



IN3120/4120 - Search technology

Group 2, Second session
TA: Markus Sverdvik Heiervang



Today's plan

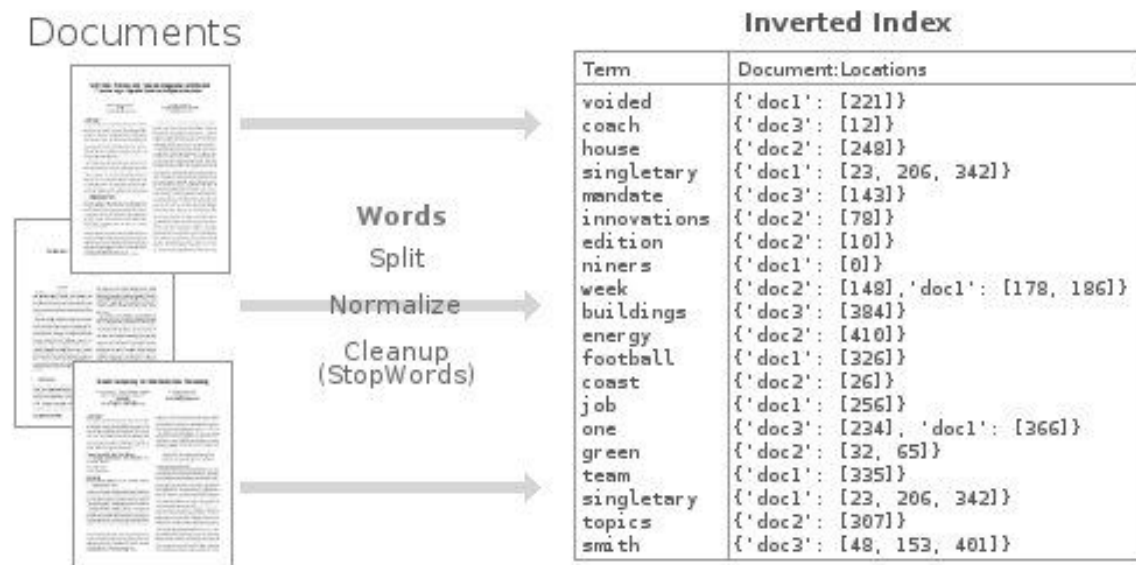
- In depth coverage of inverted indices
 - Inverted index walkthrough in jupyter notebook
 - Merging of postings
 - Index compression
-
- Q&A, eventually topics you want me to cover



Inverted index: what is it?

- A data structure which lets you effectively find every document of a corpus that contains a given term
- Posting: typically a collection containing document id, and optionally additional information such as term frequency, term position in the document, etc.

Inverted index: what is it?



Inverted index as opposed to forward index

docID			geo-scopeID	
1			Europe	
2			Europe	
3			France	
4			England	
5			Portugal	
6			Quebec	
7			Europe	
8			Spain	

Forward Index

geo-scopeID			docID		
Europe			1	2	7
France			3		
Portugal			5		
England			4		
Quebec			6		
Spain			8		

Inverted Index



Examples and how to build

Postings merging

```
INTERSECT( $p_1, p_2$ )
1   $answer \leftarrow \langle \rangle$ 
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4      then  $\text{ADD}(answer, \text{docID}(p_1))$ 
5           $p_1 \leftarrow \text{next}(p_1)$ 
6           $p_2 \leftarrow \text{next}(p_2)$ 
7      else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
8          then  $p_1 \leftarrow \text{next}(p_1)$ 
9          else  $p_2 \leftarrow \text{next}(p_2)$ 
10 return  $answer$ 
```

Necessary when
searching for documents
containing both term a
and b, using the query:
a AND b

Example:
You've heard rumors that your favorite
rapper Tupac is dead. To find out if it is true
or not, you search using an intersectional
query:
Tupac AND Dead

Figure 2: pseudocode from the book

Postings merging

Problem: what if you want an intersection of $n > 2$ terms?

Solution:

INTERSECT($\langle t_1, \dots, t_n \rangle$)

```
1  terms  $\leftarrow$  SORTBYINCREASINGFREQUENCY( $\langle t_1, \dots, t_n \rangle$ )
2  result  $\leftarrow$  postings(first(terms))
3  terms  $\leftarrow$  rest(terms)
4  while terms  $\neq$  NIL and result  $\neq$  NIL
5  do result  $\leftarrow$  INTERSECT(result, postings(first(terms)))
6    terms  $\leftarrow$  rest(terms)
7  return result
```

the query:

a and b and c and d

will be evaluated as such:

((a and b) and c) and d)

Figure 3: pseudocode from the book

Index compression



- When working with big data and limited memory, compression can be very useful
- Typical when serializing an inverted index (i.e. storing the inverted index on a disk)

Say we have a posting list for a given term t , containing only document ids, in an index of a relatively large corpus. There are many words that occur rarely, such as t

$\text{Postings}(t) = [7, 40, 7000, 43021, 140236, 2773002]$

Large numbers take up much space: $7 == 111$, $2773002 == 1010100101000000001010$.

The idea is to keep the first id, and then only keep the differences compressed

$\text{Differences}(\text{Postings}(t)) = [7-0, 40-7, 7000-40, 43021-7000, \dots]$

Then we apply a number-compression algorithm on these differences. You will learn more in the lecture about compression



What more do we need for a search engine?

- A ranker
- A filter for cheats
- A query processor
 - e.g. if the query is “Tupac dead”
 - search (Tupac AND dead), then (Tupac OR dead)
 - NLP/NLU
- etc ...



Tomorrows lecture

In tomorrows lecture, Aleksander will cover most of the string algorithms in this course

important algorithms for assignment B:

- Suffix arrays
- Tries and aho-corasick algorithm

I will go further in depth on some of these topics on the next group session