

INF5300

Linear feature transforms

- Linear feature transforms
- Principal component analysis (PCA)
- Fisher's linear discriminant analysis

Curriculum: See links to pdfs on course page.

2015.02.04

INF 5300

1

Linear feature transforms

- We create new features by computing linear combinations of our existing features, x_0, x_1, \dots, x_{n-1} :

$$y_1 = \sum_{i=0}^{n-1} a_{i1} x_i, \quad y_2 = \sum_{i=0}^{n-1} a_{i2} x_i, \quad \dots \quad y_m = \sum_{i=0}^{n-1} a_{im} x_i$$

- In matrix notation $\mathbf{y} = \mathbf{A}^T \mathbf{x}$
- If \mathbf{y} has fewer elements than \mathbf{x} , we get a feature reduction

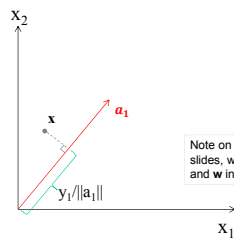
2015.02.04

INF 5300

2

Visualizing the weights in 2D/3D

$$y_1 = \sum_{i=0}^{n-1} a_{i1} x_i = \mathbf{a}_1^T \mathbf{x}$$



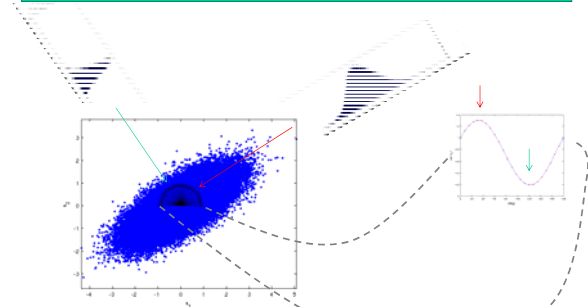
Note on naming: In the slides, we often use \mathbf{a} and \mathbf{w} interchangeably

2015.02.04

INF 5300

3

Variance of y_1



2015.02.04

INF 5300

Variance of y_1 cont.

- Assume mean of \mathbf{x} is subtracted

$$\begin{aligned}\sigma_w^2 &= \frac{1}{N} \sum_i \mathbf{w}^T \mathbf{x}_i^2 \\ &= \frac{1}{N} \sum_i (\mathbf{w}^T \mathbf{x}_i)^2 = \frac{1}{N} \sum_i \mathbf{w}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{w} = \mathbf{w}^T \left(\frac{1}{N} \sum_i \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{w} \\ &= \mathbf{w}^T \mathbf{R} \mathbf{w}\end{aligned}$$

The sample covariance matrix, \mathbf{R}

Called σ_w^2 on some slides

2015.02.04

INF 5300

y_1

5

Variance and projection residuals

Single sample

Projection onto \mathbf{w} , assuming $\|\mathbf{w}\|=1$

$$\begin{aligned}\| \mathbf{x}_i - (\mathbf{w} \cdot \mathbf{x}_i) \mathbf{w} \|^2 &= (\mathbf{x}_i - (\mathbf{w} \cdot \mathbf{x}_i) \mathbf{w}) \cdot (\mathbf{x}_i - (\mathbf{w} \cdot \mathbf{x}_i) \mathbf{w}) \\ &= \mathbf{x}_i \cdot \mathbf{x}_i - \mathbf{x}_i \cdot (\mathbf{w} \cdot \mathbf{x}_i) \mathbf{w} \\ &\quad - (\mathbf{w} \cdot \mathbf{x}_i) \mathbf{w} \cdot \mathbf{x}_i + (\mathbf{w} \cdot \mathbf{x}_i) \mathbf{w} \cdot (\mathbf{w} \cdot \mathbf{x}_i) \mathbf{w} \\ &= \|\mathbf{x}_i\|^2 - 2(\mathbf{w} \cdot \mathbf{x}_i)^2 + (\mathbf{w} \cdot \mathbf{x}_i)^2 \mathbf{w} \cdot \mathbf{w} \\ &= \mathbf{x}_i \cdot \mathbf{x}_i - (\mathbf{w} \cdot \mathbf{x}_i)^2\end{aligned}$$

All n samples (not dimensions)

$$MSE(\hat{\mathbf{w}}) = \frac{1}{n} \left(\sum_{i=1}^n \|\mathbf{x}_i\|^2 - \sum_{i=1}^n (\mathbf{w} \cdot \mathbf{x}_i)^2 \right)$$

Note: Max variance \leftrightarrow min projection residuals!

Indie of \mathbf{w}

σ_w^2

$\mathbf{w} \cdot \mathbf{w} = 1$

2015.02.04

INF 5300

6

Criterion function

- Goal: Find transform minimizing representation error
- We start with a single weight-vector, \mathbf{w} , giving us a single feature, y_1
- Let $J(\mathbf{w}) = \sigma_w^2$
- Now, let's find $\max_{\mathbf{w}} J(\mathbf{w})$ s.t. $\|\mathbf{w}\|=1$

As we learned on the previous slide, maximizing this is equivalent to minimizing representation error

2015.02.04

INF 5300

7

Maximizing variance of y_1

$$\begin{aligned}\mathcal{L}(\mathbf{w}, \lambda) &\equiv \sigma_w^2 - \lambda(\mathbf{w}^T \mathbf{w} - 1) \\ \frac{\partial \mathcal{L}}{\partial \lambda} &= \mathbf{w}^T \mathbf{w} - 1 \\ \frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= 2\mathbf{R}\mathbf{w} - 2\lambda\mathbf{w}\end{aligned}$$

Equating zero

$$\begin{aligned}\mathbf{w}^T \mathbf{w} &= 1 \\ \mathbf{R}\mathbf{w} &= \lambda\mathbf{w}\end{aligned}$$

The maximizing \mathbf{w} is an eigenvector of \mathbf{R} !

And $\sigma_w^2 = \lambda$! [Why?]

Lagrangian function for maximizing σ_w^2 with the constraint $\mathbf{w}^T \mathbf{w} = 1$

Unfamiliar with Lagrangian multipliers? See link to short note on the course's web page - very useful!

2015.02.04

INF 5300

8

Eigenvectors of covariance matrices

Real-valued, symmetric, «n-dimensional» covariance matrix

$$\mathbf{R} = \lambda_1 \mathbf{a}_1 \mathbf{a}_1' + \lambda_2 \mathbf{a}_2 \mathbf{a}_2' + \dots + \lambda_n \mathbf{a}_n \mathbf{a}_n'$$

Eigenvalue (let's say largest) Eigenvector corresponding to λ_1 Smallest eigenvalue

$\mathbf{a}_i' \mathbf{a}_j = 0$ for $i \neq j$

Remember: $\lambda_i = \text{var of } \mathbf{x}' \mathbf{a}_i$

2015.02.04

INF 5300

9

Variance of multiple variables

$$\sigma_{y_1+y_2}^2 = \frac{1}{N} \sum_i (\mathbf{w}_1' \mathbf{x}_i + \mathbf{w}_2' \mathbf{x}_i)^2 = \dots = \mathbf{w}_1' \mathbf{R} \mathbf{w}_1 + \mathbf{w}_2' \mathbf{R} \mathbf{w}_2 + 2 \mathbf{w}_1' \mathbf{R} \mathbf{w}_2$$

That is, on the previous slide $\mathbf{a}_i' \mathbf{R} \mathbf{a}_j = 0$ for $i \neq j$

=0 if y_1 and y_2 are uncorrelated, e.g. if \mathbf{w}_1 and \mathbf{w}_2 are eigenvectors of \mathbf{R}

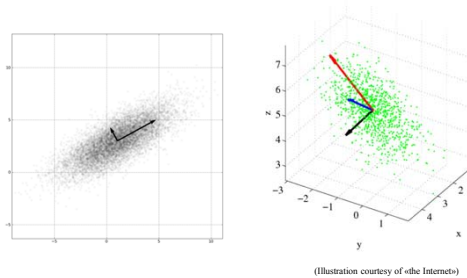
- If the weight-vectors yield uncorrelated features, their combined variance is the sum of each one's
- If \mathbf{w}_1 is the principle eigenvector of \mathbf{R} , which \mathbf{w}_2 giving an uncorrelated feature would you choose to maximize $\sigma_{y_1+y_2}^2$?
- Say \mathbf{w}_1 and \mathbf{w}_2 are the two principle eigenvectors of \mathbf{R} on the previous slide; what ratio of the total variance would they have?

2015.02.04

INF 5300

10

Example of distributions and eigenvectors



2015.02.04

INF 5300

11

Principal component transform (PCA)

- Place the m «principle» eigenvectors (the ones with the largest eigenvalues) along the columns of \mathbf{A}
- Then the transform $\mathbf{y} = \mathbf{A}^T \mathbf{x}$ gives you the m first principle components
- The m -dimensional \mathbf{y}
 - have uncorrelated elements
 - retains as much variance as possible
 - gives the best (in the mean-square sense) description of the original data (through the «image»/projection/reconstruction $\mathbf{A} \mathbf{y}$)

Note: The eigenvectors themselves can often give interesting information

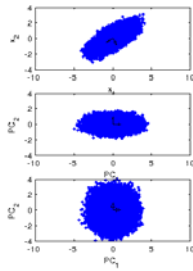
PCA is also known as Karhunen-Loeve transform

2015.02.04

INF 5300

12

PCA and rotation and «whitening»



If we use all eigenvectors in the transform, $y = Ax$, we simply rotate our data so that our new features are **uncorrelated**, i.e., $\text{cov}(y)$ is a diagonal matrix.

If we as a next step scale each feature by their σ , $y = D^{-1/2}Ax$, where D is a diagonal matrix of eigenvalues (i.e., variances), we get $\text{cov}(y)=I$. We say that we have **«whitened»** the data.

2015.02.04

INF 5300

13

PCA and multiband images

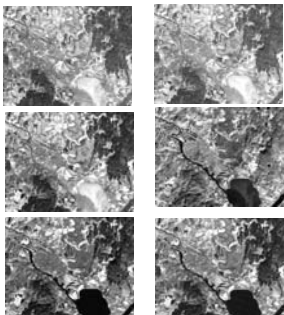
- We can compute the principal component transform for an image with n bands
- Let X be an $M \times n$ matrix having a row for each image sample
- Sample covariance matrix $R = \frac{1}{N} X^T X$
- Place the (sorted) eigenvectors along the columns of A
- $Y=XA$ will then contain the image samples, however most of the variance is in the «bands» with the lowest index (corresponding to the largest eigenvalues), and the new features are uncorrelated

2015.02.04

INF 5300

14

PCA example – original image



- Satellite image from Kjeller
- 6 spectral bands with different wavelengths

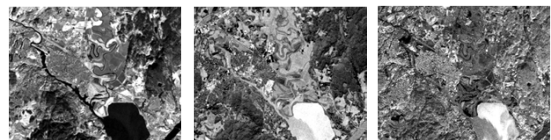
1	Blue	0.45-0.52	Max. penetration of water
2	Green	0.52-0.60	Vegetation and chlorophyll
3	Red	0.63-0.69	Vegetation type
4	Near-IR	0.76-0.90	Biomass
5	Mid-IR	1.55-1.75	Moisture/water content in vegetation/soil
7	Mid-IR	2.08-2.35	Minerals

2015.02.04

INF 5300

15

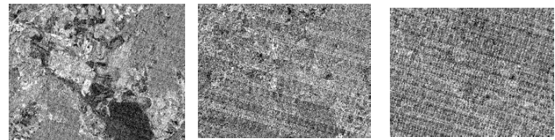
Example cont: Principal component images



Principal component 1

Principal component 2

Principal component 3



Principal component 4

Principal component 5

Principal component 6

2015.02.04

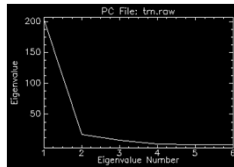
INF 5300

16

Example cont: Inspecting the eigenvalues

The mean-square representation error we get with m of the N PCA-components is given as

$$E[\|x - \hat{x}\|^2] = \sum_{i=1}^{N-1} \lambda_i - \sum_{i=1}^m \lambda_i = \sum_{i=m}^{N-1} \lambda_i$$



Plotting λ_i will give indications on how many features are needed for representation

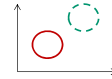
2015.02.04

INF 5300

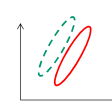
17

PCA and classification

- Reduce overfitting by detecting directions/components without any/very little variance
- Sometimes high variation means useful features for classification:



- .. and sometimes not:

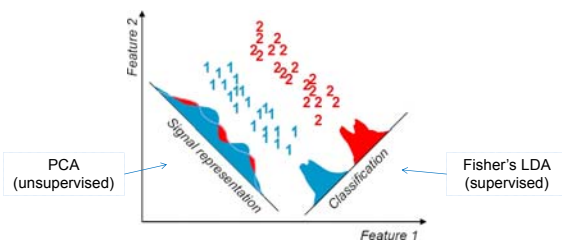


2015.02.04

INF 5300

18

Intro to Fisher's linear discriminant



2015.02.04

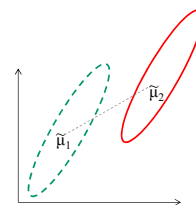
INF 5300

19

Criterion function - a first attempt

- To find a good projection vector, we need to define a measure of separation between the projections. This will be the criterion function $J(w)$
- A naive choice would be projected mean difference, $J(w) = |\tilde{\mu}_1 - \tilde{\mu}_2|^2$ s.t. $|w|=1$.

This criterion does not consider variance



w simply becomes a scaled difference in means $(\mu_1 - \mu_2)$

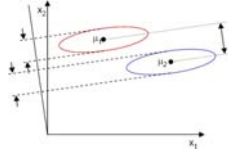
2015.02.04

INF 5300

20

A criterion function including variance

- Fisher's solution: Maximize a function that represents the difference between the means, scaled by a measure of the within-class scatter
- Define classwise scatter (scaled variance) $\tilde{s}_i^2 = \sum_{y \in \omega_i} (y - \tilde{\mu}_i)^2$
- $\tilde{s}_1^2 + \tilde{s}_2^2$ is *within class scatter*
- Fisher's criterion is then
$$J(\mathbf{w}) = \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$
- We look for a projection where examples from the same class are close to each other, while at the same time projected mean values are as far apart as possible



2015.02.04

INF 5300

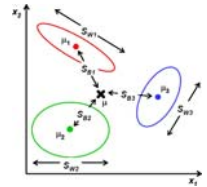
21

Scatter matrices – M classes

- Within-class scatter matrix:
$$S_w = \sum_{i=1}^M P(\omega_i) S_i$$

$$S_i = E[(x - \mu_i)(x - \mu_i)^T]$$
 Weighted average of each class' sample covariance matrix
- Between-class scatter matrix:
$$S_b = \sum_{i=1}^M P(\omega_i) (\mu_i - \mu)(\mu_i - \mu)^T$$

$$\mu = \sum_{i=1}^M \mu_i$$
 Sample covariance matrix for the means



Fisher criterion in terms of within-class and between-class scatter matrices:

$$J(\mathbf{w}) = \frac{\mathbf{w}^T S_b \mathbf{w}}{\mathbf{w}^T S_w \mathbf{w}}$$

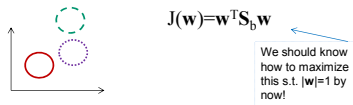
2015.02.04

INF 5300

22

Multiple classes, $S_w = \sigma^2 I$

- If $S_w = \sigma^2 I$, we can fix $\|\mathbf{w}\|=1$ and make the denominator in $J(\mathbf{w})$ independent of $\mathbf{w} \rightarrow J(\mathbf{w})$ guided by the spread of the means (S_b) only:



- Weight-vector giving maximum separability is given by principal eigenvector of S_b
 - Second best (and orthogonal to first) by next-to-principal
 - ... etc. for higher dimensional settings
 - ... until a maximum of M-1 dimensions (number of classes minus one) [If classes are «isotropically» Gaussian distributed, all discriminatory information is in this subspace!]

2015.02.04

INF 5300

23

General S_w I/II

- We saw that $S_w = I$ gave Fisher criterion independent of S_w , and only dependent on S_b
- We can get there by «whitening» the data before applying the Fisher criterion
 - Whitening data by rotation and scaling \rightarrow No general loss as distribution overlap does not change
- We must find $\mathbf{y} = \mathbf{A}^T \mathbf{x}$ that yields $S_{wy} = I$
 - We have seen that PCA gives uncorrelated data, per-feature scaling can give unit variance per feature:
 - $\mathbf{y} = \mathbf{D}^{-1/2} \mathbf{A}^T \mathbf{x}$, where \mathbf{A} has eigenvectors of S_w as columns, and \mathbf{D} is a diagonal matrix with corresponding eigenvalues

$$S_{wy} = \frac{1}{N} \sum_i (\mathbf{D}^{-1/2} \mathbf{A}^T \mathbf{x}_i) (\mathbf{D}^{-1/2} \mathbf{A}^T \mathbf{x}_i)^T = \mathbf{D}^{-1/2} \mathbf{A}^T S_w \mathbf{A} \mathbf{D}^{-1/2} = \mathbf{D}^{-1/2} \mathbf{D} \mathbf{D}^{-1/2} = \mathbf{I}$$

2015.02.04

INF 5300

24

General S_w II/II

- Let $B = D^{-1/2}A^T$ (the whitening transform)
- S_b becomes after whitening step:

$$S_{by} = BS_bB^T$$
- Ignoring the denominator (which is now independent of w), we get
 - $J_y(w) = w^T S_{by} w = w^T BS_bB^T w$, s.t. $|w|=1$
- The weight-vectors, w^* , maximizing separation are now given by the principal eigenvectors of BS_bB^T (in the whitened space)
- In the original space, $w = B^T w^* = AD^{-1/2} w^*$

Set $J_y(w^*) = J(w)$ to see this

2015.02.04

INF 5300

25

Solving Fisher more directly

- Alternatively, you can notice that

$$J(w) = \frac{w^T S_b w}{w^T S_w w}$$

- .. is a «generalized Rayleigh quotient» and look up the solution for its maximum, which is the principal eigenvector of

$$S_w^{-1} S_b$$

- The following solutions (orthogonal in S_w , i.e., $w_i^T S_w w_j = 0$, for $i \neq j$) are the next principal eigenvectors

Note that the obtained w s are identical (up to scaling) to those from the two-step procedure from the previous slides

2015.02.04

INF 5300

26

Comments on Fisher's discriminant

- In general, projection of the original feature vector to a lower dimensional space is associated with some loss of information
 - Keeping all $M-1$ dimensions gives you no reduction in classification performance for a Gaussian classifier with equal class-covariance matrices (LDA)
- Although the projection is optimal with respect to the given criterion ($J(w)$), the criterion itself might not be suitable for a given data set / classifier
- Minimizing $J(w)$ is not equivalent to minimizing the classification error

2015.02.04

INF 5300

27

Limitations of Fisher's discriminant

- It produces at most $M-1$ feature projections
- Its criterion function is based on class-wise distributions of limited complexity (all classes have a similarly-shaped Gaussian distribution)
- It will fail when the discriminatory information is not in the mean but in the variance of the data

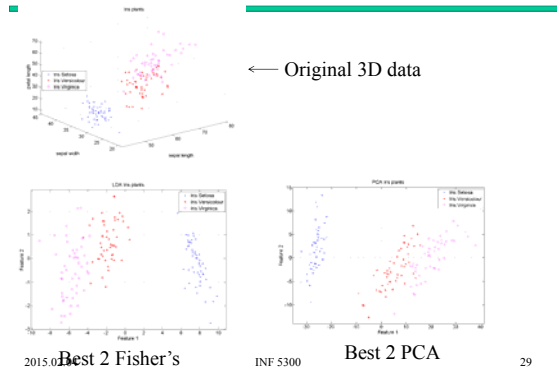


2015.02.04

INF 5300

28

Fisher's discriminant example



2015.02.04

INF 5300

29

Summary

- PCA (unsupervised)
 - Max variance \leftrightarrow min projection error
 - Eigenvectors of sample cov.mat. / scatter matrix
- Fisher's linear discriminant (supervised)
 - Maximizes spread of means while minimizing intra-class spread
 - $S_{wy} = I$ and «whitening of data»
 - Eigenvectors of $S_w^{-1}S_b$
 - At most $n_{Classes}-1$ features
 - Limitations

2015.02.04

INF 5300

30

Literature on pattern recognition

- A review on statistical pattern recognition (still good fourteen years on):
 - A. Jain, R. Duin and J. Mao: Statistical pattern recognition: a review, IEEE Trans. Pattern analysis and Machine Intelligence, vol. 22, no. 1, January 2001, pp. 4–
- Classical PR-books
 - R. Duda, P. Hart and D. Stork, Pattern Classification, 2. ed. Wiley, 2001
 - B. Ripley, Pattern Recognition and Neural Networks, Cambridge Press, 1996.
 - S. Theodoridis and K. Koutroumbas, Pattern Recognition, Academic Press, 2006.

2015.02.04

INF 5300

31