mathbf{A}^\dagger := \mathbf{V}_1 \mathbf{\Sigma}_1^{-1} \mathbf{U}_1^* \in \mathbb{C}^{n \times m}. \quad (11.7)$$

3.  $\mathbf{x} \in \mathbb{C}^n$  is a solution of the least squares problem if and only if  $\mathbf{A} \mathbf{x} = \mathbf{b}_1$ . In particular the least squares problem always has a solution.

*Proof.* By block multiplication  $\mathbf{b} = \mathbf{U} \mathbf{U}^* \mathbf{b} = [\mathbf{U}_1, \mathbf{U}_2] \begin{bmatrix} \mathbf{U}_1^* \\ \mathbf{U}_2^* \end{bmatrix} \mathbf{b} = \mathbf{b}_1 + \mathbf{b}_2$ , where  $\mathbf{b}_1 = \mathbf{U}_1 \mathbf{U}_1^* \mathbf{b} \in \mathcal{S}$  and  $\mathbf{b}_2 = \mathbf{U}_2 \mathbf{U}_2^* \mathbf{b} \in \mathcal{T}$ . Since  $\mathbf{U}_2^* \mathbf{U}_1 = \mathbf{0}$  it follows that  $\langle \mathbf{b}_1, \mathbf{b}_2 \rangle = \mathbf{b}_2^* \mathbf{b}_1 = 0$  implying Part 1. Moreover,  $\mathbf{b}_1$  is the orthogonal projection into  $\mathcal{S}$ . Since  $\mathbf{V}_1^* \mathbf{V}_1 = \mathbf{I}$  we find

$$\mathbf{A} \mathbf{A}^\dagger \mathbf{b} = (\mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^*) (\mathbf{V}_1 \mathbf{\Sigma}_1^{-1} \mathbf{U}_1^*) \mathbf{b} = \mathbf{U}_1 \mathbf{U}_1^* \mathbf{b} = \mathbf{b}_1$$

and Part 2 follows. For any  $\mathbf{x} \in \mathbb{C}^n$  we have  $\mathbf{A} \mathbf{x} - \mathbf{b}_1 \in \mathcal{S}$  and  $\mathbf{b}_2^* (\mathbf{A} \mathbf{x} - \mathbf{b}_1) = 0$  so by Pythagoras

$$\|\mathbf{b} - \mathbf{A} \mathbf{x}\|_2^2 = \|(\mathbf{b}_1 - \mathbf{A} \mathbf{x}) + \mathbf{b}_2\|_2^2 = \|\mathbf{b}_1 - \mathbf{A} \mathbf{x}\|_2^2 + \|\mathbf{b}_2\|_2^2 \geq \|\mathbf{b}_2\|_2^2.$$

We obtain the minimum value  $\|\mathbf{b}_2\|_2$  of  $\|\mathbf{b} - \mathbf{A} \mathbf{x}\|_2^2$  if and only if  $\mathbf{A} \mathbf{x} = \mathbf{b}_1$ . Since  $\mathbf{b}_1 \in \text{span}(\mathbf{A})$  we can always find such an  $\mathbf{x}$  and existence follows.  $\square$

**Example 11.11 (Projections and least squares solutions)**

We have the singular value factorization (cf. Example 6.13)

$$\mathbf{A} := \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} [2] \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \end{bmatrix}.$$

We then find

$$\mathbf{A}^\dagger = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \left[ \frac{1}{2} \right] \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 & 0 \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix} = \frac{1}{4} \mathbf{A}^T.$$

Thus,

$$\mathbf{b}_1 = \mathbf{U}_1 \mathbf{U}_1^* \mathbf{b} = \frac{1}{2} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \end{bmatrix} \mathbf{b} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = \mathbf{A} \mathbf{A}^\dagger \mathbf{b} = \begin{bmatrix} (b_1 + b_2)/2 \\ (b_1 + b_2)/2 \\ 0 \end{bmatrix}.$$

Moreover, the set of all least squares solutions is

$$\{\mathbf{x} \in \mathbb{R}^2 : \mathbf{A}\mathbf{x} = \mathbf{b}_1\} = \{x_1, x_2 \in \mathbb{R} : x_1 + x_2 = \frac{b_1 + b_2}{2}\}. \quad (11.8)$$

### Theorem 11.12 (Uniqueness)

The least squares solution is unique if and only if  $\mathbf{A}$  has linearly independent columns.

**Proof.** By what we just showed any solution  $\mathbf{x}$  of the least squares problem satisfies  $\mathbf{A}\mathbf{x} = \mathbf{b}_1$ . There is a unique such  $\mathbf{x}$  if and only if  $\text{rank}(\mathbf{A}) = n$ .  $\square$

It follows that only square or overdetermined systems can have unique solutions.

### Theorem 11.13 (Characterization)

The following is equivalent:

1.  $\mathbf{x}$  is a least squares solution,
2.  $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b} + \mathbf{z}$  for some  $\mathbf{z} \in \mathbb{C}^n$  with  $\mathbf{A}\mathbf{z} = \mathbf{0}$ ,
3.  $\mathbf{A}^* \mathbf{A}\mathbf{x} = \mathbf{A}^* \mathbf{b}$  (normal equations).

**Proof.**

1.  $\implies$  2. Let  $\mathbf{x}$  be a least squares solution, i. e.,  $\mathbf{A}\mathbf{x} = \mathbf{b}_1$ . If  $\mathbf{z} := \mathbf{x} - \mathbf{A}^\dagger \mathbf{b}$  then  $\mathbf{A}\mathbf{z} = \mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{A}^\dagger \mathbf{b} = \mathbf{b}_1 - \mathbf{b}_1 = \mathbf{0}$  and  $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b} + \mathbf{z}$ .
2.  $\implies$  3. If  $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b} + \mathbf{z}$  with  $\mathbf{A}\mathbf{z} = \mathbf{0}$  then

$$\mathbf{A}^* \mathbf{A}\mathbf{x} = \mathbf{A}^* \mathbf{A}(\mathbf{A}^\dagger \mathbf{b} + \mathbf{z}) = \mathbf{A}^*(\mathbf{A}\mathbf{A}^\dagger \mathbf{b} + \mathbf{A}\mathbf{z}) = \mathbf{A}^* \mathbf{b}_1 = \mathbf{A}^* \mathbf{b}.$$

The last equality follows since  $\mathbf{b} = \mathbf{b}_1 + \mathbf{b}_2$  and  $\mathbf{b}_2 \in \ker(\mathbf{A}^*)$ .

3.  $\implies$  1. If  $\mathbf{A}^* \mathbf{A} \mathbf{x} = \mathbf{A}^* \mathbf{b}$  then  $\mathbf{A}^*(\mathbf{A} \mathbf{x} - \mathbf{b}) = \mathbf{A}^*(\mathbf{A} \mathbf{x} - \mathbf{b}_1) = \mathbf{0}$ . But then  $\mathbf{A} \mathbf{x} - \mathbf{b}_1 \in \text{span}(\mathbf{A}) \cap \ker(\mathbf{A}^*)$  and  $\mathbf{A} \mathbf{x} - \mathbf{b}_1 = \mathbf{0}$ . It follows that  $\mathbf{x}$  is a least squares solution.

□

**Example 11.14 (Least squares solutions)**

Consider the least squares solutions  $\mathbf{x}$  in (11.8). Since  $\ker(\mathbf{A}) = \{ \begin{bmatrix} z \\ -z \end{bmatrix} : z \in \mathbb{R} \}$  it follows from Part 2 of Theorem 11.13 that

$$\mathbf{x} = \frac{1}{4} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} + \begin{bmatrix} z \\ -z \end{bmatrix} = \begin{bmatrix} \frac{b_1+b_2}{4} + z \\ \frac{b_1+b_2}{4} - z \end{bmatrix}.$$

The normal equations take the form

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \mathbf{A}^T \mathbf{b} = \begin{bmatrix} b_1 + b_2 \\ b_1 + b_2 \end{bmatrix}.$$

Thus we obtain (11.8).

**11.3.2 The generalized inverse**

Consider the matrix  $\mathbf{A}^\dagger := \mathbf{V}_1 \Sigma_1^{-1} \mathbf{U}_1^*$  in (11.7). If  $\mathbf{A}$  is square and nonsingular then  $\mathbf{A}^\dagger \mathbf{A} = \mathbf{A} \mathbf{A}^\dagger = \mathbf{I}$  and  $\mathbf{A}^\dagger$  is the usual inverse of  $\mathbf{A}$ . Thus  $\mathbf{A}^\dagger$  is a generalization of the usual inverse. The matrix  $\mathbf{A}^\dagger$  satisfies Properties (1)- (4) in Exercise 11.15 and these properties define  $\mathbf{A}^\dagger$  uniquely (cf. Exercise 11.16). The unique matrix  $\mathbf{B}$  satisfying Properties (1)- (4) in Exercise 11.15 is called the **generalized inverse** or **pseudo inverse** of  $\mathbf{A}$  and denoted  $\mathbf{A}^\dagger$ . It follows that  $\mathbf{A}^\dagger := \mathbf{V}_1 \Sigma_1^{-1} \mathbf{U}_1^*$  for any singular value factorization  $\mathbf{U}_1 \Sigma_1 \mathbf{V}_1^*$  of  $\mathbf{A}$ . We show in Exercise 11.18 that if  $\mathbf{A}$  has linearly independent columns then

$$\mathbf{A}^\dagger = (\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^*. \quad (11.9)$$

**Exercise 11.15 (The generalized inverse)**

Show that  $\mathbf{B} := \mathbf{V}_1 \Sigma_1^{-1} \mathbf{U}_1^*$  satisfies (1)  $\mathbf{A} \mathbf{B} \mathbf{A} = \mathbf{A}$ , (2)  $\mathbf{B} \mathbf{A} \mathbf{B} = \mathbf{B}$ , (3)  $(\mathbf{B} \mathbf{A})^* = \mathbf{B} \mathbf{A}$ , and (4)  $(\mathbf{A} \mathbf{B})^* = \mathbf{A} \mathbf{B}$ .

**Exercise 11.16 (Uniqueness of generalized inverse)**

Given  $\mathbf{A} \in \mathbb{C}^{m \times n}$ , and suppose  $\mathbf{B}, \mathbf{C} \in \mathbb{C}^{n \times m}$  satisfy

$$\begin{aligned} \mathbf{A} \mathbf{B} \mathbf{A} &= \mathbf{A} & (1) & & \mathbf{A} \mathbf{C} \mathbf{A} &= \mathbf{A}, \\ \mathbf{B} \mathbf{A} \mathbf{B} &= \mathbf{B} & (2) & & \mathbf{C} \mathbf{A} \mathbf{C} &= \mathbf{C}, \\ (\mathbf{A} \mathbf{B})^* &= \mathbf{A} \mathbf{B} & (3) & & (\mathbf{A} \mathbf{C})^* &= \mathbf{A} \mathbf{C}, \\ (\mathbf{B} \mathbf{A})^* &= \mathbf{B} \mathbf{A} & (4) & & (\mathbf{C} \mathbf{A})^* &= \mathbf{C} \mathbf{A}. \end{aligned}$$

Verify the following proof that  $B = C$ .

$$\begin{aligned} B &= (BA)B = (A^*)B^*B = (A^*C^*)A^*B^*B = CA(A^*B^*)B \\ &= CA(BAB) = (C)AB = C(AC)AB = CC^*A^*(AB) \\ &= CC^*(A^*B^*A^*) = C(C^*A^*) = CAC = C. \end{aligned}$$

**Exercise 11.17 (Verify that a matrix is a generalized inverse)**

Show that the matrices  $A = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 0 & 0 \end{bmatrix}$  and  $B = \frac{1}{4} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix}$  satisfy the axioms in Exercise 11.15. Thus we can conclude that  $B = A^\dagger$  without computing the singular value decomposition of  $A$ .

**Exercise 11.18 (Linearly independent columns and generalized inverse)**

Suppose  $A \in \mathbb{C}^{m \times n}$  has linearly independent columns. Show that  $A^*A$  is nonsingular and  $A^\dagger = (A^*A)^{-1}A^*$ . If  $A$  has linearly independent rows, then show that  $AA^*$  is nonsingular and  $A^\dagger = A^*(AA^*)^{-1}$ .

**Exercise 11.19 (The generalized inverse of a vector)**

Show that  $\mathbf{u}^\dagger = (\mathbf{u}^*\mathbf{u})^{-1}\mathbf{u}^*$  if  $\mathbf{u} \in \mathbb{C}^{n,1}$  is nonzero.

**Exercise 11.20 (The generalized inverse of an outer product)**

If  $A = \mathbf{u}\mathbf{v}^*$  where  $\mathbf{u} \in \mathbb{C}^m$ ,  $\mathbf{v} \in \mathbb{C}^n$  are nonzero, show that

$$A^\dagger = \frac{1}{\alpha}A^*, \quad \alpha = \|\mathbf{u}\|_2^2\|\mathbf{v}\|_2^2.$$

**Exercise 11.21 (The generalized inverse of a diagonal matrix)**

Show that  $\text{diag}(\lambda_1, \dots, \lambda_n)^\dagger = \text{diag}(\lambda_1^\dagger, \dots, \lambda_n^\dagger)$  where

$$\lambda_i^\dagger = \begin{cases} 1/\lambda_i, & \lambda_i \neq 0 \\ 0 & \lambda_i = 0. \end{cases}$$

**Exercise 11.22 (Properties of the generalized inverse)**

Suppose  $A \in \mathbb{C}^{m \times n}$ . Show that

a)  $(A^*)^\dagger = (A^\dagger)^*$ .

b)  $(A^\dagger)^\dagger = A$ .

c)  $(\alpha A)^\dagger = \frac{1}{\alpha}A^\dagger$ ,  $\alpha \neq 0$ .

**Exercise 11.23 (The generalized inverse of a product)**

Suppose  $k, m, n \in \mathbb{N}$ ,  $\mathbf{A} \in \mathbb{C}^{m \times n}$ ,  $\mathbf{B} \in \mathbb{C}^{n \times k}$ . Suppose  $\mathbf{A}$  has linearly independent columns and  $\mathbf{B}$  has linearly independent rows.

- a) Show that  $(\mathbf{AB})^\dagger = \mathbf{B}^\dagger \mathbf{A}^\dagger$ . Hint: Let  $\mathbf{E} = \mathbf{AF}$ ,  $\mathbf{F} = \mathbf{B}^\dagger \mathbf{A}^\dagger$ . Show by using  $\mathbf{A}^\dagger \mathbf{A} = \mathbf{BB}^\dagger = \mathbf{I}$  that  $\mathbf{F}$  is the generalized inverse of  $\mathbf{E}$ .
- b) Find  $\mathbf{A} \in \mathbb{R}^{1,2}$ ,  $\mathbf{B} \in \mathbb{R}^{2,1}$  such that  $(\mathbf{AB})^\dagger \neq \mathbf{B}^\dagger \mathbf{A}^\dagger$ .

**Exercise 11.24 (The generalized inverse of the conjugate transpose)**

Show that  $\mathbf{A}^* = \mathbf{A}^\dagger$  if and only if all singular values of  $\mathbf{A}$  are either zero or one.

**Exercise 11.25 (Linearly independent columns)**

Show that if  $\mathbf{A}$  has rank  $n$  then  $\mathbf{A}(\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^* \mathbf{b}$  is the projection of  $\mathbf{b}$  into  $\text{span}(\mathbf{A})$ . (Cf. Exercise 11.18.)

**Exercise 11.26 (Analysis of the general linear system)**

Consider the linear system  $\mathbf{Ax} = \mathbf{b}$  where  $\mathbf{A} \in \mathbb{C}^{n \times n}$  has rank  $r > 0$  and  $\mathbf{b} \in \mathbb{C}^n$ . Let

$$\mathbf{U}^* \mathbf{AV} = \begin{bmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$$

represent the singular value decomposition of  $\mathbf{A}$ .

- a) Let  $\mathbf{c} = [c_1, \dots, c_n]^T = \mathbf{U}^* \mathbf{b}$  and  $\mathbf{y} = [y_1, \dots, y_n]^T = \mathbf{V}^* \mathbf{x}$ . Show that  $\mathbf{Ax} = \mathbf{b}$  if and only if

$$\begin{bmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{y} = \mathbf{c}.$$

- b) Show that  $\mathbf{Ax} = \mathbf{b}$  has a solution  $\mathbf{x}$  if and only if  $c_{r+1} = \dots = c_n = 0$ .
- c) Deduce that a linear system  $\mathbf{Ax} = \mathbf{b}$  has either no solution, one solution or infinitely many solutions.

**Exercise 11.27 (Fredholm's alternative)**

For any  $\mathbf{A} \in \mathbb{C}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{C}^n$  show that one and only one of the following systems has a solution

$$(1) \quad \mathbf{Ax} = \mathbf{b}, \quad (2) \quad \mathbf{A}^* \mathbf{y} = \mathbf{0}, \mathbf{y}^* \mathbf{b} \neq 0.$$

In other words either  $\mathbf{b} \in \text{span}(\mathbf{A})$ , or we can find  $\mathbf{y} \in \ker(\mathbf{A}^*)$  such that  $\mathbf{y}^* \mathbf{b} \neq 0$ . This is called **Fredholm's alternative**.



## 11.4 Numerical Solution

We assume that  $m \geq n$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$ . Numerical methods can be based on normal equations, QR factorization, or Singular Value Factorization. We discuss each of these approaches in turn. Another possibility is to use an iterative method like the conjugate gradient method (cf. Exercise 9.18).

### 11.4.1 Normal equations

Suppose  $\mathbf{A}$  has linearly independent columns. The coefficient matrix  $\mathbf{B} := \mathbf{A}^T \mathbf{A}$  in the normal equations is symmetric positive definite, and we can solve these equations using the Cholesky factorization of  $\mathbf{B}$ . Consider forming the normal equations. We can use either a column oriented (inner product)- or a row oriented (outer product) approach.

$$1. \text{ inner product: } (\mathbf{A}^T \mathbf{A})_{i,j} = \sum_{k=1}^m a_{k,i} a_{k,j}, \quad i, j = 1, \dots, n,$$

$$(\mathbf{A}^T \mathbf{b})_i = \sum_{k=1}^m a_{k,i} b_k, \quad i = 1, \dots, n,$$

$$2. \text{ outer product: } \mathbf{A}^T \mathbf{A} = \sum_{k=1}^m \begin{bmatrix} a_{k1} \\ \vdots \\ a_{kn} \end{bmatrix} [a_{k1} \ \cdots \ a_{kn}], \quad \mathbf{A}^T \mathbf{b} = \sum_{k=1}^m \begin{bmatrix} a_{k1} \\ \vdots \\ a_{kn} \end{bmatrix} b_k.$$

The outer product form is suitable for large problems since it uses only one pass through the data importing one row of  $\mathbf{A}$  at a time from some separate storage.

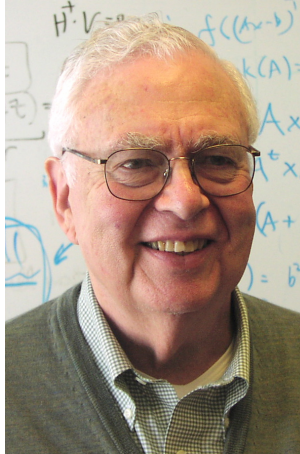
Consider the number of operations to find the least squares solution. We need  $2m$  arithmetic operations for each inner product. Since  $\mathbf{B}$  is symmetric we only need to compute  $n(n+1)/2$  such inner products. It follows that  $\mathbf{B}$  can be computed in approximately  $mn^2$  arithmetic operations. In conclusion the number of operations are  $mn^2$  to find  $\mathbf{B}$ ,  $2mn$  to find  $\mathbf{A}^T \mathbf{b}$ ,  $n^3/3$  to find  $\mathbf{R}$ ,  $n^2$  to solve  $\mathbf{R}^T \mathbf{y} = \mathbf{c}$  and  $n^2$  to solve  $\mathbf{R} \mathbf{x} = \mathbf{y}$ . If  $m \approx n$  it takes  $\frac{4}{3}n^3 = 2G_n$  arithmetic operations. If  $m$  is much bigger than  $n$  the number of operations is approximately  $mn^2$ , the work to compute  $\mathbf{B}$ .

A problem with the normal equations approach is that the linear system can be poorly conditioned. In fact the 2-norm condition number of  $\mathbf{B} := \mathbf{A}^T \mathbf{A}$  is the square of the condition number of  $\mathbf{A}$ . This follows, since the eigenvalues of  $\mathbf{B}$  are the square of the singular values of  $\mathbf{A}$  so that

$$K_2(\mathbf{B}) = \frac{\sigma_1^2}{\sigma_n^2} = \left( \frac{\sigma_1}{\sigma_n} \right)^2 = K_2(\mathbf{A})^2.$$

If  $\mathbf{A}$  is ill-conditioned, this could make the normal equations approach problematic. One difficulty which can be encountered is that the computed  $\mathbf{A}^T \mathbf{A}$  might not be positive definite. See Problem 11.36 for an example.

### 11.4.2 QR factorization



Gene Golub, 1932–2007. He pioneered use of the QR factorization to solve least square problems.

Suppose  $\mathbf{A} \in \mathbb{R}^{m \times n}$  has rank  $n$  and let  $\mathbf{b} \in \mathbb{R}^m$ . The QR factorization can be used to solve the least squares problem. Suppose  $\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1$  is a QR factorization of  $\mathbf{A}$ . Since  $\mathbf{Q}_1$  has orthonormal columns we find

$$\mathbf{A}^T \mathbf{A} = \mathbf{R}_1^T \mathbf{Q}_1^T \mathbf{Q}_1 \mathbf{R}_1 = \mathbf{R}_1^T \mathbf{R}_1, \quad \mathbf{A}^T \mathbf{b} = \mathbf{R}_1^T \mathbf{Q}_1^T \mathbf{b}.$$

Since  $\mathbf{A}$  has rank  $n$  the matrix  $\mathbf{R}_1^T$  is nonsingular and can be canceled. Thus

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b} \implies \mathbf{R}_1 \mathbf{x} = \mathbf{c}_1, \quad \mathbf{c}_1 := \mathbf{Q}_1^T \mathbf{b}.$$

We can use Householder transformations or Givens rotations to find  $\mathbf{R}_1$  and  $\mathbf{c}_1$ . Consider using the Householder triangulation algorithm Algorithm 10.8. We find  $\mathbf{R} = \mathbf{Q}^T \mathbf{A}$  and  $\mathbf{c} = \mathbf{Q}^T \mathbf{b}$ , where  $\mathbf{A} = \mathbf{Q} \mathbf{R}$  is the QR decomposition of  $\mathbf{A}$ . The matrices  $\mathbf{R}_1$  and  $\mathbf{c}_1$  are located in the first  $n$  rows of  $\mathbf{R}$  and  $\mathbf{c}$ . Using also Algorithm 2.2 we have the following method to solve the full rank least squares problem.

1. `[R, c]=housetriang(A, b) .`
2. `x=rbacksolve(R(1:n, 1:n), c(1:n), n) .`

**Example 11.28 (Solution using QR factorization)**

Consider the least squares problem with

$$\mathbf{A} = \begin{bmatrix} 1 & 3 & 1 \\ 1 & 3 & 7 \\ 1 & -1 & -4 \\ 1 & -1 & 2 \end{bmatrix} \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}.$$

This is the matrix in Example 10.11. The least squares solution  $\mathbf{x}$  is found by solving the system

$$\begin{bmatrix} 2 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

and we find  $\mathbf{x} = [1, 0, 0]^T$ .

Using Householder triangulation is a useful alternative to normal equations for solving full rank least squares problems. It can even be extended to rank deficient problems, see [2]. The 2 norm condition number for the system  $\mathbf{R}_1 \mathbf{x} = \mathbf{c}_1$  is  $K_2(\mathbf{R}_1) = K_2(\mathbf{Q}_1 \mathbf{R}_1) = K_2(\mathbf{A})$ , and as discussed in the previous section this is the square root of  $K_2(\mathbf{A}^T \mathbf{A})$ , the condition number for the normal equations. Thus if  $\mathbf{A}$  is mildly ill-conditioned the normal equations can be quite ill-conditioned and solving the normal equations can give inaccurate results. On the other hand Algorithm 10.8 is quite stable.

But using Householder transformations requires more work. The leading term in the number of arithmetic operations in Algorithm 10.8 is approximately  $2mn^2 - 2n^3/3$ , (cf. (10.5)) while the number of arithmetic operations needed to form the normal equations, taking advantage of symmetry is approximately  $mn^2$ . Thus for  $m$  much larger than  $n$  using Householder triangulation requires twice as many arithmetic operations as the approach based on the normal equations. Also, Householder triangulation have problems taking advantage of the structure in sparse problems.

### 11.4.3 Singular value factorization

This method can be used even if  $\mathbf{A}$  does not have full rank. By Theorem 11.13

$$\mathbf{x} = \mathbf{A}^\dagger \mathbf{b} + \mathbf{z},$$

where  $\mathbf{A}^\dagger$  is the generalized inverse of  $\mathbf{A}$ , is a least squares solution for any  $\mathbf{z} \in \ker(\mathbf{A})$ . If  $\mathbf{A}$  has linearly independent columns then  $\ker(\mathbf{A}) = \{\mathbf{0}\}$  and  $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}$  is the unique solution.

When  $\text{rank}(\mathbf{A})$  is less than the number of columns of  $\mathbf{A}$  then  $\ker(\mathbf{A}) \neq \{\mathbf{0}\}$ , and we have a choice of  $\mathbf{z}$ . One possible choice is  $\mathbf{z} = \mathbf{0}$  giving the solution  $\mathbf{A}^\dagger \mathbf{b}$ .

**Theorem 11.29 (Minimal solution)**

The least squares solution with minimal Euclidian norm is  $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}$  corresponding to  $\mathbf{z} = \mathbf{0}$ .

*Proof.* Suppose  $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b} + \mathbf{z}$ , with  $\mathbf{z} \in \ker(\mathbf{A})$ . Recall that if the right singular vectors of  $\mathbf{A}$  are partitioned as  $[\mathbf{v}_1, \dots, \mathbf{v}_r, \mathbf{v}_{r+1}, \dots, \mathbf{v}_n] = [\mathbf{V}_1, \mathbf{V}_2]$ , then  $\mathbf{V}_2$  is a basis for  $\ker(\mathbf{A})$ . Moreover,  $\mathbf{V}_2^* \mathbf{V}_1 = \mathbf{0}$  since  $\mathbf{V}$  has orthonormal columns. If  $\mathbf{A}^\dagger = \mathbf{V}_1 \Sigma_1^{-1} \mathbf{U}_1^*$  and  $\mathbf{z} \in \ker(\mathbf{A})$  then  $\mathbf{z} = \mathbf{V}_2 \mathbf{y}$  for some  $\mathbf{y} \in \mathbb{C}^{n-r}$  and we obtain

$$\mathbf{z}^* \mathbf{A}^\dagger \mathbf{b} = \mathbf{y}^* \mathbf{V}_2^* \mathbf{V}_1 \Sigma_1^{-1} \mathbf{U}_1^* \mathbf{b} = \mathbf{0}.$$

Thus  $\mathbf{z}$  and  $\mathbf{A}^\dagger \mathbf{b}$  are orthogonal so that by Pythagoras  $\|\mathbf{x}\|_2^2 = \|\mathbf{A}^\dagger \mathbf{b} + \mathbf{z}\|_2^2 = \|\mathbf{A}^\dagger \mathbf{b}\|_2^2 + \|\mathbf{z}\|_2^2 \geq \|\mathbf{A}^\dagger \mathbf{b}\|_2^2$  with equality for  $\mathbf{z} = \mathbf{0}$ .  $\square$

Using MATLAB a least squares solution can be found using `x=A\b` if  $\mathbf{A}$  has full rank. For rank deficient problems the function `x=lsqcov(A,b)` finds a least squares solution with a maximal number of zeros in  $\mathbf{x}$ .

**Example 11.30 (Rank deficient least squares solution)**

For  $\mathbf{A}$  as in Example 11.14 with  $\mathbf{b} = [1, 1]^T$  `lsqcov` gives the solution  $[1, 0]^T$  corresponding to  $\mathbf{z} = [1/2, -1/2]$ . The minimal norm solution is  $[1/2, 1/2]$ .

## 11.5 Perturbation Theory for Least Squares

In this section we consider what effect small changes in the data  $\mathbf{A}, \mathbf{b}$  have on the solution  $\mathbf{x}$  of the least squares problem  $\min \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$ .

If  $\mathbf{A}$  has linearly independent columns then we can write the least squares solution  $\mathbf{x}$  (the solution of  $\mathbf{A}^* \mathbf{A} \mathbf{x} = \mathbf{A}^* \mathbf{b}$ ) as

$$\mathbf{x} = \mathbf{A}^\dagger \mathbf{b} = \mathbf{A}^\dagger \mathbf{b}_1, \quad \mathbf{A}^\dagger := (\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^*,$$

where  $\mathbf{b}_1$  is the orthogonal projection of  $\mathbf{b}$  into the column space  $\text{span}(\mathbf{A})$ .

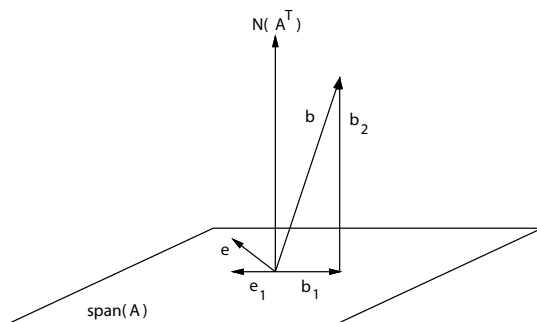
### 11.5.1 Perturbing the right hand side

Let us now consider the effect of a perturbation in  $\mathbf{b}$  on  $\mathbf{x}$ .

**Theorem 11.31 (Perturbing the right hand side)**

Suppose  $\mathbf{A} \in \mathbb{C}^{m \times n}$  has linearly independent columns, and let  $\mathbf{b}, \mathbf{e} \in \mathbb{C}^m$ . Let  $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$  be the solutions of  $\min \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$  and  $\min \|\mathbf{A}\mathbf{y} - \mathbf{b} - \mathbf{e}\|_2$ . Finally, let  $\mathbf{b}_1, \mathbf{e}_1$  be the orthogonal projections of  $\mathbf{b}$  and  $\mathbf{e}$  into  $\text{span}(\mathbf{A})$ . If  $\mathbf{b}_1 \neq \mathbf{0}$ , we have for any operator norm

$$\frac{1}{K(\mathbf{A})} \frac{\|\mathbf{e}_1\|}{\|\mathbf{b}_1\|} \leq \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(\mathbf{A}) \frac{\|\mathbf{e}_1\|}{\|\mathbf{b}_1\|}, \quad K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^\dagger\|. \quad (11.10)$$



**Figure 11.3.** Graphical interpretation of the bounds in Theorem 11.31.

**Proof.** Subtracting  $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}_1$  from  $\mathbf{y} = \mathbf{A}^\dagger \mathbf{b}_1 + \mathbf{A}^\dagger \mathbf{e}_1$  we have  $\mathbf{y} - \mathbf{x} = \mathbf{A}^\dagger \mathbf{e}_1$ . Thus  $\|\mathbf{y} - \mathbf{x}\| = \|\mathbf{A}^\dagger \mathbf{e}_1\| \leq \|\mathbf{A}^\dagger\| \|\mathbf{e}_1\|$ . Moreover,  $\|\mathbf{b}_1\| = \|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$ . Therefore  $\|\mathbf{y} - \mathbf{x}\|/\|\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{A}^\dagger\| \|\mathbf{e}_1\|/\|\mathbf{b}_1\|$  proving the rightmost inequality. From  $\mathbf{A}(\mathbf{x} - \mathbf{y}) = \mathbf{e}_1$  and  $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}_1$  we obtain the leftmost inequality.  $\square$

(11.10) is analogous to the bound (7.13) for linear systems. We see that the number  $K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^\dagger\|$  generalizes the condition number  $\|\mathbf{A}\| \|\mathbf{A}^{-1}\|$  for a square matrix. The main difference between (11.10) and (7.13) is however that  $\|\mathbf{e}\|/\|\mathbf{b}\|$  in (7.13) has been replaced by  $\|\mathbf{e}_1\|/\|\mathbf{b}_1\|$ , the orthogonal projections of  $\mathbf{e}$  and  $\mathbf{b}$  into  $\text{span}(\mathbf{A})$ . If  $\mathbf{b}$  lies almost entirely in  $\ker(\mathbf{A}^*)$ , i.e.  $\|\mathbf{b}\|/\|\mathbf{b}_1\|$  is large, then  $\|\mathbf{e}_1\|/\|\mathbf{b}_1\|$  can be much larger than  $\|\mathbf{e}\|/\|\mathbf{b}\|$ . This is illustrated in Figure 11.3. If  $\mathbf{b}$  is almost orthogonal to  $\text{span}(\mathbf{A})$ ,  $\|\mathbf{e}_1\|/\|\mathbf{b}_1\|$  will normally be much larger than  $\|\mathbf{e}\|/\|\mathbf{b}\|$ .

### Example 11.32 (Perturbing the right hand side)

Suppose

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 10^{-4} \\ 0 \\ 1 \end{bmatrix}, \quad \mathbf{e} = \begin{bmatrix} 10^{-6} \\ 0 \\ 0 \end{bmatrix}.$$

For this example we can compute  $K(\mathbf{A})$  by finding  $\mathbf{A}^\dagger$  explicitly. Indeed,

$$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}, \quad (\mathbf{A}^T \mathbf{A})^{-1} = \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix}, \quad \mathbf{A}^\dagger = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T = \begin{bmatrix} 1 & -1 & 0 \\ 0 & 1 & 0 \end{bmatrix}.$$

Thus  $K_\infty(\mathbf{A}) = \|\mathbf{A}\|_\infty \|\mathbf{A}^\dagger\|_\infty = 2 \cdot 2 = 4$  is quite small.

Consider now the projections  $\mathbf{b}_1$  and  $\mathbf{e}_1$ . We find  $\mathbf{A}\mathbf{A}^\dagger = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ . Hence

$$\mathbf{b}_1 = \mathbf{A}\mathbf{A}^\dagger \mathbf{b} = [10^{-4}, 0, 0]^T, \quad \text{and} \quad \mathbf{e}_1 = \mathbf{A}\mathbf{A}^\dagger \mathbf{e} = [10^{-6}, 0, 0]^T.$$

Thus  $\|\mathbf{e}_1\|_\infty/\|\mathbf{b}_1\|_\infty = 10^{-2}$  and (11.10) takes the form

$$\frac{1}{4}10^{-2} \leq \frac{\|\mathbf{y} - \mathbf{x}\|_\infty}{\|\mathbf{x}\|_\infty} \leq 4 \cdot 10^{-2}.$$

To verify the bounds we compute the solutions as  $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b} = [10^{-4}, 0]^T$  and  $\mathbf{y} = \mathbf{A}^\dagger (\mathbf{b} + \mathbf{e}) = [10^{-4} + 10^{-6}, 0]^T$ . Hence

$$\frac{\|\mathbf{x} - \mathbf{y}\|_\infty}{\|\mathbf{x}\|_\infty} = \frac{10^{-6}}{10^{-4}} = 10^{-2},$$

### Exercise 11.33 (Condition number)

Let

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}.$$

- Determine the projections  $\mathbf{b}_1$  and  $\mathbf{b}_2$  of  $\mathbf{b}$  on  $\text{span}(\mathbf{A})$  and  $\ker(\mathbf{A}^T)$ .
- Compute  $K(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{A}^\dagger\|_2$ .

For each  $\mathbf{A}$  we can find  $\mathbf{b}$  and  $\mathbf{e}$  so that we have equality in the upper bound in (11.10). The lower bound is best possible in a similar way.

### Exercise 11.34 (Equality in perturbation bound)

- Let  $\mathbf{A} \in \mathbb{C}^{m \times n}$ . Show that we have equality to the right in (11.10) if  $\mathbf{b} = \mathbf{A}\mathbf{y}_A$ ,  $\mathbf{e}_1 = \mathbf{y}_{A^\dagger}$  where  $\|\mathbf{A}\mathbf{y}_A\| = \|\mathbf{A}\|$ ,  $\|\mathbf{A}^\dagger \mathbf{y}_{A^\dagger}\| = \|\mathbf{A}^\dagger\|$ .
- Show that we have equality to the left if we switch  $\mathbf{b}$  and  $\mathbf{e}$  in a).
- Let  $\mathbf{A}$  be as in Example 11.32. Find extremal  $\mathbf{b}$  and  $\mathbf{e}$  when the  $l_\infty$  norm is used.

## 11.5.2 Perturbing the matrix

The analysis of the effects of a perturbation  $\mathbf{E}$  in  $\mathbf{A}$  is quite difficult. The following result is stated without proof, see [20, p. 51]. For other estimates see [2] and [28].

### Theorem 11.35 (Perturbing the matrix)

Suppose  $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{m \times n}$ ,  $m > n$ , where  $\mathbf{A}$  has linearly independent columns and  $\alpha := 1 - \|\mathbf{E}\|_2 \|\mathbf{A}^\dagger\|_2 > 0$ . Then  $\mathbf{A} + \mathbf{E}$  has linearly independent columns. Let  $\mathbf{b} = \mathbf{b}_1 + \mathbf{b}_2 \in \mathbb{C}^m$  where  $\mathbf{b}_1$  and  $\mathbf{b}_2$  are the orthogonal projections into  $\text{span}(\mathbf{A})$

and  $\ker(\mathbf{A}^*)$  respectively. Suppose  $\mathbf{b}_1 \neq \mathbf{0}$ . Let  $\mathbf{x}$  and  $\mathbf{y}$  be the solutions of  $\min\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$  and  $\min\|(\mathbf{A} + \mathbf{E})\mathbf{y} - \mathbf{b}\|_2$ . Then

$$\rho = \frac{\|\mathbf{x} - \mathbf{y}\|_2}{\|\mathbf{x}\|_2} \leq \frac{1}{\alpha} K(1 + \beta K) \frac{\|\mathbf{E}\|_2}{\|\mathbf{A}\|_2}, \quad \beta = \frac{\|\mathbf{b}_2\|_2}{\|\mathbf{b}_1\|_2}, \quad K = \|\mathbf{A}\|_2 \|\mathbf{A}^\dagger\|_2. \quad (11.11)$$

(11.11) says that the relative error in  $\mathbf{y}$  as an approximation to  $\mathbf{x}$  can be at most  $K(1 + \beta K)/\alpha$  times as large as the size  $\|\mathbf{E}\|_2/\|\mathbf{A}\|_2$  of the relative perturbation in  $\mathbf{A}$ .  $\beta$  will be small if  $\mathbf{b}$  lies almost entirely in  $\text{span}(\mathbf{A})$ , and we have approximately  $\rho \leq \frac{1}{\alpha} K \|\mathbf{E}\|_2/\|\mathbf{A}\|_2$ . This corresponds to the estimate (7.19) for linear systems. If  $\beta$  is not small, the term  $\frac{1}{\alpha} K^2 \beta \|\mathbf{E}\|_2/\|\mathbf{A}\|_2$  will dominate. In other words, the condition number is roughly  $K(\mathbf{A})$  if  $\beta$  is small and  $K(\mathbf{A})^2 \beta$  if  $\beta$  is not small. Note that  $\beta$  is large if  $\mathbf{b}$  is almost orthogonal to  $\text{span}(\mathbf{A})$  and that  $\mathbf{b}_2 = \mathbf{b} - \mathbf{A}\mathbf{x}$  is the residual of  $\mathbf{x}$ .

### Exercise 11.36 (Problem using normal equations)

Consider the least squares problems where

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1+\epsilon \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 2 \\ 3 \\ 2 \end{bmatrix}, \quad \epsilon \in \mathbb{R}.$$

- a) Find the normal equations and the exact least squares solution.
- b) Suppose  $\epsilon$  is small and we replace the  $(2, 2)$  entry  $3+2\epsilon+\epsilon^2$  in  $\mathbf{A}^T \mathbf{A}$  by  $3+2\epsilon$ . (This will be done in a computer if  $\epsilon < \sqrt{u}$ ,  $u$  being the round-off unit). For example, if  $u = 10^{-16}$  then  $\sqrt{u} = 10^{-8}$ . Solve  $\mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{A}^T \mathbf{b}$  for  $\mathbf{x}$  and compare with the  $\mathbf{x}$  found in a). (We will get a much more accurate result using the QR factorization or the singular value decomposition on this problem).

## 11.6 Perturbation Theory for Singular Values

In this section we consider what effect a small change in the matrix  $\mathbf{A}$  has on the singular values.

We recall the Hoffman-Wielandt Theorem for singular values, Theorem 6.30. If  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$  are rectangular matrices with singular values  $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n$  and  $\beta_1 \geq \beta_2 \geq \dots \geq \beta_n$ , then

$$\sum_{j=1}^n |\alpha_j - \beta_j|^2 \leq \|\mathbf{A} - \mathbf{B}\|_F^2.$$

This shows that the singular values of a matrix are well conditioned. Changing the Frobenius norm of a matrix by small amount only changes the singular values by a small amount.

Using the 2-norm we have a similar result.

**Theorem 11.37 (Perturbation of singular values)**

Let  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$  be rectangular matrices with singular values  $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n$  and  $\beta_1 \geq \beta_2 \geq \dots \geq \beta_n$ . Then

$$|\alpha_j - \beta_j| \leq \|\mathbf{A} - \mathbf{B}\|_2, \text{ for } j = 1, 2, \dots, n. \quad (11.12)$$

**Proof.** Fix  $j$  and let  $\mathcal{S}$  be the  $n - j + 1$  dimensional subspace for which the minimum in Theorem 6.29 is obtained for  $\mathbf{A}$ . Then

$$\alpha_j = \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|(\mathbf{B} + (\mathbf{A} - \mathbf{B}))\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|\mathbf{B}\mathbf{x}\|_2}{\|\mathbf{x}\|_2} + \max_{\substack{\mathbf{x} \in \mathcal{S} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\|(\mathbf{A} - \mathbf{B})\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \beta_j + \|\mathbf{A} - \mathbf{B}\|_2.$$

By symmetry we obtain  $\beta_j \leq \alpha_j + \|\mathbf{A} - \mathbf{B}\|_2$  and the proof is complete.  $\square$

The following result is an analogue of Theorem 7.31.

**Theorem 11.38 (Generalized inverse when perturbing the matrix)**

Let  $\mathbf{A}, \mathbf{E} \in \mathbb{R}^{m \times n}$  have singular values  $\alpha_1 \geq \dots \geq \alpha_n$  and  $\epsilon_1 \geq \dots \geq \epsilon_n$ . If  $\|\mathbf{A}^\dagger\|_2 \|\mathbf{E}\|_2 < 1$  then

1.  $\text{rank}(\mathbf{A} + \mathbf{E}) \geq \text{rank}(\mathbf{A})$ ,
2.  $\|(\mathbf{A} + \mathbf{E})^\dagger\|_2 \leq \frac{\|\mathbf{A}^\dagger\|_2}{1 - \|\mathbf{A}^\dagger\|_2 \|\mathbf{E}\|_2} = \frac{1}{\alpha_r - \epsilon_1}$ ,

where  $r$  is the rank of  $\mathbf{A}$ .

**Proof.** Suppose  $\mathbf{A}$  has rank  $r$  and let  $\mathbf{B} := \mathbf{A} + \mathbf{E}$  have singular values  $\beta_1 \geq \dots \geq \beta_n$ . In terms of singular values the inequality  $\|\mathbf{A}^\dagger\|_2 \|\mathbf{E}\|_2 < 1$  can be written  $\epsilon_1/\alpha_r < 1$  or  $\alpha_r > \epsilon_1$ . By Theorem 11.37 we have  $\alpha_r - \beta_r \leq \epsilon_1$ , which implies  $\beta_r \geq \alpha_r - \epsilon_1 > 0$ , and this shows that  $\text{rank}(\mathbf{A} + \mathbf{E}) > r$ . To prove 2., the inequality  $\beta_r \geq \alpha_r - \epsilon_1$  implies that

$$\|(\mathbf{A} + \mathbf{E})^\dagger\|_2 \leq \frac{1}{\beta_r} \leq \frac{1}{\alpha_r - \epsilon_1} = \frac{1/\alpha_r}{1 - \epsilon_1/\alpha_r} = \frac{\|\mathbf{A}^\dagger\|_2}{1 - \|\mathbf{A}^\dagger\|_2 \|\mathbf{E}\|_2}.$$

$\square$

## 11.7 Review Questions

**11.7.1** Do the normal equations always have a solution?

**11.7.2** When is the least squares solution unique?



- 11.7.3** Express the general least squares solution in terms of the generalized inverse.
- 11.7.4** Consider perturbing the right-hand side in a linear equation and a least squares problem. What is the main difference in the perturbation inequalities?
- 11.7.5** Why does one often prefer using QR factorization instead of normal equations for solving least squares problems.
- 11.7.6** What is an orthogonal sum?
- 11.7.7** How is an orthogonal projection defined?



## **Part V**

# **Eigenvalues and Eigenvectors**



## Chapter 12

# Numerical Eigenvalue Problems

## 12.1 Eigenpairs

Eigenpairs have applications in quantum mechanics, differential equations, elasticity in mechanics, etc, etc. Typical computational problems involve

- Finding one or a few of the eigenvalues.
- Finding one or a few of the eigenpairs.
- Finding all eigenvalues.
- Finding all eigenpairs.

In this and the next chapter we consider some numerical methods for finding one or more of the eigenvalues and eigenvectors of a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . Maybe the first method which comes to mind is to form the characteristic polynomial  $\pi_{\mathbf{A}}$  of  $\mathbf{A}$ , and then use a polynomial root finder, like Newton's method to determine one or several of the eigenvalues.

It turns out that this is not suitable as an all purpose method. One reason is that a small change in one of the coefficients of  $\pi_{\mathbf{A}}(\lambda)$  can lead to a large change in the roots of the polynomial. For example, if  $\pi_{\mathbf{A}}(\lambda) = \lambda^{16}$  and  $q(\lambda) = \lambda^{16} - 10^{-16}$  then the roots of  $\pi_{\mathbf{A}}$  are all equal to zero, while the roots of  $q$  are  $\lambda_j = 10^{-1} e^{2\pi i j / 16}$ ,  $j = 1, \dots, 16$ . The roots of  $q$  have absolute value 0.1 and a perturbation in one of the polynomial coefficients of magnitude  $10^{-16}$  has led to an error in the roots of approximately 0.1. The situation can be somewhat remedied by representing the polynomials using a different basis.

We will see that for many matrices the eigenvalues are less sensitive to perturbations in the elements of the matrix. In this text we will only consider methods

which work directly with the matrix.

## 12.2 Gerschgorin's Theorem

The following theorem is useful for locating eigenvalues of an arbitrary square matrix.

### Theorem 12.1 (Gerschgorin's circle theorem)

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . Define for  $i = 1, 2, \dots, n$

$$R_i = \{z \in \mathbb{C} : |z - a_{ii}| \leq r_i\}, \quad r_i := \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|,$$

$$C_j = \{z \in \mathbb{C} : |z - a_{jj}| \leq c_j\}, \quad c_j := \sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}|.$$

Then any eigenvalue of  $\mathbf{A}$  lies in  $R \cap C$  where  $R = R_1 \cup R_2 \cup \dots \cup R_n$  and  $C = C_1 \cup C_2 \cup \dots \cup C_n$ .

**Proof.** Suppose  $(\lambda, \mathbf{x})$  is an eigenpair for  $\mathbf{A}$ . We claim that  $\lambda \in R_i$ , where  $i$  is such that  $|x_i| = \|\mathbf{x}\|_\infty$ . Indeed,  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  implies that  $\sum_j a_{ij}x_j = \lambda x_i$  or  $(\lambda - a_{ii})x_i = \sum_{j \neq i} a_{ij}x_j$ . Dividing by  $x_i$  and taking absolute values we find

$$|\lambda - a_{ii}| = \left| \sum_{j \neq i} a_{ij}x_j/x_i \right| \leq \sum_{j \neq i} |a_{ij}| |x_j/x_i| \leq r_i$$

since  $|x_j/x_i| \leq 1$  for all  $j$ . Thus  $\lambda \in R_i$ .

Since  $\lambda$  is also an eigenvalue of  $\mathbf{A}^T$ , it must be in one of the row disks of  $\mathbf{A}^T$ . But these are the column disks  $C_j$  of  $\mathbf{A}$ . Hence  $\lambda \in C_j$  for some  $j$ .  $\square$

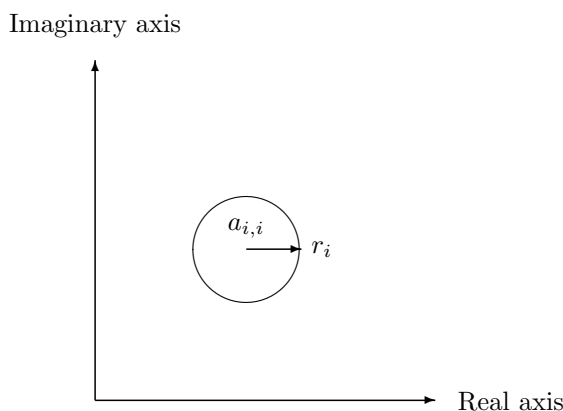
The set  $R_i$  is a subset of the complex plane consisting of all points inside a circle with center at  $a_{ii}$  and radius  $r_i$ , c.f. Figure 12.1.  $R_i$  is called a (Gerschgorin) row disk.

An eigenvalue  $\lambda$  lies in the union of the row disks  $R_1, \dots, R_n$  and also in the union of the column disks  $C_1, \dots, C_n$ . If  $\mathbf{A}$  is Hermitian then  $R_i = C_i$  for  $i = 1, 2, \dots, n$ . Moreover, in this case the eigenvalues of  $\mathbf{A}$  are real, and the Gerschgorin disks can be taken to be intervals on the real line.

### Example 12.2 (Gerschgorin)

Let  $\mathbf{T} = \text{tridiag}(-1, 2, -1) \in \mathbb{R}^{m \times m}$  be the second derivative matrix. Since  $\mathbf{A}$  is Hermitian we have  $R_i = C_i$  for all  $i$  and the eigenvalues are real. We find

$$R_1 = R_m = \{z \in \mathbb{R} : |z-2| \leq 1\}, \quad \text{and} \quad R_i = \{z \in \mathbb{R} : |z-2| \leq 2\}, \quad i = 2, 3, \dots, m-1.$$



**Figure 12.1.** The Gerschgorin disk  $R_i$ .

We conclude that  $\lambda \in [0, 4]$  for any eigenvalue  $\lambda$  of  $\mathbf{T}$ . To check this, we recall that by Lemma 3.8 the eigenvalues of  $\mathbf{T}$  are given by

$$\lambda_j = 4 \left[ \sin \frac{j\pi}{2(m+1)} \right]^2, \quad j = 1, 2, \dots, m.$$

When  $m$  is large the smallest eigenvalue  $4 \left[ \sin \frac{\pi}{2(m+1)} \right]^2$  is very close to zero and the largest eigenvalue  $4 \left[ \sin \frac{m\pi}{2(m+1)} \right]^2$  is very close to 4. Thus Gerschgorin's theorem gives a remarkably good estimate for large  $m$ .

Sometimes some of the Gerschgorin disks are distinct and we have

**Corollary 12.3 (Disjoint Gerschgorin disks)**

If  $p$  of the Gerschgorin row disks are disjoint from the others, the union of these disks contains precisely  $p$  eigenvalues. The same result holds for the column disks.

**Proof.** Consider a family of matrices

$$\mathbf{A}(t) := \mathbf{D} + t(\mathbf{A} - \mathbf{D}), \quad \mathbf{D} := \text{diag}(a_{11}, \dots, a_{nn}), \quad t \in [0, 1].$$

We have  $\mathbf{A}(0) = \mathbf{D}$  and  $\mathbf{A}(1) = \mathbf{A}$ . As a function of  $t$ , every eigenvalue of  $\mathbf{A}(t)$  is a continuous function of  $t$ . This follows from Theorem 12.8, see Exercise 12.5. The row disks  $R_i(t)$  of  $\mathbf{A}(t)$  have radius proportional to  $t$ , indeed

$$R_i(t) = \{z \in \mathbb{C} : |z - a_{ii}| \leq tr_i\}, \quad r_i := \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|.$$

Clearly  $0 \leq t_1 < t_2 \leq 1$  implies  $R_i(t_1) \subset R_i(t_2)$  and  $R_i(1)$  is a row disk of  $\mathbf{A}$  for all  $i$ . Suppose  $\bigcup_{k=1}^p R_{i_k}(1)$  are disjoint from the other disks of  $\mathbf{A}$  and set  $R^p(t) := \bigcup_{k=1}^p R_{i_k}(t)$  for  $t \in [0, 1]$ . Now  $R^p(0)$  contains only the  $p$  eigenvalues  $a_{i_1, i_1}, \dots, a_{i_p, i_p}$  of  $\mathbf{A}(0) = \mathbf{D}$ . As  $t$  increases from zero to one the set  $R^p(t)$  is disjoint from the other row disks of  $\mathbf{A}$  and by the continuity of the eigenvalues cannot lose or gain eigenvalues. It follows that  $R^p(1)$  must contain  $p$  eigenvalues of  $\mathbf{A}$ .  $\square$

**Example 12.4** Consider the matrix  $\mathbf{A} = \begin{bmatrix} 1 & \epsilon_1 & \epsilon_2 \\ \epsilon_3 & 2 & \epsilon_4 \\ \epsilon_5 & \epsilon_6 & 3 \end{bmatrix}$ , where  $|\epsilon_i| \leq 10^{-15}$  all  $i$ . By Corollary 12.3 the eigenvalues  $\lambda_1, \lambda_2, \lambda_3$  of  $\mathbf{A}$  are distinct and satisfy  $|\lambda_j - \lambda_k| \leq 2 \times 10^{-15}$  for  $j = 1, 2, 3$ .

**Exercise 12.5 (Continuity of eigenvalues)**

Suppose  $t_1, t_2 \in [0, 1]$  and that  $\mu$  is an eigenvalue of  $\mathbf{A}(t_2)$ . Show, using Theorem 12.8 with  $\mathbf{A} = \mathbf{A}(t_1)$  and  $\mathbf{E} = \mathbf{A}(t_2) - \mathbf{A}(t_1)$ , that  $\mathbf{A}(t_1)$  has an eigenvalue  $\lambda$  such that

$$|\lambda - \mu| \leq C(t_2 - t_1)^{1/n}, \text{ where } C \leq 2(\|\mathbf{D}\|_2 + \|\mathbf{A} - \mathbf{D}\|_2).$$

Thus, as a function of  $t$ , every eigenvalue of  $\mathbf{A}(t)$  is a continuous function of  $t$ .



Semyon Aranovich Gershgorin, 1901-1933 (left), Jacques Salomon Hadamard, 1865-1963 (right).

**Exercise 12.6 (Nonsingularity using Gerschgorin)**

Consider the matrix

$$\mathbf{A} = \begin{pmatrix} 4 & 1 & 0 & 0 \\ 1 & 4 & 1 & 0 \\ 0 & 1 & 4 & 1 \\ 0 & 0 & 1 & 4 \end{pmatrix}.$$

Show using Gerschgorin's theorem that  $\mathbf{A}$  is nonsingular.



**Exercise 12.7 (Gerschgorin, strictly diagonally dominant matrix)**

Show using Gerschgorin's theorem that a strictly diagonally dominant matrix  $\mathbf{A}$  ( $|a_{i,i}| > \sum_{j \neq i} |a_{i,j}|$  for all  $i$ ) is nonsingular.

**12.3 Perturbation of Eigenvalues**

In this section we study the following problem. Given matrices  $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n \times n}$ , where we think of  $\mathbf{E}$  as a perturbation of  $\mathbf{A}$ . By how much do the eigenvalues of  $\mathbf{A}$  and  $\mathbf{A} + \mathbf{E}$  differ? Not surprisingly this problem is more complicated than the corresponding problem for linear systems.

We illustrate this by considering two examples. Suppose  $\mathbf{A}_0 := \mathbf{0}$  is the zero matrix. If  $\lambda \in \sigma(\mathbf{A}_0 + \mathbf{E}) = \sigma(\mathbf{E})$ , then  $|\lambda| \leq \|\mathbf{E}\|_\infty$  by Theorem 8.28, and any zero eigenvalue of  $\mathbf{A}_0$  is perturbed by at most  $\|\mathbf{E}\|_\infty$ . On the other hand consider for  $\epsilon > 0$  the matrices

$$\mathbf{A}_1 := \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ 0 & 0 & 0 & \cdots & 0 & 0 \end{bmatrix}, \quad \mathbf{E} := \begin{bmatrix} 0 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 0 \\ \epsilon & 0 & 0 & \cdots & 0 & 0 \end{bmatrix} = \epsilon \mathbf{e}_n \mathbf{e}_1^T.$$

The characteristic polynomial of  $\mathbf{A}_1 + \mathbf{E}$  is  $\pi(\lambda) := (-1)^n(\lambda^n - \epsilon)$ , and the zero eigenvalues of  $\mathbf{A}_1$  are perturbed by the amount  $|\lambda| = \|\mathbf{E}\|_\infty^{1/n}$ . Thus, for  $n = 16$ , a perturbation of say  $\epsilon = 10^{-16}$  gives a change in eigenvalue of 0.1.

The following theorem shows that a dependence  $\|\mathbf{E}\|_\infty^{1/n}$  is the worst that can happen.

**Theorem 12.8 (Elsner's theorem(1985))**

Suppose  $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n \times n}$ . To every  $\mu \in \sigma(\mathbf{A} + \mathbf{E})$  there is a  $\lambda \in \sigma(\mathbf{A})$  such that

$$|\mu - \lambda| \leq K \|\mathbf{E}\|_2^{1/n}, \quad K = (\|\mathbf{A}\|_2 + \|\mathbf{A} + \mathbf{E}\|_2)^{1-1/n}. \quad (12.1)$$

**Proof.** Suppose  $\mathbf{A}$  has eigenvalues  $\lambda_1, \dots, \lambda_n$  and let  $\lambda_1$  be one which is closest to  $\mu$ . Let  $\mathbf{u}_1$  with  $\|\mathbf{u}_1\|_2 = 1$  be an eigenvector corresponding to  $\mu$ , and extend  $\mathbf{u}_1$  to an orthonormal basis  $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  of  $\mathbb{C}^n$ . Note that

$$\begin{aligned} \|(\mu \mathbf{I} - \mathbf{A})\mathbf{u}_1\|_2 &= \|(\mathbf{A} + \mathbf{E})\mathbf{u}_1 - \mathbf{A}\mathbf{u}_1\|_2 = \|\mathbf{E}\mathbf{u}_1\|_2 \leq \|\mathbf{E}\|_2, \\ \prod_{j=2}^n \|(\mu \mathbf{I} - \mathbf{A})\mathbf{u}_j\|_2 &\leq \prod_{j=2}^n (|\mu| + \|\mathbf{A}\mathbf{u}_j\|_2) \leq (\|\mathbf{A} + \mathbf{E}\|_2 + \|\mathbf{A}\|_2)^{n-1}. \end{aligned}$$

Using this and Hadamard's inequality (10.6) we find

$$\begin{aligned} |\mu - \lambda_1|^n &\leq \prod_{j=1}^n |\mu - \lambda_j| = |\det(\mu\mathbf{I} - \mathbf{A})| = |\det((\mu\mathbf{I} - \mathbf{A})[\mathbf{u}_1, \dots, \mathbf{u}_n])| \\ &\leq \|(\mu\mathbf{I} - \mathbf{A})\mathbf{u}_1\|_2 \prod_{j=2}^n \|(\mu\mathbf{I} - \mathbf{A})\mathbf{u}_j\|_2 \leq \|\mathbf{E}\|_2 (\|\mathbf{A} + \mathbf{E}\|_2 + \|\mathbf{A}\|_2)^{n-1}. \end{aligned}$$

The result follows by taking  $n$ th roots in this inequality.  $\square$

It follows from this theorem that the eigenvalues depend continuously on the elements of the matrix. The factor  $\|\mathbf{E}\|_2^{1/n}$  shows that this dependence is almost, but not quite, differentiable. As an example, the eigenvalues of the matrix  $\begin{bmatrix} 1 & \\ \epsilon & 1 \end{bmatrix}$  are  $1 \pm \sqrt{\epsilon}$  and this expression is not differentiable at  $\epsilon = 0$ .

Recall that a matrix is nondefective if the eigenvectors form a basis for  $\mathbb{C}^n$ . For nondefective matrices we can get rid of the annoying exponent  $1/n$  in  $\|\mathbf{E}\|_2$ . For a more general discussion than in the following theorem see [28].

**Theorem 12.9 (Absolute errors)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  has linearly independent eigenvectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  and let  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$  be the eigenvector matrix. To any  $\mu \in \mathbb{C}$  and  $\mathbf{x} \in \mathbb{C}^n$  with  $\|\mathbf{x}\|_p = 1$  we can find an eigenvalue  $\lambda$  of  $\mathbf{A}$  such that

$$|\lambda - \mu| \leq K_p(\mathbf{X}) \|\mathbf{r}\|_p, \quad 1 \leq p \leq \infty, \quad (12.2)$$

where  $\mathbf{r} := \mathbf{A}\mathbf{x} - \mu\mathbf{x}$  and  $K_p(\mathbf{X}) := \|\mathbf{X}\|_p \|\mathbf{X}^{-1}\|_p$ . If for some  $\mathbf{E} \in \mathbb{C}^{n \times n}$  it holds that  $(\mu, \mathbf{x})$  is an eigenpair for  $\mathbf{A} + \mathbf{E}$ , then we can find an eigenvalue  $\lambda$  of  $\mathbf{A}$  such that

$$|\lambda - \mu| \leq K_p(\mathbf{X}) \|\mathbf{E}\|_p, \quad 1 \leq p \leq \infty, \quad (12.3)$$

**Proof.** If  $\mu \in \sigma(\mathbf{A})$  then we can take  $\lambda = \mu$  and (12.2), (12.3) hold trivially. So assume  $\mu \notin \sigma(\mathbf{A})$ . Since  $\mathbf{A}$  is nondefective it can be diagonalized, we have  $\mathbf{A} = \mathbf{X}\mathbf{D}\mathbf{X}^{-1}$ , where  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$  and  $(\lambda_j, \mathbf{x}_j)$  are the eigenpairs of  $\mathbf{A}$  for  $j = 1, \dots, n$ . Define  $\mathbf{D}_1 := \mathbf{D} - \mu\mathbf{I}$ . Then  $\mathbf{D}_1^{-1} = \text{diag}((\lambda_1 - \mu)^{-1}, \dots, (\lambda_n - \mu)^{-1})$  exists and

$$\mathbf{X}\mathbf{D}_1^{-1}\mathbf{X}^{-1}\mathbf{r} = (\mathbf{X}(\mathbf{D} - \mu\mathbf{I})\mathbf{X}^{-1})^{-1}\mathbf{r} = (\mathbf{A} - \mu\mathbf{I})^{-1}(\mathbf{A} - \mu\mathbf{I})\mathbf{x} = \mathbf{x}.$$

Using this and Lemma 12.11 below we obtain

$$1 = \|\mathbf{x}\|_p = \|\mathbf{X}\mathbf{D}_1^{-1}\mathbf{X}^{-1}\mathbf{r}\|_p \leq \|\mathbf{D}_1^{-1}\|_p K_p(\mathbf{X}) \|\mathbf{r}\|_p = \frac{K_p(\mathbf{X}) \|\mathbf{r}\|_p}{\min_j |\lambda_j - \mu|}.$$

But then (12.2) follows. If  $(\mathbf{A} + \mathbf{E})\mathbf{x} = \mu\mathbf{x}$  then  $\mathbf{0} = \mathbf{A}\mathbf{x} - \mu\mathbf{x} + \mathbf{E}\mathbf{x} = \mathbf{r} + \mathbf{E}\mathbf{x}$ . But then  $\|\mathbf{r}\|_p = \|\mathbf{E}\mathbf{x}\|_p \leq \|\mathbf{E}\|_p$ . Inserting this in (12.2) proves (12.3).  $\square$

The equation (12.3) shows that for a nondefective matrix the absolute error can be magnified by at most  $K_p(\mathbf{X})$ , the condition number of the eigenvector matrix with respect to inversion. If  $K_p(\mathbf{X})$  is small then a small perturbation changes the eigenvalues by small amounts.

Even if we get rid of the exponent  $1/n$ , the equation (12.3) illustrates that it can be difficult or sometimes impossible to compute accurate eigenvalues and eigenvectors of matrices with almost linearly dependent eigenvectors. On the other hand the eigenvalue problem for normal matrices is better conditioned. Indeed, if  $\mathbf{A}$  is normal then it has a set of orthonormal eigenvectors and the eigenvector matrix is unitary. If we restrict attention to the 2-norm then  $K_2(\mathbf{X}) = 1$  and (12.3) implies the following result.

**Theorem 12.10 (Perturbations, normal matrix)**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is normal and let  $\mu$  be an eigenvalue of  $\mathbf{A} + \mathbf{E}$  for some  $\mathbf{E} \in \mathbb{C}^{n \times n}$ . Then we can find an eigenvalue  $\lambda$  of  $\mathbf{A}$  such that  $|\lambda - \mu| \leq \|\mathbf{E}\|_2$ .

For an even stronger result for Hermitian matrices see Corollary 5.30. We conclude that the situation for the absolute error in an eigenvalue of a Hermitian matrix is quite satisfactory. Small perturbations in the elements are not magnified in the eigenvalues.

In the proof of Theorem 12.9 we used that the  $p$ -norm of a diagonal matrix is equal to its spectral radius.

**Lemma 12.11 ( $p$ -norm of a diagonal matrix)**

If  $\mathbf{A} = \text{diag}(\lambda_1, \dots, \lambda_n)$  is a diagonal matrix then  $\|\mathbf{A}\|_p = \rho(\mathbf{A})$  for  $1 \leq p \leq \infty$ .

*Proof.* For  $p = \infty$  the proof is left as an exercise. For any  $\mathbf{x} \in \mathbb{C}^n$  and  $p < \infty$  we have

$$\|\mathbf{A}\mathbf{x}\|_p = \|[\lambda_1 x_1, \dots, \lambda_n x_n]^T\|_p = \left( \sum_{j=1}^n |\lambda_j|^p |x_j|^p \right)^{1/p} \leq \rho(\mathbf{A}) \|\mathbf{x}\|_p.$$

Thus  $\|\mathbf{A}\|_p = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|_p}{\|\mathbf{x}\|_p} \leq \rho(\mathbf{A})$ . But from Theorem 8.28 we have  $\rho(\mathbf{A}) \leq \|\mathbf{A}\|_p$  and the proof is complete.  $\square$

**Exercise 12.12 ( $\infty$ -norm of a diagonal matrix)**

Give a direct proof that  $\|\mathbf{A}\|_\infty = \rho(\mathbf{A})$  if  $\mathbf{A}$  is diagonal.

For the accuracy of an eigenvalue of small magnitude we are interested in the size of the relative error.

**Theorem 12.13 (Relative errors)**

Suppose in Theorem 12.9 that  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is nonsingular. To any  $\mu \in \mathbb{C}$  and  $\mathbf{x} \in \mathbb{C}^n$  with  $\|\mathbf{x}\|_p = 1$ , we can find an eigenvalue  $\lambda$  of  $\mathbf{A}$  such that

$$\frac{|\lambda - \mu|}{|\lambda|} \leq K_p(\mathbf{X})K_p(\mathbf{A}) \frac{\|\mathbf{r}\|_p}{\|\mathbf{A}\|_p}, \quad 1 \leq p \leq \infty, \quad (12.4)$$

where  $\mathbf{r} := \mathbf{A}\mathbf{x} - \mu\mathbf{x}$ . If for some  $\mathbf{E} \in \mathbb{C}^{n \times n}$  it holds that  $(\mu, \mathbf{x})$  is an eigenpair for  $\mathbf{A} + \mathbf{E}$ , then we can find an eigenvalue  $\lambda$  of  $\mathbf{A}$  such that

$$\frac{|\lambda - \mu|}{|\lambda|} \leq K_p(\mathbf{X})\|\mathbf{A}^{-1}\mathbf{E}\|_p \leq K_p(\mathbf{X})K_p(\mathbf{A}) \frac{\|\mathbf{E}\|_p}{\|\mathbf{A}\|_p}, \quad 1 \leq p \leq \infty, \quad (12.5)$$

**Proof.** Applying Theorem 8.28 to  $\mathbf{A}^{-1}$  we have for any  $\lambda \in \sigma(\mathbf{A})$

$$\frac{1}{\lambda} \leq \|\mathbf{A}^{-1}\|_p = \frac{K_p(\mathbf{A})}{\|\mathbf{A}\|_p}$$

and (12.4) follows from (12.2). To prove (12.5) we define the matrices  $\mathbf{B} := \mu\mathbf{A}^{-1}$  and  $\mathbf{F} := -\mathbf{A}^{-1}\mathbf{E}$ . If  $(\lambda_j, \mathbf{x})$  are the eigenpairs for  $\mathbf{A}$  then  $(\frac{\mu}{\lambda_j}, \mathbf{x})$  are the eigenpairs for  $\mathbf{B}$  for  $j = 1, \dots, n$ . Since  $(\mu, \mathbf{x})$  is an eigenpair for  $\mathbf{A} + \mathbf{E}$  we find

$$(\mathbf{B} + \mathbf{F} - \mathbf{I})\mathbf{x} = (\mu\mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{E} - \mathbf{I})\mathbf{x} = \mathbf{A}^{-1}(\mu\mathbf{I} - (\mathbf{E} + \mathbf{A}))\mathbf{x} = \mathbf{0}.$$

Thus  $(1, \mathbf{x})$  is an eigenpair for  $\mathbf{B} + \mathbf{F}$ . Applying Theorem 12.9 to this eigenvalue we can find  $\lambda \in \sigma(\mathbf{A})$  such that  $|\frac{\mu}{\lambda} - 1| \leq K_p(\mathbf{X})\|\mathbf{F}\|_p = K_p(\mathbf{X})\|\mathbf{A}^{-1}\mathbf{E}\|_p$  which proves the first estimate in (12.5). The second inequality in (12.5) follows from the submultiplicativity of the  $p$ -norm.  $\square$

## 12.4 Unitary Similarity Transformation of a Matrix into Upper Hessenberg Form

Before attempting to find eigenvalues and eigenvectors of a matrix (exceptions are made for certain sparse matrices), it is often advantageous to reduce it by similarity transformations to a simpler form. Orthogonal or unitary similarity transformations are particularly important since they are insensitive to noise in the elements of the matrix. In this section we show how this reduction can be carried out.

Recall that a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is upper Hessenberg if  $a_{i,j} = 0$  for  $j = 1, 2, \dots, i-2$ ,  $i = 3, 4, \dots, n$ . We will reduce  $\mathbf{A} \in \mathbb{C}^{n \times n}$  to upper Hessenberg form by unitary similarity transformations. Let  $\mathbf{A}_1 = \mathbf{A}$  and define  $\mathbf{A}_{k+1} = \mathbf{H}_k\mathbf{A}_k\mathbf{H}_k$  for  $k = 1, 2, \dots, n-2$ . Here  $\mathbf{H}_k$  is a Householder transformation chosen to

introduce zeros in the elements of column  $k$  of  $\mathbf{A}_k$  under the subdiagonal. The final matrix  $\mathbf{A}_{n-1}$  will be upper Hessenberg.

If  $\mathbf{A}_1 = \mathbf{A}$  is Hermitian, the matrix  $\mathbf{A}_{n-1}$  will be Hermitian and tridiagonal. For if  $\mathbf{A}_k^* = \mathbf{A}_k$  then

$$\mathbf{A}_{k+1}^* = (\mathbf{H}_k \mathbf{A}_k \mathbf{H}_k)^* = \mathbf{H}_k \mathbf{A}_k^* \mathbf{H}_k = \mathbf{A}_{k+1}.$$

Since  $\mathbf{A}_{n-1}$  is upper Hessenberg and Hermitian, it must be tridiagonal.

To describe the reduction to upper Hessenberg or tridiagonal form in more detail we partition  $\mathbf{A}_k$  as follows

$$\mathbf{A}_k = \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \\ \mathbf{D}_k & \mathbf{E}_k \end{bmatrix}.$$

Suppose  $\mathbf{B}_k \in \mathbb{C}^{k,k}$  is upper Hessenberg, and the first  $k-1$  columns of  $\mathbf{D}_k \in \mathbb{C}^{n-k,k}$  are zero, i.e.  $\mathbf{D}_k = [\mathbf{0}, \mathbf{0}, \dots, \mathbf{0}, \mathbf{d}_k]$ . Let  $\mathbf{V}_k = \mathbf{I} - \mathbf{v}_k \mathbf{v}_k^* \in \mathbb{C}^{n-k, n-k}$  be a Householder transformation such that  $\mathbf{V}_k \mathbf{d}_k = \alpha_k \mathbf{e}_1$ . Define

$$\mathbf{H}_k = \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_k \end{bmatrix} \in \mathbb{C}^{n \times n}.$$

The matrix  $\mathbf{H}_k$  is a Householder transformation, and we find

$$\begin{aligned} \mathbf{A}_{k+1} &= \mathbf{H}_k \mathbf{A}_k \mathbf{H}_k = \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_k \end{bmatrix} \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \\ \mathbf{D}_k & \mathbf{E}_k \end{bmatrix} \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_k \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{B}_k & \mathbf{C}_k \mathbf{V}_k \\ \mathbf{V}_k \mathbf{D}_k & \mathbf{V}_k \mathbf{E}_k \mathbf{V}_k \end{bmatrix}. \end{aligned}$$

Now  $\mathbf{V}_k \mathbf{D}_k = [\mathbf{V}_k \mathbf{0}, \dots, \mathbf{V}_k \mathbf{0}, \mathbf{V}_k \mathbf{d}_k] = (\mathbf{0}, \dots, \mathbf{0}, \alpha_k \mathbf{e}_1)$ . Moreover, the matrix  $\mathbf{B}_k$  is not affected by the  $\mathbf{H}_k$  transformation. Therefore the upper left  $(k+1) \times (k+1)$  corner of  $\mathbf{A}_{k+1}$  is upper Hessenberg and the reduction is carried one step further. The reduction stops with  $\mathbf{A}_{n-1}$  which is upper Hessenberg.

To find  $\mathbf{A}_{k+1}$  we use Algorithm 10.4 to find  $\mathbf{v}_k$  and  $\alpha_k$ . We store  $\mathbf{v}_k$  in the  $k$ th column of a matrix  $\mathbf{L}$  as  $\mathbf{L}(k+1 : n, k) = \mathbf{v}_k$ . This leads to the following algorithm.

**Algorithm 12.14 (Householder reduction to Hessenberg form)** This algorithm uses Householder similarity transformations to reduce a matrix  $A \in \mathbb{C}^{n \times n}$  to upper Hessenberg form. The reduced matrix  $B$  is tridiagonal if  $A$  is symmetric. Details of the transformations are stored in a lower triangular matrix  $L$ . The elements of  $L$  can be used to assemble a unitary matrix  $Q$  such that  $B = Q^* A Q$ . Algorithm 10.4 is used in each step of the reduction.

```

1 function [L,B] = hesshousegen(A)
2 n=length(A); L=zeros(n,n); B=A;
3 for k=1:n-2
4     [v,B(k+1,k)]=housegen(B(k+1:n,k));
5     L(k+1:n,k)=v; B(k+2:n,k)=zeros(n-k-1,1);
6     C=B(k+1:n,k+1:n); B(k+1:n,k+1:n)=C-v*(v'*C);
7     C=B(1:n,k+1:n); B(1:n,k+1:n)=C-(C*v)*v';
8 end

```

**Exercise 12.15 (Number of arithmetic operations, Hessenberg reduction)**

Show that the number of arithmetic operations for Algorithm 12.14 is  $\frac{10}{3}n^3 = 5G_n$ .

We can use the output of Algorithm 12.14 to assemble the matrix  $Q \in \mathbb{R}^{n \times n}$  such that  $Q$  is orthonormal and  $Q^* A Q$  is upper Hessenberg. We need to compute the product  $Q = H_1 H_2 \cdots H_{n-2}$ , where  $H_k = \begin{bmatrix} I & \mathbf{0} \\ \mathbf{0} & I - \mathbf{v}_k \mathbf{v}_k^T \end{bmatrix}$  and  $\mathbf{v}_k \in \mathbb{R}^{n-k}$ . Since  $\mathbf{v}_1 \in \mathbb{R}^{n-1}$  and  $\mathbf{v}_{n-2} \in \mathbb{R}^2$  it is most economical to assemble the product from right to left. We compute

$$Q_{n-1} = I \text{ and } Q_k = H_k Q_{k+1} \text{ for } k = n-2, n-3, \dots, 1.$$

Suppose  $Q_{k+1}$  has the form  $\begin{bmatrix} I_k & \mathbf{0} \\ \mathbf{0} & U_k \end{bmatrix}$ , where  $U_k \in \mathbb{R}^{n-k, n-k}$ . Then

$$Q_k = \begin{bmatrix} I_k & \mathbf{0} \\ \mathbf{0} & I - \mathbf{v}_k \mathbf{v}_k^T \end{bmatrix} * \begin{bmatrix} I_k & \mathbf{0} \\ \mathbf{0} & U_k \end{bmatrix} = \begin{bmatrix} I_k & \mathbf{0} \\ \mathbf{0} & U_k - \mathbf{v}_k (\mathbf{v}_k^T U_k) \end{bmatrix}.$$

This leads to the following algorithm.

**Algorithm 12.16 (Assemble Householder transformations)**

Suppose  $[L, B] = \text{hesshousegen}(A)$  is the output of Algorithm 12.14. This algorithm assembles an orthonormal matrix  $Q$  from the columns of  $L$  such that  $B = Q^* A Q$  is upper Hessenberg.

```

1 function Q = accumulateQ(L)
2 n=length(L); Q=eye(n);
3 for k=n-2:-1:1
4     v=L(k+1:n,k); C=Q(k+1:n,k+1:n);
5     Q(k+1:n,k+1:n)=C-v*(v'*C);
6 end

```

**Exercise 12.17 (Assemble Householder transformations)**

Show that the number of arithmetic operations required by Algorithm 12.16 is  $\frac{4}{3}n^3 = 2G_n$ .

**Exercise 12.18 (Tridiagonalize a symmetric matrix)**

If  $\mathbf{A}$  is real and symmetric we can modify Algorithm 12.14 as follows. To find  $\mathbf{A}_{k+1}$  from  $\mathbf{A}_k$  we have to compute  $\mathbf{V}_k \mathbf{E}_k \mathbf{V}_k$  where  $\mathbf{E}_k$  is symmetric. Dropping subscripts we have to compute a product of the form  $\mathbf{G} = (\mathbf{I} - \mathbf{v}\mathbf{v}^T)\mathbf{E}(\mathbf{I} - \mathbf{v}\mathbf{v}^T)$ . Let  $\mathbf{w} := \mathbf{E}\mathbf{v}$ ,  $\beta := \frac{1}{2}\mathbf{v}^T\mathbf{w}$  and  $\mathbf{z} := \mathbf{w} - \beta\mathbf{v}$ . Show that  $\mathbf{G} = \mathbf{E} - \mathbf{v}\mathbf{z}^T - \mathbf{z}\mathbf{v}^T$ . Since  $\mathbf{G}$  is symmetric, only the sub- or superdiagonal elements of  $\mathbf{G}$  need to be computed. Computing  $\mathbf{G}$  in this way, it can be shown that we need  $O(4n^3/3)$  operations to tridiagonalize a symmetric matrix by orthonormal similarity transformations. This is less than half the work to reduce a nonsymmetric matrix to upper Hessenberg form. We refer to [27] for a detailed algorithm.

## 12.5 Computing a Selected Eigenvalue of a Symmetric Matrix

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be symmetric with eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ . In this section we consider a method to compute an approximation to the  $m$ th eigenvalue  $\lambda_m$  for some  $1 \leq m \leq n$ . Using Householder similarity transformations as outlined in the previous section we can assume that  $\mathbf{A}$  is symmetric and tridiagonal.

$$\mathbf{A} = \begin{bmatrix} d_1 & c_1 & & & \\ c_1 & d_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & c_{n-2} & d_{n-1} & c_{n-1} \\ & & & c_{n-1} & d_n \end{bmatrix}. \quad (12.6)$$

Suppose one of the off-diagonal elements is equal to zero, say  $c_i = 0$ . We then have  $\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix}$ , where

$$\mathbf{A}_1 = \begin{bmatrix} d_1 & c_1 & & & \\ c_1 & d_2 & c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & c_{i-2} & d_{i-1} & c_{i-1} \\ & & & c_{i-1} & d_i \end{bmatrix} \quad \text{and} \quad \mathbf{A}_2 = \begin{bmatrix} d_{i+1} & c_{i+1} & & & \\ c_{i+1} & d_{i+2} & c_{i+2} & & \\ & \ddots & \ddots & \ddots & \\ & & c_{n-2} & d_{n-1} & c_{n-1} \\ & & & c_{n-1} & d_n \end{bmatrix}.$$

Thus  $\mathbf{A}$  is block diagonal and each diagonal block is tridiagonal. By 6. of Theorem 5.1 we can split the eigenvalue problem into two smaller problems involving  $\mathbf{A}_1$  and  $\mathbf{A}_2$ . We assume that this reduction has been carried out so that  $\mathbf{A}$  is irreducible, i. e.,  $c_i \neq 0$  for  $i = 1, \dots, n-1$ .

We first show that irreducibility implies that the eigenvalues are distinct.

**Lemma 12.19 (Distinct eigenvalues of a tridiagonal matrix)**

An irreducible, tridiagonal and symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  has  $n$  real and distinct eigenvalues.

**Proof.** Let  $\mathbf{A}$  be given by (12.6). By Theorem 5.23 the eigenvalues are real. Define for  $x \in \mathbb{R}$  the polynomial  $p_k(x) := \det(x\mathbf{I}_k - \mathbf{A}_k)$  for  $k = 1, \dots, n$ , where  $\mathbf{A}_k$  is the upper left  $k \times k$  corner of  $\mathbf{A}$  (the leading principal submatrix of order  $k$ ). The eigenvalues of  $\mathbf{A}$  are the roots of the polynomial  $p_n$ . Using the last column to expand for  $k \geq 2$  the determinant  $p_{k+1}(x)$  we find

$$p_{k+1}(x) = (x - d_{k+1})p_k(x) - c_k^2 p_{k-1}(x). \quad (12.7)$$

Since  $p_1(x) = x - d_1$  and  $p_2(x) = (x - d_2)(x - d_1) - c_1^2$  this also holds for  $k = 0, 1$  if we define  $p_{-1}(x) = 0$  and  $p_0(x) = 1$ . For  $M$  sufficiently large we have

$$p_2(-M) > 0, \quad p_2(d_1) < 0, \quad p_2(+M) > 0.$$

Since  $p_2$  is continuous there are  $y_1 \in (-M, d_1)$  and  $y_2 \in (d_1, M)$  such that  $p_2(y_1) = p_2(y_2) = 0$ . It follows that the root  $d_1$  of  $p_1$  separates the roots of  $p_2$ , so  $y_1$  and  $y_2$  must be distinct. Consider next

$$p_3(x) = (x - d_3)p_2(x) - c_2^2 p_1(x) = (x - d_3)(x - y_1)(x - y_2) - c_2^2(x - d_1).$$

Since  $y_1 < d_1 < y_2$  we have for  $M$  sufficiently large

$$p_3(-M) < 0, \quad p_3(y_1) > 0, \quad p_3(y_2) < 0, \quad p_3(+M) > 0.$$

Thus the roots  $x_1, x_2, x_3$  of  $p_3$  are separated by the roots  $y_1, y_2$  of  $p_2$ . In the general case suppose for  $k \geq 2$  that the roots  $z_1, \dots, z_{k-1}$  of  $p_{k-1}$  separate the roots  $y_1, \dots, y_k$  of  $p_k$ . Choose  $M$  so that  $y_0 := -M < y_1, y_{k+1} := M > y_k$ . Then

$$y_0 < y_1 < z_1 < y_2 < z_2 \cdots < z_{k-1} < y_k < y_{k+1}.$$

We claim that for  $M$  sufficiently large

$$p_{k+1}(y_j) = (-1)^{k+1-j} |p_{k+1}(y_j)| \neq 0, \quad \text{for } j = 0, 1, \dots, k+1.$$

This holds for  $j = 0, k+1$ , and for  $j = 1, \dots, k$  since

$$p_{k+1}(y_j) = -c_k^2 p_{k-1}(y_j) = -c_k^2 (y_j - z_1) \cdots (y_j - z_{k-1}).$$

It follows that the roots  $x_1, \dots, x_{k+1}$  are separated by the roots  $y_1, \dots, y_k$  of  $p_k$  and by induction the roots of  $p_n$  (the eigenvalues of  $\mathbf{A}$ ) are distinct.  $\square$



### 12.5.1 The inertia theorem

We say that two matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$  are **congruent** if  $\mathbf{A} = \mathbf{E}^* \mathbf{B} \mathbf{E}$  for some nonsingular matrix  $\mathbf{E} \in \mathbb{C}^{n \times n}$ . By Theorem 5.21 a Hermitian matrix  $\mathbf{A}$  is both congruent and similar to a diagonal matrix  $\mathbf{D}$ ,  $\mathbf{U}^* \mathbf{A} \mathbf{U} = \mathbf{D}$  where  $\mathbf{U}$  is unitary. The eigenvalues of  $\mathbf{A}$  are the diagonal elements of  $\mathbf{D}$ . Let  $\pi(\mathbf{A})$ ,  $\zeta(\mathbf{A})$  and  $\nu(\mathbf{A})$  denote the number of positive, zero and negative eigenvalues of  $\mathbf{A}$ . If  $\mathbf{A}$  is Hermitian then all eigenvalues are real and  $\pi(\mathbf{A}) + \zeta(\mathbf{A}) + \nu(\mathbf{A}) = n$ .

#### Theorem 12.20 (Sylvester's inertia theorem)

If  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$  are Hermitian and congruent then  $\pi(\mathbf{A}) = \pi(\mathbf{B})$ ,  $\zeta(\mathbf{A}) = \zeta(\mathbf{B})$  and  $\nu(\mathbf{A}) = \nu(\mathbf{B})$ .

*Proof.* Suppose  $\mathbf{A} = \mathbf{E}^* \mathbf{B} \mathbf{E}$ , where  $\mathbf{E}$  is nonsingular. Assume first that  $\mathbf{A}$  and  $\mathbf{B}$  are diagonal matrices. Suppose  $\pi(\mathbf{A}) = k$  and  $\pi(\mathbf{B}) = m < k$ . We shall show that this leads to a contradiction. Let  $\mathbf{E}_1$  be the upper left  $m \times k$  corner of  $\mathbf{E}$ . Since  $m < k$ , we can find a nonzero  $\mathbf{x}$  such that  $\mathbf{E}_1 \mathbf{x} = \mathbf{0}$  (cf. Lemma 0.32). Let  $\mathbf{y}^T = [\mathbf{x}^T, \mathbf{0}^T] \in \mathbb{C}^n$ , and  $\mathbf{z} = [z_1, \dots, z_n]^T = \mathbf{E} \mathbf{y}$ . Then  $z_i = 0$  for  $i = 1, 2, \dots, m$ . If  $\mathbf{A}$  has positive eigenvalues  $\lambda_1, \dots, \lambda_k$  and  $\mathbf{B}$  has eigenvalues  $\mu_1, \dots, \mu_n$ , where  $\mu_i \leq 0$  for  $i \geq m + 1$  then

$$\mathbf{y}^* \mathbf{A} \mathbf{y} = \sum_{i=1}^n \lambda_i |y_i|^2 = \sum_{i=1}^k \lambda_i |x_i|^2 > 0.$$

But

$$\mathbf{y}^* \mathbf{A} \mathbf{y} = \mathbf{y}^* \mathbf{E}^* \mathbf{B} \mathbf{E} \mathbf{y} = \mathbf{z}^* \mathbf{B} \mathbf{z} = \sum_{i=m+1}^n \mu_i |z_i|^2 \leq 0,$$

a contradiction.

We conclude that  $\pi(\mathbf{A}) = \pi(\mathbf{B})$  if  $\mathbf{A}$  and  $\mathbf{B}$  are diagonal. Moreover,  $\nu(\mathbf{A}) = \pi(-\mathbf{A}) = \pi(-\mathbf{B}) = \nu(\mathbf{B})$  and  $\zeta(\mathbf{A}) = n - \pi(\mathbf{A}) - \nu(\mathbf{A}) = n - \pi(\mathbf{B}) - \nu(\mathbf{B}) = \zeta(\mathbf{B})$ . This completes the proof for diagonal matrices.

Let in the general case  $\mathbf{U}_1$  and  $\mathbf{U}_2$  be unitary matrices such that  $\mathbf{U}_1^* \mathbf{A} \mathbf{U}_1 = \mathbf{D}_1$  and  $\mathbf{U}_2^* \mathbf{B} \mathbf{U}_2 = \mathbf{D}_2$  where  $\mathbf{D}_1$  and  $\mathbf{D}_2$  are diagonal matrices. Since  $\mathbf{A} = \mathbf{E}^* \mathbf{B} \mathbf{E}$ , we find  $\mathbf{D}_1 = \mathbf{F}^* \mathbf{D}_2 \mathbf{F}$  where  $\mathbf{F} = \mathbf{U}_2^* \mathbf{E} \mathbf{U}_1$  is nonsingular. Thus  $\mathbf{D}_1$  and  $\mathbf{D}_2$  are congruent diagonal matrices. But since  $\mathbf{A}$  and  $\mathbf{D}_1$ ,  $\mathbf{B}$  and  $\mathbf{D}_2$  have the same eigenvalues, we find  $\pi(\mathbf{A}) = \pi(\mathbf{D}_1) = \pi(\mathbf{D}_2) = \pi(\mathbf{B})$ . Similar results hold for  $\zeta$  and  $\nu$ .  $\square$

#### Corollary 12.21 (Counting eigenvalues using the LDLT factorization)

Suppose  $\mathbf{A} = \text{tridiag}(c_i, d_i, c_i) \in \mathbb{R}^{n \times n}$  is symmetric and that  $\alpha \in \mathbb{R}$  is such that  $\mathbf{A} - \alpha \mathbf{I}$  has an symmetric LU factorization, i.e.  $\mathbf{A} - \alpha \mathbf{I} = \mathbf{L} \mathbf{D} \mathbf{L}^T$  where  $\mathbf{L}$  is unit

lower triangular and  $\mathbf{D}$  is diagonal. Then the number of eigenvalues of  $\mathbf{A}$  strictly less than  $\alpha$  equals the number of negative diagonal elements in  $\mathbf{D}$ . The diagonal elements  $d_1(\alpha), \dots, d_n(\alpha)$  in  $\mathbf{D}$  can be computed recursively as follows

$$d_1(\alpha) = d_1 - \alpha, \quad d_k(\alpha) = d_k - \alpha - c_{k-1}^2/d_{k-1}(\alpha), \quad k = 2, 3, \dots, n. \quad (12.8)$$

**Proof.** Since the diagonal elements in  $\mathbf{R}$  in an LU factorization equal the diagonal elements in  $\mathbf{D}$  in an  $\mathbf{LDL}^T$  factorization we see that the formulas in (12.8) follows immediately from (1.4). Since  $\mathbf{L}$  is nonsingular,  $\mathbf{A} - \alpha\mathbf{I}$  and  $\mathbf{D}$  are congruent. By the previous theorem  $v(\mathbf{A} - \alpha\mathbf{I}) = v(\mathbf{D})$ , the number of negative diagonal elements in  $\mathbf{D}$ . If  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  then  $(\mathbf{A} - \alpha\mathbf{I})\mathbf{x} = (\lambda - \alpha)\mathbf{x}$ , and  $\lambda - \alpha$  is an eigenvalue of  $\mathbf{A} - \alpha\mathbf{I}$ . But then  $v(\mathbf{A} - \alpha\mathbf{I})$  equals the number of eigenvalues of  $\mathbf{A}$  which are less than  $\alpha$ .  $\square$

### Exercise 12.22 (Counting eigenvalues)

Consider the matrix in Exercise 12.6. Determine the number of eigenvalues greater than 4.5.

### Exercise 12.23 (Overflow in LDLT factorization)

Let for  $n \in \mathbb{N}$

$$\mathbf{A}_n = \begin{bmatrix} 10 & 1 & 0 & \cdots & 0 \\ 1 & 10 & 1 & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & 1 & 10 & 1 \\ 0 & \cdots & 0 & 1 & 10 \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

- Let  $d_k$  be the diagonal elements of  $\mathbf{D}$  in a symmetric factorization of  $\mathbf{A}_n$ . Show that  $5 + \sqrt{24} < d_k \leq 10$ ,  $k = 1, 2, \dots, n$ .
- Show that  $D_n := \det(\mathbf{A}_n) > (5 + \sqrt{24})^n$ . Give  $n_0 \in \mathbb{N}$  such that your computer gives an overflow when  $D_{n_0}$  is computed in floating point arithmetic.

### Exercise 12.24 (Simultaneous diagonalization)

(Simultaneous diagonalization of two symmetric matrices by a congruence transformation). Let  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$  where  $\mathbf{A}^T = \mathbf{A}$  and  $\mathbf{B}$  is symmetric positive definite. Let  $\mathbf{B} = \mathbf{U}^T \mathbf{D} \mathbf{U}$  where  $\mathbf{U}$  is orthonormal and  $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$ . Let  $\hat{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{U} \mathbf{A} \mathbf{U}^T \mathbf{D}^{-1/2}$  where

$$\mathbf{D}^{-1/2} := \text{diag}(d_1^{-1/2}, \dots, d_n^{-1/2}).$$

a) Show that  $\hat{\mathbf{A}}$  is symmetric.

Let  $\hat{\mathbf{A}} = \hat{\mathbf{U}}^T \hat{\mathbf{D}} \hat{\mathbf{U}}$  where  $\hat{\mathbf{U}}$  is orthonormal and  $\hat{\mathbf{D}}$  is diagonal. Set  $\mathbf{E} = \mathbf{U}^T \mathbf{D}^{-1/2} \hat{\mathbf{U}}^T$ .

b) Show that  $\mathbf{E}$  is nonsingular and that  $\mathbf{E}^T \mathbf{A} \mathbf{E} = \hat{\mathbf{D}}$ ,  $\mathbf{E}^T \mathbf{B} \mathbf{E} = \mathbf{I}$ .

For a more general result see Theorem 10.1 in [17].

## 12.5.2 Approximating $\lambda_m$

Corollary 12.21 can be used to determine the  $m$ th eigenvalue of  $\mathbf{A}$ , where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ . Using Gerschgorin's theorem we first find an interval  $[a, b]$ , such that  $(a, b)$  contains the eigenvalues of  $\mathbf{A}$ . Let for  $x \in [a, b]$

$$\rho(x) := \#\{k : d_k(x) > 0 \text{ for } k = 1, \dots, n\}$$

be the number of eigenvalues of  $\mathbf{A}$  which are strictly greater than  $x$ . Clearly  $\rho(a) = n$ ,  $\rho(b) = 0$ . Choosing a tolerance  $\epsilon$  and using bisection we proceed as follows:

|   |        |
|---|--------|
| <pre> h = b - a; for j = 1 : itmax     c = (a + b)/2;     if b - a &lt; eps * h         λ = (a + b)/2; return     end     k = ρ(c);     if k ≥ m a = c else b = c; end </pre> | (12.9) |
|---|--------|

We generate a sequence  $\{[a_j, b_j]\}$  of intervals, each containing  $\lambda_m$  and  $b_j - a_j = 2^{-j}(b - a)$ .

As it stands this method will fail if in (12.8) one of the  $d_k(\alpha)$  is zero. One possibility is to replace such a  $d_k(\alpha)$  by a suitable small number, say  $\delta_k = c_k \epsilon_M$ , where  $\epsilon_M$  is the Machine epsilon, typically  $2 \times 10^{-16}$  for Matlab. This replacement is done if  $|d_k(\alpha)| < |\delta_k|$ .

### Exercise 12.25 (Program code for one eigenvalue)

Suppose  $\mathbf{A} = \text{tridiag}(\mathbf{c}, \mathbf{d}, \mathbf{c})$  is symmetric and tridiagonal with elements  $d_1, \dots, d_n$  on the diagonal and  $c_1, \dots, c_{n-1}$  on the neighboring subdiagonals. Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  be the eigenvalues of  $\mathbf{A}$ . We shall write a program to compute one eigenvalue  $\lambda_m$  for a given  $m$  using bisection and the method outlined in (12.9).

- a) Write a function `k=count(c,d,x)` which for given  $x$  counts the number of eigenvalues of  $\mathbf{A}$  strictly greater than  $x$ . Use the replacement described above if one of the  $d_j(x)$  is close to zero.
- b) Write a function `lambda=findeigv(c,d,m)` which first estimates an interval  $(a,b]$  containing all eigenvalues of  $\mathbf{A}$  and then generates a sequence  $\{(a_j, b_j]\}$  of intervals each containing  $\lambda_m$ . Iterate until  $b_j - a_j \leq (b - a)\epsilon_M$ , where  $\epsilon_M$  is Matlab's machine epsilon `eps`. Typically  $\epsilon_M \approx 2.22 \times 10^{-16}$ .
- c) Test the program on  $\mathbf{T} := \text{tridiag}(-1, 2, -1)$  of size 100. Compare the exact value of  $\lambda_5$  with your result and the result obtained by using Matlab's built-in function `eig`.

### Exercise 12.26 (Determinant of upper Hessenberg matrix)

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is upper Hessenberg and  $x \in \mathbb{C}$ . We will study two algorithms to compute  $f(x) = \det(\mathbf{A} - x\mathbf{I})$ .

- a) Show that Gaussian elimination without pivoting requires  $O(n^2)$  arithmetic operations.
- b) Show that the number of arithmetic operations is the same if partial pivoting is used.
- c) Estimate the number of arithmetic operations if Givens's rotations are used.
- d) Compare the two methods discussing advantages and disadvantages.

## 12.6 Review Questions

**12.6.1** Suppose  $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n \times n}$ . To every  $\mu \in \sigma(\mathbf{A} + \mathbf{E})$  there is a  $\lambda \in \sigma(\mathbf{A})$  which is in some sense close to  $\mu$ .

- What is the general result (Elsner's theorem)?
- what if  $\mathbf{A}$  is non defective?
- what if  $\mathbf{A}$  is normal?
- what if  $\mathbf{A}$  is Hermitian?

**12.6.2** Can Gerschgorin's theorem be used to check if a matrix is nonsingular?

**12.6.3** How many arithmetic operation does it take to reduce a matrix by similarity transformations to upper Hessenberg form by Householder transformations?

**12.6.4** Give a condition ensuring that a tridiagonal symmetric matrix has real and distinct eigenvalues:

**12.6.5** What is the content of Sylvester's inertia theorem?

**12.6.6** Give an application of this theorem.



## Chapter 13

# The QR Algorithm

The QR algorithm is a method to find all eigenvalues and eigenvectors of a matrix. It is related to a simpler method called the power method and we start studying this method and its variants.

### 13.1 The Power Method and its variants

These methods can be used to compute a single eigenpair of a matrix. They also play a role in the QR algorithm.

#### 13.1.1 The power method

The **power method** in its basic form is a technique to compute the eigenvector corresponding to the largest (in absolute value) eigenvalue of a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . As a by product we can also find the corresponding eigenvalue. We define a sequence  $\{\mathbf{z}_k\}$  of vectors in  $\mathbb{C}^n$  by

$$\mathbf{z}_k := \mathbf{A}^k \mathbf{z}_0 = \mathbf{A} \mathbf{z}_{k-1}, \quad k = 1, 2, \dots \quad (13.1)$$

#### Example 13.1 (Power method)

Let

$$\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}, \quad \mathbf{z}_0 := \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

We find

$$\mathbf{z}_1 = \mathbf{A} \mathbf{z}_0 = \begin{bmatrix} 2 \\ -1 \end{bmatrix}, \quad \mathbf{z}_2 = \mathbf{A} \mathbf{z}_1 = \begin{bmatrix} 5 \\ -4 \end{bmatrix}, \quad \dots, \quad \mathbf{z}_k = \frac{1}{2} \begin{bmatrix} 1 + 3^k \\ 1 - 3^k \end{bmatrix}, \quad \dots$$

It follows that  $2\mathbf{z}_k/3^k$  converges to the eigenvector  $[1, -1]$  corresponding to the dominant eigenvalue  $\lambda = 3$ . The sequence of Rayleigh quotients  $\{\mathbf{z}_k^T \mathbf{A} \mathbf{z}_k / \mathbf{z}_k^T \mathbf{z}_k\}$  will converge to the dominant eigenvalue  $\lambda = 3$ .

To understand better what happens we expand  $\mathbf{z}_0$  in terms of the eigenvectors

$$\mathbf{z}_0 = \frac{1}{2} \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2.$$

Since  $\mathbf{A}^k$  has eigenpairs  $(\lambda_j^k, \mathbf{v}_j)$ ,  $j = 1, 2$  we find

$$\mathbf{z}_k = c_1 \lambda_1^k \mathbf{v}_1 + c_2 \lambda_2^k \mathbf{v}_2 = c_1 3^k \mathbf{v}_1 + c_2 1^k \mathbf{v}_2.$$

Thus  $3^{-k} \mathbf{z}_k = c_1 \mathbf{v}_1 + 3^{-k} c_2 \mathbf{v}_2 \rightarrow c_1 \mathbf{v}_1$ . Since  $c_1 \neq 0$  the result is convergence to the dominant eigenvector.

Let  $\mathbf{A} \in \mathbb{C}^{n \times n}$  have eigenpairs  $(\lambda_j, \mathbf{v}_j)$ ,  $j = 1, \dots, n$  with  $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$ .

Given  $\mathbf{z}_0 \in \mathbb{C}^n$  we assume that

- (i)  $|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n|$ ,
  - (ii)  $\mathbf{z}_0^T \mathbf{v}_1 \neq 0$
  - (iii)  $\mathbf{A}$  has linearly independent eigenvectors.
- (13.2)

The first assumption means that  $\mathbf{A}$  has a dominant eigenvalue  $\lambda_1$  of algebraic multiplicity one. The second assumption says that  $\mathbf{z}_0$  has a component in the direction  $\mathbf{v}_1$ . The third assumption is not necessary, but is included in order to simplify the analysis.

To see what happens let  $\mathbf{z}_0 = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \dots + c_n \mathbf{v}_n$ , where by assumption (ii) of (13.2) we have  $c_1 \neq 0$ . Since  $\mathbf{A}^k \mathbf{v}_j = \lambda_j^k \mathbf{v}_j$  for all  $j$  we see that

$$\mathbf{z}_k = c_1 \lambda_1^k \mathbf{v}_1 + c_2 \lambda_2^k \mathbf{v}_2 + \dots + c_n \lambda_n^k \mathbf{v}_n, \quad k = 0, 1, 2, \dots \quad (13.3)$$

Dividing by  $\lambda_1^k$  we find

$$\frac{\mathbf{z}_k}{\lambda_1^k} = c_1 \mathbf{v}_1 + c_2 \left(\frac{\lambda_2}{\lambda_1}\right)^k \mathbf{v}_2 + \dots + c_n \left(\frac{\lambda_n}{\lambda_1}\right)^k \mathbf{v}_n, \quad k = 0, 1, 2, \dots \quad (13.4)$$

Assumption (i) of (13.2) implies that  $(\lambda_j/\lambda_1)^k \rightarrow 0$  as  $k \rightarrow \infty$  for all  $j \geq 2$  and we obtain

$$\lim_{k \rightarrow \infty} \frac{\mathbf{z}_k}{\lambda_1^k} = c_1 \mathbf{v}_1, \quad (13.5)$$

the dominant eigenvector of  $\mathbf{A}$ . It can be shown that this also holds for defective matrices as long as (i) and (ii) of (13.2) hold, see for example page 58 of [27].



In practice we need to scale the iterates  $\mathbf{z}_k$  somehow and we normally do not know  $\lambda_1$ . Instead we choose a norm on  $\mathbb{C}^n$ , set  $\mathbf{x}_0 = \mathbf{z}_0/\|\mathbf{z}_0\|$  and generate for  $k = 1, 2, \dots$  unit vectors as follows:

$$\begin{aligned} (i) \quad & \mathbf{y}_k = \mathbf{A}\mathbf{x}_{k-1} \\ (ii) \quad & \mathbf{x}_k = \mathbf{y}_k/\|\mathbf{y}_k\|. \end{aligned} \tag{13.6}$$

**Lemma 13.2 (Convergence of the power method)**

Suppose (13.2) holds. Then

$$\lim_{k \rightarrow \infty} \left( \frac{|\lambda_1|}{\lambda_1} \right)^k \mathbf{x}_k = \frac{c_1}{|c_1|} \frac{\mathbf{v}_1}{\|\mathbf{v}_1\|}.$$

In particular, if  $\lambda_1 > 0$  and  $c_1 > 0$  then the sequence  $\{\mathbf{x}_k\}$  will converge to the eigenvector  $\mathbf{u}_1 := \mathbf{v}_1/\|\mathbf{v}_1\|$  of unit length.

**Proof.** By induction on  $k$  it follows that  $\mathbf{x}_k = \mathbf{z}_k/\|\mathbf{z}_k\|$  for all  $k \geq 0$ , where  $\mathbf{z}_k = \mathbf{A}^k \mathbf{z}_0$ . Indeed, this holds for  $k = 1$ , and if it holds for  $k - 1$  then  $\mathbf{y}_k = \mathbf{A}\mathbf{x}_{k-1} = \mathbf{A}\mathbf{z}_{k-1}/\|\mathbf{z}_{k-1}\| = \mathbf{z}_k/\|\mathbf{z}_{k-1}\|$  and  $\mathbf{x}_k = (\mathbf{z}_k/\|\mathbf{z}_{k-1}\|)(\|\mathbf{z}_{k-1}\|/\|\mathbf{z}_k\|) = \mathbf{z}_k/\|\mathbf{z}_k\|$ . But then

$$\mathbf{x}_k = \frac{\mathbf{z}_k}{\|\mathbf{z}_k\|} = \frac{c_1 \lambda_1^k}{|c_1 \lambda_1^k|} \frac{\mathbf{v}_1 + \frac{c_2}{c_1} \left( \frac{\lambda_2}{\lambda_1} \right)^k \mathbf{v}_2 + \dots + \frac{c_n}{c_1} \left( \frac{\lambda_n}{\lambda_1} \right)^k \mathbf{v}_n}{\left\| \mathbf{v}_1 + \frac{c_2}{c_1} \left( \frac{\lambda_2}{\lambda_1} \right)^k \mathbf{v}_2 + \dots + \frac{c_n}{c_1} \left( \frac{\lambda_n}{\lambda_1} \right)^k \mathbf{v}_n \right\|}, \quad k = 0, 1, 2, \dots,$$

and this implies the lemma.  $\square$

Suppose we know an approximate eigenvector  $\mathbf{u}$  of  $\mathbf{A}$ , but not the corresponding eigenvalue  $\mu$ . One way of estimating  $\mu$  is to minimize the Euclidian norm of the residual  $r(\lambda) := \mathbf{A}\mathbf{u} - \lambda\mathbf{u}$ .

**Theorem 13.3 (The Rayleigh quotient minimizes the residual)**

Let  $\mathbf{A} \in \mathbb{C}^{n \times n}$ ,  $\mathbf{u} \in \mathbb{C}^n \setminus \{\mathbf{0}\}$ , and let  $\rho: \mathbb{C} \rightarrow \mathbb{R}$  be given by  $\rho(\lambda) = \|\mathbf{A}\mathbf{u} - \lambda\mathbf{u}\|_2$ . Then  $\rho$  is minimized when  $\lambda := \frac{\mathbf{u}^* \mathbf{A} \mathbf{u}}{\mathbf{u}^* \mathbf{u}}$ , the Rayleigh quotient for  $\mathbf{A}$ .

**Proof.** Assume  $\mathbf{u}^* \mathbf{u} = 1$  and extend  $\mathbf{u}$  to an orthonormal basis  $\{\mathbf{u}, \mathbf{U}\}$  for  $\mathbb{C}^n$ . Then  $\mathbf{U}^* \mathbf{u} = \mathbf{0}$  and

$$\begin{bmatrix} \mathbf{u}^* \\ \mathbf{U}^* \end{bmatrix} (\mathbf{A}\mathbf{u} - \lambda\mathbf{u}) = \begin{bmatrix} \mathbf{u}^* \mathbf{A}\mathbf{u} - \lambda \mathbf{u}^* \mathbf{u} \\ \mathbf{U}^* \mathbf{A}\mathbf{u} - \lambda \mathbf{U}^* \mathbf{u} \end{bmatrix} = \begin{bmatrix} \mathbf{u}^* \mathbf{A}\mathbf{u} - \lambda \\ \mathbf{U}^* \mathbf{A}\mathbf{u} \end{bmatrix}.$$

By unitary invariance of the Euclidian norm

$$\rho(\lambda)^2 = |\mathbf{u}^* \mathbf{A}\mathbf{u} - \lambda|^2 + \|\mathbf{U}^* \mathbf{A}\mathbf{u}\|_2^2,$$

and  $\rho$  has a global minimum at  $\lambda = \mathbf{u}^* \mathbf{A} \mathbf{u}$ .  $\square$

**Exercise 13.4 (Orthogonal vectors)**

Show that  $\mathbf{u}$  and  $\mathbf{A} \mathbf{u} - \lambda \mathbf{u}$  are orthogonal when  $\lambda = \frac{\mathbf{u}^* \mathbf{A} \mathbf{u}}{\mathbf{u}^* \mathbf{u}}$ .

Using Rayleigh quotients we can incorporate the calculation of the eigenvalue into the power iteration. We can then compute the residual and stop the iteration when the residual is sufficiently small. But what does it mean to be sufficiently small? Recall that if  $\mathbf{A}$  is nonsingular with a nonsingular eigenvector matrix  $\mathbf{X}$  and  $(\mu, \mathbf{u})$  is an approximate eigenpair with  $\|\mathbf{u}\|_2 = 1$ , then by (12.4) we can find an eigenvalue  $\lambda$  of  $\mathbf{A}$  such that

$$\frac{|\lambda - \mu|}{|\lambda|} \leq K_2(\mathbf{X})K_2(\mathbf{A}) \frac{\|\mathbf{A} \mathbf{u} - \mu \mathbf{u}\|_2}{\|\mathbf{A}\|_2}.$$

Thus if the relative residual is small and both  $\mathbf{A}$  and  $\mathbf{X}$  are well conditioned then the relative error in the eigenvalue will be small.

This discussion leads to the power method with Rayleigh quotient computation. Given  $\mathbf{A} \in \mathbb{C}^{n \times n}$ , a starting vector  $\mathbf{z} \in \mathbb{C}^n$ , a maximum number  $K$  of iterations, and a convergence tolerance  $tol$ . The power method combined with a Rayleigh quotient estimate for the eigenvalue is used to compute a dominant eigenpair  $(l, \mathbf{x})$  of  $\mathbf{A}$  with  $\|\mathbf{x}\|_2 = 1$ . The integer  $it$  returns the number of iterations needed in order for  $\|\mathbf{A} \mathbf{x} - l \mathbf{x}\|_2 / \|\mathbf{A}\|_F < tol$ . If no such eigenpair is found in  $K$  iterations the value  $it = K + 1$  is returned.

**Algorithm 13.5 (The power method)**

```

1 function [l, x, it]=powerit(A, z, K, tol)
2 af=norm(A, 'fro'); x=z/norm(z);
3 for k=1:K
4     y=A*x; l=x'*y;
5     if norm(y-l*x)/af<tol
6         it=k; x=y/norm(y); return
7     end
8     x=y/norm(y);
9 end
10 it=K+1;
```

**Example 13.6 (Power method)**

We try powerit on the three matrices

$$\mathbf{A}_1 := \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, \quad \mathbf{A}_2 := \begin{bmatrix} 1.7 & -0.4 \\ 0.15 & 2.2 \end{bmatrix}, \quad \text{and} \quad \mathbf{A}_3 = \begin{bmatrix} 1 & 2 \\ -3 & 4 \end{bmatrix}.$$

In each case we start with the random vector  $\mathbf{z} = [0.6602, 0.3420]$  and  $\text{tol} = 10^{-6}$ . For  $\mathbf{A}_1$  we get convergence in 7 iterations, for  $\mathbf{A}_2$  it takes 174 iterations, and for  $\mathbf{A}_3$  we do not get convergence.

The matrix  $\mathbf{A}_3$  does not have a dominant eigenvalue since the two eigenvalues are complex conjugate of each other. Thus the basic condition (i) of (13.2) is not satisfied and the power method diverges. The enormous difference in the rate of convergence for  $\mathbf{A}_1$  and  $\mathbf{A}_2$  can be explained by looking at (13.4). The rate of convergence depends on the ratio  $\frac{|\lambda_2|}{|\lambda_1|}$ . If this ratio is small then the convergence is fast, while it can be quite slow if the ratio is close to one. The eigenvalues of  $\mathbf{A}_1$  are  $\lambda_1 = 5.3723$  and  $\lambda_2 = -0.3723$  giving a quite small ratio of 0.07 and the convergence is fast. On the other hand the eigenvalues of  $\mathbf{A}_2$  are  $\lambda_1 = 2$  and  $\lambda_2 = 1.9$  and the corresponding ratio is 0.95 resulting in slow convergence.

A variant of the power method is the **shifted power method**. In this method we choose a number  $s$  and apply the power method to the matrix  $\mathbf{A} - s\mathbf{I}$ . The number  $s$  is called a shift since it shifts an eigenvalue  $\lambda$  of  $\mathbf{A}$  to  $\lambda - s$  of  $\mathbf{A} - s\mathbf{I}$ . Sometimes the convergence can be faster if the shift is chosen intelligently. For example, if we apply the shifted power method to  $\mathbf{A}_2$  in Example 13.6 with shift 1.8, then with the same starting vector and  $\text{tol}$  as above, we get convergence in 17 iterations instead of 174 for the unshifted algorithm.

### 13.1.2 The inverse power method

Another variant of the power method with Rayleigh quotient is the **inverse power method**. This method can be used to determine any eigenpair  $(\lambda, \mathbf{x})$  of  $\mathbf{A}$  as long as  $\lambda$  has algebraic multiplicity one. In the inverse power method we apply the power method to the inverse matrix  $(\mathbf{A} - s\mathbf{I})^{-1}$ , where  $s$  is a shift. If  $\mathbf{A}$  has eigenvalues  $\lambda_1, \dots, \lambda_n$  in no particular order then  $(\mathbf{A} - s\mathbf{I})^{-1}$  has eigenvalues

$$\mu_1(s) = (\lambda_1 - s)^{-1}, \mu_2(s) = (\lambda_2 - s)^{-1}, \dots, \mu_n(s) = (\lambda_n - s)^{-1}.$$

Suppose  $\lambda_1$  is a simple eigenvalue of  $\mathbf{A}$ . Then  $\lim_{s \rightarrow \lambda_1} |\mu_1(s)| = \infty$ , while  $\lim_{s \rightarrow \lambda_1} \mu_j(s) = (\lambda_j - \lambda_1)^{-1} < \infty$  for  $j = 2, \dots, n$ . Hence, by choosing  $s$  sufficiently close to  $\lambda_1$  the inverse power method will converge to that eigenvalue.

For the inverse power method (13.6) is replaced by

$$\begin{aligned} (i) \quad & (\mathbf{A} - s\mathbf{I})\mathbf{y}_k = \mathbf{x}_{k-1} \\ (ii) \quad & \mathbf{x}_k = \mathbf{y}_k / \|\mathbf{y}_k\|. \end{aligned} \tag{13.7}$$

Note that we solve the linear system rather than computing the inverse matrix. Normally the PLU factorization of  $\mathbf{A} - s\mathbf{I}$  is precomputed in order to speed up the computation.

### 13.1.3 Rayleigh quotient iteration

A variant of the inverse power method is known simply as **Rayleigh quotient iteration**. In this method we change the shift from iteration to iteration, using the previous Rayleigh quotient  $s_{k-1}$  as the current shift. In each iteration we need to compute the following quantities

$$\begin{aligned} (i) \quad & (\mathbf{A} - s_{k-1}\mathbf{I})\mathbf{y}_k = \mathbf{x}_{k-1}, \\ (ii) \quad & \mathbf{x}_k = \mathbf{y}_k / \|\mathbf{y}_k\|, \\ (iii) \quad & s_k = \mathbf{x}_k^* \mathbf{A} \mathbf{x}_k, \\ (iv) \quad & \mathbf{r}_k = \mathbf{A} \mathbf{x}_k - s_k \mathbf{x}_k. \end{aligned}$$

We can avoid the calculation of  $\mathbf{A} \mathbf{x}_k$  in (iii) and (iv). Let

$$\rho_k := \frac{\mathbf{y}_k^* \mathbf{x}_{k-1}}{\mathbf{y}_k^* \mathbf{y}_k}, \quad \mathbf{w}_k := \frac{\mathbf{x}_{k-1}}{\|\mathbf{y}_k\|_2}.$$

Then

$$\begin{aligned} s_k &= \frac{\mathbf{y}_k^* \mathbf{A} \mathbf{y}_k}{\mathbf{y}_k^* \mathbf{y}_k} = s_{k-1} + \frac{\mathbf{y}_k^* (\mathbf{A} - s_{k-1}\mathbf{I}) \mathbf{y}_k}{\mathbf{y}_k^* \mathbf{y}_k} = s_{k-1} + \frac{\mathbf{y}_k^* \mathbf{x}_{k-1}}{\mathbf{y}_k^* \mathbf{y}_k} = s_{k-1} + \rho_k, \\ \mathbf{r}_k &= \mathbf{A} \mathbf{x}_k - s_k \mathbf{x}_k = \frac{\mathbf{A} \mathbf{y}_k - (s_{k-1} + \rho_k) \mathbf{y}_k}{\|\mathbf{y}_k\|_2} = \frac{\mathbf{x}_{k-1} - \rho_k \mathbf{y}_k}{\|\mathbf{y}_k\|_2} = \mathbf{w}_k - \rho_k \mathbf{x}_k. \end{aligned}$$

Another problem is that the linear system in *i*) becomes closer and closer to singular as  $s_k$  converges to the eigenvalue. Thus the system becomes more and more ill-conditioned and we can expect large errors in the computed  $\mathbf{y}_k$ . This is indeed true, but we are lucky. Most of the error occurs in the direction of the eigenvector and this error disappears when we normalize  $\mathbf{y}_k$  in *ii*). Miraculously, the normalized eigenvector will be quite accurate.

Given an approximation  $(s, \mathbf{x})$  to an eigenpair  $(\lambda, \mathbf{v})$  of a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . The following algorithm computes a hopefully better approximation to  $(\lambda, \mathbf{v})$  by doing one Rayleigh quotient iteration. The length  $nr$  of the new residual is also returned

#### Algorithm 13.7 (Rayleigh quotient iteration)

```

1 function [x, s, nr]=rayleight(A, x, s)
2 n=length(x);
3 y=(A-s*eye(n,n))\x;
4 yn=norm(y);
5 w=x/yn;
6 x=y/yn;
7 rho=x'*w;
8 s=s+rho;
9 nr=norm(w-rho*x);

```

| $k$                | 1        | 2         | 3         | 4         | 5         |
|--------------------|----------|-----------|-----------|-----------|-----------|
| $\ \mathbf{r}\ _2$ | 1.0e+000 | 7.7e-002  | 1.6e-004  | 8.2e-010  | 2.0e-020  |
| $ s - \lambda_1 $  | 3.7e-001 | -1.2e-002 | -2.9e-005 | -1.4e-010 | -2.2e-016 |

**Table 13.9.** Quadratic convergence of Rayleigh quotient iteration.

Since the shift changes from iteration to iteration the computation of  $\mathbf{y}$  in `rayleighit` will require  $O(n^3)$  arithmetic operations for a full matrix. For such a matrix it might pay to reduce it to an upper Hessenberg form or tridiagonal form before starting the iteration. However, if we have a good approximation to an eigenpair then only a few iterations are necessary to obtain close to machine accuracy.

If Rayleigh quotient iteration converges the convergence will be quadratic and sometimes even cubic. We illustrate this with an example.

**Example 13.8 (Rayleigh quotient iteration)**

The smallest eigenvalue of the matrix  $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$  is  $\lambda_1 = (5 - \sqrt{33})/2 \approx -0.37$ . Starting with  $\mathbf{x} = [1, 1]^T$  and  $s = 0$  `rayleighit` converges to this eigenvalue and corresponding eigenvector. In Table 13.9 we show the rate of convergence by iterating `rayleighit` 5 times. The errors are approximately squared in each iteration indicating quadratic convergence.

## 13.2 The basic QR Algorithm

The QR algorithm is an iterative method to compute all eigenvalues and eigenvectors of a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . The matrix is reduced to triangular form by a sequence of unitary similarity transformations computed from the QR factorization of  $\mathbf{A}$ . Recall that for a square matrix the QR factorization and the QR decomposition are the same. If  $\mathbf{A} = \mathbf{QR}$  is a QR factorization then  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  is unitary,  $\mathbf{Q}^* \mathbf{Q} = \mathbf{I}$  and  $\mathbf{R} \in \mathbb{C}^{n \times n}$  is upper triangular.

The basic QR algorithm takes the following form:

$$\begin{array}{l}
 \mathbf{A}_1 = \mathbf{A} \\
 \text{for } k = 1, 2, \dots \\
 \quad \mathbf{Q}_k \mathbf{R}_k = \mathbf{A}_k \quad (\text{QR factorization of } \mathbf{A}_k) \\
 \quad \mathbf{A}_{k+1} = \mathbf{R}_k \mathbf{Q}_k. \\
 \text{end}
 \end{array} \tag{13.8}$$

The determination of the QR factorization of  $\mathbf{A}_k$  and the computation of  $\mathbf{R}_k \mathbf{Q}_k$  is called a QR step. It is not at all clear that a QR step does anything

useful. At this point, since  $\mathbf{R}_k = \mathbf{Q}_k^* \mathbf{A}_k$  we find

$$\mathbf{A}_{k+1} = \mathbf{R}_k \mathbf{Q}_k = \mathbf{Q}_k^* \mathbf{A}_k \mathbf{Q}_k, \quad (13.9)$$

so  $\mathbf{A}_{k+1}$  is unitary similar to  $\mathbf{A}_k$ . By induction  $\mathbf{A}_{k+1}$  is unitary similar to  $\mathbf{A}$ . Thus, each  $\mathbf{A}_k$  has the same eigenvalues as  $\mathbf{A}$ . We shall see that the basic QR algorithm is related to the power method.

Here are two examples to illustrate what happens.

**Example 13.10 (QR iteration; real eigenvalues)**

We start with

$$\mathbf{A}_1 = \mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} = \left( \frac{1}{\sqrt{5}} \begin{bmatrix} 2 & -1 \\ 1 & 2 \end{bmatrix} \right) * \left( \frac{1}{\sqrt{5}} \begin{bmatrix} 5 & 4 \\ 0 & 3 \end{bmatrix} \right) = \mathbf{Q}_1 \mathbf{R}_1$$

and obtain

$$\mathbf{A}_2 = \mathbf{R}_1 \mathbf{Q}_1 = \frac{1}{5} \begin{bmatrix} 5 & 4 \\ 0 & 3 \end{bmatrix} * \begin{bmatrix} 2 & -1 \\ 1 & 2 \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 14 & 3 \\ 3 & 6 \end{bmatrix} = \begin{bmatrix} 2.8 & 0.6 \\ 0.6 & 1.2 \end{bmatrix}.$$

Continuing we find

$$\mathbf{A}_4 \approx \begin{bmatrix} 2.997 & -0.074 \\ -0.074 & 1.0027 \end{bmatrix}, \quad \mathbf{A}_{10} \approx \begin{bmatrix} 3.0000 & -0.0001 \\ -0.0001 & 1.0000 \end{bmatrix}$$

$\mathbf{A}_{10}$  is almost diagonal and contains approximations to the eigenvalues  $\lambda_1 = 3$  and  $\lambda_2 = 1$  on the diagonal.

**Example 13.11 (QR iteration; complex eigenvalues)**

Applying the QR iteration (13.8) to the matrix

$$\mathbf{A}_1 = \mathbf{A} = \begin{bmatrix} 0.9501 & 0.8913 & 0.8214 & 0.9218 \\ 0.2311 & 0.7621 & 0.4447 & 0.7382 \\ 0.6068 & 0.4565 & 0.6154 & 0.1763 \\ 0.4860 & 0.0185 & 0.7919 & 0.4057 \end{bmatrix}$$

we obtain

$$\mathbf{A}_{14} = \left[ \begin{array}{c|c|c|c} 2.323 & 0.047223 & -0.39232 & -0.65056 \\ \hline -2.1e-10 & 0.13029 & 0.36125 & 0.15946 \\ \hline -4.1e-10 & -0.58622 & 0.052576 & -0.25774 \\ \hline 1.2e-14 & 3.3e-05 & -1.1e-05 & 0.22746 \end{array} \right].$$

This matrix is almost quasi-triangular and estimates for the eigenvalues  $\lambda_1, \dots, \lambda_4$  of  $\mathbf{A}$  can now easily be determined from the diagonal blocks of  $\mathbf{A}_{14}$ . The  $1 \times 1$  blocks

give us two real eigenvalues  $\lambda_1 \approx 2.323$  and  $\lambda_4 \approx 0.2275$ . The middle  $2 \times 2$  block has complex eigenvalues resulting in  $\lambda_2 \approx 0.0914 + 0.4586i$  and  $\lambda_3 \approx 0.0914 - 0.4586i$ . From Gerschgorin's circle theorem 12.1 and Corollary 12.3 it follows that the approximations to the real eigenvalues are quite accurate. We would also expect the complex eigenvalues to have small absolute errors.

These two examples illustrate what happens in general. The sequence  $(\mathbf{A}_k)_k$  converges to the triangular Schur form (Cf. Theorem 5.13) if all the eigenvalues are real or the quasi-triangular Schur form (Cf. Definition 5.17) if some of the eigenvalues are complex.

### 13.2.1 Relation to the power method

Let us show that the basic QR algorithm is related to the power method. We obtain the QR factorization of the powers  $\mathbf{A}^k$  as follows:

#### Theorem 13.12 (QR and power)

For  $k = 1, 2, 3, \dots$ , the QR factorization of  $\mathbf{A}^k$  is  $\mathbf{A}^k = \tilde{\mathbf{Q}}_k \tilde{\mathbf{R}}_k$ , where

$$\tilde{\mathbf{Q}}_k := \mathbf{Q}_1 \cdots \mathbf{Q}_k \text{ and } \tilde{\mathbf{R}}_k := \mathbf{R}_k \cdots \mathbf{R}_1, \quad (13.10)$$

and  $\mathbf{Q}_1, \dots, \mathbf{Q}_k, \mathbf{R}_1, \dots, \mathbf{R}_k$  are the matrices generated by the basic QR algorithm (13.8).

**Proof.** By (13.9)

$$\mathbf{A}_k = \mathbf{Q}_{k-1}^* \mathbf{A}_{k-1} \mathbf{Q}_{k-1} = \mathbf{Q}_{k-1}^* \mathbf{Q}_{k-2}^* \mathbf{A}_{k-2} \mathbf{Q}_{k-2} \mathbf{Q}_{k-1} = \cdots = \tilde{\mathbf{Q}}_{k-1}^* \mathbf{A} \tilde{\mathbf{Q}}_{k-1}. \quad (13.11)$$

The proof is by induction on  $k$ . Clearly  $\tilde{\mathbf{Q}}_1 \tilde{\mathbf{R}}_1 = \mathbf{Q}_1 \mathbf{R}_1 = \mathbf{A}_1$ . Suppose  $\tilde{\mathbf{Q}}_{k-1} \tilde{\mathbf{R}}_{k-1} = \mathbf{A}^{k-1}$  for some  $k \geq 2$ . Since  $\mathbf{Q}_k \mathbf{R}_k = \mathbf{A}_k$  and using (13.11)

$$\tilde{\mathbf{Q}}_k \tilde{\mathbf{R}}_k = \tilde{\mathbf{Q}}_{k-1} (\mathbf{Q}_k \mathbf{R}_k) \tilde{\mathbf{R}}_{k-1} = \tilde{\mathbf{Q}}_{k-1} \mathbf{A}_k \tilde{\mathbf{R}}_{k-1} = (\tilde{\mathbf{Q}}_{k-1} \tilde{\mathbf{Q}}_{k-1}^*) \mathbf{A} \tilde{\mathbf{Q}}_{k-1} \tilde{\mathbf{R}}_{k-1} = \mathbf{A}^k.$$

□

Since  $\tilde{\mathbf{R}}_k$  is upper triangular, its first column is a multiple of  $\mathbf{e}_1$  so that

$$\mathbf{A}^k \mathbf{e}_1 = \tilde{\mathbf{Q}}_k \tilde{\mathbf{R}}_k \mathbf{e}_1 = \tilde{r}_{11}^{(k)} \tilde{\mathbf{Q}}_k \mathbf{e}_1 \text{ or } \tilde{\mathbf{q}}_1^{(k)} := \tilde{\mathbf{Q}}_k \mathbf{e}_1 = \frac{1}{\tilde{r}_{11}^{(k)}} \mathbf{A}^k \mathbf{e}_1.$$

Since  $\|\tilde{\mathbf{q}}_1^{(k)}\|_2 = 1$  the first column of  $\tilde{\mathbf{Q}}_k$  is the result of applying the normalized power iteration (13.6) to the starting vector  $\mathbf{x}_0 = \mathbf{e}_1$ . If this iteration converges we conclude that the first column of  $\tilde{\mathbf{Q}}_k$  must converge to a dominant eigenvector of  $\mathbf{A}$ . It can be shown that the first column of  $\mathbf{A}_k$  must then converge to  $\lambda_1 \mathbf{e}_1$ , where  $\lambda_1$  is a dominant eigenvalue of  $\mathbf{A}$ . This is clearly what happens in Examples 13.10 and 13.11. Indeed, what is observed in practice is that the sequence  $(\tilde{\mathbf{Q}}_k^* \mathbf{A} \tilde{\mathbf{Q}}_k)_k$  converges to a (quasi-triangular) Schur form of  $\mathbf{A}$ .

$$\mathbf{A} = \begin{bmatrix} x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & 0 & x \end{bmatrix} \xrightarrow{\mathbf{P}_{12}^*} \begin{bmatrix} x & x & x & x \\ \mathbf{x} & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & 0 & x \end{bmatrix} \xrightarrow{\mathbf{P}_{23}^*} \begin{bmatrix} x & x & x & x \\ x & x & x & x \\ 0 & \mathbf{x} & x & x \\ 0 & 0 & 0 & x \end{bmatrix} \xrightarrow{\mathbf{P}_{34}^*} \begin{bmatrix} x & x & x & x \\ x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & \mathbf{x} & x \end{bmatrix}.$$

Figure 13.1. Post multiplication in a QR step.

### 13.2.2 Invariance of the Hessenberg form

One QR step requires  $O(n^3)$  arithmetic operations for a matrix  $\mathbf{A}$  of order  $n$ . By an initial reduction of  $\mathbf{A}$  to upper Hessenberg form  $\mathbf{H}_1$  using Algorithm 12.14, the cost of a QR step can be reduced to  $O(n^2)$ . Consider a QR step on  $\mathbf{H}_1$ . We first determine plane rotations  $\mathbf{P}_{i,i+1}$ ,  $i = 1, \dots, n-1$  so that  $\mathbf{P}_{n-1,n} \cdots \mathbf{P}_{1,2} \mathbf{H}_1 = \mathbf{R}_1$  is upper triangular. The details were described in Section 10.4. Thus  $\mathbf{H}_1 = \mathbf{Q}_1 \mathbf{R}_1$ , where  $\mathbf{Q}_1 = \mathbf{P}_{1,2}^* \cdots \mathbf{P}_{n-1,n}^*$  is a QR factorization of  $\mathbf{H}_1$ . To finish the QR step we compute  $\mathbf{R}_1 \mathbf{Q}_1 = \mathbf{R}_1 \mathbf{P}_{1,2}^* \cdots \mathbf{P}_{n-1,n}^*$ . This postmultiplication step is illustrated by the Wilkinson diagram in Figure 13.1.

The postmultiplication by  $\mathbf{P}_{i,i+1}$  introduces a nonzero in position  $(i+1, i)$  leaving the other elements marked by a zero in Figure 13.1 unchanged. Thus the final matrix  $\mathbf{R} \mathbf{P}_{1,2}^* \cdots \mathbf{P}_{n-1,n}^*$  is upper Hessenberg and a QR step leaves the Hessenberg form invariant.

In conclusion, to compute  $\mathbf{A}_{k+1}$  from  $\mathbf{A}_k$  requires  $O(n^2)$  arithmetic operations if  $\mathbf{A}_k$  is upper Hessenberg and  $O(n)$  arithmetic operations if  $\mathbf{A}_k$  is tridiagonal.

### 13.2.3 Deflation

If a subdiagonal element  $a_{i+1,i}$  of an upper Hessenberg matrix  $\mathbf{A}$  is equal to zero, then the eigenvalues of  $\mathbf{A}$  are the union of the eigenvalues of the two smaller matrices  $A(1:i, 1:i)$  and  $A(i+1:n, i+1:n)$ . Thus if during the iteration the  $(i+1, i)$  element of  $\mathbf{A}_k$  is sufficiently small then we can continue the iteration on the two smaller submatrices separately.

To see what effect this can have on the eigenvalues of  $\mathbf{A}$  suppose  $|a_{i+1,i}^{(k)}| \leq \epsilon$ . Let  $\hat{\mathbf{A}}_k := \mathbf{A}_k - a_{i+1,i}^{(k)} \mathbf{e}_{i+1} \mathbf{e}_i^T$  be the matrix obtained from  $\mathbf{A}_k$  by setting the  $(i+1, i)$  element equal to zero. Since  $\mathbf{A}_k = \tilde{\mathbf{Q}}_{k-1}^* \mathbf{A} \tilde{\mathbf{Q}}_{k-1}$  we have

$$\hat{\mathbf{A}}_k = \tilde{\mathbf{Q}}_{k-1}^* (\mathbf{A} + \mathbf{E}) \tilde{\mathbf{Q}}_{k-1}, \quad \mathbf{E} = \tilde{\mathbf{Q}}_{k-1} (a_{i+1,i}^{(k)} \mathbf{e}_{i+1} \mathbf{e}_i^T) \tilde{\mathbf{Q}}_{k-1}^*.$$

Since  $\tilde{\mathbf{Q}}_{k-1}$  is unitary,  $\|\mathbf{E}\|_F = \|a_{i+1,i}^{(k)} \mathbf{e}_{i+1} \mathbf{e}_i^T\|_F = |a_{i+1,i}^{(k)}| \leq \epsilon$  and setting  $a_{i+1,i}^{(k)} = 0$  amounts to a perturbation in the original  $\mathbf{A}$  of at most  $\epsilon$ . For how to choose  $\epsilon$  see the discussion on page 94-95 in [27].

This deflation occurs often in practice and can with a proper implementation reduce the computation time considerably. It should be noted that to find the eigenvectors of the original matrix one has to continue with some care, see [27].



## 13.3 The Shifted QR Algorithms

Like in the inverse power method it is possible to speed up the convergence by introducing shifts. The **explicitly shifted QR algorithm** works as follows:

```

 $\mathbf{A}_1 = \mathbf{A}$ 
for  $k = 1, 2, \dots$ 
  Choose a shift  $s_k$ 
   $\mathbf{Q}_k \mathbf{R}_k = \mathbf{A}_k - s_k \mathbf{I}$    (QR factorization of  $\mathbf{A}_k - s_k \mathbf{I}$ )
   $\mathbf{A}_{k+1} = \mathbf{R}_k \mathbf{Q}_k + s_k \mathbf{I}$ .
end

```

Since  $\mathbf{R}_k = \mathbf{Q}_k^* (\mathbf{A}_k - s_k \mathbf{I})$  we find

$$\mathbf{A}_{k+1} = \mathbf{Q}_k^* (\mathbf{A}_k - s_k \mathbf{I}) \mathbf{Q}_k + s_k \mathbf{I} = \mathbf{Q}_k^* \mathbf{A}_k \mathbf{Q}_k$$

and  $\mathbf{A}_{k+1}$  and  $\mathbf{A}_k$  are unitary similar.

The shifted QR algorithm is related to the power method with shift, cf. Theorem 13.12 and also the inverse power method. In fact the last column of  $\mathbf{Q}_k$  is the result of one iteration of the inverse power method to  $\mathbf{A}^*$  with shift  $s_k$ . Indeed, since  $\mathbf{A} - s_k \mathbf{I} = \mathbf{Q}_k \mathbf{R}_k$  we have  $(\mathbf{A} - s_k \mathbf{I})^* = \mathbf{R}_k^* \mathbf{Q}_k^*$  and  $(\mathbf{A} - s_k \mathbf{I})^* \mathbf{Q}_k = \mathbf{R}_k^*$ . Thus, since  $\mathbf{R}_k^*$  is lower triangular with  $n, n$  element  $\bar{r}_{nn}^{(k)}$  we find  $(\mathbf{A} - s_k \mathbf{I})^* \mathbf{Q}_k \mathbf{e}_n = \mathbf{R}_k^* \mathbf{e}_n = \bar{r}_{nn}^{(k)} \mathbf{e}_n$  from which the conclusion follows.

The shift  $s_k := \mathbf{e}_n^T \mathbf{A}_k \mathbf{e}_n$  is called the **Rayleigh quotient shift**, while the eigenvalue of the lower right  $2 \times 2$  corner of  $\mathbf{A}_k$  closest to the  $n, n$  element of  $\mathbf{A}_k$  is called the **Wilkinson shift**. This shift can be used to find complex eigenvalues of a real matrix. The convergence is very fast and at least quadratic both for the Rayleigh quotient shift and the Wilkinson shift.

By doing two QR iterations at a time it is possible to find both real and complex eigenvalues without using complex arithmetic. The corresponding algorithm is called the **implicitly shifted QR algorithm**

After having computed the eigenvalues we can compute the eigenvectors in steps. First we find the eigenvectors of the triangular or quasi-triangular matrix. We then compute the eigenvectors of the upper Hessenberg matrix and finally we get the eigenvectors of  $\mathbf{A}$ .

Practical experience indicates that only  $O(n)$  iterations are needed to find all eigenvalues of  $\mathbf{A}$ . Thus both the explicit- and implicit shift QR algorithms are normally  $O(n^3)$  algorithms.

For further remarks and detailed algorithms see [27].

## 13.4 A Convergence Theorem

There is no theorem which proves convergence of the QR algorithm in general. The following theorem shows convergence of the basic QR algorithm under somewhat restrictive assumptions.

### Theorem 13.13 (Convergence of basis QR)

Suppose in the basic QR algorithm (13.8) that

1.  $\mathbf{A} \in \mathbb{R}^{n \times n}$  can be diagonalized,  $\mathbf{X}^{-1}\mathbf{A}\mathbf{X} = \mathbf{\Lambda} := \text{diag}(\lambda_1, \dots, \lambda_n)$ .
2. The eigenvalues  $\lambda_1, \dots, \lambda_n$  are real with  $|\lambda_1| > |\lambda_2| > \dots > |\lambda_n| > 0$ .
3. The inverse of the eigenvector matrix has an LU factorization  $\mathbf{X}^{-1} = \mathbf{L}\mathbf{R}$ .

Let  $\tilde{\mathbf{Q}}_k = \mathbf{Q}_1 \dots \mathbf{Q}_k$  for  $k \geq 1$ . Then there is a diagonal matrix  $\mathbf{D}_k$  with diagonal elements  $\pm 1$  such that  $\tilde{\mathbf{Q}}_k \mathbf{D}_k \rightarrow \mathbf{Q}$ , where  $\mathbf{Q}^T \mathbf{A} \mathbf{Q}$  is triangular and  $\mathbf{Q}$  is the  $Q$ -factor in the QR factorization of the eigenvector matrix  $\mathbf{X}$ .

**Proof.** In this proof we assume that every QR factorization has an  $\mathbf{R}$  with positive diagonal elements so that the factorization is unique. Let  $\mathbf{X} = \mathbf{Q}\mathbf{R}$  be the QR factorization of  $\mathbf{X}$ . We observe that  $\mathbf{Q}^T \mathbf{A} \mathbf{Q}$  is upper triangular. For since  $\mathbf{X}^{-1} \mathbf{A} \mathbf{X} = \mathbf{\Lambda}$  we have  $\mathbf{R}^{-1} \mathbf{Q}^T \mathbf{A} \mathbf{Q} \mathbf{R} = \mathbf{\Lambda}$  so that  $\mathbf{Q}^T \mathbf{A} \mathbf{Q} = \mathbf{R} \mathbf{\Lambda} \mathbf{R}^{-1}$  is upper triangular. Since  $\mathbf{A}_{k+1} = \tilde{\mathbf{Q}}_k^T \mathbf{A} \tilde{\mathbf{Q}}_k$ , it is enough to show that  $\tilde{\mathbf{Q}}_k \mathbf{D}_k \rightarrow \mathbf{Q}$  for some diagonal matrix  $\mathbf{D}_k$  with diagonal elements  $\pm 1$ .

We define the nonsingular matrices

$$\mathbf{F}_k := \mathbf{R} \mathbf{\Lambda}^k \mathbf{L} \mathbf{\Lambda}^{-k} \mathbf{R}^{-1} = \hat{\mathbf{Q}}_k \hat{\mathbf{R}}_k, \quad \mathbf{G}_k := \hat{\mathbf{R}}_k \mathbf{R} \mathbf{\Lambda}^k \mathbf{R}, \quad \mathbf{D}_k := \text{diag} \left( \frac{\delta_1}{|\delta_1|}, \dots, \frac{\delta_n}{|\delta_n|} \right),$$

where  $\delta_1, \dots, \delta_n$  are the diagonal elements in the upper triangular matrix  $\mathbf{G}_k$  and  $\mathbf{F}_k = \hat{\mathbf{Q}}_k \hat{\mathbf{R}}_k$  is the QR factorization of  $\mathbf{F}_k$ . Then

$$\begin{aligned} \mathbf{A}^k &= \mathbf{X} \mathbf{\Lambda}^k \mathbf{X}^{-1} = \mathbf{Q} \mathbf{R} \mathbf{\Lambda}^k \mathbf{L} \mathbf{R} = \mathbf{Q} (\mathbf{R} \mathbf{\Lambda}^k \mathbf{L} \mathbf{\Lambda}^{-k} \mathbf{R}^{-1}) (\mathbf{R} \mathbf{\Lambda}^k \mathbf{R}) \\ &= \mathbf{Q} \mathbf{F}_k (\mathbf{R} \mathbf{\Lambda}^k \mathbf{R}) = \mathbf{Q} \hat{\mathbf{Q}}_k \hat{\mathbf{R}}_k (\mathbf{R} \mathbf{\Lambda}^k \mathbf{R}) = (\mathbf{Q} \hat{\mathbf{Q}}_k \mathbf{D}_k^{-1}) (\mathbf{D}_k \mathbf{G}_k), \end{aligned}$$

and this is the QR factorization of  $\mathbf{A}^k$ . Indeed,  $\mathbf{Q} \hat{\mathbf{Q}}_k \mathbf{D}_k^{-1}$  is a product of orthonormal matrices and therefore orthonormal. Moreover  $\mathbf{D}_k \mathbf{G}_k$  is a product of upper triangular matrices and therefore upper triangular. Note that  $\mathbf{D}_k$  is chosen so that this matrix has positive diagonal elements. By Theorem 13.12  $\mathbf{A}^k = \tilde{\mathbf{Q}}_k \tilde{\mathbf{R}}_k$  is also the QR factorization of  $\mathbf{A}^k$ , and we must have  $\tilde{\mathbf{Q}}_k = \mathbf{Q} \hat{\mathbf{Q}}_k \mathbf{D}_k^{-1}$  or  $\tilde{\mathbf{Q}}_k \mathbf{D}_k = \mathbf{Q} \hat{\mathbf{Q}}_k$ . The theorem will follow if we can show that  $\hat{\mathbf{Q}}_k \rightarrow \mathbf{I}$ .

The matrix  $\mathbf{\Lambda}^k \mathbf{L} \mathbf{\Lambda}^{-k}$  is lower triangular with elements  $(\frac{\lambda_i}{\lambda_j})^k l_{ij}$  on and under the diagonal. Thus for  $n = 3$

$$\mathbf{\Lambda}^k \mathbf{L} \mathbf{\Lambda}^{-k} = \begin{bmatrix} 1 & 0 & 0 \\ (\frac{\lambda_2}{\lambda_1})^k l_{21} & 1 & 0 \\ (\frac{\lambda_3}{\lambda_1})^k l_{31} & (\frac{\lambda_3}{\lambda_2})^k l_{32} & 1 \end{bmatrix}.$$

By Assumption 2. it follows that  $\mathbf{\Lambda}^k \mathbf{L} \mathbf{\Lambda}^{-k} \rightarrow \mathbf{I}$ , and hence  $\mathbf{F}_k \rightarrow \mathbf{I}$ . Since  $\hat{\mathbf{R}}_k^T \hat{\mathbf{R}}_k$  is the Cholesky factorization of  $\mathbf{F}_k^T \mathbf{F}_k$  it follows that  $\hat{\mathbf{R}}_k^T \hat{\mathbf{R}}_k \rightarrow \mathbf{I}$ . By the continuity of the Cholesky factorization it holds  $\hat{\mathbf{R}}_k \rightarrow \mathbf{I}$  and hence  $\hat{\mathbf{R}}_k^{-1} \rightarrow \mathbf{I}$ . But then  $\hat{\mathbf{Q}}_k = \mathbf{F}_k \hat{\mathbf{R}}_k^{-1} \rightarrow \mathbf{I}$ .  $\square$

### Exercise 13.14 (QR convergence detail)

Use Theorem 7.31 to show that  $\hat{\mathbf{R}}_k \rightarrow \mathbf{I}$  implies  $\hat{\mathbf{R}}_k^{-1} \rightarrow \mathbf{I}$ .

## 13.5 Review Questions

**13.5.1** What is the main use of the power method?

**13.5.2** Can the QR method be used to find all eigenvectors of a matrix?

**13.5.3** Can the power method be used to find an eigenvalue?

**13.5.4** Do the power method converge to an eigenvector corresponding to a complex eigenvalue?

**13.5.5** What is the inverse power method?

**13.5.6** Give a relation between the QR algorithm and the power method.

**13.5.7** How can we make the basic QR algorithm converge faster?



**Part VI**

**Appendix**



## Appendix A

# Determinants

The first systematic treatment of determinants was given by Cauchy in 1812. He adopted the word “determinant”. The first use of determinants was made by Leibniz in 1693 in a letter to De L’Hôpital. By the beginning of the 20th century the theory of determinants filled four volumes of almost 2000 pages (Muir, 1906–1923. Historic references can be found in this work). The main use of determinants in this text will be to study the characteristic polynomial of a matrix.

In this section we prove the elementary properties of determinants that we need.

### A.1 Permutations

For  $n \in \mathbb{N}$ , let  $N_n = \{1, 2, \dots, n\}$ . A *permutation* is a function  $\sigma : N_n \rightarrow N_n$  which is one-to-one and onto. That is,  $\{\sigma(1), \sigma(2), \dots, \sigma(n)\}$  is a rearrangement of  $\{1, 2, \dots, n\}$ . If  $n = 2$ , there are two permutations  $\{1, 2\}$  and  $\{2, 1\}$ , while for  $n = 3$  we have six permutations  $\{1, 2, 3\}$ ,  $\{1, 3, 2\}$ ,  $\{2, 1, 3\}$ ,  $\{2, 3, 1\}$ ,  $\{3, 1, 2\}$  and  $\{3, 2, 1\}$ . We denote the set of all permutations on  $N_n$  by  $S_n$ . There are  $n!$  elements in  $S_n$ .

If  $\sigma, \tau$  are two permutations in  $S_n$ , we can define their product  $\sigma\tau$  as

$$\sigma\tau = \{\sigma(\tau(1)), \sigma(\tau(2)), \dots, \sigma(\tau(n))\}.$$

For example if  $\sigma = \{1, 3, 2\}$  and  $\tau = \{3, 2, 1\}$ , then  $\sigma\tau = \{\sigma(3), \sigma(2), \sigma(1)\} = \{2, 3, 1\}$ , while  $\tau\sigma = \{\tau(1), \tau(3), \tau(2)\} = \{3, 1, 2\}$ . Thus in general  $\sigma\tau \neq \tau\sigma$ . It is easily shown that the product of two permutations  $\sigma, \tau$  is a permutation, i.e.  $\sigma\tau : N_n \rightarrow N_n$  is one-to-one and onto.

The permutation  $\epsilon = \{1, 2, \dots, n\}$  is called the *identity permutation* in  $S_n$ .

We have  $\epsilon\sigma = \sigma\epsilon = \sigma$  for all  $\sigma \in S_n$ .

Since each  $\sigma \in S_n$  is one-to-one and onto, it has a unique inverse  $\sigma^{-1}$ . To define  $\sigma^{-1}(j)$  for  $j \in N_n$ , we find the unique  $i$  such that  $\sigma(i) = j$ . Then  $\sigma^{-1}(j) = i$ . We have  $\sigma^{-1}\sigma = \sigma\sigma^{-1} = \epsilon$ . As an example, if  $\sigma = \{2, 3, 1\}$  then  $\sigma^{-1} = \{3, 1, 2\}$ , and  $\sigma^{-1}\sigma = \sigma\sigma^{-1} = \{1, 2, 3\} = \epsilon$ .

With each  $\sigma \in S_n$  we can associate a + or - sign. We define

$$\text{sign}(\sigma) = \frac{g(\sigma)}{|g(\sigma)|},$$

where

$$g(\sigma) = \prod_{i=2}^n (\sigma(i) - \sigma(1))(\sigma(i) - \sigma(2)) \cdots (\sigma(i) - \sigma(i-1)).$$

For example if  $\epsilon = \{1, 2, 3, 4\}$  and  $\sigma = \{4, 3, 1, 2\}$ , then

$$\begin{aligned} g(\epsilon) &= (2-1)(3-1)(3-2)(4-1)(4-2)(4-3) = 1! \cdot 2! \cdot 3! > 0, \\ g(\sigma) &= (3-4)(1-4)(1-3)(2-4)(2-3)(2-1) \\ &= (-1)(-3)(-2)(-2)(-1) \cdot 1 = -1! \cdot 2! \cdot 3! < 0. \end{aligned}$$

Thus  $\text{sign}(\epsilon) = +1$  and  $\text{sign}(\sigma) = -1$ .

$g(\sigma)$  contains one positive factor  $(2-1)$  and five negative ones. The negative factors are called *inversions*. The number of inversions equals the number of times a bigger integer precedes a smaller one in  $\sigma$ . That is, in  $\{4, 3, 1, 2\}$  4 precedes 3, 1 and 2 (three inversions corresponding to the negative factors  $(3-4)$ ,  $(1-4)$  and  $(2-4)$  in  $g(\sigma)$ ), and 3 precedes 1 and 2 ( $(1-3)$  and  $(2-3)$  in  $g(\sigma)$ ). This makes it possible to compute  $\text{sign}(\sigma)$  without actually writing down  $g(\sigma)$ .

In general, the sign function has the following properties

1.  $\text{sign}(\epsilon) = 1$ .
2.  $\text{sign}(\sigma\tau) = \text{sign}(\sigma)\text{sign}(\tau)$  for  $\sigma, \tau \in S_n$ .
3.  $\text{sign}(\sigma^{-1}) = \text{sign}(\sigma)$  for  $\sigma \in S_n$ .

Since all factors in  $g(\epsilon)$  are positive, we have  $g(\epsilon) = |g(\epsilon)|$  and  $\text{sign}(\epsilon) = 1$ . This proves 1. To prove 2 we first note that for any  $S_n$

$$\text{sign}(\sigma) = \frac{g(\sigma)}{g(\epsilon)}.$$

Since  $g(\sigma)$  and  $g(\epsilon)$  contain the same factors apart from signs and  $g(\epsilon) > 0$ , we have  $|g(\sigma)| = g(\epsilon)$ . Now

$$\text{sign}(\sigma\tau) = \frac{g(\sigma\tau)}{g(\epsilon)} = \frac{g(\sigma\tau)}{g(\tau)} \frac{g(\tau)}{g(\epsilon)} = \frac{g(\sigma\tau)}{g(\tau)} \text{sign}(\tau).$$



We have to show that  $g(\sigma\tau)/g(\tau) = g(\sigma)/g(\epsilon)$ . We write  $g(\sigma)/g(\epsilon)$  in the form

$$\frac{g(\sigma)}{g(\epsilon)} = \prod_{i=2}^n \prod_{j=1}^{i-1} r_{\sigma}(i, j), \quad r_{\sigma}(i, j) = \frac{\sigma(i) - \sigma(j)}{i - j}.$$

Now

$$\frac{g(\sigma\tau)}{g(\tau)} = \frac{\prod_{i=2}^n (\sigma(\tau(i)) - \sigma(\tau(1))) \cdots (\sigma(\tau(i)) - \sigma(\tau(i-1)))}{\prod_{i=2}^n (\tau(i) - \tau(1)) \cdots (\tau(i) - \tau(i-1))} = \prod_{i=2}^n \prod_{j=1}^{i-1} r_{\sigma}(\tau(i), \tau(j)).$$

$\tau$  is a permutation so  $g(\sigma)/g(\epsilon)$  and  $g(\sigma\tau)/g(\tau)$  contain the same factors. Moreover, the sign of the factors are the same since  $r(i, j) = r(j, i)$  for all  $i \neq j$ . Thus  $g(\sigma)/g(\epsilon) = g(\sigma\tau)/g(\tau)$ , and 2 is proved. Finally, 3 follows from 1 and 2;  $1 = \text{sign}(\epsilon) = \text{sign}(\sigma\sigma^{-1}) = \text{sign}(\sigma)\text{sign}(\sigma^{-1})$  so that  $\sigma$  and  $\sigma^{-1}$  have the same sign.

### Example A.1 (Properties of permutations)

It can be shown that  $\rho(\sigma\tau) = (\rho\sigma)\tau$  for  $\rho, \sigma, \tau \in S_n$ , i.e. multiplication of permutations is associative. (In fact, we have

1. Multiplication is associative.
2. There exists an identity permutation  $\epsilon$ .
3. Every permutation has an inverse.

Thus the set  $S_n$  of permutations is a group with respect to multiplication.  $S_n$  is called the symmetric group of degree  $n$ ).

## A.2 Basic Properties of Determinants

For any  $A \in \mathbb{C}^{n \times n}$  the determinant of  $A$  is defined the number

$$\det(A) = \sum_{\sigma \in S_n} \text{sign}(\sigma) a_{\sigma(1),1} a_{\sigma(2),2} \cdots a_{\sigma(n),n}. \quad (\text{A.1})$$

This sum ranges of all  $n!$  permutations of  $\{1, 2, \dots, n\}$ . We also denote the determinant by (Cayley, 1841)

$$\begin{vmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{vmatrix}.$$

From the definition we have

$$\begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{21}a_{12}.$$

The first term on the right corresponds to the identity permutation  $\epsilon$  given by  $\epsilon(i) = i$ ,  $i = 1, 2$ . The second term comes from the permutation  $\sigma = \{2, 1\}$ . For  $n = 3$

$$\begin{vmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{vmatrix} = a_{11}a_{22}a_{33} - a_{11}a_{32}a_{23} - a_{21}a_{12}a_{33} \\ + a_{21}a_{32}a_{13} + a_{31}a_{12}a_{23} - a_{31}a_{22}a_{13}.$$

The following is a list of properties of determinants.

1. **Triangular matrix** The determinant of a triangular matrix is the product of the diagonal elements.  $\det(A) = a_{11}a_{22} \cdots a_{nn}$ . In particular  $\det(I) = 1$ .
2. **Transpose**  $\det(A^T) = \det(A)$ .
3. **Homogeneity** For any  $\beta_i \in \mathbb{C}$ ,  $i = 1, 2, \dots, n$ , we have

$$\det([\beta_1 \mathbf{a}_1, \beta_2 \mathbf{a}_2, \dots, \beta_n \mathbf{a}_n]) = \beta_1 \beta_2 \cdots \beta_n \det([\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]).$$

4. **Permutation of columns** If  $\tau \in S_n$  then

$$\det(\mathbf{B}) := \det([\mathbf{a}_{\tau(1)}, \mathbf{a}_{\tau(2)}, \dots, \mathbf{a}_{\tau(n)}]) = \text{sign}(\tau) \det([\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]).$$

5. **Additivity**

$$\det([\mathbf{a}_1, \dots, \mathbf{a}_{k-1}, \mathbf{a}_k + \mathbf{a}'_k, \mathbf{a}_{k+1}, \dots, \mathbf{a}_n]) \\ = \det([\mathbf{a}_1, \dots, \mathbf{a}_n]) + \det([\mathbf{a}_1, \dots, \mathbf{a}'_k, \dots, \mathbf{a}_n]).$$

6. **Singular matrix**  $\det(A) = 0$  if and only if  $A$  is singular.
7. **Product rule** If  $A, B \in \mathbb{C}^{n \times n}$  then  $\det(AB) = \det(A) \det(B)$ .
8. **Block triangular** If  $A$  is block triangular with diagonal blocks  $B$  and  $C$  then  $\det(A) = \det(B) \det(C)$ .

**Proof.**

1. If  $\sigma \neq \epsilon$ , we can find distinct integers  $i$  and  $j$  such that  $\sigma(i) > i$  and  $\sigma(j) < j$ . But then  $a_{\sigma(i),i} = 0$  if  $\mathbf{A}$  is upper triangular and  $a_{\sigma(j),j} = 0$  if  $\mathbf{A}$  is lower triangular. Hence

$$\det(\mathbf{A}) = \text{sign}(\epsilon) a_{\epsilon(1),1} a_{\epsilon(2),2} \cdots a_{\epsilon(n),n} = a_{1,1} a_{2,2} \cdots a_{n,n}.$$

Since the identity matrix is triangular with all diagonal elements equal to one, we have that  $\det(\mathbf{I}) = 1$ .

2. By definition of  $\mathbf{A}^T$  and the det-function

$$\det(\mathbf{A}^T) = \sum_{\sigma \in S_n} \text{sign}(\sigma) a_{1,\sigma(1)} a_{2,\sigma(2)} \cdots a_{n,\sigma(n)}.$$

Consider an element  $a_{i,\sigma(i)}$ . If  $\sigma(i) = j$  then

$$a_{i,\sigma(i)} = a_{\sigma^{-1}(j),j}.$$

Since  $\sigma(1), \sigma(2), \dots, \sigma(n)$  ranges through  $\{1, 2, \dots, n\}$ , we obtain

$$\begin{aligned} \det(\mathbf{A}^T) &= \sum_{\sigma \in S_n} \text{sign}(\sigma) a_{\sigma^{-1}(1),1} a_{\sigma^{-1}(2),2} \cdots a_{\sigma^{-1}(n),n} \\ &= \sum_{\sigma \in S_n} \text{sign}(\sigma^{-1}) a_{\sigma^{-1}(1),1} a_{\sigma^{-1}(2),2} \cdots a_{\sigma^{-1}(n),n} \\ &= \sum_{\sigma^{-1} \in S_n} \text{sign}(\sigma^{-1}) a_{\sigma^{-1}(1),1} a_{\sigma^{-1}(2),2} \cdots a_{\sigma^{-1}(n),n} \\ &= \det(\mathbf{A}). \end{aligned}$$

3. This follows immediately from the definition of  $\det[(\beta_1 \mathbf{a}_1, \beta_2 \mathbf{a}_2, \dots, \beta_n \mathbf{a}_n)]$ .

4. We have

$$\det(\mathbf{B}) = \sum_{\sigma \in S_n} \text{sign}(\sigma) a_{\sigma(1),\tau(1)} a_{\sigma(2),\tau(2)} \cdots a_{\sigma(n),\tau(n)}.$$

Fix  $i$  in  $\{1, 2, \dots, n\}$ . Let  $k = \sigma(i)$  and  $m = \tau(i)$ . Then  $\tau^{-1}(m) = i$  and  $\sigma(\tau^{-1}(m)) = k$ . Hence

$$a_{\sigma(i),\tau(i)} = a_{k,m} = a_{\sigma\tau^{-1}(m),m}.$$

Moreover,  $\text{sign}(\sigma) = \text{sign}(\tau) \text{sign}(\sigma\tau^{-1})$ . Thus

$$\det(\mathbf{B}) = \text{sign}(\tau) \sum_{\sigma \in S_n} \text{sign}(\sigma\tau^{-1}) a_{\sigma\tau^{-1}(1),1} a_{\sigma\tau^{-1}(2),2} \cdots a_{\sigma\tau^{-1}(n),n}.$$

But as  $\sigma$  ranges over  $S_n$ ,  $\sigma\tau^{-1}$  also ranges over  $S_n$ . Hence

$$\det(\mathbf{B}) = \text{sign}(\tau) \det[(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n)].$$

5. This follows at once from the definition.

6. We observe that the determinant of a matrix is equal to the product of the eigenvalues and that a matrix is singular if and only if zero is an eigenvalue (cf. Theorems 5.2, 0.54). But then the result follows.

7. To better understand the general proof, we do the  $2 \times 2$  case first. Let  $\mathbf{A} = (\mathbf{a}_1, \mathbf{a}_2)$ ,  $\mathbf{B} = (\mathbf{b}_1, \mathbf{b}_2)$ . Then

$$\mathbf{A}\mathbf{B} = (\mathbf{A}\mathbf{b}_1, \mathbf{A}\mathbf{b}_2) = (b_{1,1}\mathbf{a}_1 + b_{2,1}\mathbf{a}_2, b_{1,2}\mathbf{a}_1 + b_{2,2}\mathbf{a}_2).$$

Using the additivity, we obtain

$$\det(\mathbf{AB}) = \det(b_{1,1}\mathbf{a}_1, b_{1,2}\mathbf{a}_1) + \det(b_{2,1}\mathbf{a}_2, b_{1,2}\mathbf{a}_1) \\ + \det(b_{1,1}\mathbf{a}_1, b_{2,2}\mathbf{a}_2) + \det(b_{2,1}\mathbf{a}_2, b_{2,2}\mathbf{a}_2).$$

Next we have by homogeneity

$$\det(\mathbf{AB}) = b_{1,1}b_{1,2} \det(\mathbf{a}_1, \mathbf{a}_1) + b_{2,1}b_{1,2} \det(\mathbf{a}_2, \mathbf{a}_1) \\ + b_{1,1}b_{2,2} \det(\mathbf{a}_1, \mathbf{a}_2) + b_{2,1}b_{2,2} \det(\mathbf{a}_2, \mathbf{a}_2).$$

Property 6 implies that  $\det(\mathbf{a}_1, \mathbf{a}_1) = \det(\mathbf{a}_2, \mathbf{a}_2) = 0$ . Using Property 4, we obtain  $\det(\mathbf{a}_2, \mathbf{a}_1) = -\det(\mathbf{a}_1, \mathbf{a}_2)$  and

$$\det(\mathbf{AB}) = (b_{1,1}b_{2,2} - b_{2,1}b_{1,2}) \det(\mathbf{a}_1, \mathbf{a}_2) = \det(\mathbf{B}) \det(\mathbf{A}).$$

The proof for  $n > 2$  follows the  $n = 2$  case step by step. Let  $\mathbf{C} = (\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n) = \mathbf{AB}$ . Then

$$\mathbf{c}_i = \mathbf{A}\mathbf{b}_i = b_{1,i}\mathbf{a}_1 + b_{2,i}\mathbf{a}_2 + \dots + b_{n,i}\mathbf{a}_n, \quad i = 1, 2, \dots, n.$$

Using the additivity, we obtain

$$\det(\mathbf{AB}) = \sum_{i_1=1}^n \sum_{i_2=1}^n \dots \sum_{i_n=1}^n \det[(b_{i_1,1}\mathbf{a}_{i_1}, b_{i_2,2}\mathbf{a}_{i_2}, \dots, b_{i_n,n}\mathbf{a}_{i_n})].$$

Next we have by homogeneity

$$\det(\mathbf{AB}) = \sum_{i_1=1}^n \sum_{i_2=1}^n \dots \sum_{i_n=1}^n b_{i_1,1}b_{i_2,2} \dots b_{i_n,n} \det[(\mathbf{a}_{i_1}, \mathbf{a}_{i_2}, \dots, \mathbf{a}_{i_n})].$$

Property 6 implies that  $\det[(\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_n})] = 0$  if any two of the indices  $i_1, \dots, i_n$  are equal. Therefore we only get a contribution to the sum whenever  $i_1, \dots, i_n$  is a permutation of  $\{1, 2, \dots, n\}$ . Thus

$$\det(\mathbf{AB}) = \sum_{\sigma \in S_n} b_{\sigma(1),1} \dots b_{\sigma(n),n} \det[(\mathbf{a}_{\sigma(1)}, \dots, \mathbf{a}_{\sigma(n)})].$$

By Property 4 we obtain

$$\det(\mathbf{AB}) = \sum_{\sigma \in S_n} \text{sign}(\tau) b_{\sigma(1),1} \dots b_{\sigma(n),n} \det[(\mathbf{a}_1, \dots, \mathbf{a}_n)].$$

According to the definition of  $\det(\mathbf{B})$  this is equal to  $\det(\mathbf{B}) \det(\mathbf{A})$ .

8. Suppose  $\mathbf{A}$  is block upper triangular. Let

$$S_{n,k} = \{\sigma \in S_n : \sigma(i) \leq k \text{ if } i \leq k, \text{ and } \sigma(i) \geq k+1 \text{ if } i \geq k+1\}.$$

We claim that  $a_{\sigma(1),1} \cdots a_{\sigma(n),n} = 0$  if  $\sigma \notin S_{n,k}$ , because if  $\sigma(i) > k$  for some  $i \leq k$  then  $a_{\sigma(i),i} = 0$  since it lies in the zero part of  $\mathbf{A}$ . If  $\sigma(i) \leq k$  for some  $i \geq k+1$ , we must have  $\sigma(j) > k$  for some  $j \leq k$  to make “room” for  $\sigma(i)$ , and  $a_{\sigma(j),j} = 0$ . It follows that

$$\det(\mathbf{A}) = \sum_{\sigma \in S_{n,k}} \text{sign}(\sigma) a_{\sigma(1),1} \cdots a_{\sigma(n),n}.$$

Define

$$\rho(i) = \begin{cases} \sigma(i) & i = 1, \dots, k \\ i & i = k+1, \dots, n, \end{cases} \quad \tau(i) = \begin{cases} i & i = 1, \dots, k \\ \sigma(i) & i = k+1, \dots, n. \end{cases}$$

If  $\sigma \in S_{n,k}$ ,  $\rho$  and  $\tau$  will be permutations. Moreover,  $\sigma = \rho\tau$ . Define  $\hat{\rho}$  and  $\hat{\tau}$  in  $S_k$  and  $S_{n-k}$  respectively by  $\hat{\rho}(i) = \rho(i)$ ,  $i = 1, \dots, k$ , and  $\hat{\tau}(i) = \tau(i+k) - k$  for  $i = 1, \dots, n-k$ . As  $\sigma$  ranges over  $S_{n,k}$ ,  $\hat{\rho}$  and  $\hat{\tau}$  will take on all values in  $S_k$  and  $S_{n-k}$  respectively. Since  $\text{sign}(\hat{\rho}) = \text{sign}(\rho)$  and  $\text{sign}(\hat{\tau}) = \text{sign}(\tau)$ , we find

$$\text{sign}(\sigma) = \text{sign}(\rho)\text{sign}(\tau) = \text{sign}(\hat{\rho})\text{sign}(\hat{\tau}).$$

Then

$$\begin{aligned} \det(\mathbf{A}) &= \sum_{\hat{\rho} \in S_k} \sum_{\hat{\tau} \in S_{n-k}} \text{sign}(\hat{\rho})\text{sign}(\hat{\tau}) b_{\hat{\rho}(1),1} \cdots b_{\hat{\rho}(k),k} d_{\hat{\tau}(1),1} \cdots d_{\hat{\tau}(n-k),n-k} \\ &= \det(\mathbf{B}) \det(\mathbf{D}). \end{aligned}$$

□

## A.3 The Adjoint Matrix and Cofactor Expansion

We start with a useful formula for the solution of a linear system.

Let  $\mathbf{A}_j(\mathbf{b})$  denote the matrix obtained from  $\mathbf{A}$  by replacing the  $j$ th column of  $\mathbf{A}$  by  $\mathbf{b}$ . For example,

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}, & \mathbf{b} &= \begin{bmatrix} 3 \\ 6 \end{bmatrix}, & \mathbf{A}_1(\mathbf{b}) &= \begin{bmatrix} 3 & 2 \\ 6 & 1 \end{bmatrix}, & \mathbf{A}_2(\mathbf{b}) &= \begin{bmatrix} 1 & 3 \\ 2 & 6 \end{bmatrix}, \\ \mathbf{I} &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, & \mathbf{x} &= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, & \mathbf{I}_1(\mathbf{x}) &= \begin{bmatrix} x_1 & 0 \\ x_2 & 1 \end{bmatrix}, & \mathbf{I}_2(\mathbf{x}) &= \begin{bmatrix} 1 & x_1 \\ 0 & x_2 \end{bmatrix}. \end{aligned}$$

**Theorem A.2 (Cramer's rule (1750))**

Suppose  $\mathbf{A} \in \mathbb{C}^{n \times n}$  with  $\det(\mathbf{A}) \neq 0$  and  $\mathbf{b} \in \mathbb{C}^n$ . Let  $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$  be the unique solution of  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Then

$$x_j = \frac{\det(\mathbf{A}_j(\mathbf{b}))}{\det(\mathbf{A})}, \quad j = 1, 2, \dots, n.$$

**Proof.** Since  $1 = \det(\mathbf{I}) = \det(\mathbf{A}\mathbf{A}^{-1}) = \det(\mathbf{A})\det(\mathbf{A}^{-1})$  we have  $\det(\mathbf{A}^{-1}) = 1/\det(\mathbf{A})$ . Then

$$\begin{aligned} \frac{\det(\mathbf{A}_j(\mathbf{b}))}{\det(\mathbf{A})} &= \det(\mathbf{A}^{-1}\mathbf{A}_j(\mathbf{b})) \\ &= \det([\mathbf{A}^{-1}\mathbf{a}_1, \dots, \mathbf{A}^{-1}\mathbf{a}_{j-1}, \mathbf{A}^{-1}\mathbf{b}, \mathbf{A}^{-1}\mathbf{a}_{j+1}, \dots, \mathbf{A}^{-1}\mathbf{a}_n]) \\ &= \det([\mathbf{e}_1, \dots, \mathbf{e}_{j-1}, \mathbf{x}, \mathbf{e}_{j+1}, \dots, \mathbf{e}_n]) = x_j, \end{aligned}$$

where we used Property 8 for the last equality.  $\square$

Let  $\mathbf{A}_{i,j}$  denote the submatrix of  $\mathbf{A}$  obtained by deleting the  $i$ th row and  $j$ th column of  $\mathbf{A}$ . For example,

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}, & \mathbf{A}_{1,1} &= \begin{bmatrix} 5 & 6 \\ 8 & 9 \end{bmatrix}, & \mathbf{A}_{1,2} &= \begin{bmatrix} 4 & 6 \\ 7 & 9 \end{bmatrix}, \\ \mathbf{A}_{2,1} &= \begin{bmatrix} 2 & 3 \\ 8 & 9 \end{bmatrix}, & \mathbf{A}_{2,2} &= \begin{bmatrix} 1 & 3 \\ 7 & 9 \end{bmatrix}, & & \text{etc.} \end{aligned}$$

**Definition A.3 (Cofactor and Adjoint)**

For  $\mathbf{A} \in \mathbb{C}^{n \times n}$  and  $1 \leq i, j \leq n$  the determinant  $\det(\mathbf{A}_{i,j})$  is called the **cofactor** of  $a_{ij}$ . The matrix  $\text{adj}(\mathbf{A}) \in \mathbb{C}^{n \times n}$  with elements  $(-1)^{i+j} \det(\mathbf{A}_{j,i})$  is called the **adjoint** of  $\mathbf{A}$ .

**Theorem A.4 (The inverse as an adjoint)**

If  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is nonsingular then

$$\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \text{adj}(\mathbf{A}).$$

**Proof.** Let  $\mathbf{A}^{-1} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ , where  $\mathbf{x}_j = [x_{1j}, \dots, x_{nj}]^T$ . The equation  $\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$  implies that  $\mathbf{A}\mathbf{x}_j = \mathbf{e}_j$  for  $j = 1, \dots, n$  and by Cramer's rule

$$x_{ij} = \frac{\det(\mathbf{A}_i(\mathbf{e}_j))}{\det(\mathbf{A})} = (-1)^{i+j} \frac{\det(\mathbf{A}_{ji})}{\det(\mathbf{A})}, \quad j = 1, 2, \dots, n.$$

For the last equality we first interchange the first and  $i$ th column of  $\mathbf{A}_i(\mathbf{e}_j)$ . By Property 4 it follows that  $\det(\mathbf{A}_i(\mathbf{e}_j)) = (-1)^{i-1} \det([\mathbf{e}_j, \mathbf{a}_1, \dots, \mathbf{a}_{i-1}, \mathbf{a}_{i+1}, \dots, \mathbf{a}_n])$ . We then interchange row  $j$  and row 1. Using Property 8 we obtain

$$\det(\mathbf{A}_i(\mathbf{e}_j)) = (-1)^{i+j-2} \det(\mathbf{A}_{ji}) = (-1)^{i+j} \det(\mathbf{A}_{ji}).$$

□

### Corollary A.5 (The adjoint and the inverse)

For any  $\mathbf{A} \in \mathbb{C}^{n \times n}$  we have

$$\mathbf{A} \operatorname{adj}(\mathbf{A}) = \operatorname{adj}(\mathbf{A})\mathbf{A} = \det(\mathbf{A})\mathbf{I}. \quad (\text{A.2})$$

**Proof.** If  $\mathbf{A}$  is nonsingular then (A.2) follows from Theorem A.4. We simply multiply by  $\mathbf{A}$  from the left and from the right. Suppose next that  $\mathbf{A}$  is singular with  $m$  zero eigenvalues  $\lambda_1, \dots, \lambda_m$  and nonzero eigenvalues  $\lambda_{m+1}, \dots, \lambda_n$ . We define  $\epsilon_0 := \min_{m+1 \leq j \leq n} |\lambda_j|$ . For any  $\epsilon \in (0, \epsilon_0)$  the matrix  $\mathbf{A} + \epsilon\mathbf{I}$  has nonzero eigenvalues  $\epsilon, \dots, \epsilon, \lambda_{m+1} + \epsilon, \dots, \lambda_n + \epsilon$  and hence is nonsingular. By what we have proved

$$(\mathbf{A} + \epsilon\mathbf{I}) \operatorname{adj}(\mathbf{A} + \epsilon\mathbf{I}) = \operatorname{adj}(\mathbf{A} + \epsilon\mathbf{I})(\mathbf{A} + \epsilon\mathbf{I}) = \det(\mathbf{A} + \epsilon\mathbf{I})\mathbf{I}. \quad (\text{A.3})$$

Since the elements in  $\mathbf{A} + \epsilon\mathbf{I}$  and  $\operatorname{adj}(\mathbf{A} + \epsilon\mathbf{I})$  depend continuously on  $\epsilon$  we can take limits in (A.3) to obtain (A.2). □

### Corollary A.6 (Cofactor expansion)

For any  $\mathbf{A} \in \mathbb{C}^{n \times n}$  we have

$$\det(\mathbf{A}) = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{ij}) \text{ for } i = 1, \dots, n, \quad (\text{A.4})$$

$$\det(\mathbf{A}) = \sum_{i=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{ij}) \text{ for } j = 1, \dots, n. \quad (\text{A.5})$$

**Proof.** By (A.2) we have  $\mathbf{A} \operatorname{adj}(\mathbf{A}) = \det(\mathbf{A})\mathbf{I}$ . But then  $\det(\mathbf{A}) = \mathbf{e}_i^T \mathbf{A} \operatorname{adj}(\mathbf{A}) \mathbf{e}_i = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det(\mathbf{A}_{ij})$  which is (A.4). Applying this row expansion to  $\mathbf{A}^T$  we find  $\det(\mathbf{A}^T) = \sum_{j=1}^n (-1)^{i+j} a_{ji} \det(\mathbf{A}_{ji})$ . Switching the roles of  $i$  and  $j$  proves (A.5). □

## A.4 Computing Determinants

A determinant of an  $n$ -by- $n$  matrix computed from the definition can contain up to  $n!$  terms and we need other methods to compute determinants.

A matrix can be reduced to upper triangular form using elementary row operations. We can then use Property 1. to compute the determinant. The elementary operations using either rows or columns are

1. Interchanging two rows(columns).
2. Multiply a row(column) by a scalar  $\alpha$ .
3. Add a constant multiple of one row(column) to another row(column).

Let  $\mathbf{B}$  be the result of performing an elementary operation on  $\mathbf{A}$ . For the three elementary operations the numbers  $\det(\mathbf{A})$  and  $\det(\mathbf{B})$  are related as follows.

1.  $\det(\mathbf{B}) = -\det(\mathbf{A})$  (from Property 4.)
2.  $\det(\mathbf{B}) = \alpha \det(\mathbf{A})$  (from Property 3.)
3.  $\det(\mathbf{B}) = \det(\mathbf{A})$  (from Properties 5., 7.)

It follows from Property 2. that it is enough to show this for column operations. The proof of 1. and 2. are immediate. For 3. suppose we add  $\alpha$  times column  $k$  to column  $i$  for some  $k \neq i$ . Then using Properties 5. and 7. we find

$$\begin{aligned} \det(\mathbf{B}) &= \det \left( [\mathbf{a}_1, \dots, \mathbf{a}_{i-1}, \mathbf{a}_i + \alpha \mathbf{a}_k, \mathbf{a}_{i+1}, \dots, \mathbf{a}_n] \right) \\ &\stackrel{5.}{=} \det(\mathbf{A}) + \det \left( [\mathbf{a}_1, \dots, \mathbf{a}_{i-1}, \alpha \mathbf{a}_k, \mathbf{a}_{i+1}, \dots, \mathbf{a}_n] \right) \stackrel{7.}{=} \det(\mathbf{A}) \end{aligned}$$

## A.5 Some Useful Determinant Formulas

Suppose  $\mathbf{A} \in \mathbb{C}^{m \times n}$  and suppose for an integer  $r \leq \min\{m, n\}$  that  $\mathbf{i} = \{i_1, \dots, i_r\}$  and  $\mathbf{j} = \{j_1, \dots, j_r\}$  are integers with  $1 \leq i_1 < i_2 < \dots < i_r \leq m$  and  $1 \leq j_1 < j_2 < \dots < j_r$ . We let

$$\mathbf{A}(\mathbf{i}, \mathbf{j}) = \begin{bmatrix} a_{i_1, j_1} & \cdots & a_{i_1, j_r} \\ \vdots & & \vdots \\ a_{i_r, j_1} & \cdots & a_{i_r, j_r} \end{bmatrix}$$

be the submatrix of  $\mathbf{A}$  consisting of rows  $i_1, \dots, i_r$  and columns  $j_1, \dots, j_r$ . The following formula bears a strong resemblance to the formula for matrix multiplication.

### Theorem A.7 (Cauchy-Binet formula)

Let  $\mathbf{A} \in \mathbb{C}^{m \times p}$ ,  $\mathbf{B} \in \mathbb{C}^{p \times n}$  and  $\mathbf{C} = \mathbf{AB}$ . Suppose  $1 \leq r \leq \min\{m, n, p\}$  and let



$\mathbf{i} = \{i_1, \dots, i_r\}$  and  $\mathbf{j} = \{j_1, \dots, j_r\}$  be integers with  $1 \leq i_1 < i_2 < \dots < i_r \leq m$  and  $1 \leq j_1 < j_2 < \dots < j_r \leq n$ . Then

$$\det(\mathbf{C}(\mathbf{i}, \mathbf{j})) = \sum_{\mathbf{k}} \det(\mathbf{A}(\mathbf{i}, \mathbf{k})) \det(\mathbf{B}(\mathbf{k}, \mathbf{j})), \quad (\text{A.6})$$

where we sum over all  $\mathbf{k} = \{k_1, \dots, k_r\}$  with  $1 \leq k_1 < k_2 < \dots < k_r \leq p$ .



## Appendix B

# Computer Arithmetic

### B.1 Absolute and Relative Errors

Suppose  $a$  and  $b$  are real or complex scalars. If  $b$  is an approximation to  $a$  then there are different ways of measuring the error in  $b$ .

**Definition B.1 (Absolute error)**

*The absolute error in  $b$  as an approximation to  $a$  is the number  $\epsilon := |a - b|$ . The number  $e := b - a$  is called the error in  $b$  as an approximation to  $a$ . This is what we have to add to  $a$  to get  $b$ .*

Note that the absolute error is symmetric in  $a$  and  $b$ , so that  $\epsilon$  is also the absolute error in  $a$  as an approximation to  $b$

**Definition B.2 (Relative error)** *If  $a \neq 0$  then the relative error in  $b$  as an approximation to  $a$  is defined by*

$$\rho = \rho_b := \frac{|b - a|}{|a|}.$$

*We say that  $a$  and  $b$  agree to approximately  $-\log_{10} \rho$  digits.*

As an example, if  $a := 31415.9265$  and  $b := 31415.8951$ , then  $\rho = 0.999493 * 10^{-6}$  and  $a$  and  $b$  agree to approximately 6 digits.

We have  $b = a(1 + r)$  for some  $r$  if and only if  $\rho = |r|$ .

We can also consider the relative error  $\rho_a := |a - b|/|a|$  in  $a$  as an approximation to  $b$ .

**Lemma B.3 (Relative errors)**

*If  $a, b \neq 0$  and  $\rho_b < 1$  then  $\rho_a \leq \rho_b / (1 - \rho_b)$ .*

**Proof.** Since  $|a|\rho_b = |b - a| \geq |a| - |b|$  we obtain  $|b| \geq |a| - |a - b| = (1 - \rho_b)|a|$ . Then

$$\rho_a = \frac{|b - a|}{|b|} \leq \frac{|b - a|}{(1 - \rho_b)|a|} = \frac{\rho_b}{1 - \rho_b}.$$

□

If  $\rho_b$  is small then  $\rho_a$  is small and it does not matter whether we choose  $\rho_a$  or  $\rho_b$  to discuss relative error.

## B.2 Floating Point Numbers

We shall assume that the reader is familiar with different number systems (binary, octal, decimal, hexadecimal) and how to convert from one number system to another. We use  $(x)_\beta$  to indicate a number written to the base  $\beta$ . If no parenthesis and subscript are used, the base 10 is understood. For instance,

$$\begin{aligned}(100)_2 &= 4, \\ (0.1)_2 &= 0.5, \\ 0.1 &= (0.1)_{10} = (0.0001100110011001\dots)_2.\end{aligned}$$

In general,

$$x = (c_m c_{m-1} \dots c_0 . d_1 d_2 \dots d_n)_\beta$$

means

$$x = \sum_{i=0}^m c_i \beta^i + \sum_{i=1}^n d_i \beta^{-i}, \quad 0 \leq c_i, d_i \leq \beta - 1.$$

We can move the decimal point by adding an exponent:

$$y = x \cdot \beta^e,$$

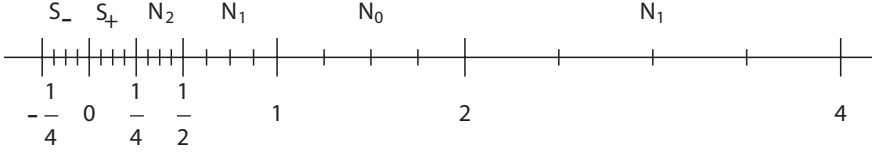
for example

$$(0.1)_{10} = (1.100110011001\dots)_2 \cdot 2^{-4}.$$

We turn now to a description of the floating-point numbers. We will only describe a **standard system**, namely the binary IEEE floating-point standard. Although it is not used by all systems, it has been widely adopted and is used in MATLAB. For a more complete introduction to the subject see [12],[26].

We denote the real numbers which are represented in our computer by  $\mathcal{F}$ . The set  $\mathcal{F}$  are characterized by three integers  $t$ , and  $\underline{e}, \bar{e}$ . We define

$$\epsilon_M := 2^{-t}, \quad \text{machine epsilon,} \quad (\text{B.1})$$



**Figure B.1.** *Distribution of some positive floating-point numbers*

and

$$\begin{aligned}
 \mathcal{F} &:= \{0\} \cup \mathcal{S} \cup \mathcal{N}, \text{ where} \\
 \mathcal{N} &:= \mathcal{N}_+ \cup \mathcal{N}_-, \quad \mathcal{N}_+ := \cup_{e=\underline{e}}^{\bar{e}} \mathcal{N}_e, \quad \mathcal{N}_- := -\mathcal{N}_+, \\
 \mathcal{N}_e &:= \{(1.d_1d_2 \cdots d_t)_2\} * 2^e = \{1, 1 + \epsilon_M, 1 + 2\epsilon_M, \dots, 2 - \epsilon_M\} * 2^e, \\
 \mathcal{S} &:= \mathcal{S}_+ \cup \mathcal{S}_-, \quad \mathcal{S}_+ := \{\epsilon_M, 2\epsilon_M, 3\epsilon_M, \dots, 1 - \epsilon_M\} * 2^{\underline{e}}, \quad \mathcal{S}_- := -\mathcal{S}_+.
 \end{aligned}
 \tag{B.2}$$

**Example B.4 (Floating numbers)**

Suppose  $t := 2$ ,  $\bar{e} = 3$  and  $\underline{e} := -2$ . Then  $\epsilon_M = 1/4$  and we find

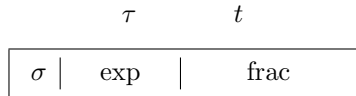
$$\begin{aligned}
 \mathcal{N}_{-2} &= \left\{ \frac{1}{4}, \frac{5}{16}, \frac{3}{8}, \frac{7}{16} \right\}, \quad \mathcal{N}_{-1} = \left\{ \frac{1}{2}, \frac{5}{8}, \frac{3}{4}, \frac{7}{8} \right\}, \quad \mathcal{N}_0 = \left\{ 1, \frac{5}{4}, \frac{3}{2}, \frac{7}{4} \right\}, \\
 \mathcal{N}_1 &= \left\{ 2, \frac{5}{2}, 3, \frac{7}{2} \right\}, \quad \mathcal{N}_2 = \{4, 5, 6, 7\}, \quad \mathcal{N}_3 = \{8, 10, 12, 14\}, \\
 \mathcal{S}_+ &= \left\{ \frac{1}{16}, \frac{1}{8}, \frac{3}{16} \right\}, \quad \mathcal{S}_- = \left\{ -\frac{3}{16}, -\frac{1}{8}, -\frac{1}{16} \right\}.
 \end{aligned}$$

The position of some of these sets on the real line is shown in Figure B.1

1. The elements of  $\mathcal{N}$  are called **normalized (floating-point) numbers**. They consists of three parts, the sign +1 or -1, the **mantissa**  $(1.d_1d_2 \cdots d_t)_2$ , and the **exponent part**  $2^e$ .
2. the elements in  $\mathcal{N}_+$  has the sign +1 indicated by the bit  $\sigma = 0$  and the elements in  $\mathcal{N}_-$  has the sign bit  $\sigma = 1$ . Thus the sign of a number is  $(-1)^\sigma$ . The standard system has two zeros +0 and -0.
3. The mantissa is a number between 1 and 2. It consists of  $t + 1$  binary digits.
4. The number  $e$  in the exponent part is restricted to the range  $\underline{e} \leq e \leq \bar{e}$ .
5. The positive normalized numbers are located in the interval  $[r_m, r_M]$ , where

$$r_m := 2^{\underline{e}}, \quad r_M := (2 - \epsilon_M) * 2^{\bar{e}}.
 \tag{B.3}$$

6. The elements in  $\mathcal{S}$  are called **subnormal** or **denormalized**. As for normalized numbers they consists of three parts, but the mantissa is less than one in size. The main use of subnormal numbers is to soften the effect of underflow. If a number is in the range  $(0, (1 - \epsilon_M/2) * 2^e)$ , then it is rounded to the nearest subnormal number or to zero.
7. Two additional symbols "Inf" and "NaN" are used for special purposes.
8. The symbol **Inf** is used to represent numbers outside the interval  $[-r_M, r_M]$  (**overflow**), and results of arithmetic operations of the form  $x/0$ , where  $x \in \mathcal{N}$ . Inf has a sign, +Inf and -Inf.
9. The symbol **NaN** stands for "not a number". a NaN results from illegal operations of the form  $0/0, 0 * \text{Inf}, \text{Inf}/\text{Inf}, \text{Inf} - \text{Inf}$  and so on.
10. The choices of  $t, \bar{e}$ , and  $\underline{e}$  are to some extent determined by the architecture of the computer. A floating-point number, say  $x$ , occupies  $n := 1 + \tau + t$  bits, where 1 bit is used for the sign,  $\tau$  bits for the exponent, and  $t$  bits for the fractional part of the mantissa.



Here  $\sigma = 0$  if  $x > 0$  and  $\sigma = 1$  if  $x < 0$ , and  $\text{exp} \in \{0, 1, 2, 3, \dots, 2^\tau - 1\}$  is an integer. The integer  $\text{frac}$  is the fractional part  $d_1 d_2 \dots d_t$  of the mantissa. The value of a normalized number in the standard system is

$$x = (-1)^\sigma * (1.\text{frac})_2 * 2^{\text{exp}-b}, \text{ where } b := 2^{\tau-1} - 1. \tag{B.4}$$

The integer  $b$  is called the **bias**.

11. To explain the choice of  $b$  we note that the extreme values  $\text{exp} = 0$  and  $\text{exp} = 2^\tau - 1$  are used for special purposes. The value  $\text{exp} = 0$  is used for the number zero and the subnormal numbers, while  $\text{exp} = 2^\tau - 1$  is used for Inf and NaN. Since  $2b = 2^\tau - 2$ , the remaining numbers of  $\text{exp}$ , i. e.,  $\text{exp} \in \{1, 2, \dots, 2^\tau - 2\}$  correspond to  $e$  in the set  $\{1 - b, 2 - b, \dots, b\}$ . Thus in a standard system we have

$$\underline{e} = 1 - b, \quad \bar{e} = b := 2^{\tau-1} - 1. \tag{B.5}$$

12. The most common choices of  $\tau$  and  $t$  are shown in the following table

| precision | $\tau$ | $t$ | $b$   | $\epsilon_M = 2^{-t}$ | $r_m = 2^{1-b}$         | $r_M$                  |
|-----------|--------|-----|-------|-----------------------|-------------------------|------------------------|
| half      | 5      | 10  | 15    | $9.8 \times 10^{-4}$  | $6.1 \times 10^{-5}$    | $6.6 \times 10^4$      |
| single    | 8      | 23  | 127   | $1.2 \times 10^{-7}$  | $1.2 \times 10^{-38}$   | $3.4 \times 10^{38}$   |
| double    | 11     | 52  | 1023  | $2.2 \times 10^{-16}$ | $2.2 \times 10^{-308}$  | $1.8 \times 10^{308}$  |
| quad      | 15     | 112 | 16383 | $1.9 \times 10^{-34}$ | $3.4 \times 10^{-4932}$ | $1.2 \times 10^{4932}$ |

Here  $b$  is given by (B.5) and  $r_M$  by (B.3). The various lines correspond to a normalized number occupying **half** a word of 32 bits, one word (**single precision**), two words (**double precision**), and 4 words (**quad precision**).

## B.3 Rounding and Arithmetic Operations

The standard system is a closed system. Every  $x \in \mathbb{R}$  has a representation as either a floating-point number, or Inf or NaN, and every arithmetic operation produces a result. We denote the computer representation of a real number  $x$  by  $\text{fl}(x)$ .

### B.3.1 Rounding

To represent a real number  $x$  there are three cases.

$$\text{fl}(x) = \begin{cases} \text{Inf}, & \text{if } x > r_M, \\ -\text{Inf}, & \text{if } x < -r_M, \\ \text{round to zero}, & \text{otherwise.} \end{cases}$$

To represent a real number with  $|x| \leq r_M$  the system chooses a machine number  $\text{fl}(x)$  closest to  $x$ . This is known as **rounding**. When  $x$  is midway between two numbers in  $\mathcal{F}$  we can either choose the one of larger magnitude (**round away from zero**), or pick the one with a zero last bit (**round to zero**). The standard system uses round to zero. As an example, if  $x = 1 + \epsilon_M/2$ , then  $x$  is midway between 1 and  $1 + \epsilon_M$ . Therefore  $\text{fl}(x) = 1 + \epsilon_M$  if round away from zero is used, while  $\text{fl}(x) = 1$  if  $x$  is rounded to zero. This is because the machine representation of 1 has  $\text{frac} = 0$ .

The following lemma gives a bound for the relative error in rounding.

#### Theorem B.5 (Relative error in rounding)

If  $r_m \leq |x| \leq r_M$  then

$$\text{fl}(x) = x(1 + \delta), \quad |\delta| \leq u_M := \frac{1}{2}\epsilon_M = 2^{-t-1}.$$

**Proof.** Suppose  $2^e < x < 2^{e+1}$ . Then  $\text{fl}(x) \in \{1, 1 + \epsilon_M, 1 + 2\epsilon_M, \dots, 2 - \epsilon_M\} * 2^e$ . These numbers are uniformly spaced with spacing  $\epsilon_M * 2^e$  and therefore  $|\text{fl}(x) - x| \leq \frac{1}{2}\epsilon_M 2^e \leq \frac{1}{2}\epsilon_M * |x|$ . The proof for a negative  $x$  is similar.  $\square$

The number  $u_M$  is called the **rounding unit**.

### B.3.2 Arithmetic operations

Suppose  $x, y \in \mathcal{N}$ . In a standard system we have

$$\text{fl}(x \circ y) = (x \circ y)(1 + \delta), \quad |\delta| \leq u_M, \quad \circ \in \{+, -, *, /, \sqrt{\cdot}\}, \quad (\text{B.6})$$

where  $u_M$  is the rounding unit of the system. This means that the computed value is as good as the rounded exact answer. This is usually achieved by using one or several extra digits known as **guard digits** in the calculation.

## B.4 Backward Rounding-Error Analysis

The computed sum of two numbers  $\alpha_1, \alpha_2 \in \mathcal{N}$  satisfy  $\text{fl}(\alpha_1 \circ \alpha_2) = (\alpha_1 + \alpha_2)(1 + \delta)$ , where  $|\delta| \leq u_M$ , the rounding unit. If we write this as  $\text{fl}(\alpha_1 \circ \alpha_2) = \tilde{\alpha}_1 + \tilde{\alpha}_2$ , where  $\tilde{\alpha}_i := \alpha_i(1 + \delta)$  for  $i = 1, 2$ , we see that the computed sum is the exact sum of two numbers which approximate the exact summands with small relative error,  $|\delta| \leq u_M$ . The error in the addition has been boomeranged back on the data  $\alpha_1, \alpha_2$ , and in this context we call  $\delta$  the **backward error**. A similar interpretation is valid for the other arithmetic operations  $-, *, /, \sqrt{\cdot}$ , and we assume it also holds for the elementary functions  $\sin, \cos, \exp, \log$  and so on.

Suppose more generally we want to compute the value of an expression  $\phi(\alpha_1, \dots, \alpha_n)$ . Here  $\alpha_1, \dots, \alpha_n \in \mathcal{N}$  are given data, and we are using the arithmetic operations, and implementations of the standard elementary functions, in the computation. A **backward error analysis** consists of showing that the computed result is obtained as the exact result of using data  $\beta := [\beta_1, \dots, \beta_n]^T$  instead of  $\alpha := [\alpha_1, \dots, \alpha_n]$ . In symbols

$$\tilde{\phi}(\alpha_1, \dots, \alpha_n) = \phi(\beta_1, \dots, \beta_n).$$

If we can show that the relative error in  $\beta$  as an approximation to  $\alpha$  is  $O(u_M)$  either componentwise or norm-wise in some norm, then we say that the algorithm to compute  $\phi(\alpha_1, \dots, \alpha_n)$  is **backward stable**. Normally the constant  $K$  in the  $O(u_M)$  term will grow with  $n$ . Typically  $K = p(n)$  for some polynomial  $p$  is acceptable, while an exponential growth of  $K$  can be problematic.

### B.4.1 Computing a sum

We illustrate this discussion by computing the backward error in the sum of  $n$  numbers  $s := \alpha_1 + \dots + \alpha_n$ , where  $\alpha_i \in \mathcal{N}$  for all  $i$ . We have the following



algorithm.

```

s1 := α1
for k = 2 : n
    sk := fl(sk-1 + αk)
end
s̃ := sn

```

Using a standard system we obtain for  $n = 3$

$$\begin{aligned}
 s_2 &= \text{fl}(\alpha_1 + \alpha_2) = \alpha_1(1 + \delta_2) + \alpha_2(1 + \delta_2), \\
 s_3 &= \text{fl}(s_2 + \alpha_3) = s_2(1 + \delta_3) + \alpha_3(1 + \delta_3) = \alpha_1(1 + \eta_1) + \alpha_2(1 + \eta_2) + \alpha_3(1 + \eta_3), \\
 \eta_1 &= \eta_2 = (1 + \delta_2)(1 + \delta_3), \quad \eta_3 = (1 + \delta_3), \quad |\delta_i| \leq u_M.
 \end{aligned}$$

In general, with  $\delta_1 := 0$ ,

$$\tilde{s} = \sum_{i=1}^n \alpha_i(1 + \eta_i). \quad \eta_i = (1 + \delta_i) \dots (1 + \delta_n), \quad |\delta_i| \leq u_M, \quad i = 1, \dots, n. \quad (\text{B.7})$$

With  $\phi(\alpha_1, \dots, \alpha_n) := \alpha_1 + \dots + \alpha_n$  this shows that

$$\tilde{s} = \tilde{\phi}(\alpha_1, \dots, \alpha_n) = \phi(\beta_1, \dots, \beta_n), \quad \beta_i = \alpha_i(1 + \eta_i). \quad (\text{B.8})$$

The following lemma gives a convenient bound on the  $\eta$  factors.

**Lemma B.6 (Bound on factors)**

Suppose for integers  $k, m$  with  $0 \leq m \leq k$  and  $k \geq 1$  that

$$1 + \eta_k := \frac{(1 + \delta_1) \dots (1 + \delta_m)}{(1 + \delta_{m+1}) \dots (1 + \delta_k)}, \quad |\delta_j| \leq u_M, \quad j = 1, \dots, k.$$

If  $ku_M \leq \frac{1}{11}$  then

$$|\eta_k| \leq ku'_M, \quad \text{where } u'_M := 1.1u_M. \quad (\text{B.9})$$

**Proof.** We first show that

$$ku_M \leq \alpha < 1 \implies |\eta_k| \leq k \frac{u_M}{1 - \alpha}. \quad (\text{B.10})$$

For convenience we use  $u := u_M$  in the proof. Since  $u < 1$  we have  $1/(1 - u) = 1 + u + u^2/(1 - u) > 1 + u$  and we obtain

$$(1 - u)^k \leq \frac{(1 - u)^m}{(1 + u)^{k-m}} \leq 1 + \eta_k \leq \frac{(1 + u)^m}{(1 - u)^{k-m}} \leq (1 - u)^{-k}.$$

The proof of (B.10) will be complete if we can show that

$$1 - ku \leq (1 - u)^k, \quad (1 - u)^{-k} \leq 1 + ku'.$$

The first inequality is an easy induction on  $k$ . If it holds for  $k$ , then

$$(1 - u)^{k+1} = (1 - u)^k(1 - u) \geq (1 - ku)(1 - u) = 1 - (k + 1)u + ku^2 \geq 1 - (k + 1)u.$$

The second inequality is a consequence of the first,

$$(1 - u)^{-k} \leq (1 - ku)^{-1} = 1 + \frac{ku}{1 - ku} \leq 1 + \frac{ku}{1 - \alpha} = 1 + ku'.$$

Letting  $\alpha = \frac{1}{11}$  in (B.10) we obtain (B.9).  $\square$

The number  $u'_M := 1.1u_M$ , corresponding to  $\alpha = 1/11$ , is called the **ad-justed rounding unit**. In the literature many values of  $\alpha$  can be found. [26] uses  $\alpha = 1/10$  giving  $u'_M = 1.12u_M$ , while in [12] the value  $\alpha = 0.01$  can be found. In the classical work [34] one finds  $1/(1 - \alpha) = 1.06$ .

Let us return to the backward error (B.8) in a sum of  $n$  numbers. Since  $\delta_1 = 0$  we see that

$$|\eta_1| \leq (n - 1)u'_M, \quad |\eta_i| \leq (n - i + 1)u'_M, \quad \text{for } i = 2, \dots, n.$$

or more simply

$$|\eta_i| \leq (n - 1)u'_M, \quad \text{for } i = 1, \dots, n. \quad (\text{B.11})$$

This shows that the algorithm for computing a sum is backward stable.

The bounds from a backward rounding-error analysis can be used together with a condition number to bound the actual error in the computed result. To see this for the sum, we subtract the exact sum  $s = \alpha_1 + \dots + \alpha_n$  from the computed sum  $\tilde{s} = \alpha_1(1 + \eta_1) + \dots + \alpha_n(1 + \eta_n)$ , to get

$$|\tilde{s} - s| = |\alpha_1\eta_1 + \dots + \alpha_n\eta_n| \leq (|\alpha_1| + \dots + |\alpha_n|)(n - 1)u'_M.$$

Thus the relative error in the computed sum of  $n$  numbers is bounded as follows

$$\left| \frac{\tilde{s} - s}{s} \right| \leq \kappa(n - 1)u'_M, \quad \text{where } \kappa := \frac{|\alpha_1| + \dots + |\alpha_n|}{\alpha_1 + \dots + \alpha_n}. \quad (\text{B.12})$$

This bound shows that the backward error can be magnified by at most  $\kappa$ . The number  $\kappa$  is called the **condition number** for the sum.

The condition number measures how much a relative error in each of the components in a sum can be magnified in the final sum. The backward error shows how large these relative perturbations can be in the actual algorithm we used to compute the sum. Using backward error analysis and condition number

separates the process of estimating the error in the final result into two distinct jobs.

A problem where small relative changes in the data leads to large relative changes in the exact result is called **ill conditioned**. We see that computing a sum can be ill-conditioned if the exact value of the sum is close to zero and some of the individual terms have large absolute values with opposite signs.

### B.4.2 Computing an inner product

Computing an inner product  $p := \alpha_1\gamma_1 + \dots + \alpha_n\gamma_n$  is also backward stable using the standard algorithm

```

 $p_1 := \text{fl}(\alpha_1\gamma_1)$ 
for  $k = 2 : n$ 
     $p_k := \text{fl}(p_{k-1} + \text{fl}(\alpha_k\gamma_k))$ 
end
 $\tilde{p} := p_n$ 

```

For a backward error analysis of this algorithm we only need to modify (B.7) slightly. All we have to do is to add terms  $\text{fl}(\alpha_k\gamma_k) = \alpha_k\gamma_k(1 + \pi_k)$  to the terms of the sum. The result is

$$\tilde{p} = \sum_{k=1}^n \alpha_k\gamma_k(1 + \eta_k), \quad \eta_k = (1 + \pi_k)(1 + \delta_k) \cdots (1 + \delta_n), \quad k = 1, \dots, n,$$

where  $\delta_1 = 0$ . Thus for the inner product of  $n$  terms we obtain

$$\left| \frac{\tilde{p} - p}{p} \right| \leq \kappa n u_M, \quad \kappa := \frac{|\alpha_1\gamma_1| + \dots + |\alpha_n\gamma_n|}{|\alpha_1\gamma_1 + \dots + \alpha_n\gamma_n|}. \quad (\text{B.13})$$

The computation can be ill conditioned if the exact value is close to zero and some of the components are large in absolute value.

### B.4.3 Computing a matrix product

Using matrix norms we can bound the backward error in matrix algorithms. Suppose we want to compute the matrix product  $\mathbf{C} = \mathbf{A} * \mathbf{B}$ . Let  $n$  be the number of columns of  $\mathbf{A}$  and the number of rows of  $\mathbf{B}$ . Each element in  $\mathbf{C}$  is the inner product of a row of  $\mathbf{A}$  and a column of  $\mathbf{B}$ . Thus if  $\tilde{\mathbf{C}}$  is the computed product then from (B.13)

$$\left| \frac{\tilde{c}_{ij} - c_{ij}}{c_{ij}} \right| \leq \kappa_{ij} n u'_M, \quad \kappa_{ij} := \frac{|a_1 b_1| + \dots + |a_n b_n|}{|a_1 b_1 + \dots + a_n b_n|}, \quad \text{all } i, j. \quad (\text{B.14})$$

We write this as  $|\tilde{c}_{ij} - c_{ij}| \leq \kappa_{ij}|c_{ij}|nu'_M$ . Using the infinity matrix norm we find

$$\sum_j |\tilde{c}_{ij} - c_{ij}| \leq nu'_M \sum_j \kappa_{ij}|c_{ij}| \leq \kappa nu'_M \sum_j |c_{ij}| \leq \kappa nu'_M \|\mathbf{C}\|_\infty, \text{ all } i,$$

where  $\kappa := \max_{ij} \kappa_{ij}$ . Maximizing over  $i$  we obtain

$$\frac{\|\tilde{\mathbf{C}} - \mathbf{C}\|_\infty}{\|\mathbf{C}\|_\infty} \leq \kappa nu'_M. \tag{B.15}$$

The calculation of a matrix product can be ill conditioned if one or more of the product elements are small and the corresponding inner products have large terms of opposite signs.

## Appendix C

# Differentiation of Vector Functions

For any sufficiently differentiable  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  we recall that the partial derivative with respect to the  $i$ th variable of  $f$  is defined by

$$D_i f(\mathbf{x}) := \frac{\partial f(\mathbf{x})}{\partial x_i} := \lim_{h \rightarrow \mathbf{0}} \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h}, \quad \mathbf{x} \in \mathbb{R}^n,$$

where  $\mathbf{e}_i$  is the  $i$ th unit vector in  $\mathbb{R}^n$ . For each  $\mathbf{x} \in \mathbb{R}^n$  we define the **gradient**  $\nabla f(\mathbf{x}) \in \mathbb{R}^n$ , and the **hessian**  $\mathbf{H}f = \nabla \nabla^T f(\mathbf{x}) \in \mathbb{R}^{n,n}$  of  $f$  by

$$\nabla f := \begin{bmatrix} D_1 f \\ \vdots \\ D_n f \end{bmatrix}, \quad \mathbf{H}f := \nabla \nabla^T f := \begin{bmatrix} D_1 D_1 f & \cdots & D_1 D_n f \\ \vdots & & \vdots \\ D_n D_1 & \cdots & D_n D_n f \end{bmatrix}, \quad (\text{C.1})$$

where  $\nabla^T f := (\nabla f)^T$  is the row vector gradient. The operators  $\nabla \nabla^T$  and  $\nabla^T \nabla$  are quite different. Indeed,  $\nabla^T \nabla f = D_1^2 f + \cdots + D_n^2 f =: \nabla^2$  the **Laplacian** of  $f$ , while  $\nabla \nabla^T$  can be thought of as an outer product resulting in a matrix.

### Lemma C.1 (Product rules)

For  $f, g : \mathbb{R}^n \rightarrow \mathbb{R}$  we have the product rules

1.  $\nabla(fg) = f\nabla g + g\nabla f, \quad \nabla^T(fg) = f\nabla^T g + g\nabla^T f,$
2.  $\nabla \nabla^T(fg) = \nabla f \nabla^T g + \nabla g \nabla^T f + f \nabla \nabla^T g + g \nabla \nabla^T f.$
3.  $\nabla^2(fg) = 2\nabla^T f \nabla g + f \nabla^2 g + g \nabla^2 f.$

We define the **Jacobian** of a vector function  $\mathbf{f} = [f_1, \dots, f_m]^T : \mathbb{R}^n \rightarrow \mathbb{R}^m$  as the  $m, n$  matrix

$$\nabla^T \mathbf{f} := \begin{bmatrix} D_1 f_1 & \cdots & D_n f_1 \\ \vdots & & \vdots \\ D_1 f_m & \cdots & D_n f_m \end{bmatrix}.$$

As an example, if  $f(\mathbf{x}) = f(x, y) = x^2 - xy + y^2$  and  $\mathbf{g}(x, y) := [f(x, y), x - y]^T$  then

$$\begin{aligned} \nabla f(x, y) &= \begin{bmatrix} 2x - y \\ -x + 2y \end{bmatrix}, & \nabla^T \mathbf{g}(x, y) &= \begin{bmatrix} 2x - y & -x + 2y \\ 1 & -1 \end{bmatrix}, \\ \mathbf{H}f(x, y) &= \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}. \end{aligned}$$

The second order Taylor expansion in  $n$  variables can be expressed in terms of the gradient and the hessian.

**Lemma C.2 (Second order Taylor expansion)**

Suppose  $f \in C^2(\Omega)$ , where  $\Omega \in \mathbb{R}^n$  contains two points  $\mathbf{x}, \mathbf{x} + \mathbf{h} \in \Omega$ , such that the line segment  $L := \{\mathbf{x} + t\mathbf{h} : t \in (0, 1)\} \subset \Omega$ . Then

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \mathbf{h}^T \nabla f(\mathbf{x}) + \frac{1}{2} \mathbf{h}^T \nabla \nabla^T f(\mathbf{c}) \mathbf{h}, \text{ for some } \mathbf{c} \in L. \quad (\text{C.2})$$

**Proof.** Let  $g : [0, 1] \rightarrow \mathbb{R}$  be defined by  $g(t) := f(\mathbf{x} + t\mathbf{h})$ . Then  $g \in C^2[0, 1]$  and by the chain rule

$$\begin{aligned} g(0) &= f(\mathbf{x}) & g(1) &= f(\mathbf{x} + \mathbf{h}), \\ g'(t) &= \sum_{i=1}^n h_i \frac{\partial f(\mathbf{x} + t\mathbf{h})}{\partial x_i} = \mathbf{h}^T \nabla f(\mathbf{x} + t\mathbf{h}), \\ g''(t) &= \sum_{i=1}^n \sum_{j=1}^n h_i h_j \frac{\partial^2 f(\mathbf{x} + t\mathbf{h})}{\partial x_i \partial x_j} = \mathbf{h}^T \nabla \nabla^T f(\mathbf{x} + t\mathbf{h}) \mathbf{h}. \end{aligned}$$

Inserting these expressions in the second order Taylor expansion

$$g(1) = g(0) + g'(0) + \frac{1}{2} g''(u), \text{ for some } u \in (0, 1),$$

we obtain (C.2) with  $\mathbf{c} = \mathbf{x} + u\mathbf{h}$ .  $\square$

The gradient and hessian of some functions involving matrices can be found from the following lemma.

**Lemma C.3 (Functions involving matrices)**

For any  $m, n \in \mathbb{N}$ ,  $\mathbf{B} \in \mathbb{R}^{n,n}$ ,  $\mathbf{C} \in \mathbb{R}^{m,n}$ , and  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{y} \in \mathbb{R}^m$  we have

1.  $\nabla(\mathbf{y}^T \mathbf{C}) = \nabla^T(\mathbf{C}\mathbf{x}) = \mathbf{C}$ ,
2.  $\nabla(\mathbf{x}^T \mathbf{B}\mathbf{x}) = (\mathbf{B} + \mathbf{B}^T)\mathbf{x}$ ,  $\nabla^T(\mathbf{x}^T \mathbf{B}\mathbf{x}) = \mathbf{x}^T(\mathbf{B} + \mathbf{B}^T)$ ,
3.  $\nabla \nabla^T(\mathbf{x}^T \mathbf{B}\mathbf{x}) = \mathbf{B} + \mathbf{B}^T$ .

**Proof.**

1. We find  $D_i(\mathbf{y}^T \mathbf{C}) = \lim_{h \rightarrow 0} \frac{1}{h} ((\mathbf{y} + h\mathbf{e}_i)^T \mathbf{C} - \mathbf{y}^T \mathbf{C}) = \mathbf{e}_i^T \mathbf{C}$  and  $D_i(\mathbf{C}\mathbf{x}) = \lim_{h \rightarrow 0} \frac{1}{h} (\mathbf{C}(\mathbf{x} + h\mathbf{e}_i) - \mathbf{C}\mathbf{x}) = \mathbf{C}\mathbf{e}_i$  and 1. follows.

2. Here we find

$$\begin{aligned} D_i(\mathbf{x}^T \mathbf{B}\mathbf{x}) &= \lim_{h \rightarrow 0} \frac{1}{h} ((\mathbf{x} + h\mathbf{e}_i)^T \mathbf{B}(\mathbf{x} + h\mathbf{e}_i) - \mathbf{x}^T \mathbf{B}\mathbf{x}) \\ &= \lim_{h \rightarrow 0} (\mathbf{e}_i^T \mathbf{B}\mathbf{x} + \mathbf{x}^T \mathbf{B}\mathbf{e}_i + h\mathbf{e}_i^T \mathbf{e}_i) = \mathbf{e}_i^T (\mathbf{B} + \mathbf{B}^T)\mathbf{x}, \end{aligned}$$

and the first part of 2. follows. Taking transpose we obtain the second part.

3. Combining 1. and 2. we obtain 3.

□





# Bibliography

- [1] Beckenbach, E. F, and R. Bellman, *Inequalities*, Springer Verlag, Berlin, Fourth Printing, 1983.
- [2] Björck, Åke, *Numerical Methods for Least Squares Problems*, SIAM, Philadelphia, 1995.
- [3] E. Cohen, R. F. Riesenfeld, G. Elber, *Geometric Modeling with Splines: An Introduction*, A.K. Peters, Ltd., 2001,
- [4] Demmel, J. W., *Applied Numerical Linear Algebra*, SIAM, Philadelphia, 1997.
- [5] Golub, G. H., and C. F. Van Loan, *Matrix Computations*, John Hopkins University Press, Baltimore, MD, third edition, 1996.
- [6] Grcar, Joseph F., Mathematicians of Gaussian elimination, *Notices of the AMS*, **58** (2011), 782–792.
- [7] Greenbaum, Anne, *Iterative Methods for Solving Linear Systems*, SIAM, Philadelphia, 1997.
- [8] Hackbush, Wolfgang, *Iterative Solution of Large Sparse Systems of Equations*, Springer-Verlag, Berlin, 1994.
- [9] Hall, C. A. and W. W. Meyer, Optimal error bounds for cubic spline interpolation. *J. Approx. Theory*, **16** (1976), 105122.
- [10] Hestenes, Magnus, *Conjugate Direction Methods in Optimization*, Springer-Verlag, Berlin, 1980.
- [11] Hestenes, M. and E. Stiefel, *Methods of conjugate gradients for solving linear systems*, *Journal of Research of the National Bureau of Standards* **29**(1952), 409–439.

- 
- [12] Higham, N. J., *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, 1996.
- [13] Horn, Roger A. and Charles R. Johnson, *Matrix Analysis*, Cambridge University Press, Cambridge, UK, 1985.
- [14] Horn, Roger A. and Charles R. Johnson, *Topics in Matrix Analysis*, Cambridge University Press, Cambridge, UK, 1991.
- [15] Ipsen, I.C.F., *Numerical Matrix Analysis: Linear Systems and Least Squares*, SIAM, Philadelphia, 2009.
- [16] Kato, *Perturbation Theory for Linear Operators*, Springer.
- [17] Lancaster, P., and Rodman, L., "Canonical forms for hermitian matrix pairs under strict equivalence and congruence", *SIAM Review*, vol. 47, 2005, 407-443.
- [18] Laub, A. J., *Matrix Analysis for Scientists and Engineers*, SIAM Philadelphia, 2005.
- [19] Laub, A. J., *Computational Matrix Analysis*, SIAM Philadelphia, 2012.
- [20] Lawson, C.L. and R. J. Hanson, *Solving Least Squares Problems*, Prentice-Hall, Englewood Cliffs, N.J, 1974.
- [21] Lax, Peter D., *Linear Algebra*, John Wiley & Sons, New York, 1997.
- [22] Lay, D.C: *Linear algebra and its applications*, 2012. Addison Wesley / Pearson. Fourth edition.
- [23] Leon, Steven J., *Linear Algebra with Applications*, Prentice Hall, NJ, Seventh Edition, 2006.
- [24] Meyer, Carl D., *Matrix Analysis and Applied Linear Algebra*, Siam Philadelphia, 2000.
- [25] Steel, J. Michael, *The Cauchy-Schwarz Master Class*, Cambridge University Press, Cambridge, UK, 2004.
- [26] Stewart, G. W., *Matrix Algorithms Volume I: Basic Decompositions*, Siam Philadelphia, 1998.
- [27] Stewart, G. W., *Matrix Algorithms Volume II: Eigensystems*, Siam Philadelphia, 2001.
- [28] Stewart, G. W. and Ji-guang Sun, *Matrix Perturbation Theory*, Academic Press, San Diego, 1990.

- 
- [29] Stewart, G. W., *Introduction to Matrix Computations*, Academic press, New York, 1973.
  - [30] Trefethen, Lloyd N., and David Bau III, *Numerical Linear Algebra*, Siam Philadelphia, 1997.
  - [31] Tveito, A., and R. Winther, *Partial Differential Equations*, Springer, Berlin.
  - [32] Van Loan, Charles, *Computational Frameworks for the Fast Fourier Transform*, Siam Philadelphia, 1992.
  - [33] Varga, R. S., *Matrix Iterative Analysis/ 2nd Edn.*, Springer Verlag, New York, 2000.
  - [34] Wilkinson, J. H., *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford, 1965.
  - [35] Young, D. M., *Iterative Solution of Large Linear Systems*, Academic Press, New York, 1971.
  - [36] Zhang, F., *Matrix Theory*, Springer, New York, 1999.



# List of Figures

|      |  |     |
|------|--|-----|
| 1    | The construction of $\mathbf{v}_1$ and $\mathbf{v}_2$ in Gram-Schmidt. The constant $c$ is given by $c := \langle \mathbf{s}_2, \mathbf{v}_1 \rangle / \langle \mathbf{v}_1, \mathbf{v}_1 \rangle$ . . . . . | 19  |
| 2    | The triangle $T$ defined by the three points $P_1, P_2$ and $P_3$ . . . . .  | 27  |
| 2.1  | Lower triangular $5 \times 5$ band matrices: $d = 1$ (left) and $d = 2$ (right). . . . .   | 52  |
| 2.2  | Gaussian elimination . . . . .   | 77  |
| 3.1  | Numbering of grid points . . . . .   | 91  |
| 3.2  | The 5-point stencil . . . . .  | 92  |
| 3.3  | Band structure of the 2D test matrix, $n = 9, n = 25, n = 100$ . . . . .   | 93  |
| 4.1  | Fill-in in the Cholesky factor of the Poisson matrix ( $n = 100$ ). . . . .  | 104 |
| 6.1  | The ellipse $y_1^2/9 + y_2^2 = 1$ (left) and the rotated ellipse $\mathbf{A}\mathbf{S}$ (right). . . . .   | 158 |
| 7.1  | A convex function. . . . .   | 180 |
| 8.1  | $\rho(\mathbf{G}_\omega)$ with $\omega \in [0, 2]$ for $n = 100$ , (lower curve) and $n = 2500$ (upper curve). . . . .   | 203 |
| 9.1  | Level curves for $Q(x, y)$ given by (9.2). Also shown is a steepest descent iteration (left) and a conjugate gradient iteration (right) to find the minimum of $Q$ . (cf Examples 9.3,9.6) . . . . .         | 216 |
| 9.2  | The orthogonal projection of $\mathbf{x} - \mathbf{x}_0$ into $\mathbb{W}_k$ . . . . .   | 231 |
| 9.3  | This is an illustration of the proof of Theorem 9.27 for $k = 3$ . $f \equiv Q - Q^*$ has a double zero at $\mu_1$ and one zero between $\mu_2$ and $\mu_3$ . . . . .  | 235 |
| 10.1 | The Householder transformation in Exercise 10.2 . . . . .  | 250 |
| 10.2 | A plane rotation. . . . .  | 261 |

---

|      |  |     |
|------|--|-----|
| 11.1 | A least squares fit to data. . . . .                                   | 267 |
| 11.2 | The orthogonal projection of $\mathbf{b}$ into $\mathcal{S}$ . . . . . | 271 |
| 11.3 | Graphical interpretation of the bounds in Theorem 11.31. . . . .       | 281 |
| 12.1 | The Gerschgorin disk $R_i$ . . . . .                                   | 291 |
| 13.1 | Post multiplication in a QR step. . . . .                              | 316 |
| B.1  | Distribution of some positive floating-point numbers . . . . .         | 337 |

# List of Tables

|      |   |     |
|------|---|-----|
| 8.1  | The number of iterations $k_n$ to solve the discrete Poisson problem with $n$ unknowns using the methods of Jacobi, Gauss-Seidel, and SOR (see text) with a tolerance $10^{-8}$ . . . . . | 194 |
| 8.2  | Spectral radii for $\mathbf{G}_J$ , $\mathbf{G}_1$ , $\mathbf{G}_{\omega^*}$ and the smallest integer $k_n$ such that $\rho(\mathbf{G})^{k_n} \leq 10^{-8}$ . . . . .                     | 204 |
| 9.12 | The number of iterations $K$ for the averaging problem on a $\sqrt{n} \times \sqrt{n}$ grid for various $n$ . . . . .   | 223 |
| 9.14 | The number of iterations $K$ for the Poisson problem on a $\sqrt{n} \times \sqrt{n}$ grid for various $n$ . . . . .   | 224 |
| 9.33 | The number of iterations $K$ (no preconditioning) and $K_{pre}$ (with preconditioning) for the problem (9.50) using the discrete Poisson problem as a preconditioner. . . . .             | 244 |
| 13.9 | Quadratic convergence of Rayleigh quotient iteration. . . . .   | 313 |





# List of Algorithms

|       |  |     |
|-------|--|-----|
| 1.3   | trifactor . . . . .                                | 40  |
| 1.4   | trisolve . . . . .                                 | 40  |
| 2.1   | forwardsolve (row oriented) . . . . .              | 52  |
| 2.2   | backsolve (row oriented) . . . . .                 | 53  |
| 2.4   | Forward solve (column oriented) . . . . .          | 54  |
| 2.5   | Backsolve (column oriented) . . . . .              | 54  |
| 2.47  | bandcholesky . . . . .                             | 72  |
| 2.53  | bandsemi-cholesky . . . . .                        | 76  |
| 4.1   | Fast Poisson solver . . . . .                      | 107 |
| 4.4   | Recursive FFT . . . . .                            | 112 |
| 8.3   | Jacobi . . . . .                                   | 194 |
| 8.4   | SOR . . . . .                                      | 195 |
| 9.11  | Conjugate gradient iteration . . . . .             | 222 |
| 9.13  | Testing conjugate gradient . . . . .               | 224 |
| 9.30  | Preconditioned conjugate gradient . . . . .        | 239 |
| 10.4  | Generate a Householder transformation . . . . .    | 252 |
| 10.8  | Householder triangulation . . . . .                | 255 |
| 10.24 | Upper Hessenberg linear system . . . . .           | 263 |
| 12.14 | Householder reduction to Hessenberg form . . . . . | 298 |
| 12.16 | Assemble Householder transformations . . . . .     | 298 |
| 13.5  | The power method . . . . .                         | 310 |
| 13.7  | Rayleigh quotient iteration . . . . .              | 312 |



# List of Exercises

|      |   |    |
|------|---|----|
| 0.24 | The $\mathbf{A}^T \mathbf{A}$ inner product . . . . .       | 17 |
| 0.25 | Angle between vectors in complex case . . . . .             | 17 |
| 0.41 | The inverse of a general $2 \times 2$ matrix . . . . .      | 24 |
| 0.42 | The inverse of a special $2 \times 2$ matrix . . . . .      | 24 |
| 0.43 | Sherman-Morrison formula . . . . .                          | 24 |
| 0.44 | Cramer's rule; special case . . . . .                       | 25 |
| 0.45 | Adjoint matrix; special case . . . . .                      | 26 |
| 0.47 | Determinant equation for a plane . . . . .                  | 26 |
| 0.48 | Signed area of a triangle . . . . .                         | 27 |
| 0.49 | Vandermonde matrix . . . . .                                | 27 |
| 0.50 | Cauchy determinant (1842) . . . . .                         | 28 |
| 0.51 | Inverse of the Hilbert matrix . . . . .                     | 28 |
| 1.2  | Gaussian elimination example . . . . .                      | 37 |
| 1.8  | Strict diagonal dominance . . . . .                         | 42 |
| 1.9  | LU factorization of 2. derivative matrix . . . . .          | 42 |
| 1.10 | Inverse of 2. derivative matrix . . . . .                   | 43 |
| 1.11 | Central difference approximation of 2. derivative . . . . . | 43 |
| 1.12 | Two point boundary value problem . . . . .                  | 43 |
| 1.13 | Two point boundary value problem; computation . . . . .     | 44 |
| 1.14 | Matrix element as a quadratic form . . . . .                | 47 |
| 1.15 | Outer product expansion of a matrix . . . . .               | 47 |
| 1.16 | The product $\mathbf{A}^T \mathbf{A}$ . . . . .             | 47 |
| 1.17 | Outer product expansion . . . . .                           | 47 |
| 1.18 | System with many right hand sides; compact form . . . . .   | 47 |
| 1.19 | Block multiplication example . . . . .                      | 47 |
| 1.20 | Another block multiplication example . . . . .              | 47 |
| 2.3  | Column oriented forward- and backsolve . . . . .            | 53 |
| 2.6  | Computing the inverse of a triangular matrix . . . . .      | 54 |
| 2.15 | Row interchange . . . . .                                   | 58 |
| 2.16 | LU of singular matrix . . . . .                             | 59 |
| 2.17 | LU and determinant . . . . .                                | 59 |

|      |   |     |
|------|---|-----|
| 2.18 | Diagonal elements in U . . . . .  | 59  |
| 2.20 | Finite sums of integers . . . . .   | 60  |
| 2.21 | Operations . . . . .  | 60  |
| 2.22 | Multiplying triangular matrices . . . . .   | 61  |
| 2.30 | Making block LU into LU . . . . .   | 64  |
| 2.39 | Positive definite characterizations . . . . .                                       | 68  |
| 2.61 | Using PLU of $\mathbf{A}$ to solve $\mathbf{A}^T \mathbf{x} = \mathbf{b}$ . . . . . | 84  |
| 2.62 | Using PLU to compute the determinant . . . . .                                      | 84  |
| 2.63 | Using PLU to compute the inverse . . . . .  | 84  |
| 2.67 | Direct proof of Theorem 2.64 . . . . .  | 87  |
| 3.2  | $4 \times 4$ Poisson matrix . . . . .   | 92  |
| 3.5  | Properties of Kronecker products . . . . .  | 95  |
| 3.11 | 2. derivative matrix is positive definite . . . . .                                 | 99  |
| 3.12 | 1D test matrix is positive definite? . . . . .                                      | 100 |
| 3.13 | Eigenvalues for 2D test matrix of order 4 . . . . .                                 | 100 |
| 3.14 | Nine point scheme for Poisson problem . . . . .                                     | 100 |
| 3.15 | Matrix equation for nine point scheme . . . . .                                     | 100 |
| 3.16 | Biharmonic equation . . . . .   | 101 |
| 4.5  | Fourier matrix . . . . .  | 113 |
| 4.6  | Sine transform as Fourier transform . . . . .                                       | 113 |
| 4.7  | Explicit solution of the discrete Poisson equation . . . . .                        | 113 |
| 4.8  | Improved version of Algorithm 4.1 . . . . .   | 114 |
| 4.9  | Fast solution of 9 point scheme . . . . .   | 114 |
| 4.10 | Algorithm for fast solution of 9 point scheme . . . . .                             | 114 |
| 4.11 | Fast solution of biharmonic equation . . . . .                                      | 115 |
| 4.12 | Algorithm for fast solution of biharmonic equation . . . . .                        | 115 |
| 4.13 | Check algorithm for fast solution of biharmonic equation . . . . .                  | 115 |
| 4.14 | Fast solution of biharmonic equation using 9 point rule . . . . .                   | 115 |
| 5.3  | Eigenvalues of an idempotent matrix . . . . .                                       | 121 |
| 5.4  | Eigenvalues of a nilpotent matrix . . . . .   | 121 |
| 5.5  | Eigenvalues of a unitary matrix . . . . .   | 121 |
| 5.6  | Nonsingular approximation of a singular matrix . . . . .                            | 121 |
| 5.7  | Companion matrix . . . . .  | 121 |
| 5.16 | Schur decomposition example . . . . .   | 126 |
| 5.19 | Skew-Hermitian matrix . . . . .   | 127 |
| 5.20 | Eigenvalues of a skew-Hermitian matrix . . . . .                                    | 127 |
| 5.31 | Eigenvalue perturbation for Hermitian matrices . . . . .                            | 132 |
| 5.33 | Hoffman-Wielandt . . . . .  | 132 |
| 5.42 | Find eigenpair example . . . . .  | 135 |
| 5.45 | Jordan example . . . . .  | 138 |
| 5.46 | Big Jordan example . . . . .  | 138 |
| 5.49 | Properties of the Jordan form . . . . .   | 139 |

|      |  |     |
|------|--|-----|
| 5.50 | Powers of a Jordan block . . . . .   | 139 |
| 5.52 | Minimal polynomial example . . . . .   | 140 |
| 5.53 | Similar matrix polynomials . . . . .   | 140 |
| 5.54 | Minimal polynomial of a diagonalizable matrix . . . . .                        | 141 |
| 5.59 | Biorthogonal expansion . . . . .   | 142 |
| 5.60 | Generalized Rayleigh quotient . . . . .  | 143 |
| 6.14 | SVD examples . . . . .   | 155 |
| 6.15 | More SVD examples . . . . .  | 155 |
| 6.17 | Counting dimensions of fundamental subspaces . . . . .                         | 156 |
| 6.18 | Rank and nullity relations . . . . .   | 156 |
| 6.19 | Orthonormal bases example . . . . .  | 156 |
| 6.20 | Some spanning sets . . . . .   | 156 |
| 6.21 | Singular values and eigenpair of composite matrix . . . . .                    | 157 |
| 6.27 | Rank example . . . . .   | 160 |
| 6.28 | Another rank example . . . . .   | 161 |
| 7.4  | Consistency of sum norm? . . . . .   | 166 |
| 7.5  | Consistency of max norm? . . . . .   | 167 |
| 7.6  | Consistency of modified max norm . . . . .                                     | 167 |
| 7.8  | What is the sum norm subordinate to? . . . . .                                 | 167 |
| 7.9  | What is the max norm subordinate to? . . . . .                                 | 167 |
| 7.16 | Spectral norm . . . . .  | 171 |
| 7.17 | Spectral norm of the inverse . . . . .   | 171 |
| 7.18 | $p$ -norm example . . . . .  | 172 |
| 7.21 | Unitary invariance of the spectral norm . . . . .                              | 172 |
| 7.22 | $\ \mathbf{A}\mathbf{U}\ _2$ rectangular $\mathbf{A}$ . . . . .                | 172 |
| 7.23 | $p$ -norm of diagonal matrix . . . . .   | 173 |
| 7.24 | spectral norm of a column vector . . . . .                                     | 173 |
| 7.25 | Norm of absolute value matrix . . . . .  | 173 |
| 7.32 | Sharpness of perturbation bounds . . . . .                                     | 178 |
| 7.33 | Condition number of 2. derivative matrix . . . . .                             | 178 |
| 7.44 | When is a complex norm an inner product norm? . . . . .                        | 184 |
| 7.45 | $p$ norm for $p = 1$ and $p = \infty$ . . . . .                                | 184 |
| 7.46 | The $p$ - norm unit sphere . . . . .   | 184 |
| 7.47 | Sharpness of $p$ -norm inequality . . . . .                                    | 184 |
| 7.48 | $p$ -norm inequalities for arbitrary $p$ . . . . .                             | 184 |
| 8.2  | Richardson and Jacobi . . . . .  | 193 |
| 8.13 | Convergence of the R-method when eigenvalues have positive real part . . . . . | 200 |
| 8.16 | Example: GS converges, J diverges . . . . .                                    | 201 |
| 8.17 | Divergence example for J and GS . . . . .                                      | 202 |
| 8.18 | Strictly diagonally dominance; The J method . . . . .                          | 202 |
| 8.19 | Strictly diagonally dominance; The GS method . . . . .                         | 202 |

|       |  |     |
|-------|--|-----|
| 8.23  | Convergence example for fix point iteration . . . . .                        | 205 |
| 8.24  | Estimate in Lemma 8.22 can be exact . . . . .                                | 205 |
| 8.25  | Slow spectral radius convergence . . . . .                                   | 206 |
| 8.31  | A special norm . . . . .   | 209 |
| 8.33  | When is $\mathbf{A} + \mathbf{E}$ nonsingular? . . . . .                     | 210 |
| 9.1   | Paraboloid . . . . .   | 216 |
| 9.4   | Steepest descent iteration . . . . .   | 219 |
| 9.7   | Conjugate gradient iteration, II . . . . .                                   | 221 |
| 9.8   | Conjugate gradient iteration, III . . . . .                                  | 221 |
| 9.9   | The cg step length is optimal . . . . .                                      | 221 |
| 9.10  | Starting value in cg . . . . .   | 222 |
| 9.15  | The $\mathbf{A}$ -inner product . . . . .                                    | 225 |
| 9.17  | Program code for testing steepest descent . . . . .                          | 226 |
| 9.18  | Using cg to solve normal equations . . . . .                                 | 227 |
| 9.20  | Maximum of a convex function . . . . .                                       | 228 |
| 9.25  | Krylov space and cg iterations . . . . .                                     | 232 |
| 9.28  | Another explicit formula for the Chebyshev polynomial . . . . .              | 235 |
| 10.2  | Reflector . . . . .  | 250 |
| 10.5  | What does algorithm housegen do when $\mathbf{x} = \mathbf{e}_1$ ? . . . . . | 252 |
| 10.6  | Examples of Householder transformations . . . . .                            | 252 |
| 10.7  | $2 \times 2$ Householder transformation . . . . .                            | 253 |
| 10.16 | QR decomposition . . . . .   | 258 |
| 10.17 | Householder triangulation . . . . .  | 259 |
| 10.20 | QR using Gram-Schmidt, II . . . . .  | 260 |
| 10.22 | Plane rotation . . . . .   | 261 |
| 10.23 | Solving upper Hessenberg system using rotations . . . . .                    | 262 |
| 11.7  | Straight line fit (linear regression) . . . . .                              | 269 |
| 11.8  | Straight line fit using shifted power form . . . . .                         | 269 |
| 11.9  | Fitting a circle to points . . . . .   | 270 |
| 11.15 | The generalized inverse . . . . .  | 274 |
| 11.16 | Uniqueness of generalized inverse . . . . .                                  | 274 |
| 11.17 | Verify that a matrix is a generalized inverse . . . . .                      | 275 |
| 11.18 | Linearly independent columns and generalized inverse . . . . .               | 275 |
| 11.19 | The generalized inverse of a vector . . . . .                                | 275 |
| 11.20 | The generalized inverse of an outer product . . . . .                        | 275 |
| 11.21 | The generalized inverse of a diagonal matrix . . . . .                       | 275 |
| 11.22 | Properties of the generalized inverse . . . . .                              | 275 |
| 11.23 | The generalized inverse of a product . . . . .                               | 276 |
| 11.24 | The generalized inverse of the conjugate transpose . . . . .                 | 276 |
| 11.25 | Linearly independent columns . . . . .                                       | 276 |
| 11.26 | Analysis of the general linear system . . . . .                              | 276 |
| 11.27 | Fredholm's alternative . . . . .   | 276 |

---

|       |   |     |
|-------|---|-----|
| 11.33 | Condition number . . . . .                                      | 282 |
| 11.34 | Equality in perturbation bound . . . . .                        | 282 |
| 11.36 | Problem using normal equations . . . . .                        | 283 |
| 12.5  | Continuity of eigenvalues . . . . .                             | 292 |
| 12.6  | Nonsingularity using Gerschgorin . . . . .                      | 292 |
| 12.7  | Gerschgorin, strictly diagonally dominant matrix . . . . .      | 293 |
| 12.12 | $\infty$ -norm of a diagonal matrix . . . . .                   | 295 |
| 12.15 | Number of arithmetic operations, Hessenberg reduction . . . . . | 298 |
| 12.17 | Assemble Householder transformations . . . . .                  | 298 |
| 12.18 | Tridiagonalize a symmetric matrix . . . . .                     | 299 |
| 12.22 | Counting eigenvalues . . . . .                                  | 302 |
| 12.23 | Overflow in LDLT factorization . . . . .                        | 302 |
| 12.24 | Simultaneous diagonalization . . . . .                          | 302 |
| 12.25 | Program code for one eigenvalue . . . . .                       | 303 |
| 12.26 | Determinant of upper Hessenberg matrix . . . . .                | 304 |
| 13.4  | Orthogonal vectors . . . . .                                    | 310 |
| 13.14 | QR convergence detail . . . . .                                 | 319 |





# List of Theorems

|      |  |    |
|------|--|----|
| 0.7  | Basis subset of a spanning set . . . . .                   | 7  |
| 0.9  | Dimension of a vector space . . . . .                      | 8  |
| 0.10 | Enlarging vectors to a basis . . . . .                     | 8  |
| 0.13 | Dimension formula for sums of subspaces . . . . .          | 9  |
| 0.19 | Basic properties of vector norms . . . . .                 | 13 |
| 0.22 | Cauchy-Schwarz inequality . . . . .                        | 15 |
| 0.23 | Inner product norm . . . . .                               | 16 |
| 0.27 | Pythagoras . . . . .                                       | 17 |
| 0.29 | Gram-Schmidt . . . . .                                     | 18 |
| 0.30 | Orthogonal Extension of basis . . . . .                    | 19 |
| 0.34 | Linear systems; existence and uniqueness . . . . .         | 21 |
| 0.37 | Complex linear system; existence and uniqueness . . . . .  | 22 |
| 0.38 | Product of nonsingular matrices . . . . .                  | 22 |
| 0.39 | When is a square matrix invertible? . . . . .              | 23 |
| 0.54 | Zero eigenvalue . . . . .                                  | 30 |
| 0.55 | Eigenvalues of a triangular matrix . . . . .               | 30 |
| 1.6  | Strict diagonal dominance . . . . .                        | 41 |
| 1.7  | Weak diagonal dominance . . . . .                          | 41 |
| 2.12 | LU theorem . . . . .                                       | 57 |
| 2.26 | Symmetric LU theorem . . . . .                             | 62 |
| 2.27 | Block LU theorem . . . . .                                 | 63 |
| 2.35 | A general criterium . . . . .                              | 66 |
| 2.37 | The nonsymmetric case . . . . .                            | 67 |
| 2.38 | Symmetric positive definite characterization . . . . .     | 68 |
| 2.44 | Cholesky . . . . .   | 70 |
| 2.49 | Positive symmetric semidefinite characterization . . . . . | 73 |
| 2.51 | Characterization, semi-Cholesky factorization . . . . .    | 74 |
| 2.52 | Bandwidth semi-Cholesky factor . . . . .                   | 75 |
| 2.59 | Gaussian elimination is well defined . . . . .             | 83 |
| 2.60 | PLU theorem . . . . .                                      | 83 |
| 2.64 | LU factorization . . . . .                                 | 85 |

|      |  |     |
|------|--|-----|
| 2.65 | Nonzero pivots . . . . .   | 85  |
| 3.7  | Properties of Kronecker products . . . . .   | 95  |
| 3.10 | Eigenpairs of 2D test matrix . . . . .   | 99  |
| 4.3  | Fast Fourier transform . . . . .   | 111 |
| 5.1  | Transformations of eigenpairs . . . . .  | 119 |
| 5.2  | Sums and products of eigenvalues; trace . . . . .                                    | 120 |
| 5.9  | Eigenpairs of similar matrices . . . . .   | 122 |
| 5.12 | Unitary matrix . . . . .   | 123 |
| 5.13 | Schur decomposition . . . . .  | 124 |
| 5.14 | Schur form, real eigenvalues . . . . .   | 125 |
| 5.21 | Orthonormal eigenpairs characterization . . . . .                                    | 127 |
| 5.27 | Minmax . . . . .   | 129 |
| 5.28 | Maxmin . . . . .   | 130 |
| 5.30 | Eigenvalue perturbation for Hermitian matrices . . . . .                             | 131 |
| 5.32 | Hoffman-Wielandt theorem . . . . .   | 132 |
| 5.35 | Eigenvectors of diagonalizable matrices . . . . .                                    | 133 |
| 5.36 | Distinct eigenvalues . . . . .   | 133 |
| 5.41 | Geometric multiplicity of similar matrices . . . . .                                 | 135 |
| 5.44 | The Jordan form of a matrix . . . . .  | 136 |
| 5.56 | Biorthogonality . . . . .  | 141 |
| 5.57 | Simple eigenvalue . . . . .  | 141 |
| 5.58 | Biorthogonal eigenvector expansion . . . . .   | 142 |
| 5.61 | The real Schur form . . . . .  | 143 |
| 6.5  | The matrices $\mathbf{A}^* \mathbf{A}$ , $\mathbf{A} \mathbf{A}^*$ . . . . .         | 149 |
| 6.6  | The matrices $\mathbf{A}^* \mathbf{A}$ , $\mathbf{A} \mathbf{A}^*$ and SVD . . . . . | 150 |
| 6.7  | Existence of SVD . . . . .   | 150 |
| 6.9  | Singular values of a normal matrix . . . . .   | 152 |
| 6.16 | Singular vectors and orthonormal bases . . . . .                                     | 155 |
| 6.22 | SVF ellipse . . . . .  | 157 |
| 6.25 | Frobenius norm and singular values . . . . .   | 160 |
| 6.26 | Best low rank approximation . . . . .  | 160 |
| 6.29 | The Courant-Fischer theorem for singular values . . . . .                            | 161 |
| 6.30 | Hoffman-Wielandt theorem for singular values . . . . .                               | 162 |
| 7.2  | Matrix norm equivalence . . . . .  | 165 |
| 7.12 | onetwoinnorms . . . . .  | 169 |
| 7.14 | Spectral norm . . . . .  | 170 |
| 7.15 | Spectral norm bound . . . . .  | 171 |
| 7.20 | Unitary invariant norms . . . . .  | 172 |
| 7.26 | Perturbation in the right-hand side . . . . .  | 174 |
| 7.27 | Spectral condition number . . . . .  | 175 |
| 7.28 | Perturbation and residual . . . . .  | 175 |
| 7.29 | Nonsingularity of perturbation of identity . . . . .                                 | 175 |

---

|       |   |     |
|-------|---|-----|
| 7.30  | Nonsingularity of perturbation . . . . .                                    | 176 |
| 7.31  | Perturbation of inverse matrix . . . . .                                    | 177 |
| 7.34  | The $p$ vector norms are norms . . . . .                                    | 179 |
| 7.37  | Jensen's inequality . . . . .   | 180 |
| 7.41  | Parallelogram identity . . . . .  | 182 |
| 7.42  | When is a norm an inner product norm? . . . . .                             | 182 |
| 8.5   | Splitting matrices for R, J, and SOR . . . . .                              | 196 |
| 8.8   | Convergence of an iterative method . . . . .                                | 198 |
| 8.9   | Sufficient condition for convergence . . . . .                              | 198 |
| 8.10  | When does an iterative method converge? . . . . .                           | 199 |
| 8.11  | Convergence of Richardson's method . . . . .                                | 199 |
| 8.14  | Necessary condition for convergence of SOR . . . . .                        | 200 |
| 8.15  | SOR on positive definite matrix . . . . .                                   | 200 |
| 8.20  | The spectral radius of SOR matrix . . . . .                                 | 202 |
| 8.27  | When is $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0}$ ? . . . . . | 207 |
| 8.28  | Any consistent norm majorizes the spectral radius . . . . .                 | 208 |
| 8.29  | The spectral radius can be approximated by a norm . . . . .                 | 208 |
| 8.30  | Spectral radius convergence . . . . .                                       | 208 |
| 8.32  | Neumann series . . . . .  | 209 |
| 8.34  | The optimal $\omega$ . . . . .  | 211 |
| 9.16  | Error bound for steepest descent and conjugate gradients . . . . .          | 225 |
| 9.19  | Kantorovich inequality . . . . .  | 227 |
| 9.22  | Best approximation property . . . . .                                       | 230 |
| 9.24  | cg and best polynomial approximation . . . . .                              | 231 |
| 9.27  | A minimal norm problem . . . . .  | 234 |
| 9.29  | The error in cg is strictly decreasing . . . . .                            | 236 |
| 9.31  | Error bound preconditioned cg . . . . .                                     | 239 |
| 9.32  | Positive definite matrix . . . . .  | 242 |
| 9.34  | Eigenvalues of preconditioned matrix . . . . .                              | 244 |
| 10.3  | Zeros in vectors . . . . .  | 251 |
| 10.12 | Existence of QR decomposition . . . . .                                     | 257 |
| 10.13 | Uniqueness of QR factorization . . . . .                                    | 257 |
| 10.15 | Hadamard's inequality . . . . .   | 258 |
| 10.18 | QR and Gram-Schmidt . . . . .   | 259 |
| 11.10 | Orthogonal projection and least squares solution . . . . .                  | 272 |
| 11.12 | Uniqueness . . . . .  | 273 |
| 11.13 | Characterization . . . . .  | 273 |
| 11.29 | Minimal solution . . . . .  | 279 |
| 11.31 | Perturbing the right hand side . . . . .                                    | 280 |
| 11.35 | Perturbing the matrix . . . . .   | 282 |
| 11.37 | Perturbation of singular values . . . . .                                   | 283 |
| 11.38 | Generalized inverse when perturbing the matrix . . . . .                    | 284 |

---

|       |  |     |
|-------|--|-----|
| 12.1  | Gerschgorin's circle theorem . . . . .                 | 290 |
| 12.8  | Elsner's theorem(1985) . . . . .                       | 293 |
| 12.9  | Absolute errors . . . . .                              | 294 |
| 12.10 | Perturbations, normal matrix . . . . .                 | 295 |
| 12.13 | Relative errors . . . . .                              | 295 |
| 12.20 | Sylvester's inertia theorem . . . . .                  | 301 |
| 13.3  | The Rayleigh quotient minimizes the residual . . . . . | 309 |
| 13.12 | QR and power . . . . .                                 | 315 |
| 13.13 | Convergence of basis QR . . . . .                      | 318 |
| A.2   | Cramer's rule (1750) . . . . .                         | 329 |
| A.4   | The inverse as an adjoint . . . . .                    | 330 |
| A.7   | Cauchy-Binet formula . . . . .                         | 332 |
| B.5   | Relative error in rounding . . . . .                   | 339 |

# List of Corollaries

|       |  |     |
|-------|--|-----|
| 0.8   | Existence of a basis . . . . .   | 7   |
| 0.31  | Extending orthogonal vectors to a basis . . . . .                      | 19  |
| 0.40  | Basic properties of the inverse matrix . . . . .                       | 23  |
| 2.36  | $\mathbf{A}^T \mathbf{A}$ is symmetric positive semidefinite . . . . . | 66  |
| 2.66  | When is naive Gaussian elimination possible? . . . . .                 | 86  |
| 5.10  | Spectra of $\mathbf{AB}$ and $\mathbf{BA}$ . . . . .                   | 122 |
| 5.22  | Spectral theorem, complex form . . . . .                               | 128 |
| 5.23  | Spectral theorem (real form) . . . . .                                 | 128 |
| 5.29  | The Courant-Fischer theorem . . . . .                                  | 131 |
| 5.37  | Diagonalizable matrix . . . . .  | 134 |
| 5.47  | Geometric multiplicity . . . . .                                       | 138 |
| 7.38  | Weighted geometric/arithmetic mean inequality . . . . .                | 180 |
| 7.39  | Hölder's inequality . . . . .  | 181 |
| 7.40  | Minkowski's inequality . . . . .                                       | 182 |
| 7.43  | Are the $p$ -norms inner product norms? . . . . .                      | 184 |
| 8.12  | Rate of convergence for the R method . . . . .                         | 199 |
| 12.3  | Disjoint Gerschgorin disks . . . . .                                   | 291 |
| 12.21 | Counting eigenvalues using the LDLT factorization . . . . .            | 301 |
| A.5   | The adjoint and the inverse . . . . .                                  | 331 |
| A.6   | Cofactor expansion . . . . .   | 331 |



# List of Lemmas

|      |   |     |
|------|---|-----|
| 0.5  | Linear independence and span . . . . .                      | 7   |
| 0.14 | Change of basis matrix . . . . .                            | 10  |
| 0.32 | Underdetermined system . . . . .                            | 20  |
| 0.35 | Complex underdetermined system . . . . .                    | 21  |
| 0.52 | Characteristic equation . . . . .                           | 29  |
| 1.21 | Inverse of a block triangular matrix . . . . .              | 48  |
| 1.22 | Inverse of a triangular matrix . . . . .                    | 49  |
| 1.23 | Product of triangular matrices . . . . .                    | 49  |
| 1.24 | Unit triangular matrices . . . . .                          | 49  |
| 2.11 | LU of leading principal sub matrices . . . . .              | 57  |
| 2.25 | Symmetric LU of leading principal sub matrices . . . . .    | 62  |
| 2.31 | Quadratic form with $\mathbf{x} \in \mathbb{C}^n$ . . . . . | 65  |
| 2.33 | $\mathbf{T}$ is symmetric positive definite . . . . .       | 65  |
| 2.40 | Eigenvalues of a Hermitian matrix . . . . .                 | 68  |
| 2.41 | Symmetry and positive eigenvalues . . . . .                 | 69  |
| 2.42 | Symmetric positive definite and symmetric LU . . . . .      | 69  |
| 2.46 | Banded Cholesky factor . . . . .                            | 71  |
| 2.48 | Criteria symmetric semidefinite . . . . .                   | 72  |
| 2.58 | Gausstransformations . . . . .                              | 81  |
| 3.6  | Mixed product rule . . . . .                                | 95  |
| 3.8  | Eigenpairs of 1D test matrix . . . . .                      | 98  |
| 3.9  | Eigenpairs of a Hermitian matrix . . . . .                  | 98  |
| 4.2  | Sine transform as Fourier transform . . . . .               | 109 |
| 5.26 | Convex combination of the eigenvalues . . . . .             | 129 |
| 5.48 | A nilpotent matrix . . . . .                                | 139 |
| 6.24 | Frobenius norm properties . . . . .                         | 159 |
| 7.11 | The operator norm is a matrix norm . . . . .                | 169 |
| 7.36 | A sufficient condition for convexity . . . . .              | 179 |
| 8.22 | Number of iterations . . . . .                              | 204 |
| 8.26 | Be careful when stopping . . . . .                          | 206 |
| 9.2  | Quadratic function . . . . .                                | 217 |

---

|       |  |     |
|-------|--|-----|
| 9.5   | The residuals are orthogonal . . . . .                 | 220 |
| 9.21  | Krylov space . . . . .                                 | 229 |
| 9.23  | Krylov space and polynomials . . . . .                 | 230 |
| 9.26  | Closed forms of Chebyshev polynomials . . . . .        | 233 |
| 10.9  | Updating a Householder transformation . . . . .        | 256 |
| 11.4  | Curve fitting . . . . .                                | 268 |
| 12.11 | $p$ -norm of a diagonal matrix . . . . .               | 295 |
| 12.19 | Distinct eigenvalues of a tridiagonal matrix . . . . . | 300 |
| 13.2  | Convergence of the power method . . . . .              | 309 |
| B.3   | Relative errors . . . . .                              | 335 |
| B.6   | Bound on factors . . . . .                             | 341 |
| C.1   | Product rules . . . . .                                | 345 |
| C.2   | Second order Taylor expansion . . . . .                | 346 |
| C.3   | Functions involving matrices . . . . .                 | 346 |



# List of Definitions

|      |   |     |
|------|---|-----|
| 0.1  | Real vector space . . . . .                     | 4   |
| 0.2  | Linear combination . . . . .                    | 6   |
| 0.4  | Linear independence . . . . .                   | 6   |
| 0.6  | basis . . . . .                                 | 7   |
| 0.11 | Subspace . . . . .                              | 8   |
| 0.15 | Column space and null space . . . . .           | 11  |
| 0.16 | Vector norm . . . . .                           | 11  |
| 0.17 | Vector p-norms . . . . .                        | 12  |
| 0.18 | Equivalent norms . . . . .                      | 13  |
| 0.20 | Real inner product . . . . .                    | 14  |
| 0.21 | Complex inner product . . . . .                 | 15  |
| 0.26 | Orthogonality . . . . .                         | 17  |
| 0.28 | Orthogonal- and orthonormal bases . . . . .     | 18  |
| 0.33 | Real nonsingular or singular matrix . . . . .   | 21  |
| 0.36 | Complex nonsingular matrix . . . . .            | 21  |
| 0.53 | Characteristic polynomial of a matrix . . . . . | 30  |
| 1.5  | Diagonal dominance . . . . .                    | 41  |
| 2.9  | Principal submatrix . . . . .                   | 56  |
| 2.19 | $G_n := \frac{2}{3}n^3$ . . . . .               | 59  |
| 2.23 | Symmetric LU . . . . .                          | 61  |
| 2.43 | Cholesky . . . . .                              | 70  |
| 2.50 | Semi-Cholesky factorization . . . . .           | 74  |
| 2.55 | . . . . .                                       | 80  |
| 2.56 | Interchange matrix . . . . .                    | 80  |
| 2.57 | Gauss transformation . . . . .                  | 81  |
| 3.1  | vec operation . . . . .                         | 91  |
| 3.3  | Kronecker product . . . . .                     | 94  |
| 3.4  | Kronecker sum . . . . .                         | 94  |
| 5.8  | Similar matrices . . . . .                      | 122 |
| 5.11 | Unitary matrix . . . . .                        | 123 |
| 5.17 | Quasi-triangular matrix . . . . .               | 126 |

---

|       |  |     |
|-------|--|-----|
| 5.18  | Normal matrix . . . . .  | 127 |
| 5.25  | Rayleigh quotient . . . . .  | 129 |
| 5.34  | Diagonalizable matrix . . . . .  | 133 |
| 5.39  | Geometric multiplicity . . . . .   | 135 |
| 5.43  | Jordan block . . . . .   | 136 |
| 5.51  | Minimal polynomial of a matrix . . . . .   | 139 |
| 5.55  | Left eigenpair . . . . .   | 141 |
| 6.1   | SVD . . . . .  | 147 |
| 6.8   | SVF . . . . .  | 151 |
| 7.1   | Matrix norms . . . . .   | 165 |
| 7.3   | Consistent matrix norms . . . . .  | 166 |
| 7.7   | Subordinate matrix norms . . . . .   | 167 |
| 7.10  | Operator norm . . . . .  | 168 |
| 7.19  | Unitary invariant norm . . . . .   | 172 |
| 7.35  | Convex function . . . . .  | 179 |
| 8.7   | Convergence of $\mathbf{x}_{k+1} := \mathbf{G}\mathbf{x}_k + \mathbf{c}$ . . . . . | 198 |
| 10.1  | Householder transformation . . . . .   | 250 |
| 10.10 | QR decomposition . . . . .   | 256 |
| 10.21 | Givens rotation, plane rotation . . . . .  | 261 |
| 11.1  | Least squares problem . . . . .  | 265 |
| A.3   | Cofactor and Adjoint . . . . .   | 330 |
| B.1   | Absolute error . . . . .   | 335 |
| B.2   | Relative error . . . . .   | 335 |

# Index

- 1D test matrix, 93
- 2D test matrix, 93
  
- convex combinations, 228
- eigenvector expansion, 133
- singular values, 148
  
- absolute error, 174, 335
- adjoint matrix, 330
- adjusted rounding unit, 342
- algebraic multiplicity, 135
- algorithms
  - assemble Householder transformations, 298
  - backsolve, 53
  - backsolve column oriented, 54
  - bandcholesky, 72
  - cg, 222
  - fastpoisson, 107
  - forwardsolve, 52
  - forwardsolve column oriented, 54
  - housegen, 252
  - Householder reduction to Hessenberg form, 298
  - Householder triangulation, 255
  - Jacobi, 194
  - preconditioned cg, 239
  - Rayleigh quotient iteration, 312
  - SOR, 195
  - testing conjugate gradient, 224
  - the power method, 310
  - trifactor, 40
  - trisolve, 40
  - upper Hessenberg linear system, 263
  
- back substitution, 37
- back substitution, 52
- backward error, 340
- backward stable, 340
- banded matrix, 3
  - symmetric LU factorization, 72
- banded symmetric LU factorization, 72
- bandsemi-cholesky, 76
- biharmonic equation, 101
  - fast solution method, 115
  - nine point rule, 115
- block LU theorem, 63
  
- Cauchy determinant, 28
- Cauchy-Binet formula, 332
- Cauchy-Schwarz inequality, 16
- Cayley Hamilton Theorem, 140
- central difference, 43
- central difference approximation
  - second derivative, 43
- change of basis matrix, 10
- characteristic equation, 30
- characteristic polynomial, 30, 119
- Chebyshev polynomial, 233
- Cholesky factor, 70

- Cholesky factorization, 70
- cofactor, 330
- column operations, 332
- column space (span), 11
- companion matrix, 121
- complete pivoting, 80
- complexity of an algorithm, 59
- computer arithmetic, 335
- condition number, 342
  - ill-conditioned, 174
- congruent matrices, 301
- conjugate gradient method, 215
- convergence, 225
  - derivation, 219
  - Krylov spaces, 229
  - least squares problem, 227
  - preconditioning, 237
  - preconditioning algorithm, 239
  - preconditioning convergence, 240
- convex combination, 129, 179
- convex function, 179
- Courant-Fischer theorem, 131
- Cramer's rule, 25, 330
- Crout factorization, 55
- defective matrix, 133
- deflation, 125
- determinant, 325
  - additivity, 326
  - area of a triangle, 27
  - block triangular, 326
  - Cauchy, 28
  - Cauchy-Binet, 332
  - cofactor, 330
  - cofactor expansion, 25, 331
  - homogeneity, 326
  - permutation of columns, 326
  - plane equation, 27
  - product rule, 326
  - singular matrix, 326
  - straight line equation, 26
  - transpose, 326
  - triangular matrix, 326
  - Vandermonde, 28
- diagonalizable matrix, 133
- dirac delta, 3
- direct sum, 271
- discrete cosine transform, 108
- discrete Fourier transform, 108, 109
  - Fourier matrix, 109
- discrete sine transform, 108
- double precision, 339
- eigenpair, 29, 119
  - left eigenpair, 141
- eigenvalue, 29, 119
  - algebraic multiplicity, 135
  - characteristic equation, 30
  - characteristic polynomial, 30
  - Courant-Fischer theorem, 131
  - geometric multiplicity, 135
  - Hoffman-Wielandt theorem, 132
  - location, 290
  - Rayleigh quotient, 129
  - Schur form, real, 143
  - spectral theorem, 129
  - spectrum, 29, 119
  - transformations of eigenpairs, 119
  - triangular matrix, 30
- eigenvector, 29, 119
  - left eigenvector, 141
- elementary divisors, 139
- elementary lower triangular matrix, 81
- elementary reflector, 250
- Elsner's theorem, 293
- equivalent norms, 13
- extension of basis, 19
- fast Fourier transform, 108, 110

- recursive FFT, 112
- fill-inn, 104
- finite difference method, 38
- finite dimensional vector space, 6
- fixed point form of discrete Poisson equation, 193
- fixed-point, 196
- fixed-point iteration, 196
- floating-point number
  - bias, 338
  - denormalized, 338
  - double precision, 339
  - exponent part, 337
  - guard digits, 340
  - half precision, 339
  - Inf, 338
  - mantissa, 337
  - NaN, 338
  - normalized, 337
  - overflow, 338
  - quadruple precision, 339
  - round away from zero, 339
  - round to zero, 339
  - rounding, 339
  - rounding unit, 339
  - single precision, 339
  - subnormal, 338
- forward substitution, 52
- Fourier matrix, 109
- Fredholm's alternative, 276
- Frobenius norm, 159
- Gauss transformation, 81
- Gaussian elimination
  - complete pivoting, 80
  - interchange matrix, 80
  - pivot, 78
  - pivoting, 78
  - PLU factorization, 83
  - row pivoting, 79
- generalized inverse, 274
- geometric multiplicity, 135
- Gerschgorin's theorem, 290
- Given's rotation, 261
- gradient, 66, 345
- gradient method, 218
- Gram-Schmidt, 18
- guard digits, 340
- Hölder's inequality, 13, 181
- Hadamard's inequality, 258
- half precision, 339
- hessian, 66, 345
- Hilbert matrix, 29, 269
- Hoffman-Wielandt theorem, 132
- Householder transformation, 250
- identity matrix, 3
- ill-conditioned, 343
- ill-conditioned problem, 174
- inequality
  - geometric/arithmetic mean, 181
  - Hölder, 181
  - Kantorovich, 227
  - Minkowski, 182
- Inf, 338
- inner product, 14
  - complex, 15
  - inner product norm, 14, 15
  - Pythagoras' theorem, 17
  - standard inner product in  $\mathbb{C}^n$ , 15
  - standard inner product in  $\mathbb{R}^n$ , 15
- inner product space
  - orthogonal basis, 18
  - orthonormal basis, 18
- interchange matrix, 80
- inverse power method, 311
- inverse triangle inequality, 13
- iterative method
  - convergence, 198

- Gauss-Seidel, 191
- Jacobi, 191
- SOR, 192
- SOR, convergence, 200
- SSOR, 192
- iterative methods, 189
- Jacobian, 346
- Jordan factors, 137
- Jordan form, 137
  - elementary divisors, 139
  - Jordan block, 136
  - Jordan canonical form, 136
  - principal vectors, 137
- Kronecker product, 94
  - eigenvectors, 96
  - inverse, 96
  - left product, 94
  - mixed product rule, 95
  - nonsingular, 96
  - positive definite, 96
  - propertis, 95
  - right product, 94
  - symmetry, 96
  - transpose, 96
- Kronecker sum, 94
  - nonsingular, 96
  - positive definite, 96
  - symmetry, 96
- Krylov space, 229
- Laplacian, 345
- leading principal block submatrices, 63
- leading principal minor, 57
- leading principal submatrices, 56
- least squares
  - error analysis, 280
- least squares problem, 265
- least squares solution, 265
- left eigenpair, 141
- left eigenvector, 141
- left triangular, 54
- linear combination, 6
- linear system
  - existence and uniqueness, 21, 22
  - homogeneous, 20
  - overdetermined, 20
  - residual vector, 175
  - square, 20
  - underdetermined, 20
- linearly dependent, 6
- linearly independent, 6
- LLT factorization, 70
- LU factorization, 54
  - LDLT, 61
  - symmetric, 61
  - symmetric LU, 62
- LU theorem, 57
- mantissa, 337
- matrix
  - addition, 2
  - adjoint, 25, 330
  - adjoint formula for the inverse, 25
  - block matrix, 45
  - block triangular, 48
  - blocks, 45
  - cofactor, 25
  - column space (span), 11
  - companion matrix, 121
  - computing inverse, 54
  - defective, 133
  - deflation, 125
  - diagonal, 3
  - element-by-element operations, 3
  - Hadamard product, 3
  - Hilbert, 29
  - idempotent, 121
  - ill-conditioned, 175
  - inverse, 22

- inverse Hilbert matrix, 29
- invertible, 22
- Kronecker product, 94
- leading principal minor, 56
- leading principal submatrices, 56
- left inverse, 22
- left triangular, 3
- lower Hessenberg, 3
- lower triangular, 3
- LU theorem, 57
- multiplication, 3
- negative (semi)definite, 64, 65
- Neumann series, 209
- nilpotent, 121
- nonsingular, 21
- norm, 165
- normal, 127
- null space (ker), 11
- operator norm, 168
- outer product, 47
- outer product expansion, 47
- permutation, 80
- positive definite, 64
- positive semidefinite, 64
- principal minor, 56
- principal submatrix, 56
- product of triangular matrices, 49
- quasi-triangular, 126
- right inverse, 22
- right triangular, 3
- row space, 11
- scalar multiplication, 3
- Schur product, 3
- second derivative, 39
- similar matrices, 122
- similarity transformation, 122
- singular, 21
- spectral radius, 199, 207
- strictly diagonally dominant, 41
- symmetric positive semidefinite, 65
- test matrix, 2D, 93
- test matrix, 1D, 93
- trace, 120
- triangular, 49
- tridiagonal, 3
- unit triangular, 49
- upper Hessenberg, 3
- upper trapezoidal, 253
- upper triangular, 3
- vec operation, 91
- weakly diagonally dominant, 41
- well-conditioned, 175
- matrix norm
  - consistent norm, 166
  - Frobenius norm, 159, 166
  - max norm, 166
  - operator norm, 168
  - spectral norm, 170
  - subordinate norm, 167
  - sum norm, 166
  - two-norm, 170
- minimal polynomial, 140
- Minkowski's inequality, 12, 182
- mixed product rule, 95
- naive Gaussian elimination, 85
- NaN, 338
- natural ordering, 91
- negative (semi)definite, 64, 65
- Neumann series, 209
- nilpotent matrix, 121
- nonsingular matrix, 21
- nontrivial subspaces, 9
- norm, 11
  - $l_1$ -norm, 12
  - $l_2$ -norm, 12
  - $l_\infty$ -norm, 12

- absolute norm, 173
- continuity, 13
- Euclidian norm, 12
- infinity-norm, 12
- max norm, 12
- monotone norm, 173
- one-norm, 12
- triangle inequality, 12
- two-norm, 12
- normal equations, 266
- normal matrix, 127
- null space (ker), 11
- operation count, 59
- operator norm, 168
- optimal relaxation parameter, 203
- optimal step length, 218
- orthogonal decomposition, 272
- orthogonal matrix, see orthonormal matrix, 123
- orthogonal projections, 271
- orthogonal sum, 271
- orthonormal matrix, 123
- outer product, 47
- overflow, 338
- p-norms, 12
- paraboloid, 217
- parallelogram identity, 182
- partial pivoting, 78
- permutation, 323
  - identity, 323
  - inversion, 324
  - sign, 324
  - symmetric group, 325
- permutation matrix, 80
- perpendicular vectors, 17
- pivot row, 78
- pivot vector, 80
- plane rotation, 261
- PLU factorization, 58, 77
- Poisson matrix, 92
- Poisson problem, 90
  - five point stencil, 91
  - nine point scheme, 100
  - Poisson matrix, 92
  - variable coefficients, 240
- Poisson problem (1D), 38
- polarization identity, 184
- positive definite, 64
- positive semidefinite, 64
- power method, 307
  - inverse, 311
  - Rayleigh quotient iteration, 312
  - shifted, 311
- preconditioned conjugate gradient method, 215
- preconditioning, 237
- preconditioning matrix, 196
- principal minor, 57
- principal submatrix, 56
- principal vectors, 137
- pseudo inverse, 274
- QR algorithm
  - implicit shift, 317
  - Rayleigh quotient shift, 317
  - shifted, 317
  - Wilkinson shift, 317
- QR decomposition, 257
- QR factorization, 257
- quadratic form, 64
- quadruple precision, 339
- rate of convergence, 204
- Rayleigh quotient, 129
  - generalized, 143
- Rayleigh quotient iteration, 312
- rectangular diagonal matrix, 147
- relative error, 174, 335
- residual vector, 175
- Richardson's method, 191
- right triangular, 54



- rotation in the  $i, j$ -plane, 261
- rounding unit, 339
- rounding-error analysis
  - adjusted rounding unit, 342
  - backward error, 340
  - backward stable, 340
  - condition number, 342
  - ill-conditioned, 343
- row operations, 332
- row space, 11
  
- scalar product, 14
- scaled partial pivoting, 79
- Schur factors, 124
- Schur form, real, 143
- second derivative matrix, 39
- semi-Cholesky factorization, 74
- Sherman-Morrison formula, 24
- shifted power method, 311
- similar matrices, 122
- similarity transformation, 122
- single precision, 339
- singular value decomposition, 147
- singular value factorization, 151
- singular values
  - Courant-Fischer theorem, 161
  - error analysis, 283
  - Hoffman-Wielandt theorem, 162
- singular vectors, 148
- span, 6
- spectral radius, 199, 207
- spectral theorem, 129
- spectrum, 29, 119
- splitting matrices for  $R$ ,  $J$ , and  $SOR$ , 197
- splitting matrix, 196
- standard inner product in  $\mathbb{C}^n$ , 123
- steepest descent, 218
- stencil, 91
- sum of subspaces, 271
  
- sums of integers, 60
- SVD, 147
- SVF, 151
- Sylvester's inertia theorem, 301
- symmetric positive semidefinite, 65
  
- trace, 120
- triangle inequality, 12
- triangular matrix
  - left triangular, 54
  - right triangular, 54
  - unit triangular, 55
- trivial subspace, 9
- two point boundary value problem, 38
  
- unit triangular, 55
- unit vectors, 3
- unitary matrix, 123
- upper trapezoidal matrix, 253
  
- vector
  - angle, 17
  - linearly dependent, 6
  - linearly independent, 6
  - nontrivial subspaces, 9
  - orthogonal, 17
  - orthonormal, 17
- vector norm, 11
- vector space
  - basis, 7
  - change of basis matrix, 10
  - complementary, 9
  - complex inner product space, 15
  - dimension, 8
  - dimension formula for sums of subspaces, 9
  - direct sum, 9
  - enlarging vectors to a basis, 8
  - examples of subspaces, 8

- existence of basis, 8
- intersection, 9
- normed, 12
- orthogonal vectors, 17
- real, 4
- real inner product space, 14
- span, 7
- subspace, 8
- sum, 9
- union, 9
- vectorization, 91

weights, 268