

## Information Retrieval

Books

(Available on the web)

Information Retrieval, von Rijsbergen Univ. of Glasgow

x Information Retrieval : Algorithms &amp; Data Structure

Ch 8?

W. Frakes &amp; R. Baeza-Yates

Modern Information Retrieval : R. Baeza-Yates

N. - ?

⇒ acm.org / dl

Conf. ACM SIGIR  
Info Ret.ACM SIGKM  
Information & Knowledge  
management

ACM Years on Information Systems

" " " Database Systems

Journal of the American Society for Information Science & Technology  
(JASIST)

Information Processing &amp; Management (IPM)

New Event Detection &amp; Tracking

⇒ Efficient &amp; Effective

+  
in terms of time & space

collect

↳ 400 000 diff. news → determine new events

2007 → 9 new resources ⇒ thirty million --

(Collection)

Midterm ⇒ 25 March 2008 ⇒ Tuesday

Final ⇒ 21<sup>st</sup> May - 31<sup>st</sup> Mayad hoc queries ⇒ word .. word ..  
queries ↗they do similarity comparison  
rank from most similar to the less similar

abundance problem

calculate similarity between queries

Efficient  
effective system

→ relevant documents at the top

2

## System Evaluation:

### Recall & Precision

TREC = Text Retrieval Conference  
 trec eval package  
 provides performance

bpref : binary preference

there is a part of data which has not been evaluated by human being & assumed as irrelevant.

## Clustering & Cluster Validation:

grouping.

$tf \cdot idf$   
 term frequency inverse document frequency

a word appears in all of the documents  $\rightarrow$  binary  
 if a term appears smaller number of times  $\rightarrow$  can be used for differentiation

( fundamental file structures )

information  $\rightarrow$   $d_{10}, d_{11}, d_{30}, d_5$

posting list

retrieval  $\rightarrow d_1, d_{10}, d_{51}$

## Query Processing:

ACM Computing Surveys, Alistair Moffat, Justin Zobel

Inverted Files for text search engines

### Signature Files:

For each document we have a bit array referred to as signature:

Document	Information	Retrieval	Performance	evaluation
	[ ] [ ] [ ] [ ] [ ]	[ ] [ ] [ ] [ ] [ ]	[ ] [ ] [ ] [ ] [ ]	[ ] [ ] [ ] [ ] [ ]

f : size of the array  
 → set random locations 1  
 y not always random  
 Random Number Generator  
 with 3 seed

then we superimpose this  $\Rightarrow$  order this column wise  $\Rightarrow$  0111

Query signature.  
 Document "

if ( $O_1 \& O_2$ ) =  $O_3$

then document is retrieved.

{ if (0's & 1's) = 0's  
then document is retrieved }

assume you're dealing with huge doc

↳ false drop resolution

↳ so has to make this more efficient

### N-gram

1-gram  
bigram  
trigram

Information: in, int, fo, or, on, m, ma, at, ti, is, on

in \* on → word begins with in & end with on.

so \* m \* t

PAT TREES: (Patricia Trees)

cm  
snider

(Suffix Trees)

trie

good for  
intelligence  
purpose

similar to priority queue.  
record everything in telephone conv  
If you construct pat tree  
you can get anything  
any string

for movies → part tree structure  
you can easily find it (Not so easy to construct)

### Compression:

Assume we have posting list

posting list  
information → [10, 15, 2, 5, 17, 22, 25, 32, 36, 42, 47] → the #s of the doc's that keyword appears  
retrieval → [10, 15, 22, 30, 47] (skip some)  
→ [27, 32, 42] → then

But the numbers will increase → the digits will ↑



of gap content  
↳ document  
⇒ by doing so  
we make the #s  
↳ smaller

$$10 + \text{gap} = 15$$

$$15 + \text{gap} = 17$$

say for some reason we begin with 27

with step

first begins with 27

with  
& the stopping (information) ⇒ we know we can  
skip first 5  
of them since

$$25 < 27$$

Self-Indexing  
Morphology & Label

ACM, TOIS, 1986

④

## INFORMATION RETRIEVAL SYSTEMS OVERVIEW

### Information Systems ?

DBms      IR  
structured    unstructured  
data          data  
Precise Answers      ~ answers

### Traditional IR Systems ?

MEDLARS provided by National Institutes of Health

MEDLINE online version of MEDLARS

LEXIS/NEXIS law

STAIRS provided by IBM

1993-1994  
Altavista

Google

SMART

research IR system

Cornell Univ.  
Gerald Salton

### Components of An IR System

User Interface

Indexing System

Storage System (the file system that we want to use)

Query Processing System

Zobel & Moffat

Inverted index for text search engines

### An Indexing Example:

Doc 1 : Information Retrieval by Parallel Document Ranking

Doc 2 : An analysis of parallel text retrieval systems

Doc 3 : Information retrieval in the law office : an overview

Stopwords : Frequent words (not good for distinguishing)

documents from each other

stopword list

~200

{the  
and  
a  
in}

Cornell

SIGIR forum (christ for)

say we use  
Doc 3

Stopword list:

an  
by  
in  
at  
overview  
systems  
the

Indexing Text  
Information  
Retrieval  
law  
office

\* uncontrolled environment

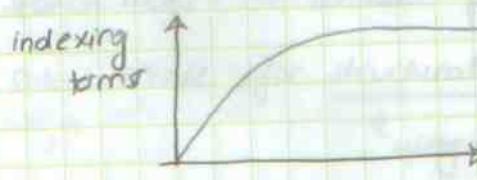
: include any word that is in the document or long as it doesn't appear in the stopword list

### \* Controlled Indexing

the indexing terms will be provided before hand

Moth : Medical search Hierarchical ?

✓ after a while any word encountered may have been used before



Indexing Terms in Alphabetical Order:

- 1) analysis
- 2) document
- 3) information
- 4) law
- 5) office
- 6) parallel
- 7) ranking
- 8) retrieval
- 9) tax

$n = 9$  (no of words)

$m = 3$  (no of Documents)

(full of 0's  
\* of 1's are 1's)

(sparse matrix)  
binary matrix

$$D = \begin{bmatrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 & t_7 & t_8 & t_9 \\ d_1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 1 & 0 \\ d_2 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 \\ d_3 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$n \times m$

$9 \times 3$

$$d_{ij} = \begin{cases} 1 & , \text{ if } t_j \text{ appears in } d_i \\ 0 & , \text{ otherwise} \end{cases}$$

SNSPEC database (online)

Entry: David Lewis

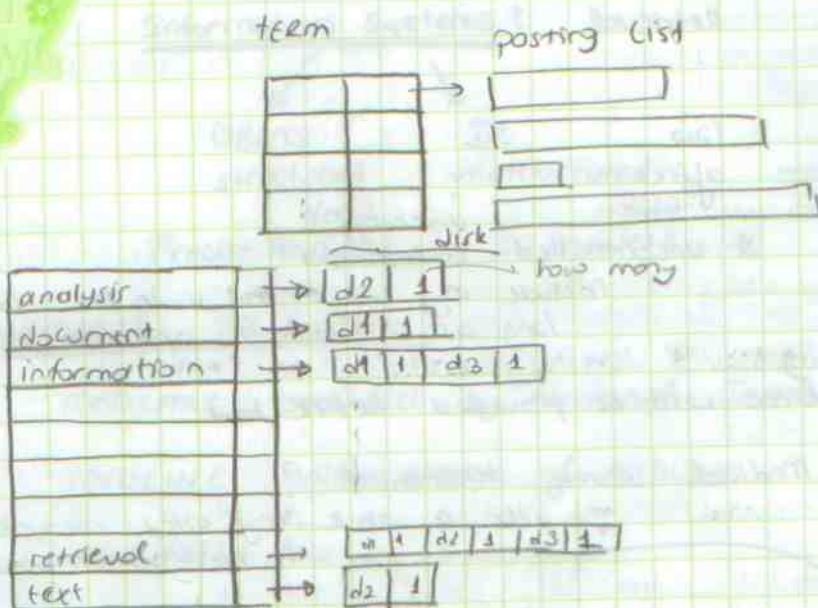
$m = 12,684$   
 $n = 14573$

D matrix of SNSPEC : density of 1's  $\approx 2\%$

ACM,  
time  
MPL

6

## Inverted Index for the Example Collection:



Assume a Boolean Query Environment:

(Q : information & retrieval)

$$(d_1 \ d_2) \cap (d_1 \ d_2 \ d_3) = (d_1 \ d_2)$$

19.10.2008

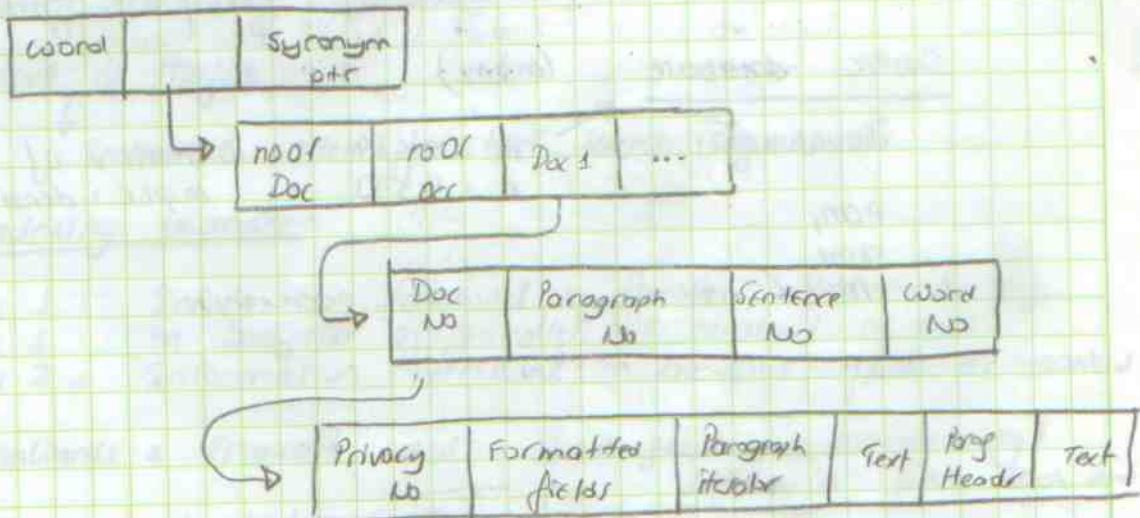
### Example IR Systems:

STAIRS : Storage Information Retrieval Systems

Designed for mainframes, IBM

File System / Structure

#### Dictionary



## Formatted Field:

Author Name

Journal Name

Page No

STAIRS has two programs:

1. utility programs : for database initialization & maintenance
2. Query & Retrieval Utility Systems (AQUARIUS)

## Modes Of Operations:

Search Mode: for textual IR

Select Mode: for structural IR

using formatted fields

## Queries:

HEART

HEART or DISEASE

HEART & DISEASE \$3

WITH

SAME

in the  
same sentence  
in the same  
paragraph

- Find matching documents using a boolean query
- Rank matching documents according to their significance

Value of a query term: flabberg

a: the frequency of a word in the document

b:

retrieved set

c: the no. of documents in the retrieved set in which the term occurs

$$\text{Query Term Value} : \frac{a \times b}{c}$$

$$\text{Score of a document} = \sum \text{value of all query terms which appear in this document}$$

## Evaluation of IR Systems:

users point of view

implementors point of view ) ← System Resources

► effectiveness

efficiency : response time

system resource requirement X

- coverage

- links should be alive

- ease of use

## How to measure effectiveness?

TREC : Text Retrieval Conference [1992, ...]

[www.nist.gov/trec/](http://www.nist.gov/trec/) ? (twik)  $\downarrow$

$\hookrightarrow$  national institute standards of technology

trec-eval package

TREC 3

Appendix A

### Precision (P)

✓
X
X
X
V

$$\frac{2}{5} = 0.40$$

p10  
p20

precision @ 10 or 20

first one // two pages

### Recall:

$$\frac{\text{* of retrieved & relevant documents}}{\text{(total * of relevant docs in the collection)}}$$

$$\frac{2}{20}$$

$\rightarrow$  total \* of rel. doc

### Test Collection in IR

A set of documents

A set of queries & their relevant documents  $\rightarrow 50$

Standard test collections facilitate repetition of tests

in statistical sense 50 i)

(if you get 50 relevant with)

Finding relevant documents for queries

TREC uses the pooling concept

Retrieve top 100 documents & identify the relevant ones for query

100 docs

System 1

0

System 2

0

System 10

if all distinct = 100 docs

not all distinct  $\approx$  700 docs

all other docs

they will be assumed as irrelevant

How reliable?

NIST SIGIR 1998

Justin Zobel

Documents which are not seen by the evaluators (annotators) are assumed as irrelevant.  
Actually some of them can be relevant.

A new measure (Acm SIGIR Conf. 2004, by Buckley & Voorhees)

bpref: binary preference

MAP: bpref is as reliable as MAP

$\bar{P}$  (mean Average Precision)

✓	P <sub>1</sub>
X	
X	
✓	P <sub>2</sub>
X	
X	
✓	P <sub>3</sub>

$$\text{MAP} = \frac{P_1 + P_2 + P_3 + 0}{4}$$

All together there are 4 relevant docs

top 1000

Example for Recall & Precision

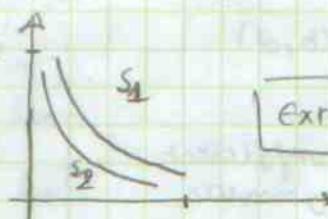
Duyarltik  
Anmol

Rank	1	2	3	4	5	6	7	8	9	10
Relevance	0	1	0	1	1	1	1	0	0	0
Precision	0/1	1/2	1/3	2/4	3/5	4/6	5/7	6/8	7/9	8/10
Recall	0/10	1/10	1/10	2/10	3/10	4/10	5/10	6/10	7/10	8/10

Total no of relevant docs = 10

more relevant precision ↗

S<sub>1</sub> is better than S<sub>2</sub>



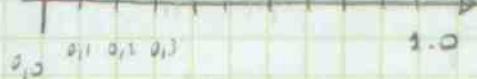
Extrapolation ↗ to prevent big gaps

Makes the curve nice...

TREC



11 point



t-test → available

student t-test

2-tail

1-tail

not only looks at avg  
compares the individual  
values  
gives a "p" value

→ p = 0,05

If p value is p < 0,05  
significant

Is the difference statistically significant?

Q.No	MAP <sub>1</sub>	MAP <sub>2</sub>
1		
2		
3		
40		

0.92 0.55

10

## Similarity Calculation:

MOTIVATION :- For ranking documents according to their similarity to submitted query

browsing

- Cluster documents according to similarity to each other then use these clusters to find additional relevant documents like your favorite document

\* There are several similarity coefficients

Most of them is symmetric  $\Rightarrow s(a,b) = s(b,a)$

Van Rijsbergen, Information Retrieval, Univ. of Glasgow

### Similarity Coefficient

Binary

Weighted

- \* Dot product
- \* Cosine

$$\cos \theta$$

if they overlap  $\Rightarrow \cos \theta = 1$

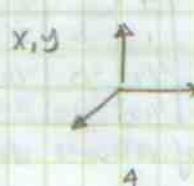
- \* Dice Coeff

- \* Jaccard

$$s(a,b) = s(b,a)$$

~~dx~~

over coefficient  
asymmetr.



$$D = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_m \end{bmatrix}$$

n-dimensions

### Similarity Coefficient

Inner Product  
(Dot Product)

### Binary

$$x \cap y$$

### Weighted

$$\sum x_i y_i$$

Cosine

$$\frac{|x \cap y|}{|x|^{1/2} |y|^{1/2}}$$

$$\frac{\sum x_i y_i}{\sqrt{\sum x_i^2 \sum y_i^2}}$$

Dice

$$\frac{2|x \cap y|}{|x| + |y|}$$

$$\frac{2 \sum x_i y_i}{\sum x_i^2 + \sum y_i^2}$$

Jaccard

$$\frac{|x \cap y|}{(|x| + |y|) - |x \cup y|}$$

$$\frac{\sum x_i y_i}{\sum x_i^2 + \sum y_i^2 - \sum x_i y_i}$$

$$X = (1 \ 0 \ 1 \ 1 \ 1)$$

$$|X| = 4$$

$$Y = (1 \ 1 \ 0 \ 1 \ 0)$$

$$|Y| = 5$$

Inner Product = 2

$$\text{Dice} = \frac{2 \cdot 2}{4+3} = \frac{4}{7}$$

Jaccard

$$\frac{2}{4+3-2} = \frac{2}{5}$$

$$\text{Cosine} = \frac{2}{\sqrt{4} \sqrt{5}}$$

$\Rightarrow$  tf idf (we assign higher values to words which appear more frequently in documents  $\Rightarrow$  like stopwords)

+  
inverse document frequency

Salton Buckley Term Weighting Approaches Information processing and management. Pg 61

$$X = (2 \ 0 \ 1 \ 3 \ 2)$$

$$Y = (1 \ 0 \ 2 \ 1 \ 5)$$

$$\text{Dice} = \frac{2(2+0+2+3+1)}{(4+0+1+9+4) + (1+0+4+1+25)} = 0,69$$

$$\text{Cosine} = \frac{17}{(18, 31)^{1/2}} = \frac{17}{23,6} \approx 0,72$$

Clustery

## # How To CALCULATE SIMILARITY AMONG DOCUMENTS #

$$D = d_1 \begin{bmatrix} t_1 & t_2 & \dots & t_n \end{bmatrix}$$

$$S_{ij} = S_{ji}$$

Brute force approach  
1) straight forward approach

$$S = \begin{bmatrix} 1 & 2 & 3 & \dots & m \\ 1 & S_{12} & S_{13} & \dots & S_{1m} \\ - & 1 & S_{23} & \dots & S_{2m} \\ & & 1 & S & \dots \\ m-1 & & & & S_{m-1m} \leftarrow m-1 \\ m & & & & 1 \end{bmatrix}$$

Total no of similar value to be calculated =  $1+2+\dots+(m-1) = \frac{m(m-1)}{2}$

(12)

2) Using the knowledge of term distributions in documents

$$D = \begin{bmatrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 \\ d_1 & 1 & 1 & 0 & 0 & 1 & 0 \\ d_2 & 1 & 1 & 0 & 1 & 1 & 0 \\ d_3 & 0 & 0 & 0 & 0 & 0 & 1 \\ d_4 & 0 & 0 & 1 & 0 & 0 & 1 \\ d_5 & 0 & 0 & 1 & 1 & 0 & 1 \end{bmatrix}$$

$t_1 \rightarrow d_1, d_2$      $t_4 \rightarrow d_2, d_5$   
 $t_2 \rightarrow d_1, d_2$      $t_5 \rightarrow d_1, d_2$   
 $t_3 \rightarrow d_4, d_5$      $t_6 \rightarrow d_3, d_4, d_5$

Consider  $d_1$ : $t_1 \quad t_2 \quad t_5$ 

$$\left( \begin{array}{c} d_1 \\ d_2 \end{array} \right) \cup \left( \begin{array}{c} d_1 \\ d_2 \end{array} \right) \cup \left( \begin{array}{c} d_1 \\ d_2 \end{array} \right) = \left( \begin{array}{c} d_1 \\ d_2 \end{array} \right)$$

$S_{12}$

$$\left[ \begin{array}{c} 1 & S_{12} & S_{13} & S_{14} & S_{15} \\ 1 & S_{23} & S_{24} & S_{25} & \\ 1 & S_{34} & S_{35} & & \\ 1 & S_{45} & & & \\ 1 & & & & \end{array} \right]$$

 $d_2$ :

$$\left( \begin{array}{c} d_1 \\ d_2 \end{array} \right) \cup \left( \begin{array}{c} t_1 \\ d_2 \end{array} \right) \cup \left( \begin{array}{c} t_4 \\ d_2 \end{array} \right) \cup \left( \begin{array}{c} t_5 \\ d_2 \end{array} \right) = \left( \begin{array}{c} d_1 \\ d_2 \\ d_5 \end{array} \right)$$

$S_{25}$

 $d_3$ :

$$t_6 \cup \left( \begin{array}{c} d_3 \\ d_4 \\ d_5 \end{array} \right) \Rightarrow S_{34}, S_{35}$$

 $d_4$ :

$$\left( \begin{array}{c} d_4 \\ d_5 \end{array} \right) \cup \left( \begin{array}{c} d_3 \\ d_4 \\ d_5 \end{array} \right) = \left( \begin{array}{c} d_3 \\ d_4 \\ d_5 \end{array} \right)$$

$S_{45}$

There are 10 similarity values to be calculated, but we need to calculate only 5 of them.

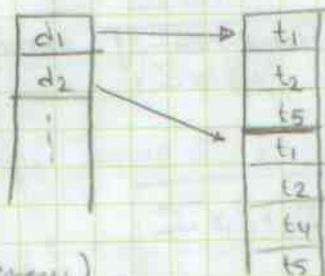
$5/10 \Rightarrow 50\% \text{ savings}$

3) Using the Inverted File for Term

$$d_1 \rightarrow [t_1 \quad t_2 \quad | \quad t_5]$$

$$d_2 \rightarrow [t_1 \quad | \quad t_2 \quad | \quad t_4 \quad | \quad t_5]$$

use pointers to first two

similarity  $\rightarrow$  (say common elements)

in order to  
get rid of  
internal  
pointer

$t_1, t_2$  are the same  $\checkmark$  increment counter  
 $t_2, t_2$  " " " "  
 $t_5, t_4$  " " " - increment pointer  
 $t_5, t_5$  " " " " counter

	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	$t_6$	docum freq
$d_1$	1	1	0	0	1	0	$t_1 \rightarrow \langle 1, 1 \rangle \langle 2, 1 \rangle$
$d_2$	1	1	0	1	1	0	$t_2 \rightarrow \langle 1, 1 \rangle \langle 2, 1 \rangle$
$d_3$	0	0	0	0	0	1	$t_5 \rightarrow \langle 1, 1 \rangle \langle 2, 1 \rangle$
$d_4$	0	0	1	0	0	1	$t_3 \rightarrow \langle 4, 1 \rangle \langle 5, 1 \rangle$
$d_5$	0	0	1	1	0	1	$t_6 \rightarrow \langle 3, 1 \rangle \langle 4, 1 \rangle \langle 5, 1 \rangle$

Dice Coef :  $\frac{21 \times 0.4}{1 \times (1 + 14)}$

$d_1 \quad d_2 \quad d_3 \quad d_4 \quad d_5$

3 4 1 2 3

document length array

$$S = \begin{bmatrix} 1 & S_{12} & S_{13} & S_{14} & S_{15} \\ -1 & S_{23} & S_{24} & S_{25} \end{bmatrix}$$

1  
2  
3  
4  
5

Consider  $d_1$ :

Similarity Array:

$S_{11} \quad S_{12} \quad S_{13} \quad S_{14} \quad S_{15}$  mail boxes

X	Ø	0	0	0
---	---	---	---	---

$t_1, t_2, t_5$

1  
2  
3

$$\frac{2 \times 3}{3+4} = \frac{6}{7}$$

Consider  $d_2$ :

$S_{21} \quad S_{22} \quad S_{23} \quad S_{24} \quad S_{25}$

X	X	0	0	Ø
---	---	---	---	---

1

$t_1, t_2, t_4, t_5$

$$\frac{2 \times 1}{3+4} = \frac{2}{7}$$

$1,1 \quad 1,1 \quad 2,1 \quad 1,1$  ignore  $(1,1)$  & increment others  
 $2,1 \quad 2,1 \quad 5,1 \quad 2,1$

XX Consider the Computation Requirements?

m

$x_d$ : depth of indexing (no of avg. terms / doc)

$t_g$ : term generality (avg. posting list length, avg. no of docs / term)

$O(m \times_d t_g)$

$t_4 \rightarrow \langle 2, 1 \rangle \langle 5, 1 \rangle \longleftrightarrow \langle 5, 1 \rangle \langle 2, 1 \rangle$

try to reduce the comparison looking at this order

in best stable rep X

Kendisiye gelme 13  
kader random looking around

\* Doing the calculations faster: organize the posting lists in reverse order: first higher numbered documents

$$x_d \cdot \frac{m-1}{m} + x_d \cdot \frac{m-2}{m} + \dots$$

$$x_d \cdot \frac{1}{m} \left( \frac{m(m-1)}{2} \right) = x_d \cdot \frac{m}{2}$$

CLUSTERING: Jain, Murty, Flynn Data Clustering  
Acm computing surveys, Sept 1995

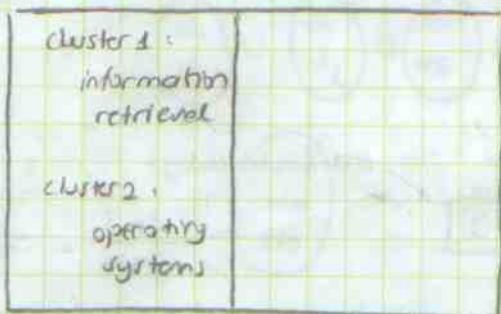
Clusty  
Sophia



Read (link on the course web page)

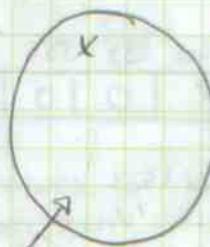
How to use clusters for IR?

1- For clustering the search results



2- For browsing:

X → interesting document



look at the other members of this document's cluster.

3- Cluster-based retrieval

Saltton & Studenti

C.T Yu

- First choose the best matching clusters

- Then match your query with the documents of these selected clusters

Clusters → cluster centroids

Documents → document vector

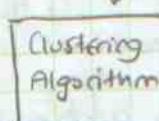
## Classification of Clustering Algorithms

- | Categorization : Groups are defined before categorization ← supervised classification
- | Clustering : Groups are undefined when we begin ← unsupervised

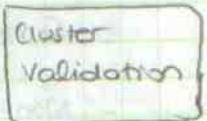
According to the generated clustering structure

According to its working principles

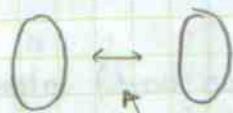
$$D = [ ] \rightarrow$$



change parameters

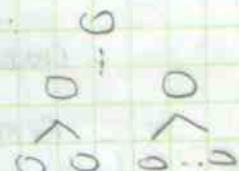


Algorithm Results



Human Results

check the consistency



- partitioning :  $C_i \cap C_j = \emptyset$ ,  $i \neq j$  - hierarchical

- overlapping :  $C_i \cap C_j \neq \emptyset$ ,  $i \neq j$

- single-pass :
- multi-pass :
- graph-theoretical
- based on user queries

omiecienski : ACM TODS, 1990

Peter J. Denning  $\Rightarrow$  (working set model) (akademik omic)

IT Professionel olige bir column'ı var  
Matematik the mass

Abstraction  $\hookrightarrow$  kavramlar ustaki ustaki (ingiliz once yazmis)

Flow  $\Rightarrow$  sevgilimiz birisi yaparken nasil de olsun gecisler joran  
mitaly cikis+minali

Communication  
of April  
ACM 2007

## \* cluster hypothesis

(Van Rijsbergen), (1972)

Functional Description: what ← easy

Operational Description: how ← difficult

Paul Erdős  
Jain & Duben  
Clustering Algorithm  
1988

### 1 \* Single-Pass Algorithms:

#### a) SEED oriented

- Some documents are selected as cluster initiates
- Non-seed documents are assigned to clusters initiated by seeds

? How many number of clusters ( $n_c$ )?

#### \* SEED-selection

- Random Selection (not too bad) although it doesn't seem so it may be good for some applications



- Choose first  $n_c$
- Generate  $n_c$  synthetic seed documents
- Use posting list (Peter Willett)

take the first make a cluster  
get the second if similar  
put it in the first one's cluster  
or it creates its own cluster  
take third if it similar to any  
of them ... and so on.

Anderberg : 1973

Jain & Dubes : 1988

Jain, Flynn, Murkin : 1999

### Single-Pass Algorithms : (continued)

#### a) Seed oriented

How to choose seed.

$$n_C = \frac{m \times n}{t}$$

no of non-zero elements in  $D$ .

$$D = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}_{dm}$$

$$D = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \quad \frac{3 \times 4}{12} = 1$$

$$D = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad \frac{4 \times 4}{4} = 4$$



#### b) Heuristic Approach :

Von Rijssbergen

used in news event detection & tracking

temporal order  
time

1. Process documents (objects) serially (one by one)

2. The first document becomes a cluster by itself.

3. Consider the next document

if it is not similar to existing clusters then  
it starts its own cluster

else

Join to the most similar cluster(s)

#### Questions:

which similarity measure

↑  
as we change the cluster  
also change their  
centroids

similarity threshold  
order of processing

}



15

20.11.2011 Nisan SIU.

GB  
19.11.11

How to represent clusters (cluster centroids)?

## 2. K-Means Algorithms

A typical approach

use a seed oriented approach  
"the heuristic approach."1- Obtain the initial clusters using an efficient algorithm.2- repeat

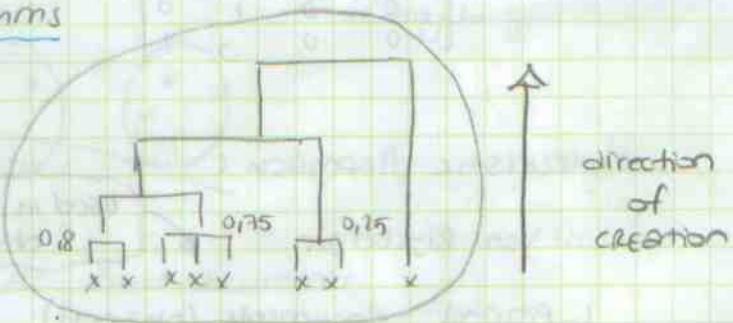
- Generate cluster centroids
- Reassign objects to the clusters according to their similarity to the cluster centroids

until all documents stay in their previous cluster OR100% stability  
90% "no of iterations = limit  
10 iterations?

## Graph Theoretical Algorithms

Agglomerative

\* bottom-up

This cluster is referred as a  $\Rightarrow$  DENDROGRAM

Individual similarities are used as a starting point, and a gluing process collects similar items, or group, into larger groups.

Single-link  
Complete-link  
Average-link

SPSS  
SAS  
Matlab?

Ellen Voorhees  
1965  
Sailor

Cornell



**Single-Link:** the similarity between a pair of clusters is taken to be the similarity between the most similar pair of documents, one of which appears in each cluster; thus each cluster member will be more similar to at least one member in that same cluster than to any member of another cluster.

04/03/2008

	A	B	C	D
A	1.0	0.3	0.5	0.6
B	-	1.0	0.4	0.5
C	-	-	1.0	0.3
D	-	-	-	1.0

## CLUSTER EXE

(including sources)

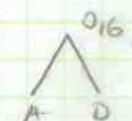
C++

OOR (Object Windows Library)

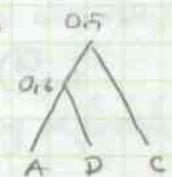
STEP	PAIR	SIMILARITY VALUE
1	AD	0.6
2	AC	0.5
3	BD	0.5
4	BC	0.4
5	AB	0.3
6	CD	0.3

## STEP Sim Pair

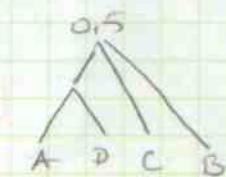
1 AD, 0.6



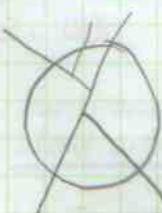
2 AC, 0.5



3 BD, 0.5



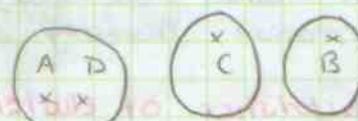
Agglomeration



Divisive



Dendrogram



	A	B	C	D
A	1.0	0.5	0.5	0.6
B	-	1.0	0.4	0.5
C	-	-	1.0	0.3
D	-	-	-	1.0

Similarity matrix implied by the dendrogram.

S vs Si Product moment correlation

[-1, +1]

+1 → total agreement  
-1 → total disagreement

⇒

If the clustering structure reflects the nature of the original data then there will be a high agreement of the values of  $S$  &  $S_1$  matrices

0.8 → good no for agreement

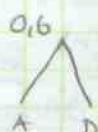
R Jain + ?

Advances in Computers

1985

### Complete Algorithm:

① AD 0.6

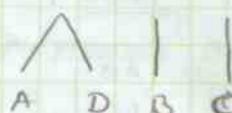


② AC, 0.6  
AC, CD

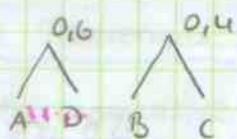
the sim between C-D not known  
thus we cannot join this one

③ BD 0.5

AB?

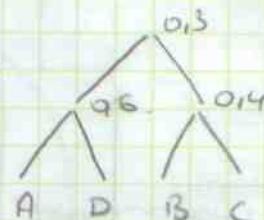


④ BC 0.4



⑤ AB : 0.3

AB 0.3      BD 0.5      AC 0.5      CD 0.3



### Average Link:

Ellen Voorhees, 1985 Cornell

Survey: Peter Willett on hierarchical clustering  
Information Processing & Management 1985

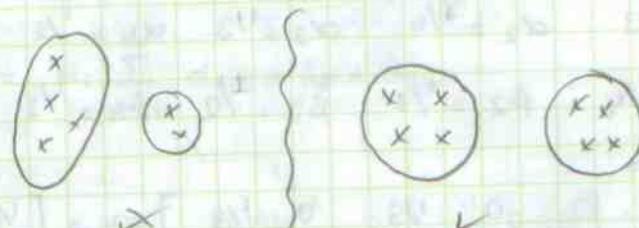
### Desirable Characteristics of Clustering Algorithms:

**EFFECTIVE:** generates a meaningful clustering structure + provides an effective IR environment

**EFFICIENT:** Time & Space

## Guidelines for partition-based clustering

- 1- Order Independence
- 2- Small errors made during indexing should not effect the clustering structure
- 3- Clustering structure should be maintainable (notice that usually we work in a dynamic environment)
- 4- Uniform distribution of objects among the clusters.



- 5- Clustering structure should be stable (addition of new items does not cause huge changes)
- 6- Small number of parameters  
(to guess this will make the use of clustering algorithm easier)

Cover Coefficient-based clustering methodology (C3m)

ACM Transactions on Database Systems (ACM TODS)

DEC. 1990

### Partitioning Seed Oriented

1- Determine no of clusters ( $n_c$ )

2- find cluster seeds

3- Assign non-seed, to clusters initiated by seed documents

$$\begin{array}{l} m = \text{no of docs} \\ n = \text{no of terms} \\ n_c = \text{no of clusters} \end{array}$$

$$\frac{(m-n_c) * n_c}{\text{non-seeds} \quad \text{seeds}}$$

## Relationship between indexing & clustering:

Example:

$$D = \begin{bmatrix} d_1 & t_1 & t_2 & t_3 & t_4 & t_5 & t_6 \\ d_2 & 1 & 0 & 0 & 1 & 0 & 1 \\ d_3 & 1 & 1 & 1 & 1 & 0 & 0 \\ d_4 & 1 & 0 & 0 & 0 & 1 & 1 \\ d_5 & 0 & 0 & 0 & 0 & 1 & 1 \\ d_6 & 1 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$$

→

$$\alpha_1 = 1/3 \quad \alpha_2 = 1/4 \quad \alpha_3 = 1/3 \quad \alpha_4 = 1/2 \quad \alpha_5 = 1/2 \quad \text{length of } 1$$

$$\downarrow \beta_1 = 1/4 \quad \beta_2 = 1/1 \quad \beta_3 = 1/2 \quad \beta_4 = 1/2 \quad \beta_5 = 1/2 \quad \beta_6 = 1/3$$

$$S = \begin{bmatrix} 1/3 & 0 & 0 & 1/3 & 0 & 1/3 \\ 1/4 & 1/4 & 1/4 & 1/4 & 0 & 0 \\ 1/3 & 0 & 0 & 0 & 1/3 & 1/3 \\ 0 & 0 & 0 & 0 & 1/2 & 1/2 \\ 1/2 & 0 & 1/2 & 0 & 0 & 0 \end{bmatrix}_{5 \times 6}$$

$$S' = \begin{bmatrix} 1/4 & 0 & 0 & 1/2 & 0 & 1/3 \\ 1/4 & 1 & 1/2 & 1/2 & 0 & 0 \\ 1/4 & 0 & 0 & 0 & 1/2 & 1/3 \\ 0 & 0 & 0 & 0 & 1/2 & 1/3 \\ 1/4 & 0 & 1/2 & 0 & 0 & 0 \end{bmatrix}_{5 \times 6}$$

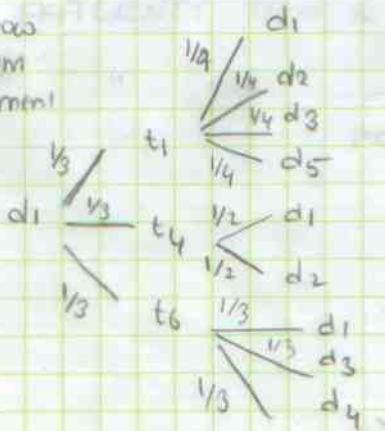
$$S'^T = \begin{bmatrix} 1/4 & 1/4 & 1/4 & 0 & 1/4 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1/2 & 0 & 0 & 1/2 \\ 1/2 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 1/2 & 1/2 & 0 \\ 1/3 & 0 & 1/3 & 1/3 & 0 \end{bmatrix}_{5 \times 6}$$

$$C = S S'^T = \text{result}$$

$$c_{11} = 1/3 * 1/4 + 0 + 0 + 1/3 * 1/2 + 0 + 1/3 * 1/3 = 0.363$$

$$C = \left[ \quad \right]_{n \times n}$$

SEE how  
to go from  
one document  
to another  
document



$$1 \leq i, j \leq m$$

$$c_{ij} = \alpha_i \sum_{k=1}^n d_{ik} \times \beta_k \times d_{jk}$$

Probability of choosing  
any term of  $d_j$  from  $d_i$ .

5 5 5 / 5

$$C_{11} = \alpha_1 \sum_{k=1}^6 d_{1k} * \beta_k * d_{1k}$$

$$m=5 \quad n=6 \quad = \alpha_1 \sum_{k=1}^6 d_{1k}^2 * \beta_k$$

$$= \alpha_1 (d_{11}^2 \beta_1 + d_{14}^2 \beta_4 + d_{16}^2 \beta_6)$$

$$= \frac{1}{3} \left( \frac{1}{4} + \frac{1}{2} + \frac{1}{3} \right) = 0.361$$

$$C_{12} = \alpha_1 \sum_{k=1}^6 d_{1k} * \beta_k * d_{2k}$$

$$C = \begin{bmatrix} 0.361 & 0.250 & 0.194 & 0.111 & 0.083 \\ 0.188 & 0.563 & 0.065 & 0.000 & 0.166 \\ 0.194 & 0.083 & 0.361 & 0.177 & 0.063 \\ 0.167 & 0.000 & 0.417 & 0.417 & 0.000 \\ 0.125 & 0.375 & 0.125 & 0.000 & 0.375 \end{bmatrix}$$

1-  $C_{ij} = 0 \Leftrightarrow C_{ji} = 0$

2-  $C_{ij} > 0 \Leftrightarrow C_{ji} > 0$

3- Row sum = 1

4-  $C_{ij}$  and  $C_{ji}$  may be non-zero and not equal to each other

$C_{ij}$  &  $C_{ji}$   
can be different  
when they are  
non-zero

$\hookrightarrow C_{12}, C_{21}$

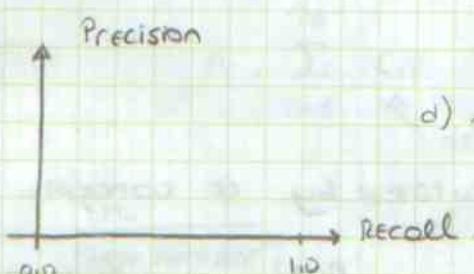
$$0.250, 0.188 > 0 \quad C_{12} \neq C_{21}$$

Two documents are identical  $\Rightarrow$  what would happen?

5-  $d_i \equiv d_j \Rightarrow C_{ii} = C_{jj} = C_{ij} = C_{ji}$

### ABOUT HOMEWORK

①



TREC 6 Appendix A.

d) find tree eval package  
bpref

$\Rightarrow$  11 point representation

24

$$\text{di} \xrightarrow{\frac{1}{1/3}} \frac{1}{1/3 + 1} \text{di} = \frac{1}{\frac{4}{3}} \text{di} = \frac{3}{4} \text{di}$$

$$\frac{1}{3} + \frac{1}{3} + \frac{1}{3} = 1$$

If  $\text{di}$  is unique, i.e., its terms do not appear in any other doc  $c_{ii} = 1$

Collection  $\Rightarrow$  all docs are unique

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

$$n_c = m = \sum_{i=1}^m c_{ii} = \sum_i 1 = m$$

$\Rightarrow$  all identical

$$\begin{bmatrix} 1/m & & \\ & 1/m & \\ & & 1/m \end{bmatrix}$$

$$n_c = 1 = \sum_{i=1}^m \frac{1}{m} = 1$$

$$\overbrace{\text{all identical}}$$

$$\overbrace{\text{all unique}}$$

11/03/2008

### \* C<sup>3</sup>M (Cover-Coefficient-Based Clustering Algorithm) (continued ...)

Partitioning (flat)

Seed based

$n_c$ : no. of clusters is calculated by  $cc$  concept.

$$C = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

$$n_c = m$$

all docs  
are unique

$$n_c = 1$$

$$C = \begin{bmatrix} 1/m & & \\ & 1/m & \\ & & 1/m \end{bmatrix}$$

all docs  
are identical

$$D = \begin{bmatrix} d_1 & & d_n \\ \vdots & \ddots & \vdots \\ d_m & & d_m \end{bmatrix}_{m \times n}$$

$$c_{ij} = \alpha_i \sum_{k=1}^n d_{ik} \times p_k \times d_{jk}$$

$$D = d_1 \begin{bmatrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 \end{bmatrix} \rightarrow d_1 = 1/3$$

$$\begin{bmatrix} 1 & 0 & 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$$

$$\beta_1 = 1/4$$

$$C = \begin{bmatrix} 0,361 & 0,250 & 0,194 & 0,111 & 0,082 \\ 0,188 & 0,563 & 0,065 & 0,000 & 0,188 \\ 0,194 & 0,083 & 0,361 & 0,277 & 0,083 \\ 0,167 & 0,000 & 0,417 & 0,417 & 0,000 \\ 0,125 & 0,375 & 0,125 & 0,000 & 0,375 \end{bmatrix}$$

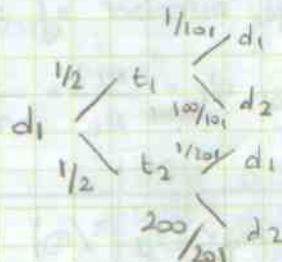
diagonal values  
gets ↘ as the collection  
gets ↗

$c_{ii} \rightarrow$  small no

$$c_{ii} > c_{ij}$$

for a weighted  $D$  matrix  $\Rightarrow c_{ii}$  can be  $< c_{ij}$

$$\begin{bmatrix} 1 & 1 & 0 & \dots & 0 \\ 100 & 200 & 0 & 0 & \dots & 0 \end{bmatrix}$$



- $D$  is binary

$$\rightarrow \min(c_{ii}) = 1/m$$

$$c_{ii} = c_{jj} = c_{ij} = c_{ji} \Rightarrow d_i = d_j \leftarrow \text{true for both binary and weighted cases.}$$

$$n_c = \sum_{i=1}^m c_{ii} = \sum s_i$$

$\nwarrow$  decoupling coefficient

$$\frac{n_c}{\text{large number}}$$

$$2,54 * 10^4 \rightarrow 1,2 * 10^{-2}$$

$$2,54 * 10^4, \quad 000001.2 \Rightarrow \underline{0,0000012} * 10^4$$

$$\frac{0,00 * 10^4}{2,54 * 10^4}$$

$\delta$  = average decoupling coefficient

$$\delta = \frac{1}{m} \sum_{i=1}^m \delta_i = \frac{1}{m} \sum_{i=1}^m c_{ii}$$

$$n_c = \sum_{i=1}^m \delta_i = m * \delta \rightarrow m * \frac{1}{m} \sum_{i=1}^m \delta_i$$

Coupling Coefficient  $\Psi_i = 1 - \delta_i$

avg. coupling coef.  $\Rightarrow \frac{1}{m} \sum_{i=1}^m \Psi_i$

$$n_c = m * \delta$$

\* Average number of docs per cluster

$$d_c = \frac{m}{n_c} = \frac{m}{m * \delta} = \frac{1}{\delta}$$

$$\max(1, m/n) \leq d_c \leq m$$

$$\max(1, n/m) \leq \delta_i \leq n$$

$$C = \left[ \begin{array}{c|ccccc} & & & & & n_c \\ \hline & \diagdown & & & & \\ & & \ddots & & & \\ & & & \diagdown & & \\ & & & & \ddots & \\ & & & & & 1 \end{array} \right]_{m \times m}$$

for terms

$$C' = \left[ \begin{array}{c|ccccc} & & & & & n'_c \\ \hline & \diagdown & & & & \\ & & \ddots & & & \\ & & & \diagdown & & \\ & & & & \ddots & \\ & & & & & n \end{array} \right]_{n \times n}$$

$n_c = n'_c$

15/5/09

20<sup>th</sup> March  $\Rightarrow$  form your group  
Thursday & select the project.

Elements of Style  
by Strunk & White

27

## Seed Selection:

cluster SEED power

docs → highest seed power

$\{ \begin{matrix} X \\ X \\ X \end{matrix} \}$  choose top  $n_c$  docs  
as the seeds

$$\left[ \begin{matrix} 1 & 1 \\ 1 & 1 \\ \dots & \dots \\ 0 & 0 \end{matrix} \right]$$

7 possible

X doc with the lowest seed power

$$P_i = S_i * \Psi_i * \underbrace{X_{di}}$$

if its higher provides connection with the collection of docs

no of terms in  $d_i$

(depth of indexing)

the rest of the docs

13/03/2008

$$D = \begin{bmatrix} d_1 & t_1 & t_2 & t_3 & t_4 & t_5 & t_6 \\ d_2 & 1 & 1 & 1 & 1 & 0 & 0 \\ d_3 & 1 & 0 & 0 & 0 & 1 & 1 \\ d_4 & 0 & 0 & 0 & 0 & 1 & 1 \\ d_5 & 1 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$$

$$C = \begin{bmatrix} 0,361 & 0,250 & 0,194 & 0,111 & 0,033 \\ 0,166 & 0,563 & 0,065 & 0,033 & 0,111 \\ 0,194 & 0,063 & 0,0361 & 0,227 & 0,123 \\ 0,167 & 0,000 & 0,419 & 0,419 & 0 \\ 0,115 & 0,395 & 0,115 & 0 & 0,349 \end{bmatrix}$$

$$P_i = S_i * \Psi_i * \underbrace{X_{di}}$$

seed power

no of terms in  $d_i$

$(\alpha_i^{-1})$

$$S_i = C_{ii}$$

$$\Psi_i = 1 - S_i$$

$$\frac{P_i}{d_2 \rightarrow 0,954} \\ d_1 \rightarrow 0,692 \\ d_3 \rightarrow 0,692 \\ \vdots$$

$$C_{ii} = C_{33} + C_{33} = C_{31} ?$$

$d_1$  &  $d_3$  are not the same

$$P_1 = 0,361 * (1 - 0,361) * 3 = 0,692$$

$d_1 \rightarrow t_1, t_4, t_6$

$$P_2 = 0,563 * (1 - 0,563) * 4 = 0,984$$

choose  $d_2$  since it uses more unused terms

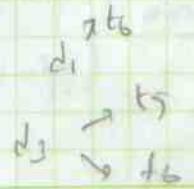
$$P_3 = 0,692$$

The previously selected seed ( $d_2$ ) contains  $t_1, t_2, t_3, t_4$ .

$$P_4 = 0,484$$

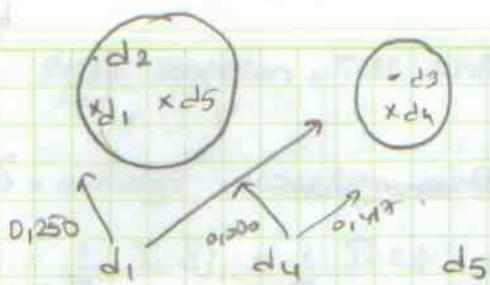
$$P_5 = 0,409$$

$$n = \sum_{i=1}^r C_{ii} \approx 2$$



27

28



non-seed  
seed  
 $3 \times 2 = 6$

$(m - n_c)$

\* No of computations Involved During Clustering?

$$m + (m - n_c) \times n_c \approx m + m \times n_c$$

↑

needed  
for finding  
 $n_c$

$m \gg n_c$

TONS = 214

INSPEC = 12,654

Inverted Index for seed Documents (IISD)

$$c_{ij} = \alpha_i \sum_{k=1}^n d_{ik} * \beta_k * d_{jk}$$

$d_1 \rightarrow [t_1 | t_6 | t_{10} | t_{15}]$

$d_2 \rightarrow [t_{10} | t_{20} | t_{30} | t_{40}]$

So doc  $d_2, d_3$  're looking

$t_1 \rightarrow (d_2, 1) < (d_3, 1)$

$t_2 \rightarrow (d_2, 1)$

$t_3 \rightarrow (d_2, 1)$

$t_4 \rightarrow (d_2, 1)$

$t_5 \rightarrow (d_3, 1)$

$t_6 \rightarrow (d_3, 1)$

⇒

cluster d1%

$c_{12}$	$c_{13}$
0	0

t4

$c_{12} = c_{21} =$

t4:  $c_{12} = c_{12} + \alpha_1 * (d_{11} * \beta_1 * d_{21})$

$= 0 + 1/3 * (1 * 1/4 * 1) = 1/n$

$c_{12} = c_{13} + \alpha_1 * (d_{11} * \beta_1 * d_{31})$

$= 0 + 1/3 * (1 * 1/4 * 1) = 1/n$

t4

$c_{12} = c_{12} + \alpha_1 * (d_{14} * \beta_4 * d_{24}) = \frac{1}{12} + \frac{1}{3} * \left(1 * \frac{1}{2} * 1\right)$

$= 0,1250$

t6  $c_{13} = c_{13} + \alpha_1 * (d_{16} * \beta_6 * d_{36}) = \frac{1}{12} + \frac{1}{3} * \left(1 * \frac{1}{3} * 1\right)$

$= 0,194$

$\frac{c_{12}}{0,1250}$

$\frac{c_{13}}{0,194}$

$c_{12} > c_{13}$  so

d1 joins the cluster of d2

Complexity:Find  $c_{ij}$  value  $\Rightarrow m * X_d$  $X_d = \text{avg no of terms/doc}$ Create IIISD  $\Rightarrow n_c * X_d$ tg: term generality,  
avg no of docs/term

Clustering

Assign non seeds  $\Rightarrow (m - n_c) * X_d * \frac{\text{no of terms}}{\text{docs}}$ R avg posting list  
in IIISD

$tg_s = \frac{2+5+1}{6} \approx 1$

(a term appears in  
many different feeds)ET O() execution time  
 $m * X_d * tg_s$ Meters  
ET

M

t

$O((m - n_c) * X_d * tg_s) = O(m * X_d * tg_s)$

$\frac{m * n}{t}$

$$0(m \cdot x_d \cdot tgs) \rightarrow \text{avg posting list IISD}$$

↑      ↙  
no of docs      avg no of terms/doc

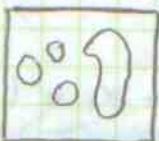
$$M - n_c \approx m$$

$$\uparrow \quad \uparrow$$

$$10^6 \quad 10^3$$

### ~~XX CLUSTER VALIDITY?~~

$$D \rightarrow [C^3 M]$$



clusters: meaningful?



Check the consistency?

"Text anchored phrase analysis"  
Klein.  
(Tubilak)

clusters generated by humans

### Test Collection

1. documents
2. queries
3. relevant docs for queries

### Relevant docs

$$q_1 \rightarrow 1, 2, 3,$$

$$q_2 \rightarrow 1, 5, 6$$



TREC

$$n_q = 50$$

A no of queries

$q_1 \rightarrow$  find the target clusters  
 $\{C_1, C_2\}$   
 $q_2 \rightarrow$   $\{C_1, C_3\}$

In a meaningful clustering environment no. of target clusters for a given query should be small?

Recall "the cluster hypothesis"

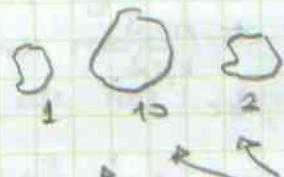
"How small is good enough?"

\* Find average number of target clusters for all queries =  $n_t$   
 (this should be a small number  
 (↳ how small?)

- Keep the clustering structure as the same

query 'loc'  
 hic baten modum  
 cluste say (5)  
 he bte clusterin  
 length'i 1000  
 tutulup dökümanları  
 random döglüyorum

$n_c$ : the same  
 no of documents in  $c_i$  ( $1 \leq c_i \leq n_c$ ) is kept the same



$x_{di}$   
 $x_{dm}$

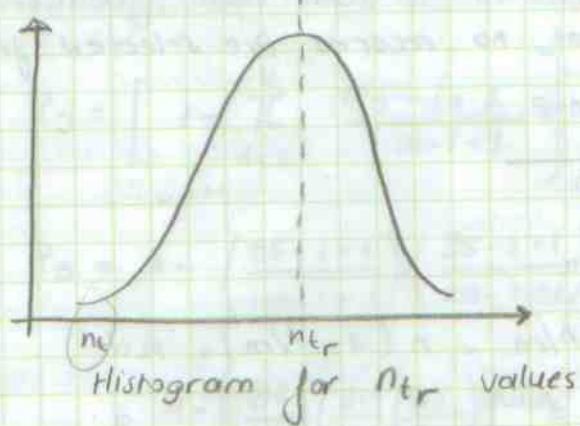
1<sup>st</sup> random  $n_t$   
 Find  $n_{tr_1}$   
 $n_{tr_2}$   
 take the  
 → average

Distribute the documents randomly

(random may not be too bad)

$\{n_t < n_{tr}\}$

\* if  $n_t > n_{tr} \Rightarrow$  INVALID CLUSTERING STRUCTURE



S.B. Yao 1977, Communications of the ACM

"Approximating block accesses in database organizations"

Blocks  
have  
some size

xx
x
xx
x
xx
x

$q \rightarrow$  relevant documents

1  
2  
3  
4  
5

Alfonso

Cardenas

using his original notations:

no. terms  $\nwarrow$   $n$ : no of records

(in each block we have the same no of docs)

$m$ : no of blocks

the probability of accessing a block

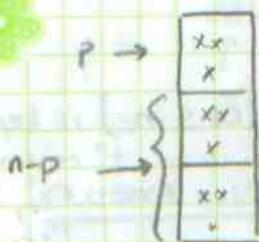
block size =  $n/m$

$k$ : no of relevant docs for  
the given query



$p = n/m \Rightarrow$  no of records in the  $j^{\text{th}}$  block

$n-p =$  no of records in the other blocks



$C_k^n =$  no of combinations that we can have if we try to select  $k$  documents out of  $n$  documents

$$\text{e.g. } b, c \in \{a, b, c\} \quad k=2$$

ab, ac, bc

$$C_k^n = \frac{n!}{k!(n-k)!}$$

$$C_2^3 = \frac{3!}{2!(3-2)!} = \frac{3!}{2!} = 3$$

#  $C_k^{n-p}$  = different values of selecting  $k$  documents from  $(n-p)$  documents.

$\Rightarrow$  The probability that no records are selected from the  $j^{\text{th}}$  block

$$\frac{C_k^{n-p}}{C_k^n}$$

$$\text{Let } d = 1 - 1/m$$

$$n-p = n - n/m = n(1 - 1/m) = nd$$

EXPECTED VALUE

$E(I_j) =$  Probability of selecting at least  $o$  record from the  $j^{\text{th}}$  block:

$$1 - C_k^{nd} / C_k^n$$

Expected no of blocks to be accessed:-

$$\sum_{j=1}^m E(I_j)$$

$$m * \left[ 1 - \frac{\frac{nd!}{k!(nd-k)!}}{\frac{n!}{k!(n-k)!}} \right]$$

$$m * \left[ 1 - \frac{nd!}{k!(nd-k)!} \cdot \frac{k!(n-k)!}{n!} \right] = m * \left[ 1 - \frac{1 \cdot 2 \dots nd}{1 \cdot 2 \dots (nd-k)} \cdot \frac{1 \cdot 2 \dots (n-k)}{1 \cdot 2 \dots n} \right]$$

$$\frac{(nd-k+1)(nd-k+2) \dots nd}{(n-k+1)(n-k+2) \dots n}$$

$$n_{tr} = m * \left[ 1 - \prod_{i=1}^k \frac{nd-i+1}{n-i+1} \right]$$

↑  
for a query

→ Buradan sonra kendi rotasyonlara imza devam ediyoruz.

\* How to use Yao's formula in a clustering environment?

$$k=3 \\ m=100 \quad (\text{no of docs})$$

$$|C_1|=5 \leftarrow \text{no of docs in } C_1$$

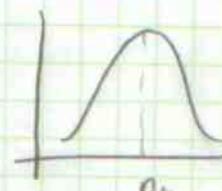
$$m_j = 100-5 = 95 \rightarrow \text{no of docs in the other clusters.}$$

- Probability that there is at least one document in this cluster.

$$P_j = \left[ 1 - \prod_{i=1}^k \frac{m_j-i+1}{m-i+1} \right]$$

$$P_5 = 1 - \left( \frac{95-1+1}{100-1+1} \right) \times \left( \frac{95-2+1}{100-2+1} \right) \times \left( \frac{95-3+1}{100-3+1} \right)$$

$$= 1 - \left[ \left( \frac{95}{100} \right) \left( \frac{94}{99} \right) \left( \frac{93}{98} \right) \right] \approx 1 - 0,86 = 0,14$$



→ ~ 10,35 indeks yoluyla (✓) doğrulandırıldı.

### Indexing:

We have two components:

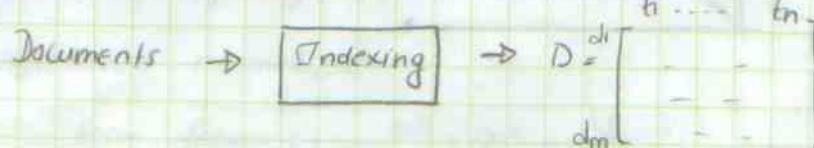
- 1- Term Selection
- 2- Assigning weights to terms.

### Database Management:

Records  $\Rightarrow$  record attributes

stNO, stName, stAge

### \* Information Retrieval:



# What's a term?

### Controlled Vocabulary



Terms are known beforehand

MESH : Medical Subject Heading  
National Institutes of Health

### Uncontrolled Vocabulary:

- ① Stoplist or no stoplist  $\Rightarrow$  use eliminates 30% of terms.
- ② Stemming vs. no stemming

if stemming then

which algorithm?

[for English  $\rightarrow$  (Porter's Algorithm published in Program)  
 $\downarrow$  journal

[for Turkish  $\rightarrow$  agglutinative language  
 $\downarrow$  words are created by using suffixes.

- simple word truncation

bilgisayar

$\downarrow$  use first 5 letters (a simple minded approach)  
f5

- word morphological structure

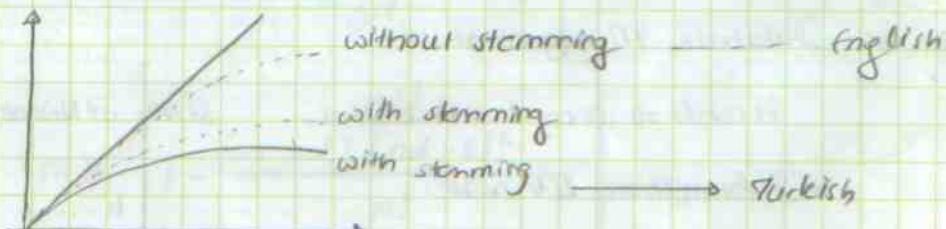
actual stem of a word

better  
best  $\rightarrow$  leme  
'good'

Its "dictionary entry"  $\rightarrow$  Lemma

### CONCERN OF THE INDEXING PROCESS:

- 1- What's the best indexing vocabulary



2. What should be the importance of a term in a document?  $tf \times idf$

3. How many terms in a certain document?

4. How many documents per term?  $\rightarrow D = [0]$

Gerard Salton AND Chris Buckley (Information Processing and Management) IPM

\* 1988, "Term Weighting Approaches in Automatic Text Retrieval"

① D

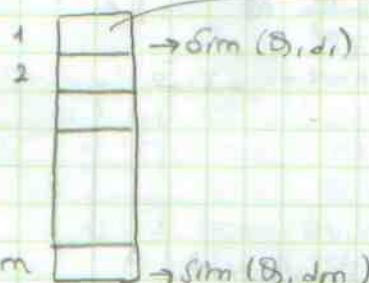
↑  
Query document matching function.

1800 diff



287 unique

Similarity



20 MAR. 2008

$$\text{Sim}(D, d_i) = \sum_{k=1}^n w_{q_k} w_{d_k}$$

$tf \times idf$

### Term Weighting Components:

Term Frequency Component (TFC)

Collection Frequency Component (CFC)

Normalization Component (NC)

d	q
TFC	TFC
CFC	CFC
NC	NC

We'll see that this one is  
not needed

(Unnecessary to normalize query components)

TFC

b: binary (1,0)

t: raw term frequency ( $tf$ )

n: augmented normalized

$$\text{term frequency} = 0.5 + 0.5 \times \frac{tf}{\max tf}$$

36

document vector  $(5 \ 0 \ 1 \ 2)$

$b: (1 \ 0 \ 1 \ 1)$

$tf: (5 \ 0 \ 1 \ 2)$

$$n: \left( 0.5 + 0.5 \frac{5}{5} \right)$$

$\max tf = 5$

$$\frac{\sqrt{m}}{\max tf}$$

$$0, 0.5 + 0.5 \cdot \frac{1}{5}, 0.5 + 0.5 \cdot \frac{2}{5}$$

CFC

X: no change use original TFC component

$$f: \ln \frac{m}{t_{gt}} + 1$$

$\rightarrow$  no of document that contain  $t_j$

$$0 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$t_{gt} = 3$

P = using 2 probabilistic inverse collection

frequency factor  $\ln \left( \frac{m - t_{gt} + 1}{t_{gt}} \right)$

$$NC \quad \gamma: \text{no of change} \quad C = \frac{1}{\sqrt{\sum w_i^2}}$$

TERM WEIGHT ASSIGNMENT

document vector

query vector

1800

$\downarrow$   
287 diff poss

3. bit.n TFC  
2. x\_ip CFC  
1. x\_ic NC

TFC 3  
CFC 3  
NC 1

$$(3 \times 3 \times 2) \times (2 \times 3 \times 1) \Rightarrow \text{we've this much possibility}$$

Example:

$$D = \begin{bmatrix} 2 & 0 & 1 & 2 & 0 \\ 0 & 2 & 1 & 3 & 1 \\ 2 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 1 \\ 2 & 1 & 0 & 1 & 0 \end{bmatrix} \quad \begin{matrix} \text{max } t_f \\ \text{min } t_f \\ \text{avg } t_f \\ \text{var } t_f \\ \text{std } t_f \end{matrix}$$

term generality

$$\begin{matrix} t_f \\ \text{tg} \\ \text{run } t_f \\ \text{avg } t_f \\ \text{var } t_f \\ \text{std } t_f \end{matrix}$$

$m=5 \quad n=5 \quad 5 \text{ docs} \quad 5 \text{ terms}$

$d_1$ :  $t_{fc}$   
 term frequency  
 $\rightarrow$  CFC  $N_C$

$$\left\{ \begin{array}{l} TFC \rightarrow t_f (2, 0, 1, 2, 0) \\ \rightarrow f \ln \frac{m}{t_f} + 1 \end{array} \right.$$

$$\left[ \begin{array}{c} \ln \frac{5}{2} + 1 \\ \ln \frac{5}{1} + 1 \\ \ln \frac{5}{2} + 1 \\ \ln \frac{5}{1} + 1 \\ \ln \frac{5}{2} + 1 \end{array} \right]$$

$$\begin{matrix} 4 & 2 & 2 & 4 & 3 \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ \ln \frac{5}{4} + 1 & \ln \frac{5}{2} + 1 & \ln \frac{5}{2} + 1 & \ln \frac{5}{4} + 1 & \ln \frac{5}{3} + 1 \\ 1.22 & 1.92 & 1.92 & 1.22 & 1.51 \end{matrix}$$

$$\Rightarrow d_1 (2 \cdot 1.22, 0 \cdot 1.92, 1.92, 2 \cdot 1.22, 0 \cdot 1.51)$$

$$\Rightarrow (2.44, 0, 1.92, 2.44, 0)$$

$$C = \frac{1}{\sqrt{\sum w_i^2}} = \frac{1}{\sqrt{(2.44)^2 + 0 + (1.92)^2 + (2.44)^2 + 0}} = \frac{1}{3.95}$$

$$d_1 \Rightarrow \left( \frac{2.44}{3.95}, 0, \frac{1.92}{3.95}, \frac{2.44}{3.95}, 0 \right)$$

normalize

$$(0.62, 0, 0.49, 0.62, 0)$$



Query weight Normalization:  $nfx$

$$Q: (1 \ 0 \ 0 \ 2 \ 0)$$

$$tg \Rightarrow (4 \ 2 \ 2 \ 4 \ 3) \leftarrow \text{coming from } D$$

$$\text{TFC: } n \Rightarrow (0,5 + 0,5 \cdot \frac{1}{2} \ 0 \ 0 \ 0,5 + 0,5 \cdot \frac{2}{2} \ 0)$$

$$(0,75 \ 0 \ 0 \ 1 \ 0)$$

$$\text{CFC} \Rightarrow m \cdot \frac{m}{k_D}$$

$$(1,22 \ 1,92 \ 1,92 \ 1,22 \ 1,51)$$

$$(0,75 \cdot 1,22 \ 0 \ 0 \ 1,22 \ 0)$$

$$\Rightarrow (0,92 \ 0 \ 0 \ 1,22 \ 0)$$

$$D = \begin{bmatrix} 0,62 & 0,00 & 0,49 & 0,62 & 0,00 \\ 0,00 & 0,46 & 0,33 & 0,63 & 0,26 \\ 0,78 & 0,00 & 0,00 & 0,33 & 0,48 \\ 0,63 & 0,00 & 0,00 & 0,00 & 0,78 \\ 0,73 & 0,56 & 0,00 & 0,37 & 0,00 \end{bmatrix}$$

$$\text{Sim}(Q, d_1) = 0,62 \cdot 0,92 + 0 \cdot 0 + 0,49 \cdot 1,22 + 0 = 1,33$$

$$\text{Sim}(Q, d_2) = 0,77$$

$$\text{Sim}(Q, d_3) = 1,20 \quad \text{Ranking: } d_1 > d_3 > d_5 > d_2 > d_4$$

$$\text{Sim}(Q, d_4) = 0,58$$

$$\text{Sim}(Q, d_5) = 1,12$$

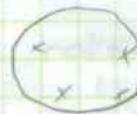
### Term Discrimination Value:

We use terms that distinguish documents from each other.

Intuition: We will be able to identify the relevant docs.

↑ effective

Document Space



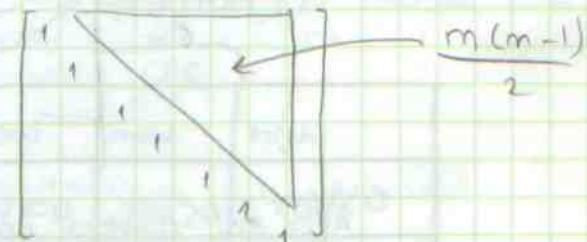
Using terms  
with better discrimination  
power

using terms with low  
discrimination  
power

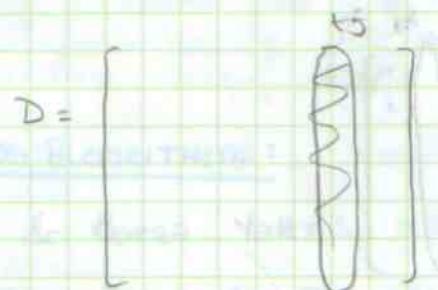
Q. space density

$$\text{Q}_S = \sum_{i=1}^m \sum_{j=1}^m S_{ij} \quad i \neq j$$

$$m(m-1)/2$$



$t_j \leftarrow$  delete this term



$\text{Q}_{S_j}$ : space density  
without  $t_j$

$t_j$  is a good discriminator. & What's its effect on  $\text{Q}_S$  (space density)

using all terms

using all terms except  $t_j$

$$\text{Q}_S < \text{Q}_{S_j}$$

4D

term characteristics  
characteristics

all terms without  $t_j$ 

Good discriminator &lt;

Bad discriminator &gt;

Indifferent discriminator ~

$$TDV_j = \frac{\partial g}{\partial t_j} - \frac{\partial g}{\partial t_j}$$

$$TDV = \frac{\partial g}{\partial t_j}$$

Theory of indexing Salton

$$d_{ij} \times TDV_j$$

How to calculate TDV using the cover coeff. concept

$$C = \begin{bmatrix} C_{11} \\ C_{21} \\ \vdots \\ C_{n1} \end{bmatrix}$$

$$n_C = \sum_{i=1}^m C_{ii} = \sum_{i=1}^m S_i$$

 $n_C$  = # of clusters using all terms

Good discriminator

 $n_{C'}$  = # of clusters without  $t_j$  $n_C > n_{C'}$ 

$$D = \begin{bmatrix} D_{11} \\ D_{21} \\ \vdots \\ D_{n1} \end{bmatrix}$$

1- Straightforward approach: requires many computations

Consider S: space density approach

Generate a collection of cluster centroid

 $t_1 \ t_2 \ t_3 \ t_4$ 

Cent.

$$D = \begin{bmatrix} 1 & 0 & 2 & 1 \\ 0 & 0 & 1 & 2 \\ 3 & 1 & 0 & 1 \end{bmatrix}$$

$m_{11} = m_{21}$   
 $t_1$

$$\text{Cent.} = \left[ \frac{4}{3} \ \frac{1}{3} \ \frac{3}{3} \ \frac{4}{3} \right]$$

$$\text{S} = \frac{\sum_{i=1}^m \text{sim}(d_i, \text{cont})}{m}$$

$$n_c = \frac{m \times n}{t}$$

$$n_{(L)} = \frac{m \times (n-1)}{(t - t_{SL})}$$

27 Mar. 2008

JASIS: T



Journal of the  
American  
Society for  
Information  
Science

S n<sub>c</sub>

t <sub>1</sub>	t <sub>1</sub>
t <sub>5</sub>	t <sub>5</sub>
t <sub>4</sub>	t <sub>3</sub>
t <sub>3</sub>	t <sub>4</sub>
1	1
1	1

Show the consistency  
of rankings

Wilcoxon (Spelling?)

(Spearman's)

TODS 214

INSPEC

(12,684 14,...)  
m n

		CFC		
		low	med	high
		TDV ≈ 0	TDV > 0	TDV < 0
		n <sub>c</sub>	TDV > 0	TDV ≈ 0
			TDV ≈ 0	TDV < 0

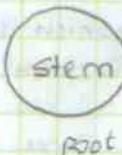
TDV's

## \* STEMMING ALGORITHMS:

Frakes & Baeza-Yates

Information Retrieval: Data Structures & Algorithms Ch 8

wanted  
want  
want



← find approximate root.

Lemma ← find the dictionary entry for a word

better  
best → good

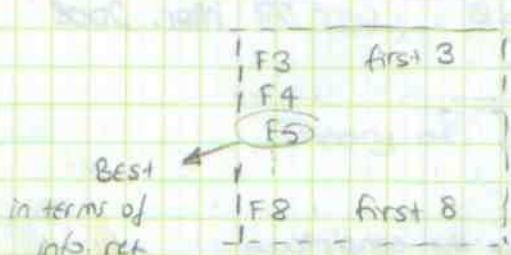
word lemma



## \* English

Porter's stemming Algorithm

### Program



↳ how this is decided

base line : no stemming

↓ by comparing the performance with the base line  $\approx 930$  increase in performance  
 $m \approx 400,000$  docs

5 years milliyet

## Stemmer vs. Lemmatizer

F5

Gemberek

NLP tools for Turkish

• Prob about Lemmatizer

→ it gives you a number of choices

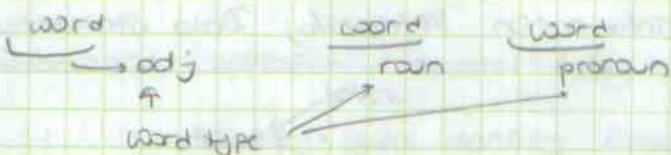
random

(how to choose the best one?)

\* How to choose one of the possible lemmas?

- Kural of Lazer

\* disambiguated text



most frequent word type

noun

verb

adjective

adverb - ?

less frequent word type

choose the one which is most frequent in Turkish!

next verb

# Performance will be better



Turkish

no stemming  
(...) F5

lemmatizer

Successor Variety

30% higher

- to measure the performance:
- bpref
- tree-eval package

Stemming (Conflation) Methods:

Manual

Automatic

Affix removal

removal

(prefix, suffix).

Successor variety

stable look-up

N-gram

Successor Variety:Root → after <sup>real</sup> root

→ many different letters

suffix → they all begin with different letters

how many different letters

AFFIX REMOVAL:

removes suffixes &amp; prefixes or both leaving a stem

SUCCESSOR VARIETY:

uses the frequencies of letter sequences in a body of text

N-GRAM:

makes decision based on the number of digram or n-gram (bigram) share

n-gram approach

bigram

trigram

n-gram

statistics

bigram

⇒ st, ta, at, ti, is, st, ti, ic, ci

statistical

⇒ st, t2, at, ti, is, st, ti, ic, ca, za

word1

word2

wordn

word1

word2

.

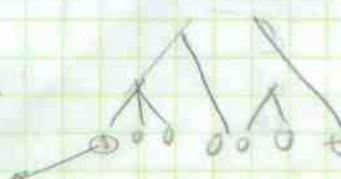
.

wordC

← find the similarity among words by using n-grams

use a clustering algorithm

words in the same group  
are assumed to have the same ~~the~~ size



## ...Stemming Algorithms:

1<sup>st</sup> April 2006  
 (alternatively 2106  
 01.00)

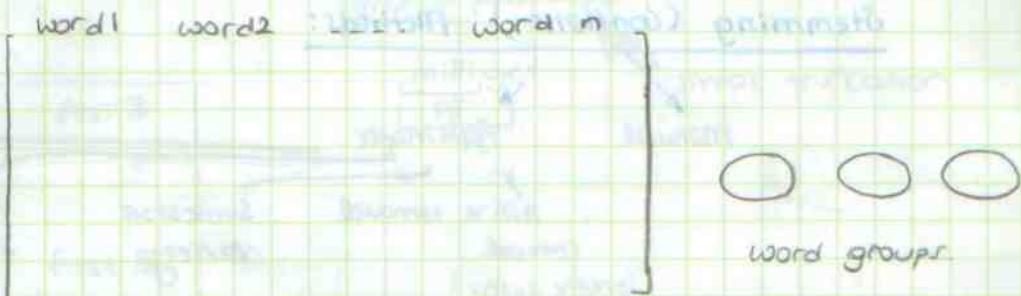
Affix Removal

n-gram

table backup

Successor Variety

$S =$   
 similarity matrix



statistical

statistics → ...

2 gran

st, ta, at, ti, is, st, th, ie, cr

 $\begin{array}{l} \text{word}_1 \rightarrow (1 - - -) \\ \text{word}_2 \rightarrow (2 - - -) \end{array}$ 
Successor Variety (SV)

SV tries to determine word stems

consider the prefixes of the word.



3-prefix  
first three  
chuctos

↑ determine the different letters  
in all words with the same prefix

CORPUS : ABLE, APE, BEATABLE, FIXABLE, READ, READABLE,  
READING, READS, RED, ROPE, RIPE

WORD : READABLE

READABLE

PREFIX	SUCCESSOR VARIETY	LETTERS
R	B	E, I, O
RE	L	A, D
REA	I	D

READABLE, READING  
READS, RED, ROPE, RIPE

PREFIX SUCCESSOR VARIETY LETTERS of word

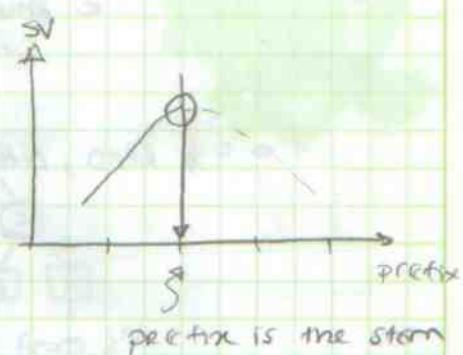
READ 3 A, I, S

READA 1 B

READAB 1 L

READABL 1 E

READABLE



## Information Storage &amp; Retrieval

W... , B...

- no stemming
- lemmatizer-based stemmer
- successor variety
- first n characters (F5, F6, ...)



matching function

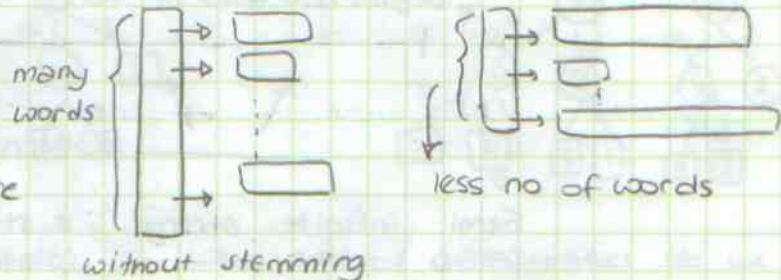
MF1, ..., MF8  
stemmers

- no stemmers: we've
- SV
- F6

24 different combinations

## # Why stemming?

- effectiveness
  - efficiency
- $\Delta$
- inverted index structure

Stemming:

Problems

Over Stemming:

□ ← (short stems)

Understemming

□ → (very long stem)

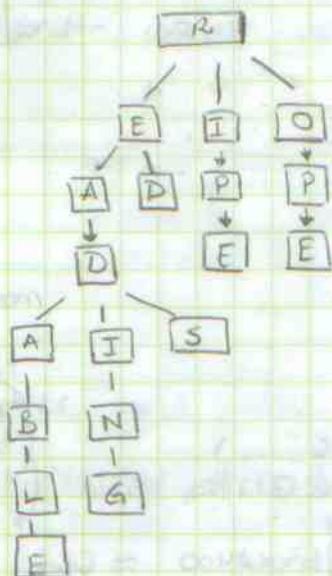


## How to Store Words?

Knuth Sorting & Searching Tree Structure

Symbol Tree Br:

\* READ, READABLE, READING, READS, RED, ROPE, RIPE



## PAT TREE CONSTRUCTION:

Patricia Tree

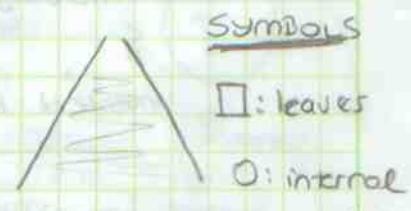
begins here goes to  $\infty$

alphabet = {0, 1, 2}

string

A      A  
pattern

Semi infinite string : s-string  
(sistring)



## Problem:

Given the text: 0 1 1 0 0 1 0 0 0 1 2 ...

Show the tree Patricia tree after the first eight sistring are inserted.

Bit position: 1 2 3 4 5 6 7 8 9 10 11  
0 1 1 0 0 1 0 0 0 1 0

sistring: 1: 0 1 1 0 0 1 0 0 ...

2: 1 1 0 0 1 0 0 ... 5: 0 1 0 0 0 1 ...

3: 1 0 0 1 0 0 ... 6: 1 0 0 0 1 0 ...

4: 0 0 1 0 0 ... 7: 0 0 0 1 ...

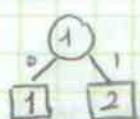
8: 0 0 1 0 ...

## # Insert sibling one by one

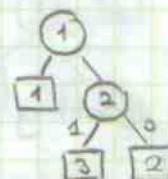
Insert 1: (i.e. sibling 1) [0 11 0 0]

1

Insert 2:



Insert 3:



1 0 0 1

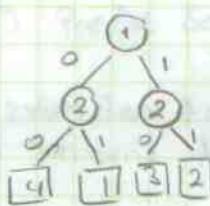
in order →  
differentiate 1<sup>st</sup> sibling from 2<sup>nd</sup>

1 → 0 11

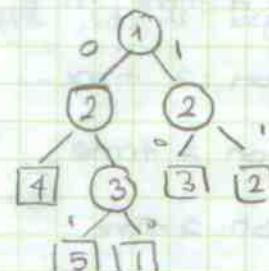
2 → 1 1

Insert 4:

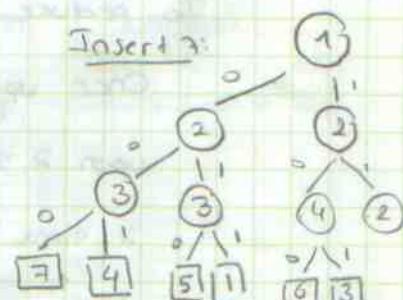
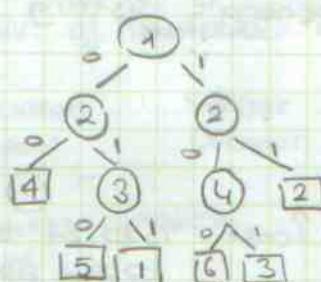
0 0 1 0 0 0



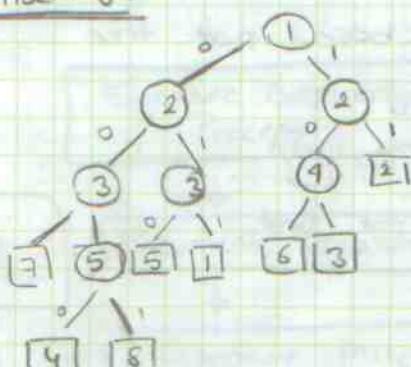
Insert 5: 0 1 0 0 0 1



Insert 6:



Insert 8:



What kind of characteristics do we have in the node of the tree structure?

7 → 0 0 0 1 -  
4 → 0 0 1 0 0 -  
8 → 0 0 1 0 1 -  
5 → 0 1 0 0 0 -  
1 → 0 1 1 0 0 1 0  
6 → 1 0 0 0 1 -  
3 → 1 0 0 1 0 -  
2 → 1 1 0 0 1 -

Increasing order



(48)

Write these positions in an array?

7 | 4 | 8 | 5 | 1 | 6 | 3 | 2      ← PAT Array (PAT)

7 < 4 < 8 < 5 < 1 < 6 < 3 < 2      in lexicographic order

We're looking for an order

Take a number search it in the array

'binary search'

08/04/08

### Another string Example

String : Once upon a time, in a far away ..

sistring1: Once upon a time ..

no of leaf nodes = n

sistring2: nice upon a time ..

no of internal nodes = n-1

sistring3: cc upon a time ..

To reduce the cost generate sistring by only using words

Once upon a time ..

upon a time

a time

← Cost = 1/320 of the character based approach

time

String: an example for a word-based pat tree

1    2    3    4    5    6    7

that is that, that is not that

sistring1: that is that, that is not that

sistring2: is that that is not

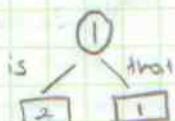
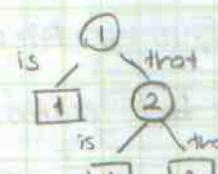
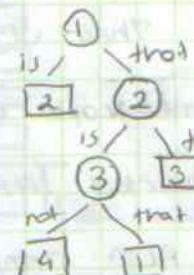
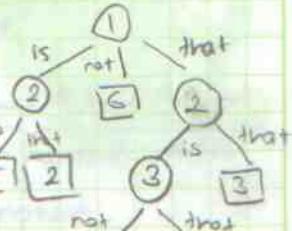
sistring3: that that is not

sistring4: that is not

sistring5: is not

Insert sis<sub>1</sub>,

1

Insert sis<sub>2</sub>,Insert sis<sub>3</sub>:Insert sis<sub>4</sub>:Insert sis<sub>5</sub>:

5 | 2 | 6 | 4 | 1 | 3

CONSTRUCTIONS OF INVERTED INDEXES WITHOUT USING THE n-GRAM CONCEPT

Chap 5: William Frakes &amp; Ricardo Baeza-Yates

use of tree  
pat

Info Ret. - Data &amp; File Structures

Prefix Search : Find the words beginning with the same prefix.

pat 1 → pat 2

n words

## \* Signature files:

1980s Univ. of Montreal must MUST SUBMITTED

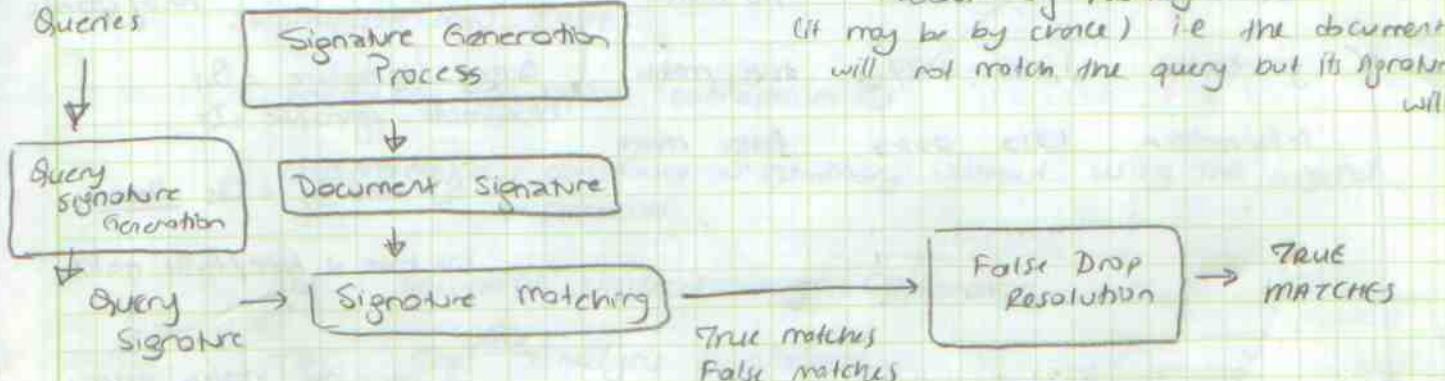
Ph.D. Supervisor : Stavros Christodoulakis  
Ph.D. Student : Christos Faloutsosbit map ← Göttlich  
concept

Document

not like inverted files

does not have exact representation  
because of hashing(it may be by choice) i.e. the document  
will not match the query but its signature  
will

Queries



false match = false drop

true match = true drop

## Why signature files?

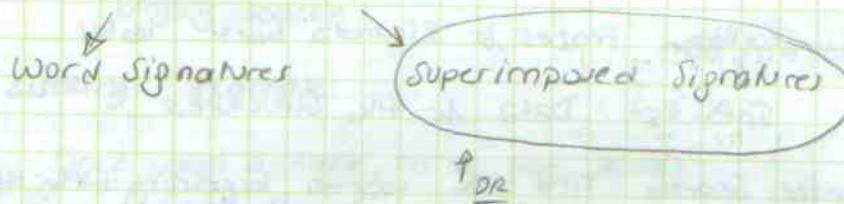
Their Size < Inverted Files

Some people question this claim?

Inverted files vs. Signature files

ACM Trans. on Database Systems.  
Justin Toben

### Signature Generation Methods



### \* Document Signature Generation Process:

1- Find term signatures

2- Superimpose (or) term signatures to obtain document signature

EXAMPLE TERM TERM SIGNATURE

object	1000	1000
signature	0010	0100
generation	1000	0100
	1010	1100

block signature

QUERY

QUERY SIGNATURE

RESULT

database 1100 0000

no match

generation

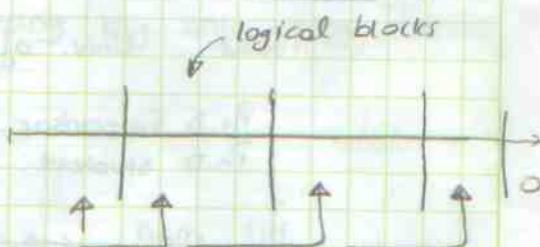
1000 0100

true match

information

1010 0000

false match



Each block contains the same no of distinct words  
(the last block may have fewer)  
10 of words

Query-Signature = Qs  
document-signature = Ds

If Qs and Ds = Qs then

We have a true/false match

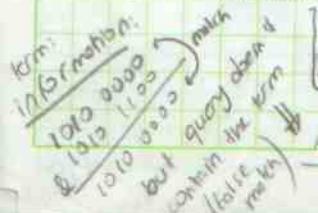
else

D does not match query

since in our query there is 1 word

here query signature = term signature

D has a match but false (Does it really contain the query term? → No).



prob of having a false match is very low

False Drop Resolution → to see whether it is a false or true match.

## 2Grams → 2GRAM SIGNATURE

Object → ob

→ [ ]

bj [ ]

je [ ]

ec [ ]

ct [ ]

+ [ ]

[ ] ← word signature

wild card characters

information  
info\*[ \* ] & [ ]  
partial match

Do not use the word itself, but generate '2grams' for the word and use its 2grams

Construction of Word Signatures  
(without using the n-gram concept)

n: signature length (F)

w: no of fs in a term signature (weight)

k: no of fs that we want to have in a signature

Let's say we want to have;

- signatures with the same weight

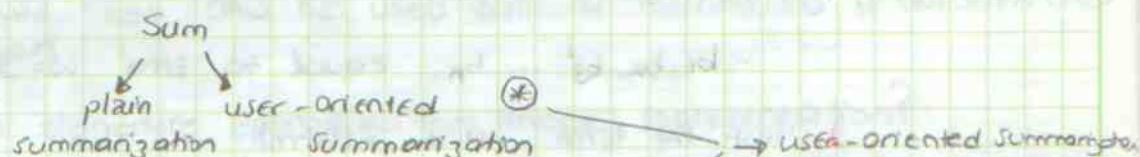
1- Initialize  $b_1 b_2 b_3 \dots b_n$ ,  $w=0$ 

Salton: (his blue book)

because many Text

10/04/2008

## User-Oriented Summarization:



## SUMMARIZATION TYPES

1- based on sentence extraction  $\oplus$ 

2. abstraction : construct a summary without using the original sentences.

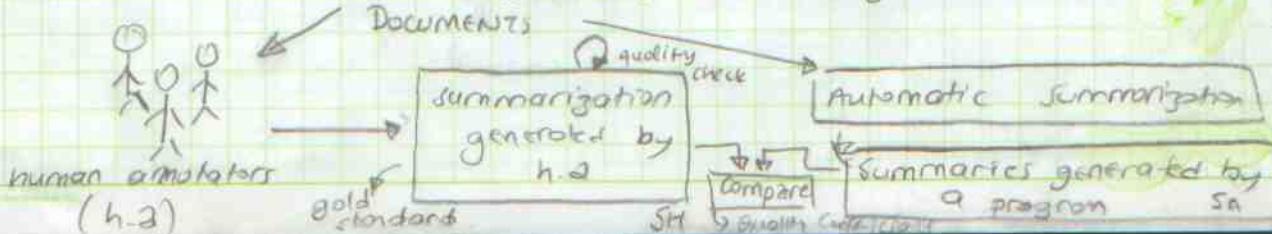
DUC : Document Understanding Conference

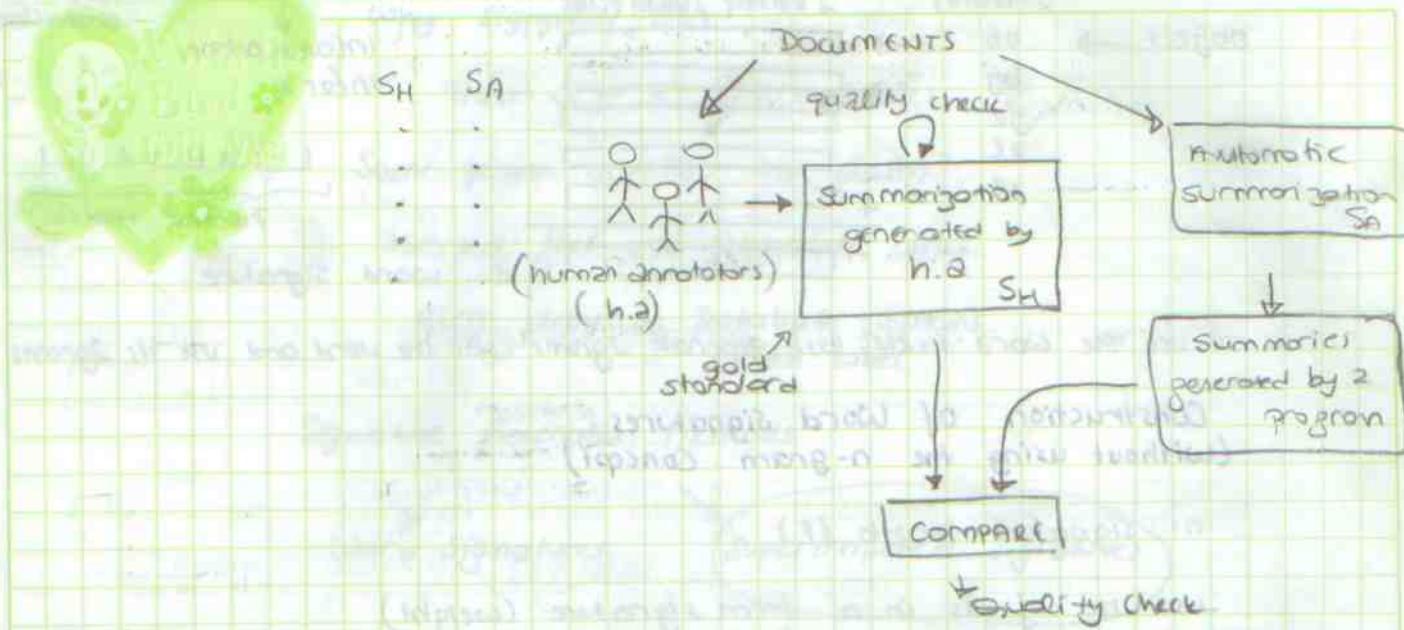
TAC : Text Analysis Conference

CNN.com → display news

meanwhile

3 sentences that summarizes the case





### ooo Signature files (cont'd)

#### Term Signature Generation:

From Salton: TEXTBOOK (Fig 7.42)

n: Signature length  
 w: no of 1's in a term signature (so far)  
 k: no of 1's we want to have in a signature ← s  
 $F \gg s$        $n \gg k$

##### 1. Initialize

$b_1, b_2, b_3, \dots, b_n$  equal to zero     $w=0$

2. Hash the term     $H \leftarrow h(\text{term})$

3. Initialize the random generator using  $H$ .

4. Determine next vector position to be set equal to 1

$$j \leftarrow \lfloor n \cdot \text{random}() \rfloor + 1$$

5. Is vector position = 1?

if  $b_j = 1$  then go to step 4

6.  $b_j \leftarrow 1$  ;  $w = w+1$

$$0 < \text{random}() < 1$$

7. If  $w < k$  then go to step 4

8. exit.

Rather than using complete words we can use n-grams of a given word and generate signature for word n-grams and superimpose them to obtain the term signature.

retrieval & 2-grams

re, et, tr, ri, ie, ev, va, al

so that we can later search terms that involve wild card characters.

Compute X  $\Rightarrow$  computer  
computation  
computational

### Use of Signature Files in Different Applications:

Paul Zelinka, et al

Dynamic Partitioning of Signature files

ACM Trans on Information Systems (ACM TOIS)

Oct 1991, pp 356-367

(LHS) Linear Hashing for superimposed signature

1. Text Retrieval

Quick Filter

2. Multimedia Representation and Retrieval

$\leftarrow$  unformatted data

3. Database Operations  $\leftarrow$  A. Thorp

$\leftarrow$  formatted data

4. In prolog database to store and access knowledge

5. CAD : Computer Aided Design.

6. OODB indexing SIGMOD

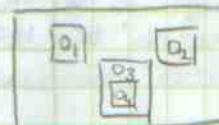
=

! SIGNATURE FILES CAN BE USED BOTH IN FORMATTED & UNFORMATTED ENVIRONMENTS

How to use signature approach for image representation?

Image  $\rightarrow$  Image Analysis Process  $\rightarrow$  Description of image in terms of the object contained

$$I = O_1, O_2(O_4, O_5(O_6)), O_3(O_9, O_{10}), O_1, O_3(O_5, O_6, O_7)$$



Object  $\Rightarrow$  Text Image

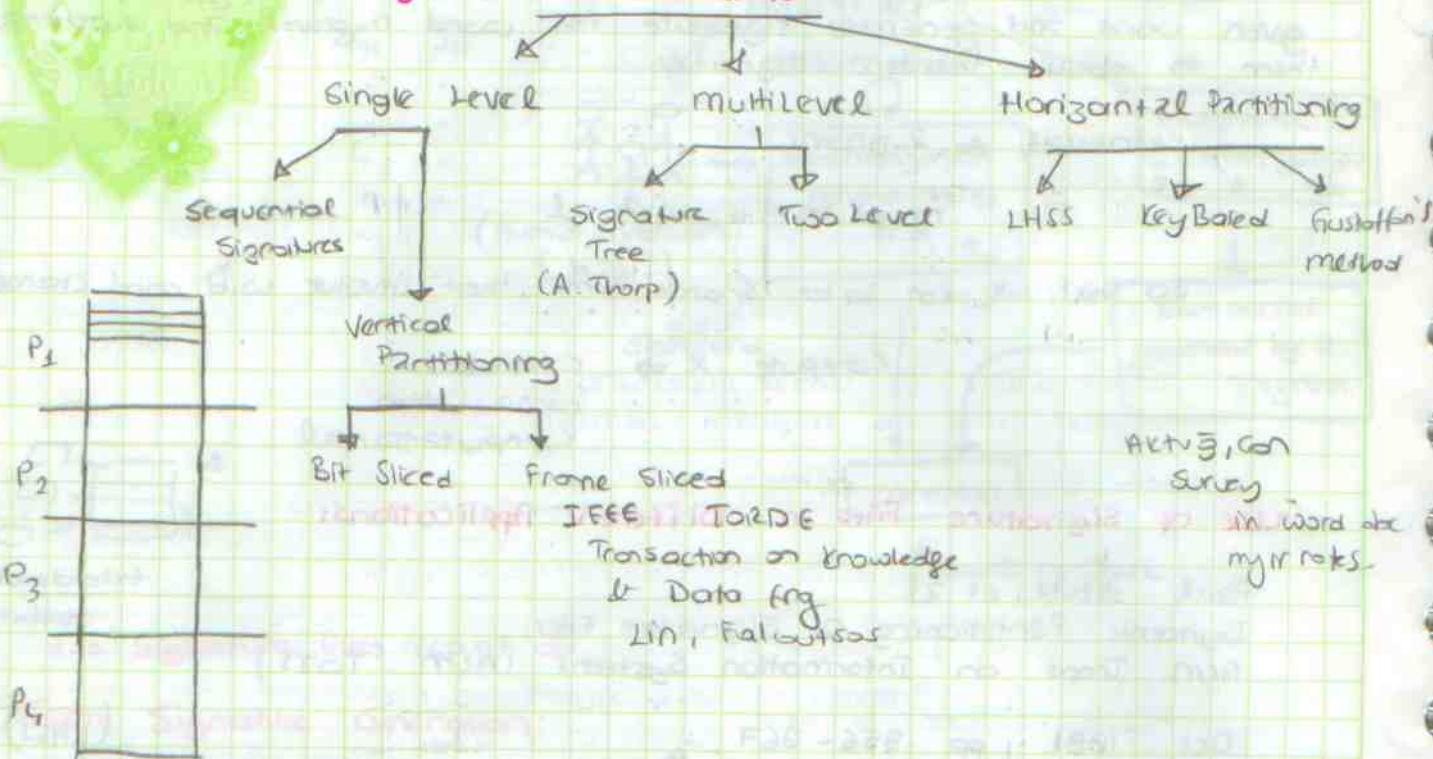


Signature  $\Rightarrow$  Text Image



Subject

## Signature File Structures:



### - Sequential Signatures (SS):

Other Names: Single Level Signatures, bit\_string approach

a bit may be 0 / 1 with equal choices		S1      0001 1110	Add new corners to the end of the file.
		S2      1100 0001	Simple Maintenance
		S3      0011 1100	Problem Search Process: $O(n)$
		S4      1100 0011	17/04/2008
		S5      0011 0110	
		S6      1100 1001	

### SSF (Sequential Signature Files)

S1	0001 1110
S2	1100 0001
S3	0011 1100
S4	1100 0011
S5	0011 0110
S6	1100 1001
⋮	⋮
Sn	⋮

F: signature size  
n: no of signature

$$\text{File size (bit)} = n * F$$

If  $(S_1 \text{ and } D_1 = S_2)$   
then access the document and perform false drop resolution

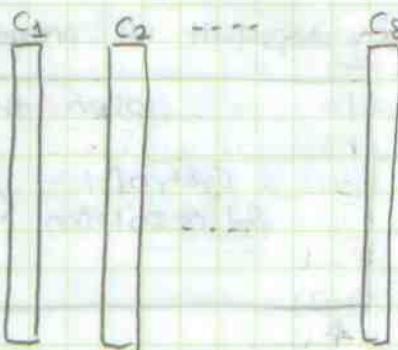
$$D_1 = 1001\ 0000$$

At  
make sure that doc contains the query terms

## Bit Sliced Signature files:

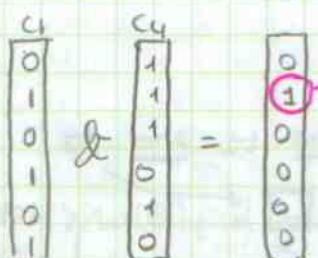
### Vertical Partitioning

C1 :	010	101
C2 :	010	101
C3 :	001	010
C4 :	111	010
C5 :	101	001
C6 :	101	010
C7 :	100	110
C8 :	010	101



Consider the same Query:

Q5 : 1001 0000



Just consider Doc2 for  
false drop resolution

How to store them

Pages in disk



Cost of query processing is proportional to  $O_{q\omega}$

(no of bits set to 1 in the  $Q_5$ )



If we assume that bit density 50%  
After processing the first bit of the  
query we eliminate 50% of the documents  
(only 50% we remain active)

After processing i no. of query bits  
only  $(\frac{1}{2})^i$  no. of documents  
remain active

$O_{q\omega}$ : query weight: no of 1s in the  
query resolution

Partial Query Matching - can also  
be used in inverted files

After a certain point rather  
than doing bit processing we may  
prefer to perform false drop  
resolution since cost of false  
drop resolution can be lower  
than bit processing  
(Assuming that we have very tall  
bit slices.)

Not using all of the bits of the query signature is referred to as partial evaluation.

When to quit bit processing?

When the

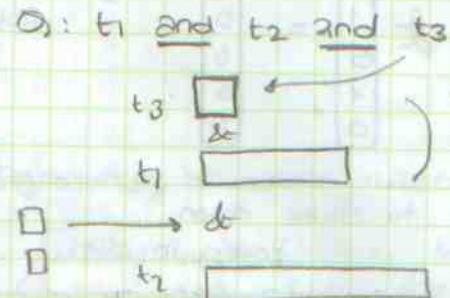
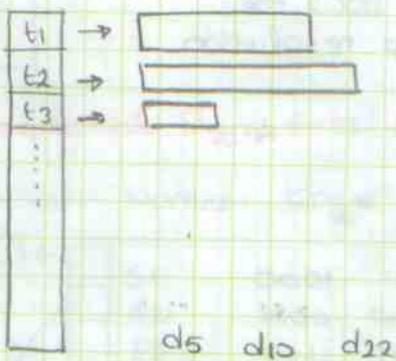
Cost of fd resolution < Cost of bit slice processing  
start fd resolution

↑

Information Processing Letters, Kacberber, con 1996  
Partial query matching

Partial Query Matching: can be used in inverted files

### Inverted Files



How about using very large signatures?  
So we'll may be able to eliminate  
not 50% but ~90%.

### How To Increase The Efficiency of Bit-Sliced File Query Processing

Use long signatures & set a small number of bit locations equal to 1

90% = 0 bits

10% = 1 bits

keep the same information like a posting list

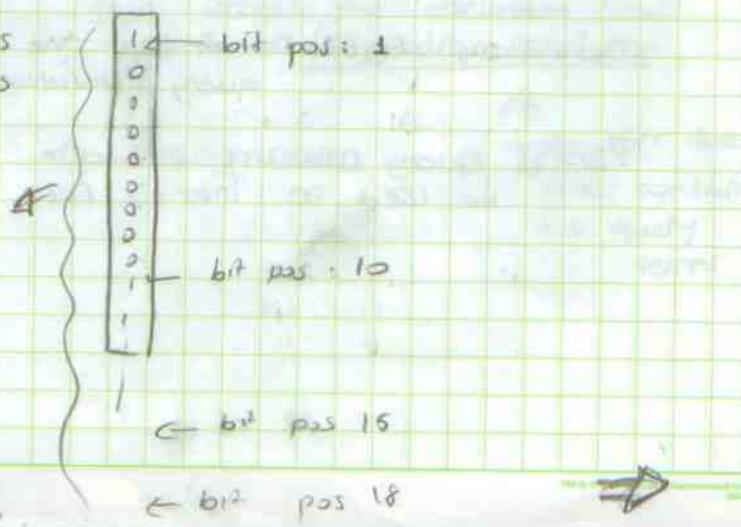
1 10 15 18 ...

use d-gap approach

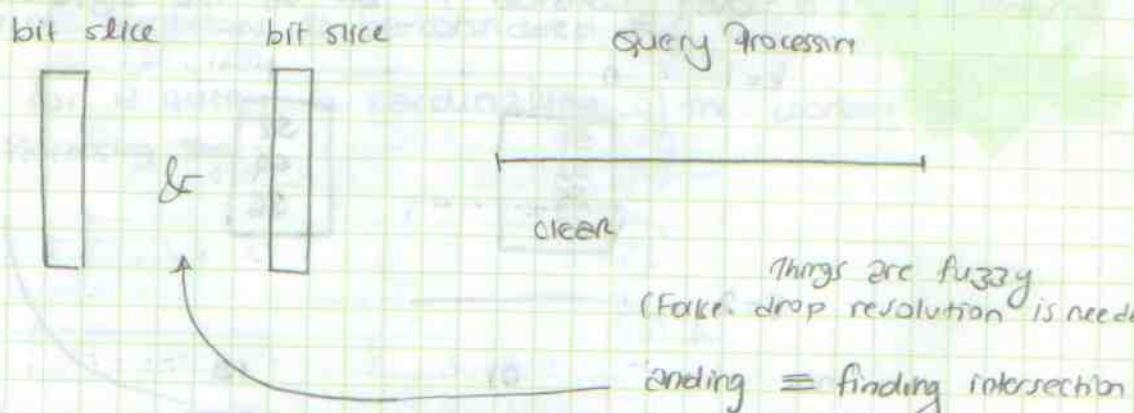
1 19 5 31 ...

(we have smaller numbers)

query processing exactly the same as inverted file approach



If we only store the locations of the 1 bits in a bit-slice then the structure becomes very similar to a posting list.



If  $i > w(S)_t = \text{no. of bits in a t-term query}$

$\downarrow$   
 no. of bits in a signature

then  $i = w(S)_t$

i.e. the volume of  $i$  cannot be greater than the no. of bits in the query.

→ Divide the Signature File into parts. In order to ↑ the efficiency of the query processing, in some way we will ignore some partitions.

### Signature File Partitioning

#### Horizontal Partitioning

It is assumed that signatures are random-bit vectors (50% - 50%)

$\uparrow$   
 $\downarrow$   
 1  
 0

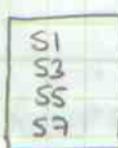
We have following set of signatures:

S1	0111	1000
S2	1000	1011
S3	0011	1100
S4	1100	0011
S5	0110	1100
S6	1011	0011
S7	0000	1111

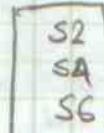
## Fixed Prefix Partitioning: (FPP)

- Just consider 1st bit of the signature

$k=1$       0



1

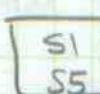


$k=2$

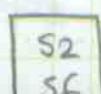
00



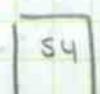
01



10



11



$k=3$

000

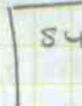
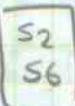
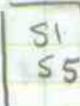
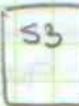
001

010

011

100

110



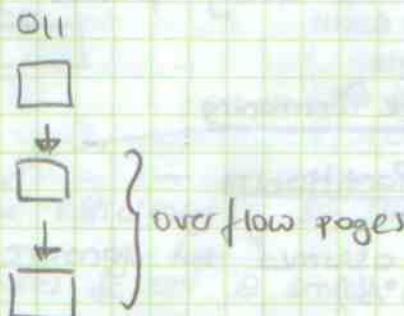
PROBLEM: We choose  $k$  value

We allocate a page for one block

We want uniform distribution among partitions since

15% - 15%

If one of the blocks overflows we will have overflow pages  $\leftarrow$  Problem



Q1 : 1000 0111

Working sets of query : a set of partitions such that (for a given query)

Q2 : 1111 0000

Q3 : 1100 0110

$$\{ p_i \mid p_i \text{key} \cap Q_k = Q_k \}$$

QUERY

$k=1$

$k=2$

$k=3$

Q1 1000

1 (1)

2 (10, 11)

4 (100, 101, 110, 111)

Q2 1111

1 (1)

2 (11)

3 (111)

Q3 1100

1 (1)

2 (11)

2 (110, 111)

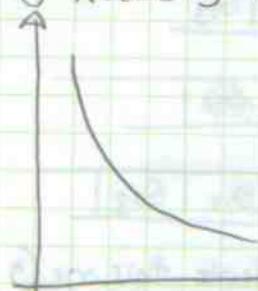
No of partitions (matching partitions)

Can we estimate the # of partitions to be considered?

# of partitions to be considered =  $2^{\# \text{ of } 0's \text{ in } k\text{-prefix} S}$

for e query = cardinality of the working set

Query Processing Time



$\rightarrow Q_w$  (query weight)

The partition that has all 1's is always selected.

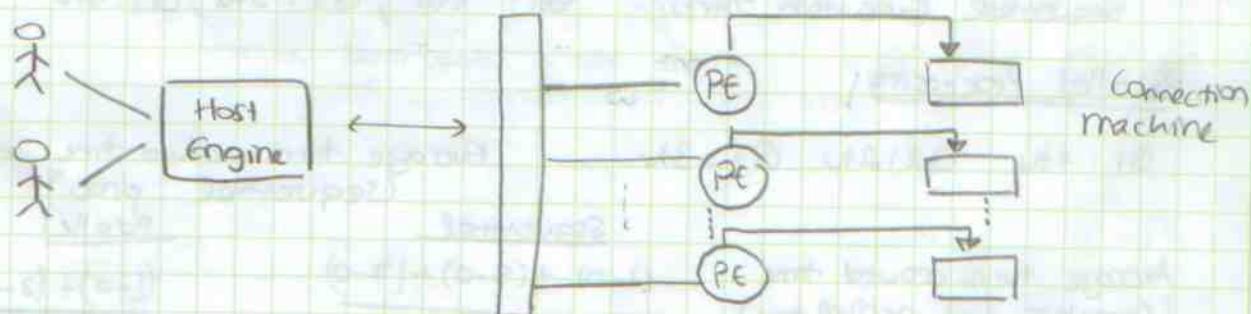
So to decrease the processing time, we would like to always keep this partition in the main memory.

- The partition with  $P_k = \text{all } 1$ 's ( $k=1 \rightarrow 1, k=2 \rightarrow 11, k=3 \rightarrow 111$ ) is always in the working set of the query.

Partition Processing can be done in - Sequential OR Parallel Partitioning.

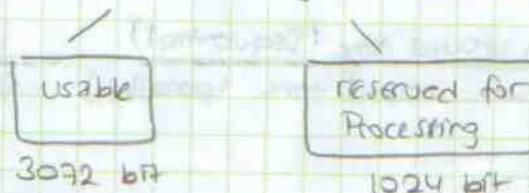
Parallel Hardware Structure:

PE : Processing Element (CPU)



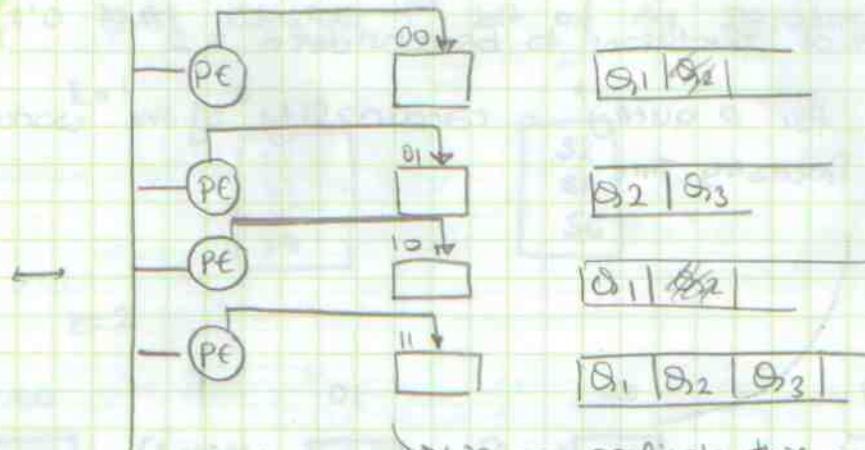
### Connection machine

Each PE has a memory of 4096 bits



60

Assume we have 4 partitions  
we want to compare sequential & parallel processing.



ASSUMPTION: All queries arrive at the same time ( $t=0$ )  
We will place queries in the queue

Another set of queries

	Query	$k=2$
$Q_1: 1110 \quad 0001$	$Q_1$	$1(11)$
$Q_2: 0000 \quad 1111$	$Q_2$	$4(00, 01, 10, 11)$
$Q_3: 0110 \quad 0011$	$Q_3$	$2(01, 11)$

### Sequential Processing

$tu$ : time unit

	$Q_1$	$Q_2$	$Q_2$	$Q_2$	$Q_3$	$Q_3$	
$t=0$	1	2	3	4	5	6	7

Sequential Execution Times:  $Q_1: 1tu ; Q_2: 5tu ; Q_3: 3tu$

### Parallel Processing

$Q_1: 1tu \quad Q_2: 2tu \quad Q_3: 3tu$

Average turn around time  
(sequential proc.)

Average turn-around time  
(complete time - arrival time)

$$\frac{(1-0)+(5-0)+(7-0)}{3} = \frac{13}{3} = 4.33 \text{ tu}$$

$$\frac{(1-0)+(2-0)+(3-0)}{3} = \frac{6}{3} = 2 \text{ tu.}$$

Throughput: # of jobs completed  
per unit time

$$3/7 = 0.43 \text{ jobs per unit time}$$

$$\text{Speed Up Ratio} = \frac{\text{Avg. Turn around time (Sequential)}}{\text{Avg. Turn around time (parallel)}} = \frac{13/3}{6/3} = \frac{13}{6} = 2.166$$

\* of p

No of partitions to be accessed for a query = 2  
no of 0's in query by

example  $Q = 1000 \text{ } 1101$

Problem: # of partitions is fixed,  
if overflow occurs it has to be  
recognized

### \* SIGNATURE FILE FOR STRUCTURED RECORDS #

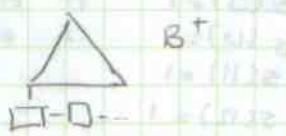
Alan Thorp File Organization & Processing p152 →

Employee (EmpNo, DeptName, Salary, No of Dependents)

Query types:

1-based on primary attribute

Display EmpName where EmpNo=100

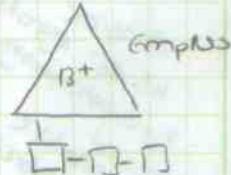


2-based on secondary attributes

Display EmpName where Salary=2000

### \* Signature Files vs. Inverted Indexes for DB Systems \*

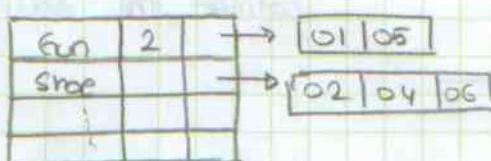
<u>EmpNo</u>	<u>EmpName</u>	<u>Dept</u>	<u>NoOfDependent</u>	<u>Salary</u>
01	Jane	Gun	1	1000
02	George	Shoe	2	1200
04	Dave	Shoe	3	1500
05	Mary	Gun	2	1800
06	Mike	Shoe	1	2000



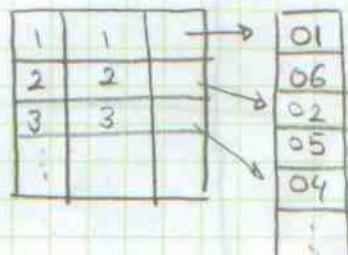
Inverted Indexes for the Secondary Attributes:

Inverted Index for Dept

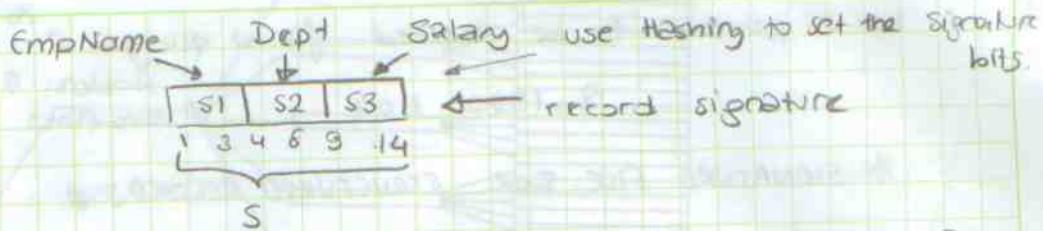
Dept



Inverted Index for NoOfDependent



62

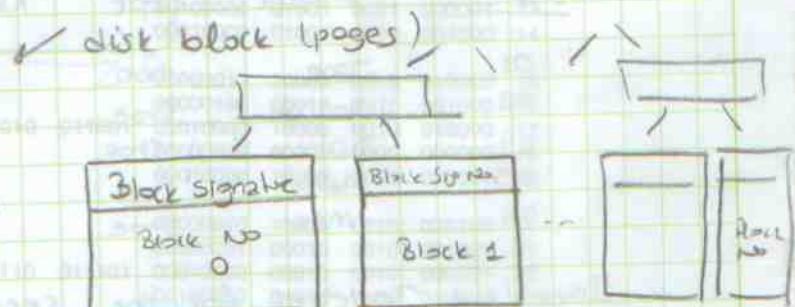
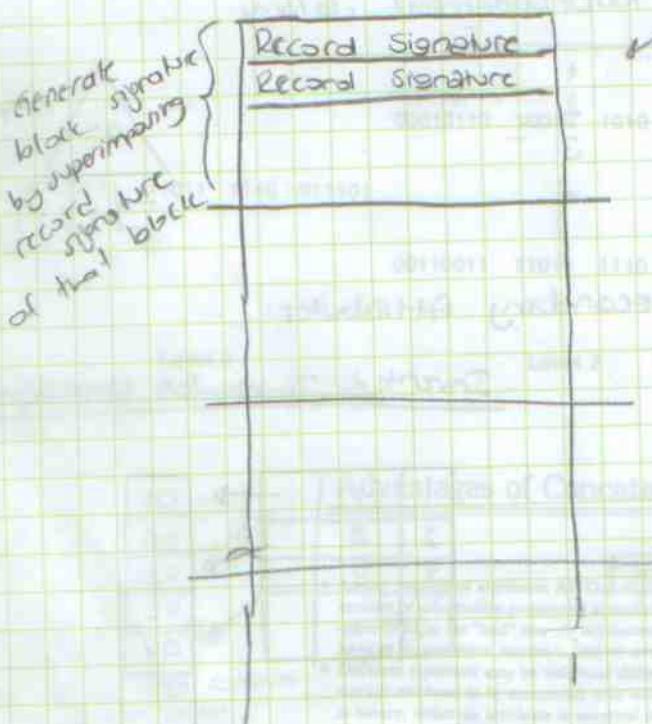


EmpName	$s(1) = 1$	if EmpName begins with [A, E]
	$s(2) = 1$	" [F, L]
	$s(3) = 1$	" [M, Z]
	$s(4) = 1$	if DeptName begins with [A, E]
	$s(5) = 1$	" [F, J]
	$s(6) = 1$	" [K, M]
	$s(7) = 1$	" [N, R]
	$s(8) = 1$	" [S, T]
	$s(9) = 1$	if EmSalary $\leq 10,000$
	$s(10) = 1$	[10000, 20000]
	$s(11) = 1$	[20000, 25000]
	$s(12) = 1$	[25000, 30000]
	$s(13) = 1$	[30000, 45000]
	$s(14) = 1$	$\gg 45000$

EXAMPLE

123 John Doe Toy 35000  
 010 00001 000010

super user  
block sig  
and compare  
our record  
with this  
one



First match every signature with super block signature if it matches consider the next level. --

How many bits to be set for a term?

to produce signatures with

50% 0 and  
50% 1

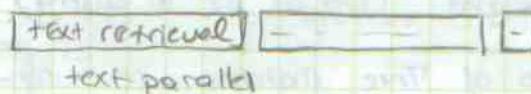
sigmod conf?  
Foliotus

$$\text{sopt(mopt)} = \frac{F \ln 2}{D} \approx \frac{0.7 F}{D}$$

F: signature size in bits

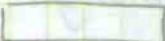
D: avg no of terms / logical block

$\propto D^{1/2}$



D=3

Bs



Another method by Foliotus & Christodoulakis

classifies terms according to their discrimination power (or more important terms  $\equiv$  terms which are more frequently used in queries) are allowed to set more no of bits.

↗ decrease % of false matches

How can we make sure that we have the same no of 0s and 1s in a block signature?

Leng & Lee 1992 ACM TODS  $\leftarrow$  Trans on Database Systems

1.  $S_1 S_2 \dots S_I = 0$  (initialization)

2. while  $sw < (F/2)$  {

take the next term & generate TS (term signature)  
 $S = S \text{ or } TS$

S = block signature

D no of distinct terms

TS<sub>1</sub>

TS<sub>2</sub>

}

TS<sub>D</sub>

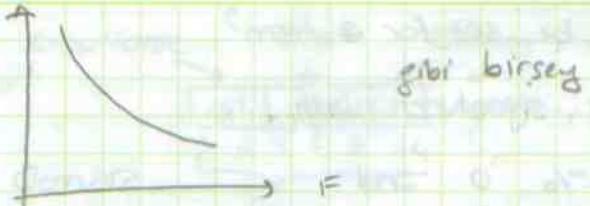
$\underline{\text{or}}$

↑  
Tries to equalize no of 0s and no of 1s in a block signature (equal weight)



64

no of False Match



### How To Calculate No. of False Drop (False Drop Probability)

F : Signature size (in bits)

S : # of bits set / term

D : # of terms / block

$Q_w$  : Query weight (No of 1s in query)

Assumption: No of True Matches << no of Terms

False Drop Probability: Probability of matching all 1s of the query by chance.



$$f_d = (1 - (1 - S/F)^D)^{Q_w}$$

$Q_w$  # of 1s we have in the query bit signature

$S/F$  : Probability of setting a bit position by a block term

$(1 - S/F)$  : Probability of not setting a bit position by a term

$(1 - S/F)^D$  : Probability of not setting a bit position by any of the terms

$(1 - (1 - S/F)^D)^{Q_w}$  : Probability of setting a bit position by any of the terms.

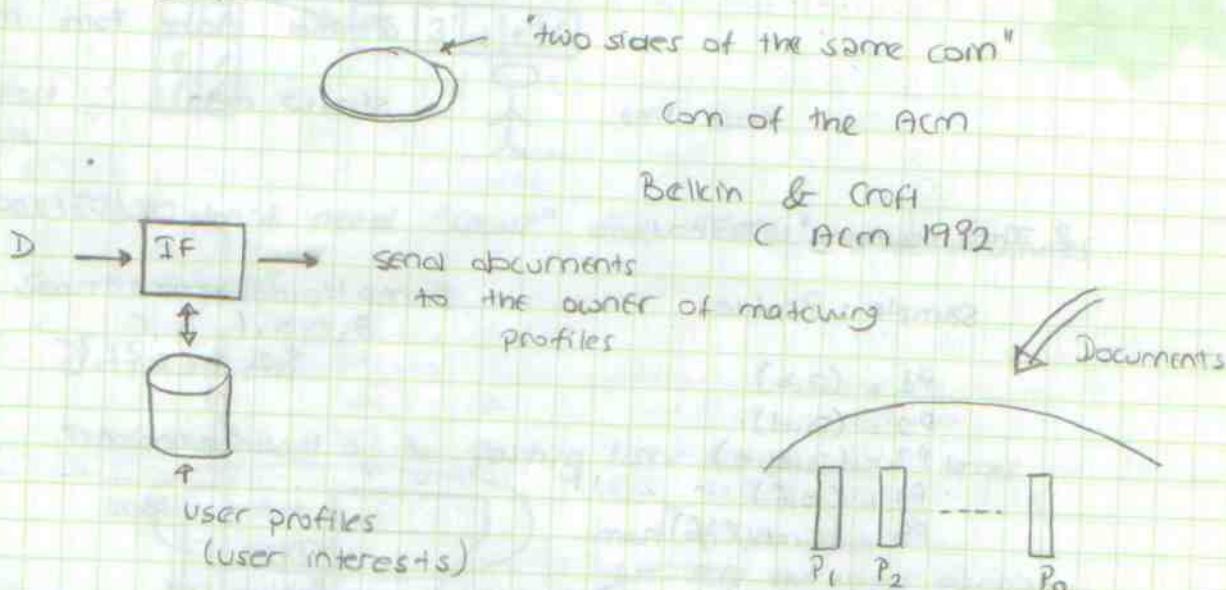
$f_d = (1 - (1 - S/F)^D)^{Q_w}$  : Probability of setting all query bit positions by chance



## INFORMATION FILTERING (IF)

(Selective Dissemination of Information)

IR / IF



Categorization : Supervised clustering

\* TDT : Topic Detection & Tracking (TREC conf.)  
TDT Conference

\* Yan T.W., Garcia molina H.

"Index Structure for Selective Dissemination of Information  
ACM TODS, 1994"

Index Structure For IF :

Assume that user profiles are in conjunctive form

t<sub>1</sub> and t<sub>2</sub> and t<sub>3</sub> ...

1. Brute Force Method :

Compare the incoming document with each profile one by one

Use an occurrence table

Start with the least frequent word in the profile, the next less frequent term so on

This makes it possible to detect non matching profiles fast.

Tan, Garcia, Molina  
Paper

no occurrence Approach table:  
Random-Brute Force Method

with occurrence table  
Ranked-Brute Force Method

base line



straw man

## 2. The Counting Method:

cl 105/2008

### Sample Profiles

$$\begin{aligned}P_1 &= (a, b) \\P_2 &= (a, d) \\P_3 &= (a, d, e) \\P_4 &= (b, f) \\P_5 &= (c, d, e, f)\end{aligned}$$

### Sample Document:

a, c, a, f, b, c  
{a, b, c, f}  
unique ones.

TOTAL

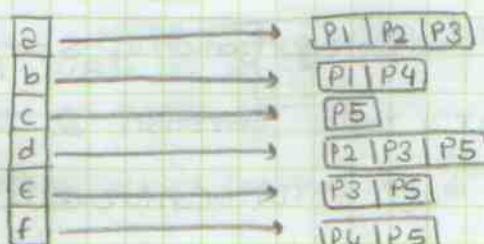
COUNT

DIRECTORY

DISK

Posting List

	2	b
P1	2	x 2
P2	2	x 1
P3	3	x 1
P4	2	0 x 2
P5	4	0 x 2



The incoming document satisfies P1 P4

For each document there is a disk access for retrieving the related (term) posting list.

### # THE KEY METHODS

In the counting Method, a profile  $(w_1, w_2, \dots, w_k)$  appears in  $k$  posting lists.

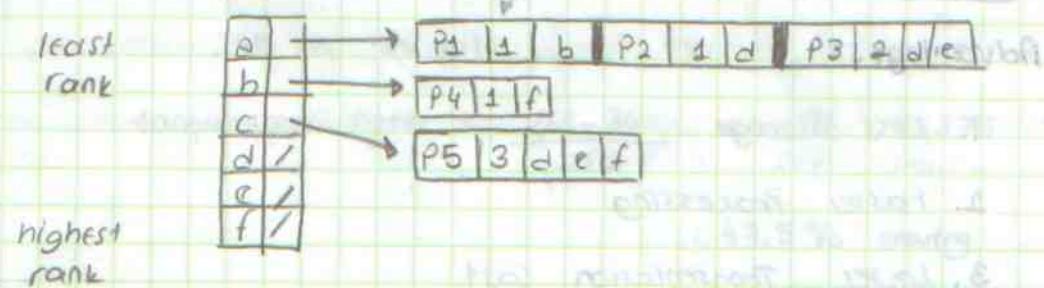
In the key methods a profile only appears in the list of one of its words. That word is called the key.

✓  
RANDOM Key METHOD  
(any word is the key)

↖  
RANKED Key METHOD  
(Store the profile in the list of the word with the lowest rank)  
lowest freq. in docs

most freq. in docs  
highest freq. in docs

300. Directory after 'd' we have 3 more nodes

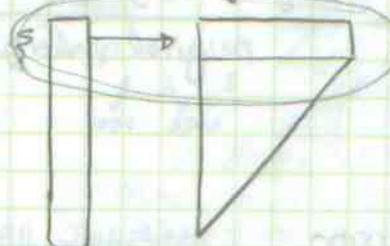


We don't need "COUNT" and "TOTAL" data structures.

### Search Algorithm

$$\mathcal{D} = \{a, b, c, f\}$$

Consider all of the posting lists for document terms

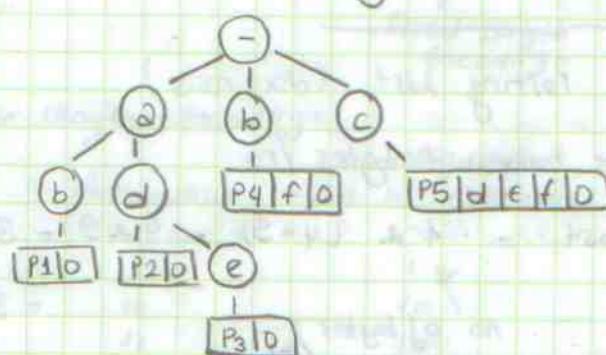


most of the time  
we skip the upper portion  
of the posting list triangle

↳ least freq. old freq. term  
deterministic geometric ratio  
(and consider terms not given by term  
to filter - likely non-existing terms, profillede, soku  
paciyar)

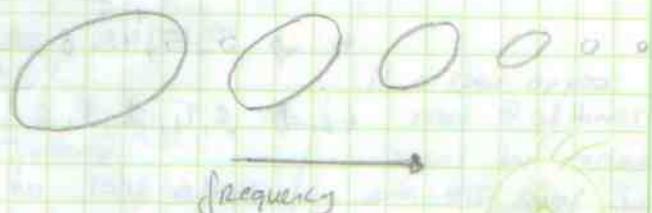
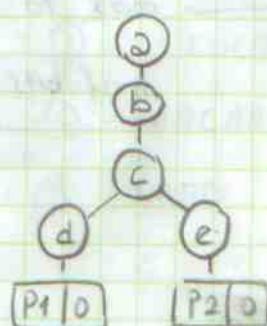
### THE TREE METHODS

Users with similar interest areas would have similar profiles.  
(their profiles would have many common terms)



Consider two profiles:

$(a, b, c, d)$   
 $(a, b, c, e)$



Compression:

08/05/2008

Advantages:

1. Less Storage
2. Faster Processing
3. Lower Transmission Cost

Disadvantages:

Compression → uncompress  
 decompression      cost

## # How to measure compression effectiveness:

$$\text{Compression Ratio} = \frac{\text{Original Length} - \text{Encoded Length}}{\text{Original Length}}$$

## # Types of compression:

Text compression

Inverted Index Compression  
(Postings Lists)

## # Inverted Index Compression:

 $t_1 \rightarrow 5, 10, 22, 30$  $t_2 \rightarrow 3, 7, 9, 10, 16$ 

→

Posting List (doc no.s)

⇒ Assume that we have 4 bytes/no

$$\text{Storage cost} = 4 * (4+5) = 4 * 9 = 36 \text{ bytes}$$

$$\begin{aligned} & \uparrow \\ & \text{no of bytes/int no} \end{aligned} = 36 * 8 = 288 \text{ bits}$$

## # Use Run Length Encoding:

 $t_1 \rightarrow 5, 5, \overbrace{12, 8}^{\text{max no 12}}$  $t_2 \rightarrow 3, 4, 2, 1, 6$ 

$$\# \text{ of bits needed} = \lceil \log_2 12 \rceil = 4$$

$$1100_2 = 12_{10}$$

Cost Using Compression:

288 vs. 86 bits

gives us much

$$\text{Compression Ratio} = \frac{288 - 86}{288} = \frac{252}{288} = 0.875$$

87.5% savings

## # Text Compression Methods:

### SPECIAL PURPOSE COMPRESSION:

#### 1. All numeric fields

123  $\Rightarrow$  F1 F2 F3  
1111 0001

EBCDIC representation

Extended Binary Coded Decimal Interchange Code

Drop F's

1 2 3  
0001 0010

50% savings

ASCII
30 31 33
0011 0010

#### 2. All alphabetic

26 letters + blank

$$\lceil \log_2 27 \rceil = 5 \quad 2^5 = 32$$

27 different characters

$$5 \text{ vs. } 8 \Rightarrow \frac{5}{8} 0.625 \approx 37.5\% \text{ savings}$$

Fixed Length Encoding

## # Variable Length Encoding:

unambiguous

Instantaneous

Char

Ambiguous

not instantaneous

(Huffman Coding)

unambiguous

Char

A

1

1

0

B

10

10

10

C

11

100

110

D

100

1000

111

11100, 111  
C D C A

①

COCA or AACAC

②

AACAAA

7 words.

③

DAAD

No look-ahead.

if a char appear more than once  
should be represented with less than  
& vice versa

for this

Huffman Coding

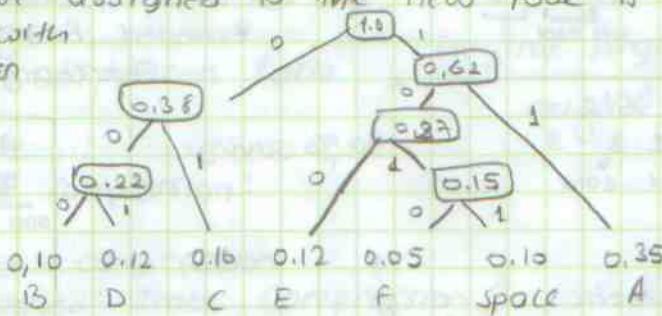
## Huffman Coding:

We have to know the relative frequency of the characters.

Alphabet	F	:	0,05	$\rightarrow$	only 5% of the characters are F
	space	:	0,10		
	B	:	0,10		
	D	:	0,12		
	E	:	0,12		
	C	:	0,16	$\rightarrow$ 2nd most frequent text if	
	A	:	0,35		

We construct a tree (Huffman Coding Tree) in an iterative manner. At each iteration we create a new node having as children two nodes with the smallest possibility values.

The value assigned to the new node is the sum of the values associated with its children.



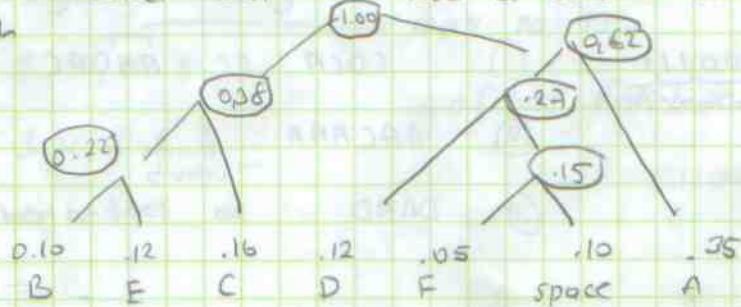
B : 000 (3)      F : 1010 (4)  
 D : 001 (3)      space : 011 (3)  
 C : 01 (2)      A : 11 (2)  
 E : 100 (3)

A    B    D    FEED  
 11    000    11

The average weighted code length:  $\sum_{i=1}^n (\text{length of char-}i) * \text{freq.}$

$$\begin{aligned}
 & \text{B} \quad \text{D} \quad \text{C} \\
 & = (3 \times .10) + (3 \times .12) + (2 \times .16) \\
 & = 2.64 \text{ bits}
 \end{aligned}$$

Huffman code is optimal, i.e. there is no other assignment of codes to select that will have a shorter average weighted code length.



NOT : Gamma Beta 'da sorumluluğuz!  
 Ders Notları 'Part-2' den ulaşılabilir.

Please read!  
 Y S will cont.  
 for Buna hazırla!