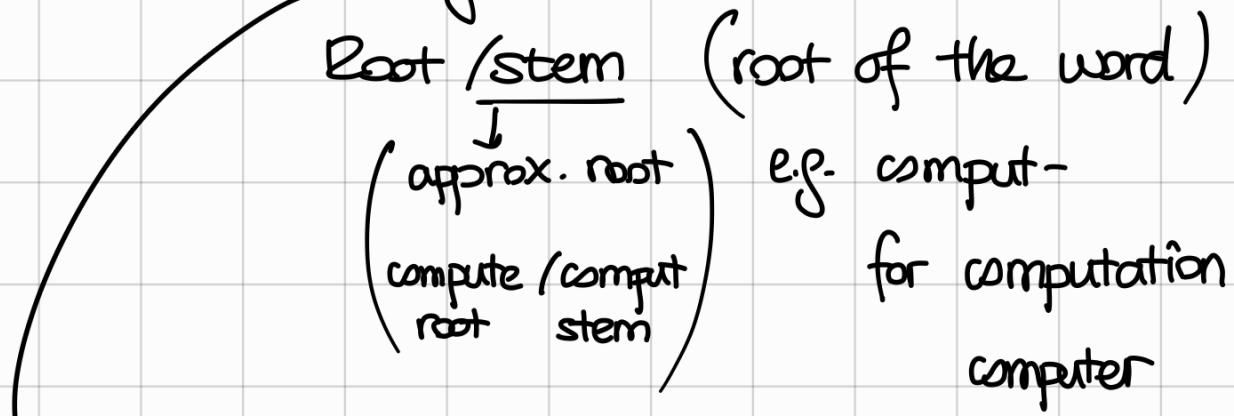


} ranking based on relevance  
 and user profile  
 } information bubble  
 (filter bubble)

Words → Indexing terms



### indexing:

- determines words (terms) to be used for the description of documents.
- assigns importance (weights) to terms.

Program: Porter's stemming algorithm

Lemmatization: indexing w/o using the orig. word

e.g. good, better, best

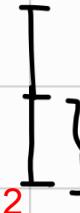
e.g. Zemberek

Stopwords: ∃ stopword lists. In those lists there are commonly used words (function words).

e.g. bir, ve, bu, ki, mi, sonra . . .



30-40%  
stopwords



50% very infreq.  
words

Text

$w_1 \ w_2 \ w_3 \ w_1 \ w_2 \ w_3$  (w:word)

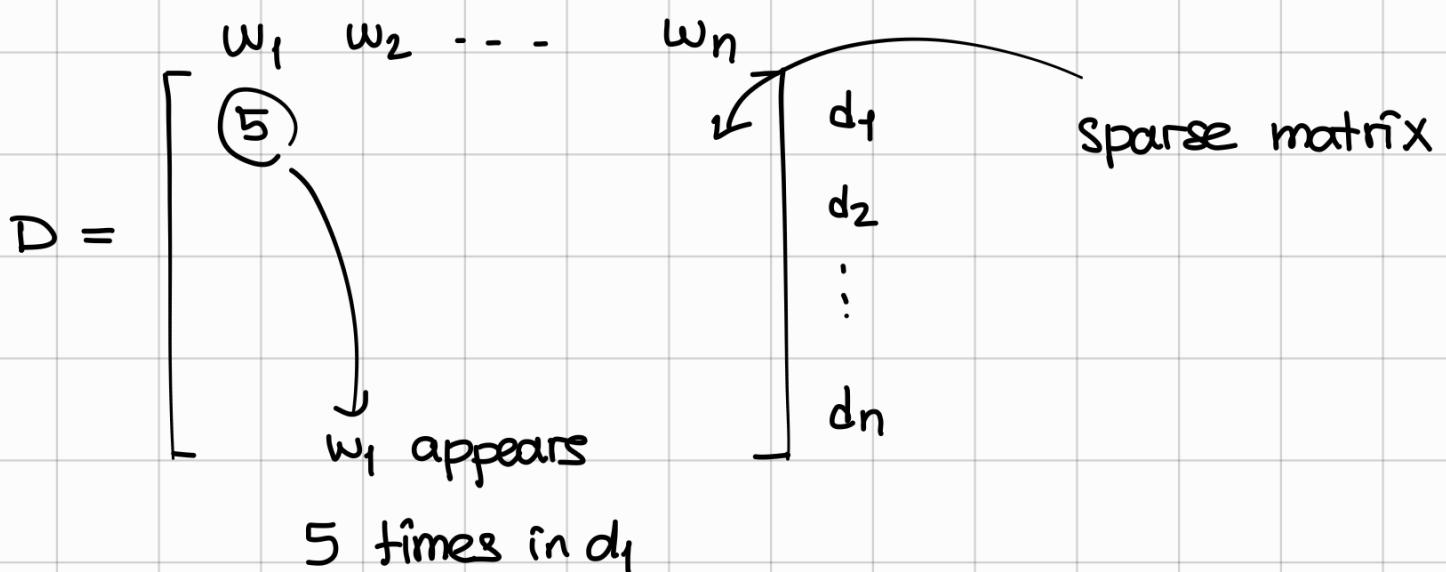
Doc 1

Doc 2

⋮

Doc n

### Document Term Matrix



TF-IDF: term freq. - inverse document freq.

Gerard Salton & Chris Buckley ?

Term weight assignment<sub>3</sub>

# Medical Collections

MEDLARS, MEDLINE



**Heap's Law:** # of docs ↑, # of terms ↑  
(but at a slower rate)

**Zipf's Law:** freq. of word  $\propto \frac{1}{\text{rank of word}}$

IR System **Evaluation**

(evaluate based on)

Effectiveness

&

Efficiency

Relevant docs

time memory

## Measuring Performance

Online

User satisfaction

Lab

Eval is based on test collection

Test Collection  $\supset$  A doc set (collection)

A query set (query collection)

Relevant docs / queries

TREC (Text Retrieval Conf.)

## Pooling

Group 1  $q_1 \rightarrow d_1, d_5, d_7$

Group 2  $q_1 \rightarrow d_5, d_{10}$

$\rightarrow q_1 \rightarrow d_1, d_5, d_7, d_{10}$

Fleiss Kappa Measure  
to judge how dependable  
the decisions are.

- Annotator, assessor

- Limitation of pooling: There might be some docs  
that are not shown to any assessor & assumed  
irrelevant.

## trec eval:

"TRECEval" refers to the software tool used for evaluating the performance of information retrieval systems in the context of TREC. It provides the means to compute various evaluation metrics, such as precision, recall, F1 score, and mean average precision (MAP), which are commonly used to assess the quality of retrieval results.

26.09.2023

## Data set

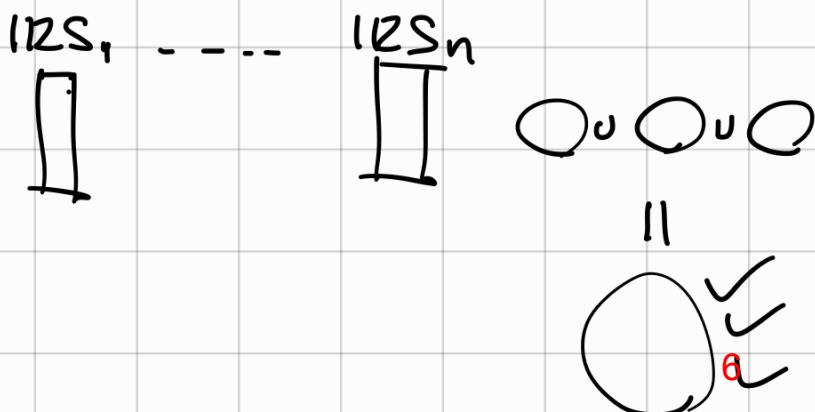
Doc. collection

query set

Rel. docs for each query

- Use test collections for reproducibility.

Pooling → TREC  
✓



Pooling in information retrieval systems is a process used to combine and rank the results obtained from multiple retrieval models or algorithms in order to generate a final ranked list of documents for a given query.

## Effectiveness Measures

- Precision, recall, f measure
- positive: deemed relevant by IRS

✓	TP (true pos.)
✗	FP (false pos.)
✓	TP
✗	
✓	
✓	

---

$$\frac{4}{6} = \frac{\# \text{ rel. docs}}{\# \text{ docs we observed}}$$

$$\text{Precision} = \frac{TP}{TP+FP}$$

$$\text{Recall} = \frac{TP}{TP+FN}$$

↓ TN  
FN  
(false neg: rel. docs that are not shown in / retrieved by the system)

R @ 10

Prec @ 5  
Prec @ 10  
Prec @ 15  
# of docs

IRS1                    IRS2  
q<sub>1</sub>                0.5                0.7

paired t-test

q<sub>750</sub>                0.3                0.4

p=0.05

Avg                0.75                0.7

ANOVA

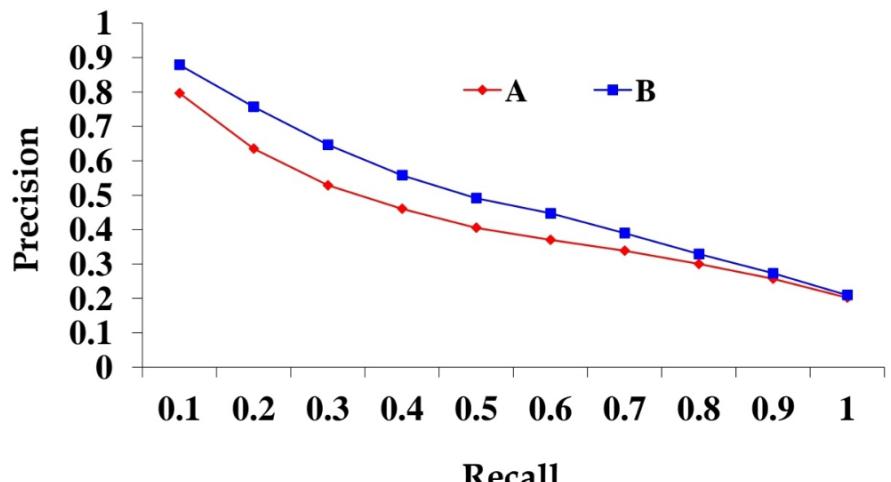
TODD. check ANOVA & MONOVA

$$F = \frac{2PR}{P+R}$$

P:precision, R:recall

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	5/8	5/9	5/10
Recall	0/10	1/10	1/10	2/10	3/10	4/10	5/10	5/10	5/10	5/10

# rel. docs = 10



→ B is better than A.

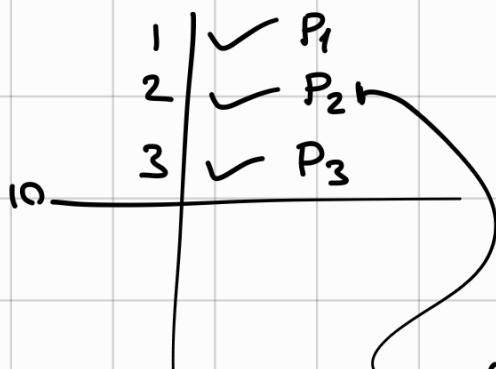
### Interpolation of Prec. Values

- The highest prec. val. after that we call pts is used as the prec. val..



## MAP (Mean Avg Precision)

- We know # rel. docs



(not necessarily  
the second doc. might  
be 7<sup>th</sup> doc.)

$$\text{MAP}_{10} = \frac{P_1 + P_2 + P_3}{\# \text{ rel. docs}}$$

the pos. where  
we cut

$$\text{MAP}_n = \frac{\sum_i P_i \text{ st } i \in \{\text{rel. docs}\} \cap \{\text{top } n \text{ docs}\}}{\# \text{ rel. docs in top } n \text{ docs}}$$

29.09.2023

- There are diff. def's of MAP.

We know the rel. docs, we haven't seen all of them.

Consider top 10 (ranked) docs

Assume that  $\exists 5$  rel. docs (3 rel. doc. in top 10)

	1	2	3	4	5	6	7	8	9	10
x		✓			✓			✓		
x										
x										

$P_2 = \frac{1}{2}$

$P_5 = \frac{2}{5}$

$P_8 = \frac{3}{8}$

$$\text{MAP}_{10} = \frac{P_2 + P_5 + P_8}{\# \text{ rel. docs for this query}}$$

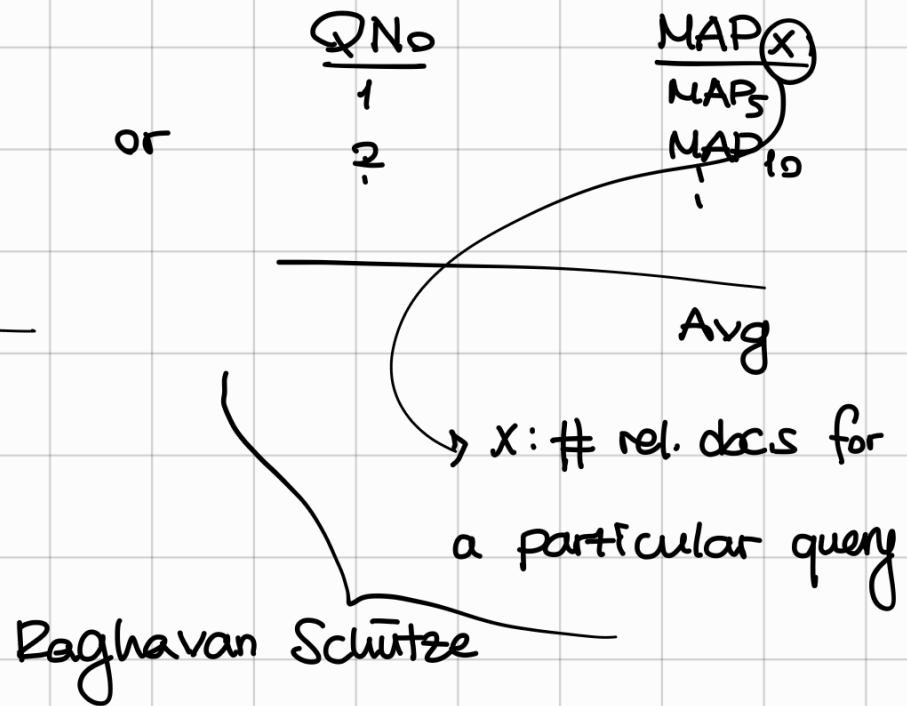
$$= \frac{0.5 + 0.4 + 0.375}{5}$$

## Performance of a System

<u>Q No</u>	<u>MAP<sub>10</sub></u>
1	...
2	...
3	...

Avg

or



- Info. Ret. - Schütze, Search Engines - Croft, Stohlmeyer

" " - van Rijnsbergen  
 $\exists$  evaluation cpt.

03.10.2023

- w/ 1000000 docs we might use only 100 docs.  $\Rightarrow$  limited accumulators
- Paper: Inverted Index Structures for Search Engines by Zohel, Moffad.

### Similarity Calculations (among the docs of a collection)

$$D = \begin{bmatrix} t_1 & \dots & t_n \\ & \vdots & \\ & \text{sparse} & \\ & d_1 & \\ & & d_m \end{bmatrix}$$

↓ doc. term matrix

$$S = \begin{bmatrix} 1 & S_{12} & S_{13} & S_{14} \\ S_{21} & 1 & S_{23} & S_{24} \\ S_{31}, S_{32} & & 1 & S_{34} \\ S_{41}, S_{42}, S_{43} & & & 1 \end{bmatrix}$$

similarity matrix

$$S(d_i, d_i) = 1$$

$S_{ij} = S_{ji}$  (symmetric) (there can be an asym. formula as well)

-  $S$  can be used in clustering docs

$$\begin{bmatrix} 1 & S_{12} & S_{13} & S_{14} \\ S_{21} & 1 & S_{23} & S_{24} \\ S_{31}, S_{32} & & 1 & S_{34} \\ S_{41}, S_{42}, S_{43} & & & 1 \end{bmatrix}$$

m-1 docs

$\frac{m \cdot (m-1)}{2}$  calculations

11

- Matrix D is weighted / binary

- Dice coefficient

- Cosine similarity

- Dot product

- Jaccard

:

Similarity Coef.

Dot product

(Inner product)

Binary

$X \cap Y$

Weighted

$\sum X_i Y_i$

Cosine

$$\frac{|X \cap Y|}{|X|^{\frac{1}{2}} \cdot |Y|^{\frac{1}{2}}}$$

$$\frac{\sum X_i Y_i}{\sqrt{\sum X_i^2 \sum Y_i^2}}$$

Dice

$$\frac{2|X \cap Y|}{|X| + |Y|}$$

$$\frac{2 \sum X_i Y_i}{\sum X_i^2 + \sum Y_i^2}$$

Jaccard

$$\frac{|X \cap Y|}{|X| + |Y| - |X \cap Y|}$$

$$\frac{\sum X_i Y_i}{\sum X_i^2 + \sum Y_i^2 - \sum X_i Y_i}$$

( i: index of the  $i^{th}$  term )

## Example

$$X = (1 0 1 1 1) \quad |X| = 4$$

$$Y = (\underline{1} \ 1 \ 0 \ \underline{1} 0) \quad |Y| = 3$$

$\downarrow$   $\downarrow$

2 common items =  $X \cap Y$

Dot product 2

Cosine  $1/\sqrt{3}$

Dice  $4/7$

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

## Example

$$X = (2 \ 0 \ 1 \ 3 \ 2) \quad \text{Dice} = \frac{2(2+2+3+10)}{(4+1+9+4)+(1+0+4+1+25)} = 0.69$$

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

Cosine  $\approx 0.72$

- When calculating the S matrix:

1. Brute force approach

$$\# \text{ of entries to be calculated} = \frac{m \cdot (m-1)}{2}$$

2. Use the knowledge of term dist. in doc

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

$m = 5$  # of docs  
 $n = 6$  # of terms



$\exists$  at least one "1" in each row & col.

represent D like this instead:

$$t_1 \rightarrow d_1, d_2 \quad t_4 \rightarrow d_2, d_5$$

$$t_2 \rightarrow d_1, d_2 \quad t_5 \rightarrow d_1, d_2$$

$$t_3 \rightarrow d_4, d_5 \quad t_6 \rightarrow d_3, d_4, d_5$$

Consider  $d_1$ :  $\exists t_1, t_2, t_5$ .

$$\begin{array}{c} d_1 \\ d_2 \end{array} \cup \begin{array}{c} d_1 \\ d_2 \end{array} \cup \begin{array}{c} d_1 \\ d_2 \end{array} = \begin{array}{c} d_1 \\ d_2 \end{array} \rightarrow \text{so only } d_2 \text{ can have non-zero similarity with } d_1.$$

Calculate  $S_{12}$ .

Similarly for  $d_2$ :  $\{d_1, d_2, d_5\} \Rightarrow$  Calc.  $S_{25}$ .

$d_3$ :  $\{d_3, d_4, d_5\} \Rightarrow$  "  $S_{34}, S_{35}$ .

$d_4$ : { " " = }  $\Rightarrow$  =  $S_{45}$ .

$d_5$ :  $\{d_2, d_3, d_4, d_5\} \Rightarrow S_{25}, S_{35}, S_{45}$  are already calculated.

5 calcs =  $S_{12}, S_{25}, S_{34}, S_{35}, S_{45}$ .

rs

$\frac{4 \cdot 5}{2} = 10$  calcs brute force.

Continuing from last lecture's example:

-  $x_d$  = depth of indexing  
(avg. # of terms / doc)

-  $t_g$  = term generability  
(avg. # of docs / term)

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

### 3. Using the inverted index file:

Inverted Index

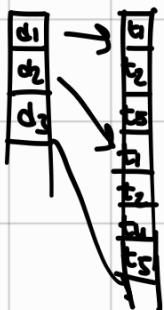
$d_1 \rightarrow t_1, t_2, t_5$

$d_2 \rightarrow t_1, t_2, t_4, t_5$

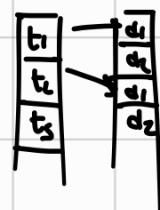
$t_1 \rightarrow d_1, d_2$        $t_4 \rightarrow d_2, d_5$

$t_2 \rightarrow d_1, d_2$        $t_5 \rightarrow d_1, d_2$

$t_3 \rightarrow d_4, d_5$        $t_6 \rightarrow d_3, d_4, d_5$



Inverted file:



posting list

$$\text{Dice coef.} = \frac{2|X \cap Y|}{|X| + |Y|}$$

Document length array: 

3	4	1	2	3
---	---	---	---	---

Consider  $d_1$ .

similarity array : 

$s_{11}$	$s_{12}$	$s_{13}$	$s_{14}$	$s_{15}$
x	0	0	0	0

$t_1 \quad x$   
 $t_2 \quad x$        $\Rightarrow s_{12} = \frac{2 \cdot 3}{3+4}$   
 $t_3 \quad 3$   
 $t_4 \quad -$   
 $t_5 \quad -$

$d_1$  contains these terms

$s_{13} = s_{14} = s_{15} = 0.$

Similarly for  $d_2$ .

$s_{21}$	$s_{22}$	$s_{23}$	$s_{24}$	$s_{25}$
x	x	0	0	0

$t_1 \quad - \quad - \quad -$   
 $t_2 \quad - \quad - \quad -$        $\Rightarrow s_{22} = s_{24} = 0.$   
 $t_3 \quad -$       1       $s_{25} = \frac{2 \cdot 1}{4+3}$   
 $t_4 \quad - \quad - \quad -$   
 $t_5 \quad - \quad - \quad -$

Similarly for  $d_3$ .

$t_6$ 

x	x	x	0	0
---	---	---	---	---

$\Rightarrow s_{34} = \frac{2 \cdot 1}{1+2}$   
 $s_{35} = \frac{2 \cdot 1}{1+3}$

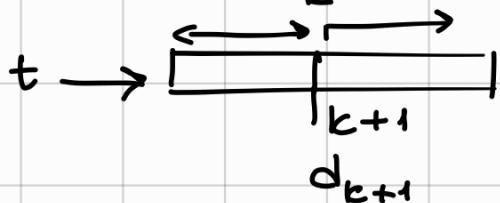
- Consider  $m$  docs

$\forall$  doc. consider  $x_d$  # terms

$\forall$  term consider  $t_g$  # docs

$$O(m \cdot x_d \cdot t_g)$$

Consider  $d_k$ . Consider a term  $t$  of this doc.



for  $d_1$ , only ignore

$d_1$

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

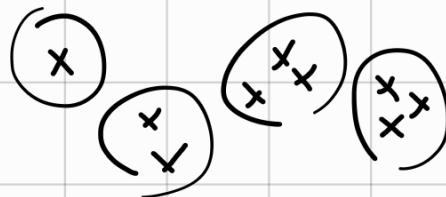
for  $d_2$ , ignore

$d_1$  and  $d_2$

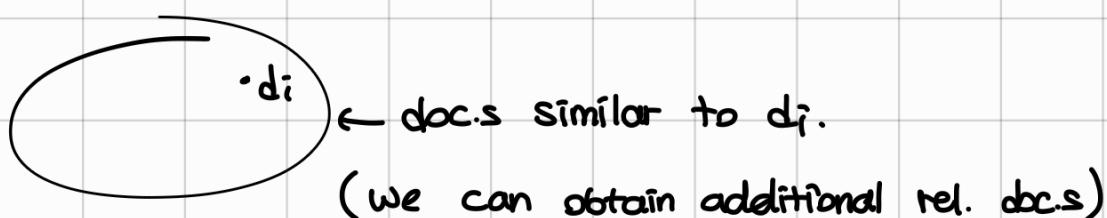
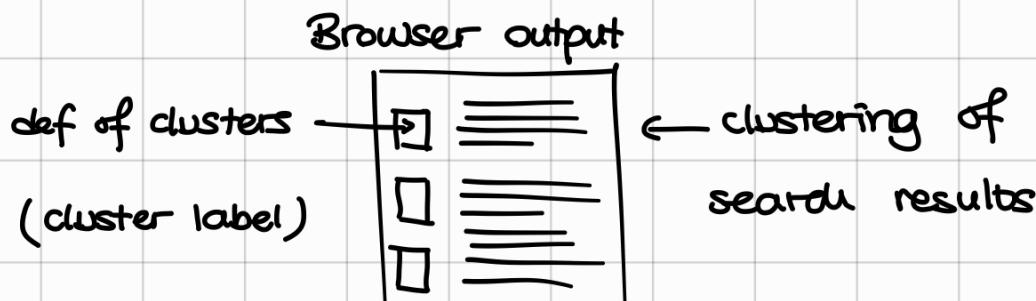
$$= x_d \cdot t_g \cdot \frac{(m-1)m \cdot 1}{2m} = x_d \cdot t_g \cdot \frac{(m-1)}{2} = O(x_d t_g m)$$

- Difference between 2 and 3 is we calculate  $|X \cap Y|$  (of all the docs that share terms with  $d_i$ ) while calculating the similarity array of  $d_i$ .

## Clustering



- How to use clustering for information retrieval?



## Cluster Based Retrieval

Identify best matching clusters

Implementation:

$$C = \begin{bmatrix} 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 \end{bmatrix} = \begin{bmatrix} \frac{1+1}{2} & 0 & \frac{1+1}{2} & \frac{1+0}{2} & 0 \\ 1 & 0 & 1 & 0.5 & 0 \end{bmatrix}$$

forall cluster obtain a centroid ←

Compare query w/ centroid

Rank clusters according to their similarity to the query

most similar cluster

least similar cluster

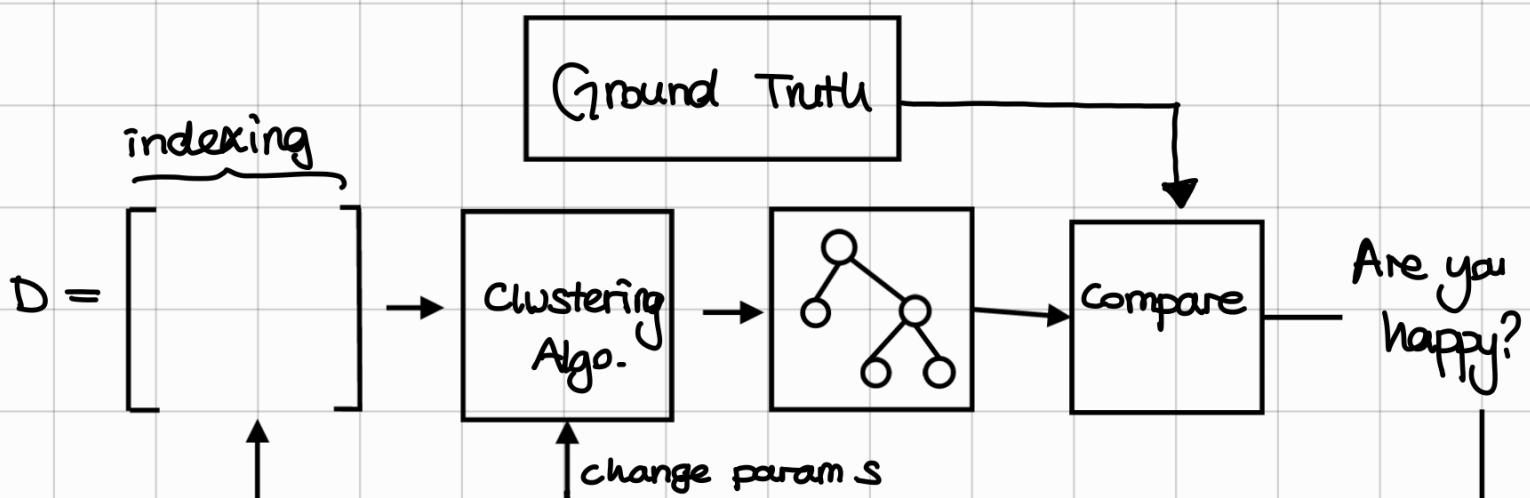
compare your query w/ cluster member

- Classification : (classes are known) supervised
- Clustering : unsupervised
  - ↳ after clustering we need to label (give summary tag) to clusters.

Cluster Hypothesis:

Docs rel. to the same query would appear in the same cluster.

- Tom Hope , Inflection in Scientific Discoveries.



Robustness      small changes in D matrix  
                   shouldn't change clusters much

## Classifying Clustering Algorithms

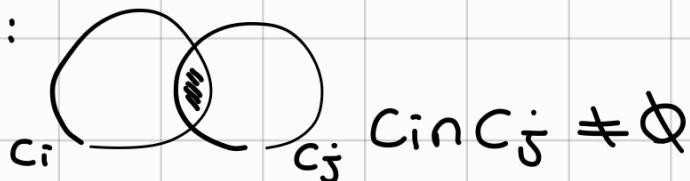
### • Structure:

According to the structure generated:

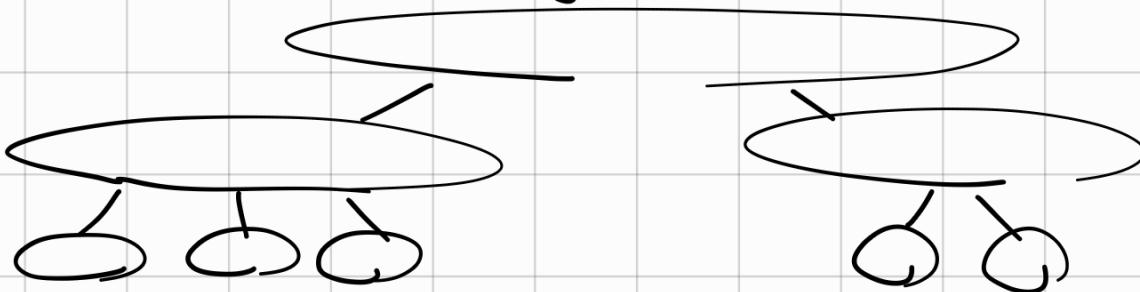
1. Partitioning: clusters do not have common items.

$$C_i \cap C_j = \emptyset \quad \forall i \neq j.$$

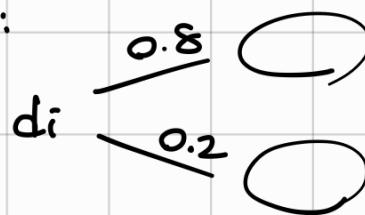
2 Overlapping:



3. Hierarchical clustering:



4. Fuzzy clustering:



10.10.2023

### • Procedure

According to how the clusters are generated:

1. Single-pass:

each item is considered only once.

- select some data item as cluster initiator (seed)

Assign non-seeds to these clusters

→ hand generated v select some of them? v random

## 2. Multi-pass:

obtain initial clusters. then try to improve them.

## 3. Heuristic algo.s:

First doc inits a cluster

Consider next one

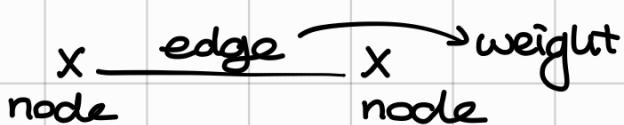
If it's diff. it starts its own cluster

- order dependent

- good for new event detection and tracking

## 4. Hierarchical:

- Graph theoretical:

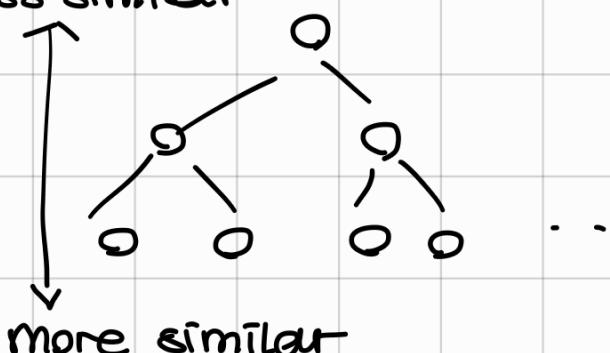


- single link

- complete link

- avg. link

less similar

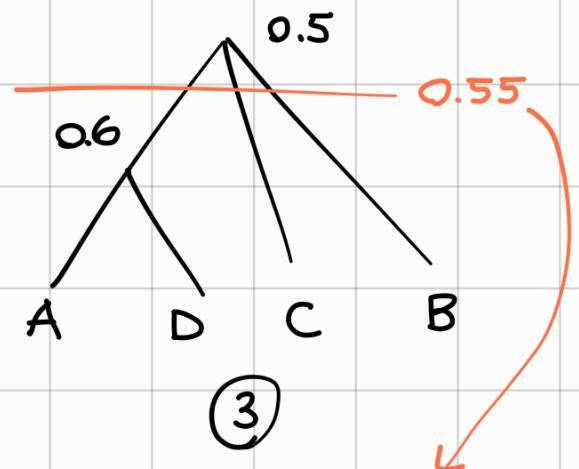
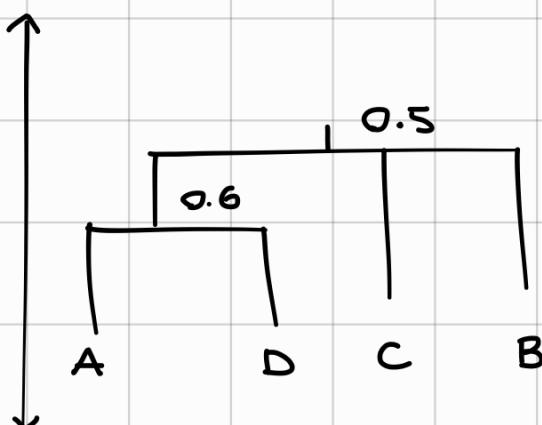
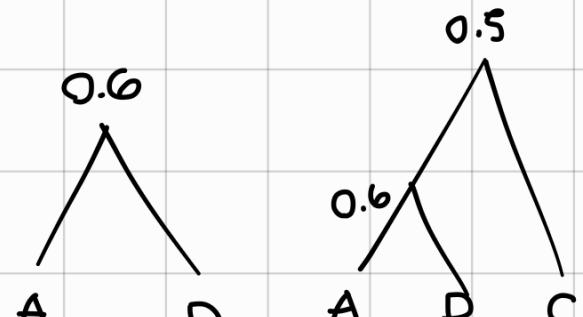


$$S = \begin{bmatrix} & A & B & C & D \\ A & 1.0 & 0.3 & 0.5 & 0.6 \\ B & 1.0 & 0.4 & 0.5 & \\ C & 1.0 & 0.3 & \\ D & 1.0 & & \end{bmatrix}$$

i. Single-link: - not order dependent

step    pair    sim. 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



partitioning threshold

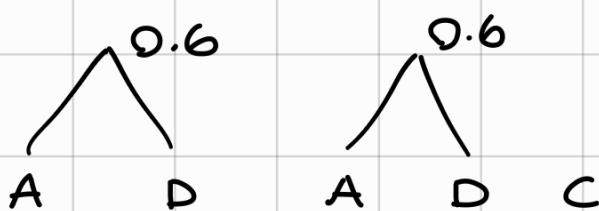
Cluster Validation:

(AD) (C) (B)

$$S' = \begin{bmatrix} A & B & C & D \\ 1.0 & 0.5 & 0.5 & 0.6 \\ 1.0 & 0.5 & 0.5 & \\ 1.0 & 0.5 & \\ 1.0 & \end{bmatrix} \quad \begin{array}{l} A \\ B \\ C \\ D \end{array}$$

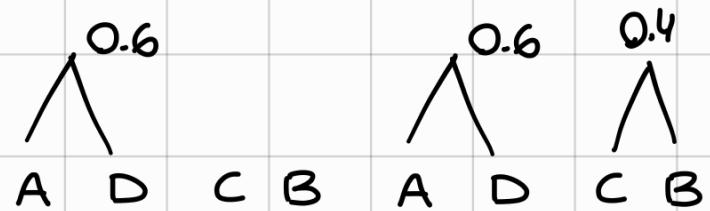
( Single-link: the similarity between a cluster and new item is taken as the most similar member of that cluster for that item. )

## ii. Complete Link:



$$C: AC = 0.5$$

①



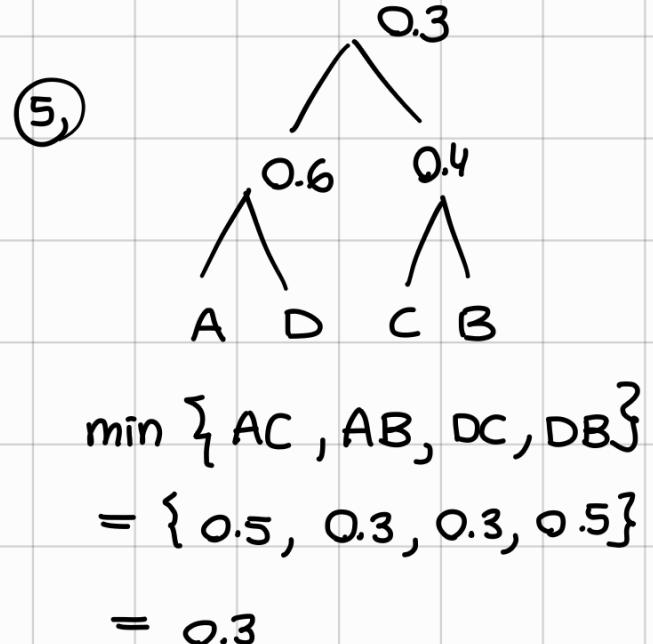
$$BD = 0.5$$

③

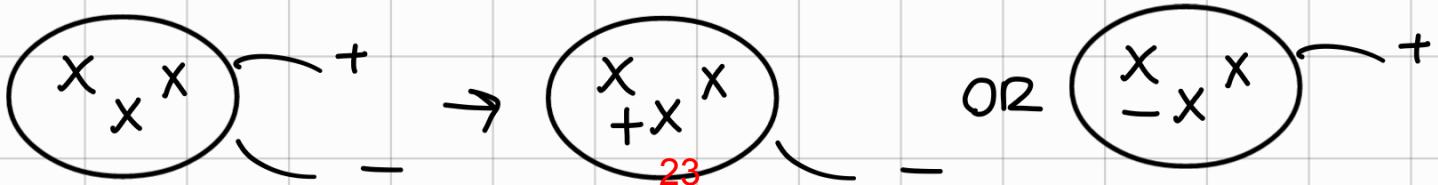
$$CB = 0.4$$

④

- AD	0.6
- AC	0.5
- BD	0.5
- BC	0.4
AB	0.3 ] intercluster
CD	0.3 ] intercluster



- final clustering structure is order dependent



## Characteristics of Desirable Clustering Algorithms

1. Effectiveness & Efficiency -
    - time
    - meaningful clusters
    - memory
  2. Maintainable: we prefer an incremental approach
    - shouldn't compute from scratch w/ each new data)
    - reclustering vs incremental clustering
  3. Order independent
  4. Robustness: small error in data shouldn't effect the clustering structure
  5. Addition of new data shouldn't destroy the existing structure
  6. Cluster sizes are comparable (partitioning)
  7. Small number of hyperparameters

e.g. Fazli hoca's paper : Cover Coef.-based Clustering  
Methodology C<sup>3</sup>M

Partitioning

Seed-oriented:

1. Find # clusters

2. Select cluster seeds

3. Assign nonseeds to clusters  
 based on their similarity to the  
 cluster seeds

C: cover coef. matrix

$$D = \begin{bmatrix} & t_1 & t_n \\ d_1 & & \\ & & \\ d_m & & \end{bmatrix}_{m \times n}$$

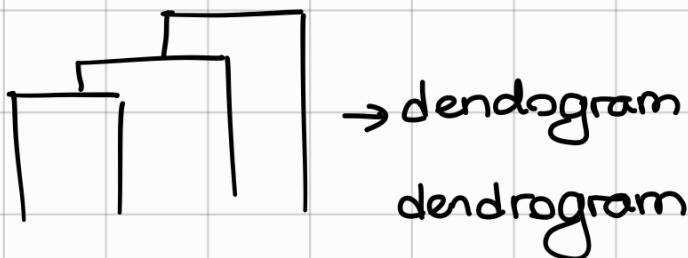
$$n_c = n_i \rightarrow C = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}_{m \times m}$$

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

$$C' = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}_{n \times n}$$

$$n'_c = \sum_{i=1}^m c'_{ii}$$

13.10.2023



Cover Coef.-based Clustering Methodology (C<sup>3</sup>M)

$$D = \begin{bmatrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 \\ \hline 1 & 0 & 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}_{m \times n}$$

$$\alpha_1 = 1/3 \quad \alpha_2 = 1/4 \quad \alpha_3 = 1/3$$

$$\alpha_4 = 1/2 \quad \alpha_5 = 1/2$$

$$\beta_1 = \frac{1}{4} \quad \beta_2 = \frac{1}{1} \quad \beta_3 = \frac{1}{2} \quad \beta_4 = \frac{1}{2} \quad \beta_5 = \frac{1}{2} \quad \beta_6 = \frac{1}{3}$$

$\alpha_i : \frac{1}{\# \text{ terms in } d_i}$

$\beta_j : \frac{1}{\# \text{ docs term } j \text{ appears in}}$

$$S = \begin{bmatrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 \\ \alpha_1 & 0 & 0 & \alpha_1 & 0 & \alpha_1 \\ \alpha_2 & \alpha_2 & \alpha_2 & \alpha_2 & 0 & 0 \\ \alpha_3 & 0 & 0 & 0 & \alpha_3 & \alpha_3 \\ 0 & 0 & 0 & 0 & \alpha_4 & \alpha_4 \\ \alpha_5 & 0 & \alpha_5 & 0 & 0 & 0 \end{bmatrix} \begin{matrix} d_1 \\ d_2 \\ d_3 \\ d_4 \\ d_5 \end{matrix} = \begin{bmatrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 \\ \frac{1}{3} & 0 & 0 & \frac{1}{3} & 0 & \frac{1}{3} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 & 0 \\ \frac{1}{3} & 0 & 0 & 0 & \frac{1}{3} & \frac{1}{3} \\ 0 & 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 & 0 \end{bmatrix} \begin{matrix} d_1 \\ d_2 \\ d_3 \\ d_4 \\ d_5 \end{matrix}$$

$\alpha$

the same as dividing  
each elem. in a row by  
the row sum. (elem  $\rightarrow \frac{\text{elem}}{\text{row sum}}$ )

$$S' = \begin{bmatrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 \\ \beta_1 & 0 & 0 & \beta_4 & 0 & \beta_6 \\ \beta_1 & \beta_2 & \beta_3 & \beta_4 & 0 & 0 \\ \beta_1 & 0 & 0 & 0 & \beta_5 & \beta_6 \\ 0 & 0 & 0 & 0 & \beta_5 & \beta_6 \\ \beta_1 & 0 & \beta_3 & 0 & 0 & 0 \end{bmatrix} \begin{matrix} d_1 \\ d_2 \\ d_3 \\ d_4 \\ d_5 \end{matrix} = \begin{bmatrix} t_1 & t_2 & t_3 & t_4 & t_5 & t_6 \\ \frac{1}{4} & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{3} \\ \frac{1}{4} & 1 & \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ \frac{1}{4} & 0 & 0 & 0 & \frac{1}{2} & \frac{1}{3} \\ 0 & 0 & 0 & 0 & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{4} & 0 & \frac{1}{2} & 0 & 0 & 0 \end{bmatrix} \begin{matrix} d_1 \\ d_2 \\ d_3 \\ d_4 \\ d_5 \end{matrix}$$

$\beta$

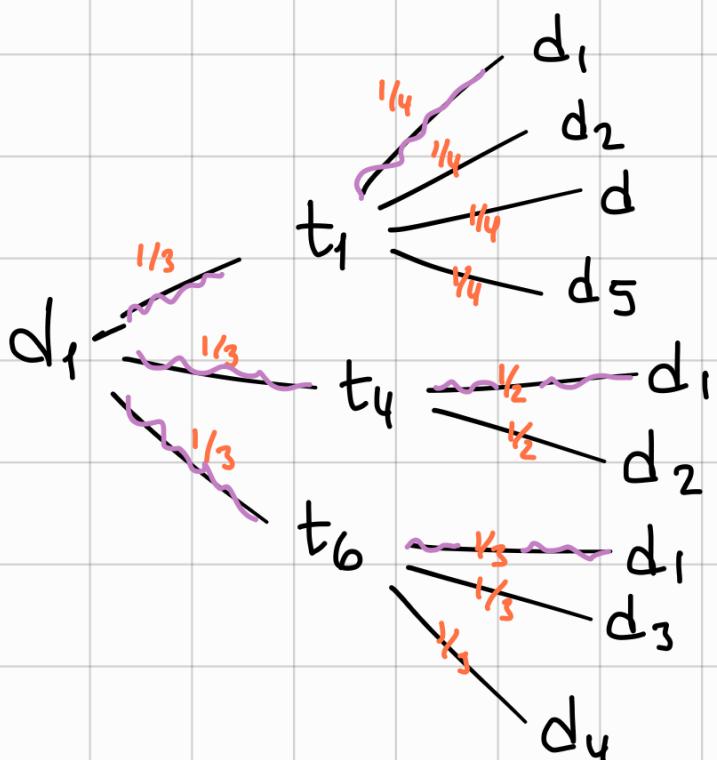
↓ (elem  $\rightarrow \frac{\text{elem}}{\text{col. sum}}$ )

$$S^T = \begin{bmatrix} \beta_1 & \beta_1 & \beta_1 & 0 & \beta_1 \\ 0 & \beta_2 & 0 & 0 & 0 \\ 0 & \beta_3 & 0 & 0 & \beta_3 \\ \beta_4 & \beta_4 & 0 & 0 & 0 \\ 0 & 0 & \beta_5 & \beta_5 & 0 \\ \beta_6 & 0 & \beta_6 & \beta_6 & 0 \end{bmatrix}$$

$$C = S \cdot S^T$$

m: #docs

$$C = \left[ C_{11} = \frac{1/3 \cdot 1/4 + 0 + 0 + 1/3 \cdot 1/2 + 0 + 1/3 \cdot 1/3}{1/3} \quad C_{12} \dots \right]$$



$$C_{ij} = \alpha_i \sum_{k=1}^n d_{ik} B_k d_{jk}$$

$$\begin{aligned}
 C_{ij} &= S_i^{\text{row}} \times S_j^{\text{col.}} = \sum_{k=1}^n S_{ik} \times S_{kj}^T \\
 &= \sum_{k=1}^n (\alpha_i \cdot d_{ik}) \cdot (\beta_k \cdot d_{kj})^T \\
 &= \sum_{k=1}^n \alpha_i \cdot d_{ik} \cdot \beta_k \cdot d_{kj} \circ
 \end{aligned}$$

- $C_{ii} \downarrow$  if  $d_i$  has a small # common terms w/ the rest of the collection.

- $n_c = \sum_{i=1}^m C_{ii} = \# \text{clusters}$

- If all does are unique:  $d_i \cap d_j = \emptyset \quad \forall i \neq j$

$$C = \begin{bmatrix} 1 & \dots & 0 \\ 0 & \ddots & 0 \\ \vdots & \vdots & \vdots \\ 0 & \dots & 1 \end{bmatrix} \quad n_c = m.$$

- If all does are identical:  $d_i \cap d_j = d_i = d_j \quad \forall i, j$

$$C = \begin{bmatrix} 1/m \\ \vdots \\ 1/m \end{bmatrix} \quad n_c = 1.$$

- decoupling coef. for doc  $i = \delta_i = C_{ii}$ .

$$\text{avg. decoupl. coef.} = \bar{\delta} = \frac{\sum C_{ii}}{m} = \frac{\sum \delta_i}{m}$$

- row  $i$ :  $C_{ii} = C_{ij}$  if  $d_i = d_j$

$$C_{ii} \geq C_{ij} \quad \text{else}$$

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

---

- For a weighted D matrix  $C_{ii}$  can be smaller than  $C_{ij}$

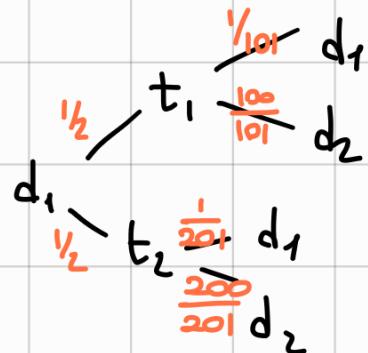
e.g. ↓

e.g.

$$C_{ii} < C_{ij}$$

$$\begin{matrix} d_1 \\ d_2 \end{matrix} \begin{bmatrix} 1 & 1 & 0 & 0 & 0 \\ 100 & 200 & 0 & 0 & 0 \end{bmatrix}$$

$$C_{11} = \frac{1}{2} \cdot \frac{1}{101} < C_{12} = \frac{1}{2} \cdot \frac{100}{101}$$



$$D \rightarrow C' = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}_{n \times n} \leftarrow \text{cover coef. matrix for terms } C_{ii}.$$

$$C'_{ij} = \beta_i \sum_{k=1}^m d_{ki} \alpha_k d_{kj}$$

(Remark.

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

$$C' = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{n \times n}$$

$$n_{C'} = \sum_{i=1}^n C'_{ii} = \sum_{i=1}^n \delta'_i$$

$$\downarrow \quad n_{C'} = n_C$$

↳ decoupling  
of terms

$$\delta' = \frac{\sum_{i=1}^n \delta'_i}{n} \rightarrow n_{C'} = \delta' \cdot n \Rightarrow \delta' \cdot n = \delta \cdot m$$

$$n_C = \delta \cdot m$$

$$\Rightarrow \delta' = \frac{\delta \cdot m}{n}$$

$$C = \begin{bmatrix} \textcircled{1} \\ \vdots \\ \textcircled{1} \end{bmatrix}_{m \times m}$$

$$C' = \begin{bmatrix} \textcircled{1} \\ \vdots \\ \textcircled{1} \end{bmatrix}_{n \times n}$$

$$\min(n_c) = 1$$

$$\max(n_c) = \min(m, n)$$

$d_c$ : avg. # docs / cluster

$d_{c'}$ : " " terms / "

$\delta_i$  = decoupling coef of  $d_i$ .

$D \rightarrow S \rightarrow n_c$

- select seeds (docs)

- assign nonseeds to seeds  $\binom{(m - n_c) \cdot n_c}{\underbrace{\text{nonseeds}}_{(m - n_c)} \underbrace{\text{seeds}}_{n_c}}$

$\Psi_i = 1 - \delta_i$  coupling coef.

- we want high decoupling  
& high seed power.

$$P_i = \text{seed power of doc } i \\ = \delta_i \cdot \psi_i \cdot \underbrace{(\# \text{ terms in doc}_i)}_{X_d}$$

Remark:  $X_{d\bar{i}} = \text{avg. } \# \text{ terms per doc.}$

$$P_1 = 0.361 \cdot (1 - 0.361) \cdot 3 = 0.692$$

$$P_2 = 0.563 \cdot (1 - 0.563) \cdot 4 = 0.884$$

$$P_3 = 0.692$$

$$P_4 = 0.484$$

$$P_5 = 0.469$$

decoupling coef

$$C = \begin{bmatrix} 0.361 & 0.250 & 0.104 & 0.111 & 0.083 \\ 0.188 & 0.563 & 0.063 & 0 & 0.188 \\ 0.194 & 0.083 & 0.361 & 0.277 & 0.088 \\ 0.167 & 0 & 0.417 & 0.417 & 0 \\ 0.125 & 0.375 & 0.125 & 0 & 0.375 \end{bmatrix} \quad m \times m$$

17.10.2023

$$C_{ii} = \delta_i$$

$$\psi_i = 1 - \delta_i$$

$$\sum_{j=1}^m C_{ij} = 1.$$

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

$$n_C = n_{C'}$$

$\Rightarrow d_4 \& d_5$  have no common term

- All docs are unique:  $C = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ \vdots & \vdots \\ 0 & 1 \end{bmatrix}$   $n_c = m$ .
- " " " identical:  $C = \begin{bmatrix} 1/m \\ 1/m \\ \vdots \\ 1/m \end{bmatrix}$   $n_c = 1$ .

$$n_c = \sum C_{ii} = n_c' = \sum C'_{ii}$$

$$C = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}_{n \times n}$$

$$\underbrace{(m - n_c) \cdot n_c}_{\text{non seeds.}}$$

## Seed Selection

$$\begin{aligned}
 p_i &= \text{seed power of doc } i \\
 &= \delta_i \cdot \psi_i \cdot (\# \text{ terms in doc}_i) \\
 &= c_{ii} \cdot (1 - c_{ii}) \cdot (\# \text{ terms in doc}_i)
 \end{aligned}$$

$$\underset{c_{ii}}{\operatorname{argmax}} p_i = 0.5$$

$$\begin{aligned}
 p_1 &= 0.692 & p_3 &= 0.692 & p_5 &= 0.469 \\
 p_2 &= 0.984 & p_4 &= 0.484 & &
 \end{aligned}$$

$d_2$  is a seed      two candidate seeds:  $d_1$  &  $d_3$ .

## False Seed Elimination

Check if  $d_1$  &  $d_3$  are identical?

(If the difference is less than a threshold, assume they are identical)

$C_{ij} = C_{ji} = C_{ii} = C_{jj} \Rightarrow d_i$  &  $d_j$  are identical.

$d_1$  and  $d_3$  are not identical.

$d_2 \rightarrow \underline{t_1}, \underline{t_2}, t_3, t_4$

$d_1 \rightarrow \underline{t_1}, \underline{t_4}, t_6$

$d_3 \rightarrow \underline{t_1}, t_5, t_6 \}$  is a seed now

bc it has less terms

in common w/  $d_2$ .

$$C = \begin{bmatrix} 0.361 & 0.250 & 0.194 & 0.111 & 0.083 \\ 0.188 & 0.563 & 0.063 & 0 & 0.188 \\ 0.194 & 0.083 & 0.361 & 0.277 & 0.088 \\ 0.167 & 0 & 0.417 & 0.417 & 0 \\ 0.125 & 0.375 & 0.125 & 0 & 0.375 \end{bmatrix}$$

↑  
the extend doc i

is covered by doc 2

For  $d_1$ :  $0.250 > 0.194$ . cluster of  $d_2$ .

$d_4$ :  $0.417 > 0$

" "  $d_3$

$d_5$ :  $0.375 > 0.125$

" "  $d_2$

The clusters are  $\{d_1, d_2, d_5\}$ ,  $\{d_3, d_4\}$ .

If we do col sums of matrix C, we get the total coverage of doc<sub>i</sub> of col. i.

Calculation of  $C_{12}$

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

$$C_{12} = \frac{1}{3} \left( d_{11} \beta_1 d_{21} + d_{12} \beta_2 d_{22} + \dots \right)$$

we are computing terms that will not contribute, a lot of times.

$$\begin{aligned} d_1 &\rightarrow t_1, t_4, t_6 \\ d_2 &\rightarrow t_1, t_2, t_3, t_4 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{d}_{12} \text{ is } 0.$$

Inverted Index for Seed Documents : IISD

Seeds  $d_2, d_3$

$t_1 \rightarrow d_2, d_3$

$t_2 \rightarrow d_2$

$t_3 \rightarrow d_2$

$t_4 \rightarrow d_2$

$t_5 \rightarrow d_3$

$t_6 \rightarrow d_3$

(For weighted matrices, we show  $t_5 \rightarrow t_3$  as  $t_5 \rightarrow \langle t_3, 1 \rangle$ )

Cluster  $d_1$ .  $d_1 \rightarrow t_1, t_4, t_6$ .

- We need to calculate  $C_{12}$  and  $C_{13}$ .

Let  $C_{12} = C_{13} = 0$ .

(Only update  $C_{12}, C_{13}$  for terms  $d_1$  has)

For  $t_1 //$

$$C_{12} = C_{12} + \alpha_1 \cdot (d_{11} \cdot B_1 \cdot d_{21}) \\ = 0 + \frac{1}{3} \cdot (1 \cdot \frac{1}{4} \cdot 1) = \frac{1}{12}$$

$$C_{13} = C_{13} + \alpha_1 \cdot (d_{11} \cdot B_1 \cdot d_{31}) \\ = 0 + \frac{1}{3} \left( 1 \cdot \frac{1}{4} \cdot 1 \right) = \frac{1}{12}$$

For  $t_4 //$

$$C_{12} = C_{12} + \alpha_1 \cdot (d_{14} \cdot B_4 \cdot d_{24}) \\ = \frac{1}{12} + \frac{1}{3} \cdot (1 \cdot \frac{1}{2} \cdot 1) = 0.250$$

For  $t_6 //$

$$C_{13} = C_{13} + \alpha_1 \cdot (d_{16} \cdot B_6 \cdot d_{36}) \\ = \frac{1}{12} + \frac{1}{3} \left( 1 \cdot \frac{1}{3} \cdot 1 \right) = 0.194$$

## IISD:

1. Find seeds according to the seed power (and then acc. to # of terms shared w/ higher seeds)

2. Compute only  $C_{ij}$   $\forall$  non seed  $i$  and seed  $j$ .

Start  $C_{ij} = 0$ .

$$C_{ij} = C_{ij} + \alpha_i \cdot d_{ik} \cdot \beta_k \cdot d_{jk} \text{ if } t_k \rightarrow d_i, d_j$$

### No IISD

$$(m - n_c)(n_c)$$

$$(m - n_c) X_d \cdot t_g$$

non seed

### IISD

$$(m - n_c) X_d t_{gs}$$

$t$  avg length of  
inverted index for seeds

$$t_{gs} \ll t_g$$

$$\mathcal{O}(m X_d t_{gs})$$

20.10.2023

Cover coef.

all unique

$$n_c = m$$

all identical

$$n_c = 1$$

$$C_{m \times m}$$

$$C'_{n \times n}$$

$$n_c = n_c'$$



$$1 < n_c \leq \min(m, n)$$

$$\text{for docs} \rightarrow \delta = \frac{\sum C_{ii}}{m} = \frac{\sum_{i=1}^m \delta_i}{m}$$

$$\delta = n_c/m \Rightarrow n_c = \delta \cdot m$$

$$\text{for terms} \rightarrow \delta' = \frac{\sum_{i=1}^n C'_{ii}}{n} = \frac{\sum_{i=1}^n \delta'_i}{n} = \frac{n'_c}{n}$$

$$\Rightarrow n'_c = \delta' \cdot n$$

$$n'_c = n_c \Rightarrow \delta' \cdot n = \delta \cdot m \Rightarrow \delta = \delta' \frac{n}{m}$$

$$\text{Observe } n_c = \frac{m \cdot n}{t} \quad ? ? ?$$

$$D_1 = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}_{2 \times 3}$$

$$D_2 = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}_{3 \times 2}$$

$$\frac{2 \cdot 3}{3} = 2$$

$$\frac{3 \cdot 2}{3} = 2$$

$$\text{Obtain } n_c = \frac{mn}{t} \text{ from } n_c = \sum_{i=1}^m \delta_i :$$

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

$$\delta_i = C_{ii} = \alpha_i \sum_{j=1}^n d_{ij}^2 \beta_j$$

$$n_c = \sum_{i=1}^m \alpha_i \sum_{j=1}^n d_{ij}^2 \beta_j = \sum_{i=1}^m \underbrace{\sum_{j=1}^n \alpha_i d_{ij}^2 \beta_j}_{\parallel}$$

$d_{ij}$   
(for  
binary matrices)

$$= \sum_{i=1}^m \sum_{j=1}^n \alpha_i d_{ij} \beta_j$$

$$\left[ \sum_{k=1}^n d_{ik} \right]^{-1} \left[ \sum_{k=1}^m d_{kj} \right]^{-1}$$

(column sum  $x_{di}$   
depth of indexing  
for  $d_{ij}$ )

(row sum  $t_{gj}$   
term generality  
for  $d_{ij}$ )

$$\Rightarrow n_c = \sum_{i=1}^m \sum_{j=1}^n d_{ij} \left[ \sum_{k=1}^n d_{ik} \right]^{-1} \left[ \sum_{k=1}^m d_{kj} \right]^{-1}$$

$$= \sum_{i=1}^m \sum_{j=1}^n d_{ij} \underbrace{\left[ \sum_{k=1}^n d_{ik} \right]}_{x_{di}} \cdot \underbrace{\left[ \sum_{k=1}^m d_{kj} \right]}_{t_{gj}}^{-1}$$

$$n_c = \sum_{i=1}^m \sum_{j=1}^n \frac{d_{ij}}{x_{di} t_{gj}}$$

$$x_d = \sum_{i=1}^m \frac{x_{di}}{m} \quad t_g = \sum_{i=1}^n \frac{t_{gi}}{n}$$

$$n_c = \sum_{i=1}^m \sum_{j=1}^n \frac{d_{ij}}{x_d t_g} = \frac{1}{x_d t_g} \sum_{i=1}^m \sum_{j=1}^n d_{ij}$$

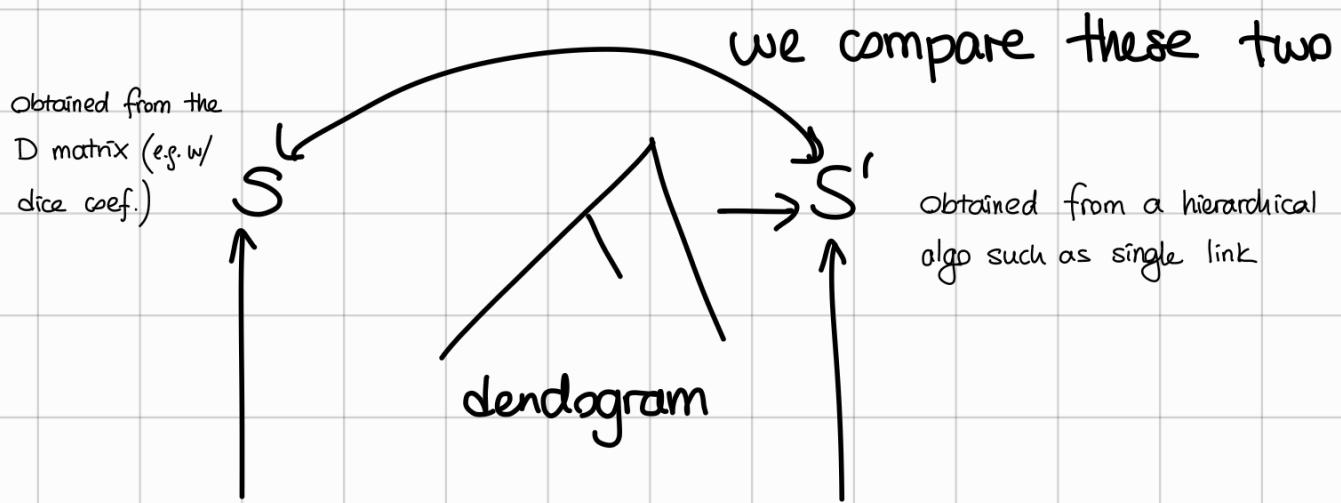
# of non zero entries in the matrix

$t =$

$$x_d = \sum_{i=1}^m \frac{x_{di}}{m} = \frac{\sum_{i=1}^m \sum_{k=1}^n d_{ik}}{m} = \frac{t}{m}$$

Similarly for  $t_g = \frac{t}{n} \Rightarrow n_c = \frac{t}{x_d t_g} = \frac{t}{\frac{t}{m} \cdot \frac{t}{n}} = \frac{mn}{t}$

$n_c = \frac{mn}{t} \quad \square.$



ab 0.6

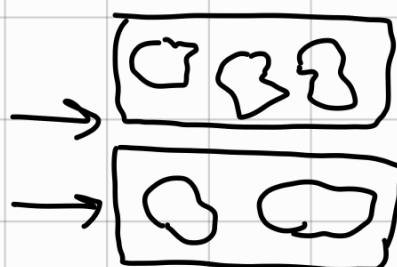
bc 0.5

⋮

Rank similarity among pairs of docs

$D \rightarrow$  clustering algo

ground truth



cluster purity = how homogenous clusters are

## Rand Index

We have cluster output (partitions)  $\{ab\} \{cd'\} \{e'f'\}$   
 " " Human based ground truth (or gold standard)  
 $\{abc\} \{d'e'f'\}$

	a	b	c	d'	e'	f'
a	-	ab	ac	ad'	ae'	af'
b	-	-	bc	bd'	be'	bf'
c						
d'						
e'				e'e'		
f'						

$$\text{Rand Index} = \frac{TP + TN}{TP + TN + FN + FP}$$

Pair	ab	ac	ad'	ae'	af'	bc	bd'	be'	bf'	cd'	ce'	cf'	de'	d'f'	e'f'
Class	TP	FN	TN	TN	TN	FN	TN	TN	TN	FP	TN	TN	FN	FN	TP

↑  
 TP : true pos.

TN  
 FP  
 FN

## Corrected Rand Coef

- it takes out matches due to randomness

## Cluster Purity

- we need to have a gold standard
- we have a partitioning structure



IR:1 OS:3 DB:0

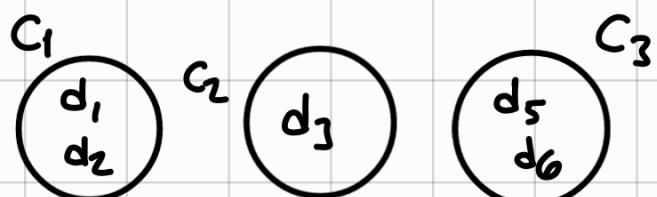


IR:3 OS:2 DB:1



IR:1 OS:2 DB:3

cluster purity  $\rightarrow \frac{3+3+3}{4+6+6} = \frac{9}{16} = \frac{\sum_{i=1}^k \max c_i}{\sum_{i=1}^k \sum_{j \in C_i}}$

 $q_1 \rightarrow 1, 2, 3$  $q_2 \rightarrow 1, 5, 6$  $n_t$ : # target cluster

$$\left. \begin{array}{ll} q_1 \rightarrow C_1, C_2 & n_{tr_1} \\ q_2 \rightarrow C_1, C_3 & n_{tr_2} \end{array} \right) \quad n_t = \frac{2+2}{2} = 2$$

avg # target cluster

## Cluster Validation

1. Keep the clustering struct. the same
2. Distribute docs randomly to clusters
3. Find the avg. # target clusters for the random query?  $n_{tr}$

For a meaningful structure

We observe:

$$n_t \ll n_{tr}$$

$n_t < n_{tr} \rightarrow$  we have hope  
we need to

$$\text{show } n_t \ll n_{tr}$$

$$n_t \geq n_{tr}$$

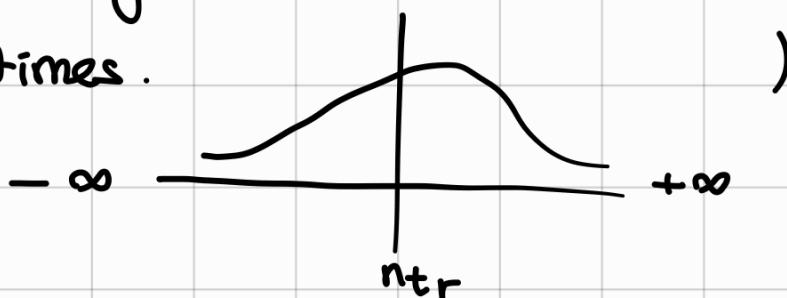
↑

no good  
don't proceed

Is there a formula for  $n_{tr}$ ?

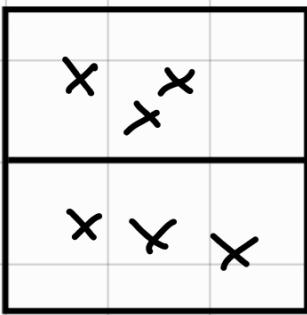
(without having to do monte Carlo simulation

1000s times.



S.B. Yao, Com. of ACM 1970

## Disk block records



**Yao's Derivation** :  $n = \# \text{ records}$

$$m = \# \text{ blocks}$$

block size =  $p = n/m \leftarrow$  # records in  $j^{\text{th}}$  block  
(since each block has same # records)

$n - p = \#$  records in the other blocks

$k = \#$  records to be accessed

$$k < n - \frac{n}{m}$$

↑                  ↓  
all              except one block

Remark.

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

18

$$a, b, c \quad C_2^3 = \frac{3!}{2! \cdot 1!} = 3$$

ab  
ac  
bc

$C_k^{n-p}$  = different ways of selecting  $k$  records from  $n-p$  records

The prob. of selecting no records from the  $i^{\text{th}}$  block:

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

$$n-p = n - \frac{n}{m} = n \underbrace{\left(1 - \frac{1}{m}\right)}_d = nd$$

$E(I_j)$  = expected value of selecting at least a record from the  $i^{\text{th}}$  block

$$= \left(1 - \frac{C_k^{nd}}{C_k^n}\right)$$

Expected number of blocks to be accessed

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

$$n_{tr} = m \times \left(1 - \frac{C_k^{nd}}{C_k^n}\right) = m \cdot \left[1 - \frac{\frac{nd!}{k! (nd-k)!}}{\frac{n!}{k! (n-k)!}}\right]$$

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

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

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

$$= m \cdot \prod_{i=1}^k \frac{nd-i+1}{n-i+1} = n_{tr} \leftarrow \begin{matrix} \text{avg \# target clusters} \\ \text{for random query} \end{matrix}$$

How to apply this for clustering?

- Using our notation

$m = \# \text{docs}$     $k = \# \text{docs we want to access}$

$m_j = \# \text{docs other than cluster } j \text{ we consider}$

let  $m=100$ ,  $k=3$ .

$|C_1| = 5$        $C_1 \rightarrow \text{cluster 1}$

$$m_1 = 100 - 5$$

$P_j = \text{Prob. that we have a doc coming from } C_j$

Remark:

$$E(I_j) = 1 - \frac{C_k^{nd}}{C_k^n} = \prod_{i=1}^k \frac{n-d-i+1}{n-i+1}$$

$$P_j = E(I_j) = \prod_{i=1}^k \frac{m_j - i + 1}{m - i + 1}$$

Expected # clusters to be accessed =  $n_{tr}$

$$n_{tr} = \sum_{j=1}^{n_c} P_j \quad (n_c : \# \text{ clusters})$$

$$P_1 = \left[ 1 - \prod_{i=1}^3 \frac{85-i+1}{100-i+1} \right]$$

27.10.2023

## Inverted File Processing

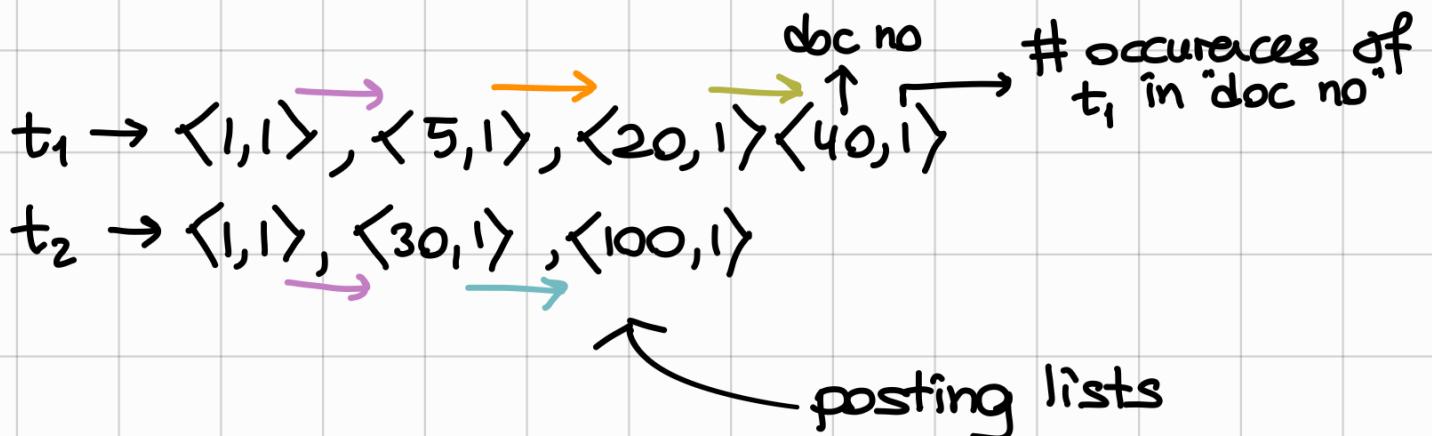
Inverted Index for Search Engines, Zobel Moffat

Idea: skips & sorts

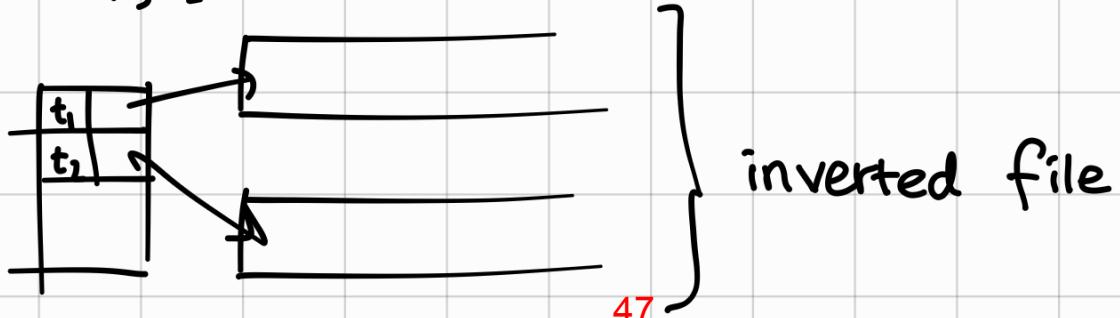
↳ (self indexing inverted files)

Boolean Retrieval

Conjunctive queries      and (can be extended to or)



$$Q = t_1, t_2$$



<u>Doc id.</u>	<u># comparisons</u>	
1	1	1 : 1 →
5	1	5 : 30 5 < 30 →
20	1	20 : 30 20 < 30
40	2	40 : 30, 40 : 100 40 > 30 → 40 > 100 →

How to decrease the size of posting list?

$\langle 5,1 \rangle \langle 8,1 \rangle \langle 12,2 \rangle \langle 13,1 \rangle \langle 15,1 \rangle \langle 18,1 \rangle \langle 23,2 \rangle \langle 28,1 \rangle \langle 29,1 \rangle$

1. Use **d-gaps** (document gaps):

use not the doc id , but the diff between the ids

$\langle 5,1 \rangle \langle 3,1 \rangle \langle 4,2 \rangle \langle 1,1 \rangle \langle 2,1 \rangle \langle 3,1 \rangle \langle 5,2 \rangle \langle 5,1 \rangle \langle 4,1 \rangle$

2. Use **skips** every 3 posting list entry :

$\langle 5,1 \rangle \langle 8,1 \rangle \langle 12,2 \rangle \langle 13,1 \rangle \langle 15,1 \rangle \langle 18,1 \rangle \langle 23,2 \rangle \langle 28,1 \rangle \langle 29,1 \rangle$

e.g. Compare w/  $[-, 0]$   $\langle 100, 1 \rangle \langle 105, 5 \rangle$

dgap - 9 comparisons  
compare 100 w/ each  
doc.  $100 > 2g$  so don't  
compare 105 at all

skips - 5 comparisons  
compare 100 w/ 13  
w/ 23  
w/ - )  
so compare w/ 28, 29 too.

### 3. Skips w/ d-gaps:

$\langle 3, 1 \rangle \langle 5, 1 \rangle \langle 3, 1 \rangle \langle 4, 2 \rangle \langle 23, 1 \rangle \langle 1, 1 \rangle \langle 2, 1 \rangle \langle 3, 1 \rangle \langle -, 0 \rangle \langle 5, 2 \rangle \langle 5, 1 \rangle \langle 4, 1 \rangle$

### 4. Frequency Ordered Inverted List (FOIL):

Reg.  $\langle 12, 2 \rangle \langle 17, 2 \rangle \langle 29, 1 \rangle \langle 32, 1 \rangle \langle 40, 6 \rangle \langle 78, 1 \rangle \langle 101, 3 \rangle \langle 106, 1 \rangle$

FOIL  $\langle 6:1:40 \rangle \langle 3:1:101 \rangle \langle 2:2:12, 17 \rangle \langle 1:4:29, 32, 78, 106 \rangle$

$f_{d,t}$ : freq. of term  $t$  in the docs

$\#$  docs w/  
this freq.

↓  
doc list

## 5. Ordering posting lists acc. to retrieval numbers

Sengör, Altingovde METU

$t_1 \rightarrow [10] \langle d_2, 2 \rangle \langle d_5, 2 \rangle [8] \langle d_3, 4 \rangle [5] \langle d_7, 2 \rangle \langle d_{10}, 8 \rangle$

$\uparrow$

$d_2$  is accessed 10 times  
in the list of retrieved docs.

### Ranked Retrieval

$f_{d,t}$  : frequency of term  $t$  in doc  $d$

$f_{q,t}$  : " " " " in query  $q$

$f_t$  : # docs containing term  $t$

$F_t$  : # term  $t$  in the collection

$N$  : # docs

$n$  : # terms

$$\text{e.g. } D = \begin{bmatrix} t_1 & t_2 & t_3 & t_4 \\ 5 & 1 & 0 & 2 \\ 1 & 0 & 0 & 0 \\ 2 & 0 & 4 & 1 \end{bmatrix} \begin{matrix} d_1 \\ d_2 \\ d_3 \end{matrix}$$

$f_{d,t}$	$ $	$f_{1,1} = 5$
$f_t$	$ $	$f_1 = 3$
$F_t$	$ $	$F_1 = 5 + 2 + 1 = 8$
$N = 3$		
$n = 4$		

# Assigning weights to doc & query terms

Assign less weight to terms w large ft (occurs in many docs)  $w_{q,t} \propto \frac{1}{ft}$

more weight is given to terms that occur frequently in a doc  $\frac{ft}{ft}, f_{d,t} \propto w_{q,t}$

less weight to docs that contain many terms

tf.idf

?

## Term weight assignment. Salton & Buckley

$$Q: w_{q,t} = \ln\left(1 + \frac{N}{ft}\right) \quad \text{weight of term } t \text{ in } q.$$

$$D: w_{d,t} = 1 + \ln f_{d,t}$$

$$w_q = \sqrt{\sum_t w_{q,t}^2}$$

$$w_d = \sqrt{\sum_t w_{d,t}^2}$$

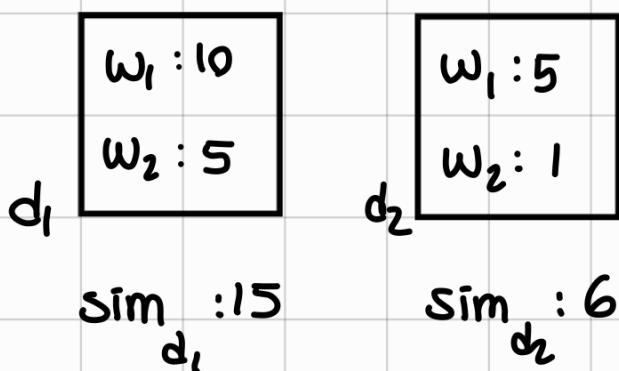
$$S_{q,d} = \frac{\sum_t w_{d,t} \cdot w_{q,t}}{w_d \cdot w_q}$$

# Salton & Buckley

- term freq. component
- doc. " "
- normalization

Assume that the similarity of a doc to a query is the total no. of these terms in the docs.

Q:  $w_1 \& w_2$



Q:  $w_1 \& w_2$

$w_1 \rightarrow \langle d_1, 5 \rangle \langle d_7, 2 \rangle \langle d_{11}, 2 \rangle \langle d_{20}, 1 \rangle$

$w_2 \rightarrow \langle d_1, 1 \rangle \langle d_4, 1 \rangle \langle d_7, 10 \rangle$

remark.  
 $w_i$  &  $t_i$   
both mean  
term 1.

Rank the entries

1		20
5		1
6		

Assume that we have  
3 accumulators

0	0	a
---	---	---

Process  $w_1$  ( $w_1$  is more important than  $w_2$ )

5	2	2
$d_1$	$d_7$	$d_{11}$

Process  $w_2$

5	2	2
6	12	↑ no change

Ranking  $d_1 : 6 \quad d_7 : 12$

$d_{11} : 2$

Organize terms according to their freq. in a doc.

$w_1 \rightarrow \langle d_1, 5 \rangle \langle d_7, 2 \rangle \langle d_{11}, 2 \rangle \langle d_{20}, 1 \rangle$

$w_2 \rightarrow \langle d_7, 10 \rangle \langle d_1, 1 \rangle \langle d_4, 1 \rangle$

0	0	0
---	---	---

$d_7$	$d_1$	$d_4$
10	1	1

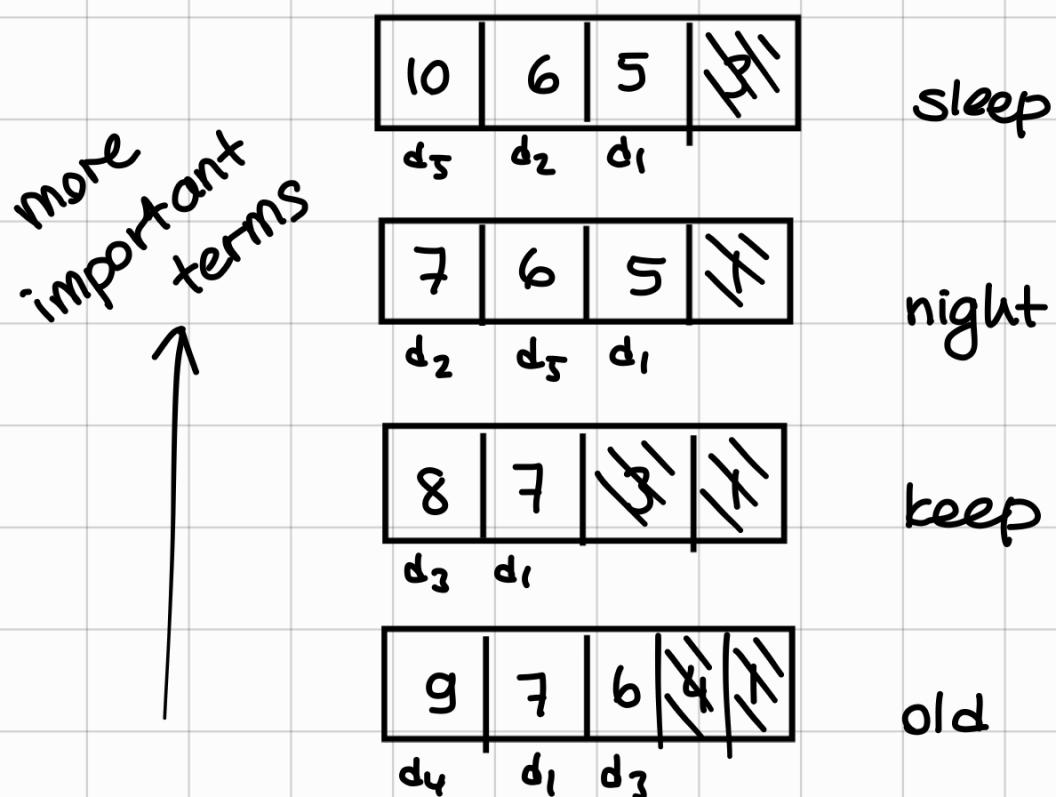
12	6	
12	2	1

process  $w_2$

process  $w_1$

$d_7 : 12, d_1 : 6, d_4 : 1$

# Interleaved processing of query terms



## Accumulators

1	2	3	4	5	6	7

--	--	--	--	--	--

conjunctive query = and  
 disjunctive " = or

<u>doc<sub>1</sub></u>
$w_1 = 5$
$w_2 = 2$

<u>doc</u>
$w_1 = 2$
$w_2 = 1$

query  $w_1$  and  $w_2$   $doc_1 = \text{total weight} \exists$   
 $doc \quad * \quad .. \quad 3$

$$w_1 \rightarrow \langle d_1, 5 \rangle \langle d_7, 2 \rangle \langle d_{11}, 2 \rangle \langle d_{20}, 1 \rangle$$

$$w_2 \rightarrow \langle d_1, 1 \rangle \langle d_4, 1 \rangle \langle d_7, 10 \rangle$$

Accumulators

6	0	0	1	.. ..
$d_1$	$d_2$	$d_3$	$d_4$	

# Limited # of Accumulators

- 3 acc.s

0	0	0
---	---	---

$d_1 \ d_2 \ d_{11}$

$d_2 : 12$

$d_1 : 6$

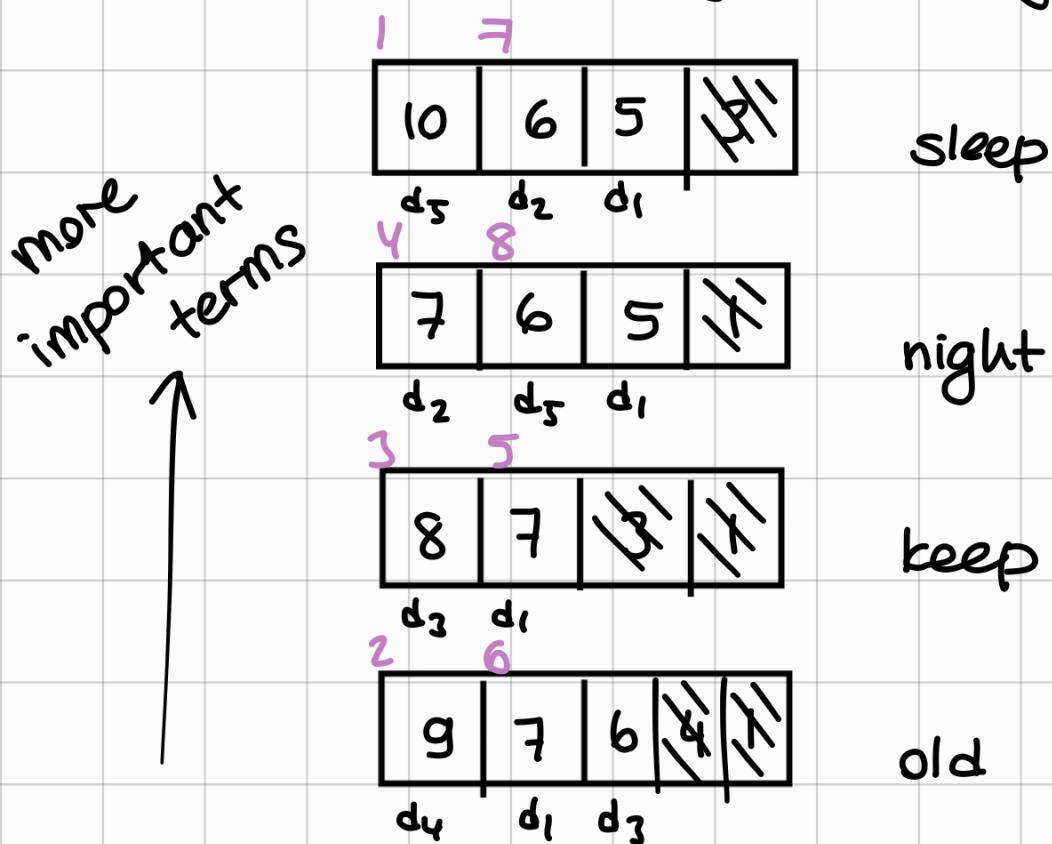
5 2 2

$d_{11} : 2$

6 12

We left lecture here last

Interleaved processing of query terms



first we care about the freq. of the term than the importance of the term

## Accumulators

1	2	3	4	5	6	7	
0	0	0	0	0	0	0	.. -.

sleep 10 / d<sub>5</sub>

10

old 9 / d<sub>4</sub>

9 10

keep 8 / d<sub>3</sub>

8 9 10

night 7 / d<sub>2</sub>

7 8 9 10

keep 7 / d<sub>1</sub>

7 7 8 9 10

: 8 more lines

24 13 14 19 16

Interleaved processing of query terms

w/ limited # accumulators

more important terms

10	6	5	X
d <sub>5</sub>	d <sub>2</sub>	d <sub>1</sub>	

sleep

7	6	5	X
d <sub>2</sub>	d <sub>5</sub>	d <sub>1</sub>	

night



	$d_5$	$d_2$
sleep	10 / $d_5$	10 / 0
night	7 / $d_2$	10 / 7
sleep	6 / $d_2$	10 / 13
night	6 / $d_5$	16 / 13

Ranking

$d_5$  1b

$d_2$  13

- - - stop since we  
don't have space for  
another doc.

03.11.2023

## Term Weighting

- Term weighting approaches in automatic text retrieval - Salton & Buckley 1988

Relational model of data 1970

No. of different weighting function (matching function)  
MF

I think it is matching  
terms w/  
docs or  
queries

# Components of Matching Functions

TFC: term freq. component

CFC: collection freq. "

NC: normalization "

• doc

• query  
(no NC needed)

## TFC

b: binary 0,1

t: raw frequency

n: augmented normalized term

$$\text{frequency} = 0.5 + \frac{tf}{\max(tf)} \cdot 0.5$$

We have 3 options to choose as  
Either b (binary TF), t (raw TF, as is)  
or n (augmented normalized TF)

applies to nonzero  $tf$ 's.

$$\underline{tf} = 0 \Rightarrow n = 0.$$

↳ term occurrence?

e.g.  $(5 \ 0 \ 1 \ 2) \Rightarrow \max(tf) = 5 = \max\{5, 0, 1, 2\}$

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

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

$$n: \left(0.5 + \frac{5 \cdot 0.5}{5} \ 0 \ \underbrace{0.5 + \frac{0.5 \cdot 1}{5}}_{0.6} \ \underbrace{0.5 + \frac{2 \cdot 0.5}{5}}_{0.7}\right)$$

## CFC

x: no change (use original TFC)

f: inverse collection frequency component

$$\ln\left(\frac{m}{tg_j}\right) + 1$$

m: # docs

(similarly for nonzero  $tf_s$ )

$tg_j$ : term generality for term j

: # unique docs containing term j

p: probabilistic inverse collection frequency factor

For CFC, either use x(TFC), f( $\frac{\text{inv. coll.}}{\text{TFC}}$ ) or p( $\frac{\text{probabilistic}}{\text{inv. coll. FC}}$ )

## NC

x: no change

2 options

c: use cosine normalization

$$\frac{1}{\sqrt{\sum w_i^2}}$$

# Different ways of assigning weights

Doc

TFC: 3

CFC: 3

NC: 2

Query

TFC: 3

CFC: 3

NC: 1

$$(3.3.2) \times (3.3.1) = 162$$

60 no normalization for query

e.g.

	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	max tf
$d_1$	2	0	1	2	0	2
$d_2$	0	2	1	3	1	3
$d_3$	2	0	0	1	1	2
$d_4$	1	0	0	0	1	1
$d_5$	2	1	0	1	0	2

$$m = 5$$

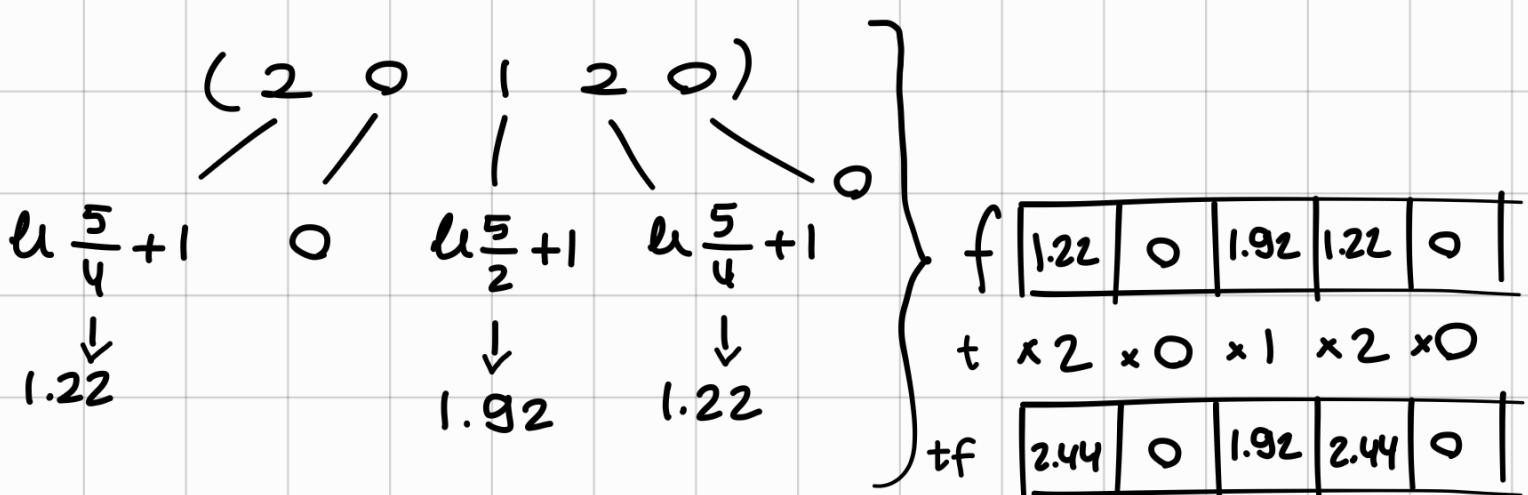
$$4 \ 2 \ 2 \ 4 \ 3 \rightarrow t_{g_j}$$

For documents use  $t_{fc}$

For doc. 1:

$$TFC : t \rightarrow (2 \ 0 \ 1 \ 2 \ 0)$$

$$CFC : f \rightarrow \ell \lfloor \frac{m}{t_{g_j}} + 1 \rfloor$$



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

$\hookrightarrow w_i = t.f$

e.g. continued

doc1, but weighted

$t_{f_C}$	$\frac{2.44}{3.95}$	0	$\frac{1.92}{3.95}$	$\frac{2.44}{3.95}$	0
	0.62	0	0.49	0.62	0

(for query normalization  $n_{fc}$  doesn't make sense because it will not change rankings)

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

	1	0	0	2	0
TFC n	$0.5 + \frac{1.05}{2}$	0	0	$0.5 + \frac{2}{2} \cdot 0.5$	0
CFC f: $f_i = \frac{m}{t_{qj}} + 1$	1.22	1.92	1.92	1.22	1.51
no normalization use nf as is	$0.75 \cdot 1.22$ $= 0.92$	0	0	$1.1 \cdot 1.22$	0

nfx: use n from TFC, f from CFC, and nothing from NC.

$$Q = (0.92 \ 0 \ 0 \ 1.22 \ 0)$$

$$\text{sim}(Q, D) = \sum w_{q_h} \cdot w_{d_h}$$

$$D = \begin{bmatrix} 0.62 & 0.00 & 0.49 & 0.62 & 0.00 \\ 0.00 & 0.66 & 0.33 & 0.63 & 0.26 \\ 0.18 & 0.00 & 0.00 & 0.39 & 0.48 \\ 0.63 & 0.00 & 0.00 & 0.00 & 0.78 \\ 0.73 & 0.58 & 0.00 & 0.37 & 0.00 \end{bmatrix}$$

$s(q, d_1) = 0.92 \cdot 0.62 + 0$   
 $+ 0 + 1.22 \cdot 0.62$   
 $+ 0 = 1.33$

$d_1 = 1.33 \quad d_4 = 0.58$   
 $d_3 = 1.29$   
 $d_5 = 1.12$   
 $d_2 = 0.77$

(remarq. no normalization for queries since it won't change the rankings)

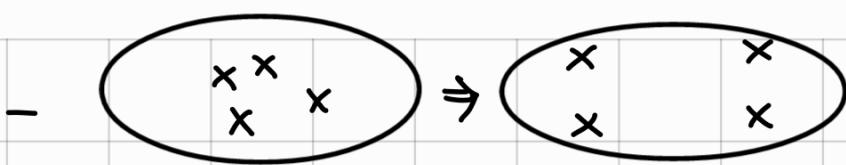
7.11.2023

### Term Discrimination Value

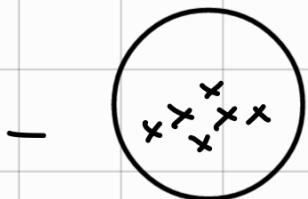
TDV: smth like importance of terms

Learning TDV, some French team, SIGIR 2020.

- Find terms w/  $TDV = 0$ .
  - Delete them from indexing vocab.
- $\Rightarrow$  IR becomes more efficient & more effective



The use of a good discriminator will separate docs from each other.



The use of a bad discriminator will make docs less distinguishable.

How to calculate TDVs:

- similarity based
- # clusters

Space density:  $Q$ : Avg. similarity val. among docs

$$S = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$$

$m \times n$   
symmetrical

$$Q = \frac{1}{\frac{m(m-1)}{2}} \cdot \sum_{i=1}^{m-1} \sum_{j=i+1}^m S_{ij}$$

$Q_j = \text{arg. similarity without } t_j$ .

$$Q \propto \text{similarity} \propto \frac{1}{\text{discrimination}} \propto \frac{1}{TDV}$$

<u>Term Type</u>	<u><math>Q:Q_j</math></u>	<u>TDV</u>	
Good	$Q_j > Q$	$> 0$	If we add $t_j$ $Q$ decreases
Bad	$Q_j < Q$	$< 0$	
Indifferent	$Q \approx Q_j$	$\approx 0$	$t_j$ is a good discriminator $\Downarrow$ $\Downarrow$ TDV of $t_j$ is high

$$TDV = Q_j - Q$$

or

$$TDV = \frac{Q_j}{Q}$$

term weight:  $TDV_e \cdot t_f$

TDV of term e

term frequency

[ASK]

of which term?  
 $t_{fe}$ ?

How to calculate  $Q$ ?

1. Obtain S matrix: expensive.

2. Obtain the collection centroid and calculate sim. of docs to the centroid.

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

$$\text{Centroid} = \left( \frac{2}{3}, \frac{1}{3}, \frac{2}{3}, \frac{2}{3} \right)$$

Salton, ... Wang - Comm. ACM.

$$Q = \frac{1}{m} \sum_{i=1}^m \text{Sim}(d_i, \text{centroid}) \quad \leftarrow \begin{matrix} \text{Vector Space} \\ \text{Model} \end{matrix}$$

TDV Calculation using Clusters

(Cover Coefficient Concept)

$n_c$  = # clusters using all terms

$n_{ce}$  = # clusters w/o using term  $t_e$

- A good discriminator should increase # clusters

Term Type	$n_c : n_{ce}$	TDV
-----------	----------------	-----

Good  $n_c > n_{ce}$   $n_c - n_{ce} > 0$

Bad  $n_c < n_{ce}$   $n_c - n_{ce} < 0$

Neutral  $n_c \approx n_{ce}$   $n_c - n_{ce} \approx 0$ .

$$n_c = \sum_{i=1}^m C_{ii} = \sum_{i=1}^m \alpha_i \cdot (d_{i1}^2 \beta_1 + d_{i2}^2 \beta_2 + \dots + d_{in}^2 \beta_n)$$

↓ cont.

JASSIST  
1989  
paper

$$TDV_e = \sum_{i=1}^{f_e} \left[ \delta_i - \alpha_i^e \left( \frac{\delta_i}{\alpha_i} - d_{ie}^2 \cdot \beta_e \right) \right]$$

$$\alpha_{ie} = \left( \alpha_i^{-1} - d_{ie} \right)^{-1}$$

↓  
original  
row sum      ↓  
subtract the  
term e

$$f_e = \{ d_i \mid d_i \in D \wedge d_{ie} \neq 0 \}$$

\* More efficient calc. of TDVs using cover coef.

concept:

$$n_c = \frac{mn}{t}$$
$$n_{ce} = \frac{m(n-1)}{(t - \underbrace{\# \text{docs that}}_{\text{contain } t_e} + \text{that contain } t_e)} = \frac{m(n-1)}{(t - t_{ge})}$$

term generality  
of term e  
 $\Rightarrow t_{ge}$

10.11.23

idf: inverse document frequency

TDV

↳ idf

**ASK**

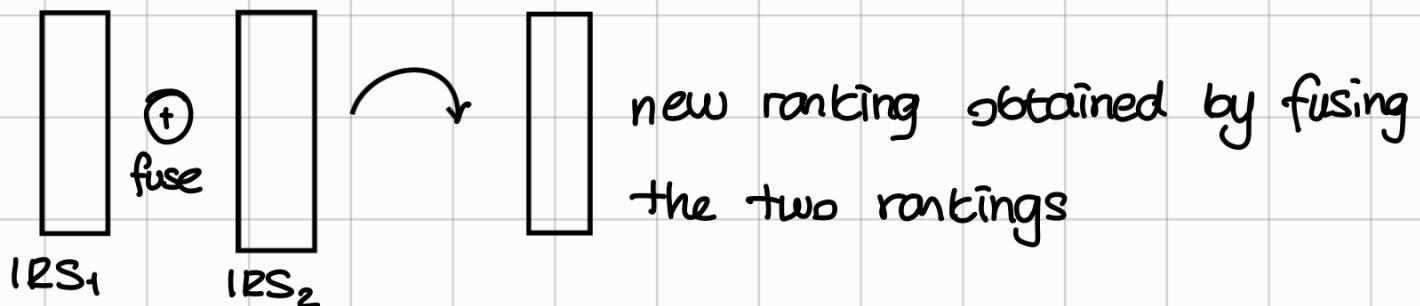
What do we mean here?

- term weighting

SIGIR 2020, Learning Term Disc Values

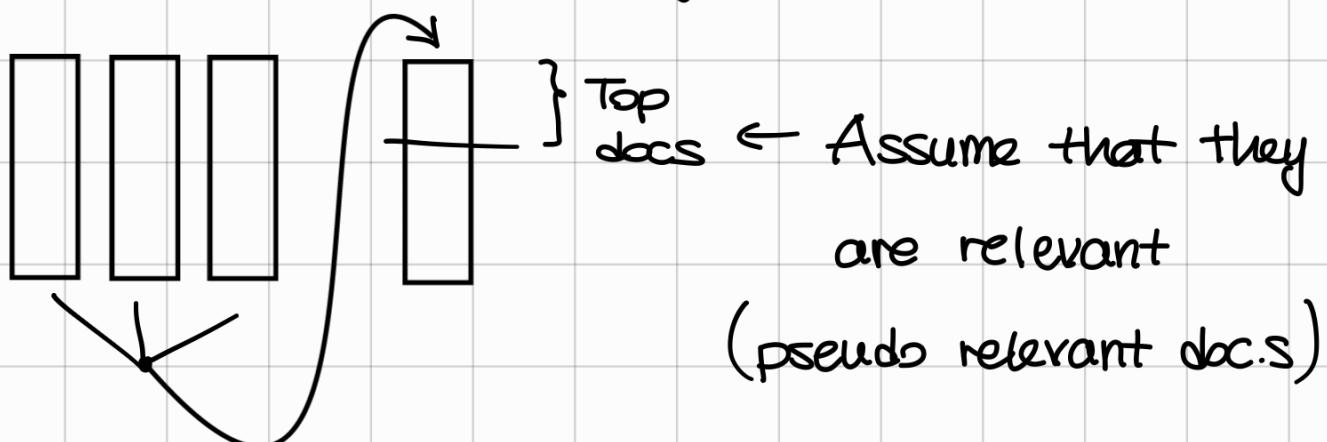
## Data Fusion

- Rank docs in diff. matching functions. Combine rankings  
⇒ more effective IRS.



- Meta search engines use data fusion

Data Fusion for Automatic Evaluation of Search Engines, IPM 2006 , Nuray and Can



- Use pseudo relevant doc.s to rank IRSs according to their effectiveness.

## Fusion Methods



1. Reciprocal rank

2. Borda Count

3. Condorcet

Bias Concept: being far away from the norm

↑ What do we use this for? For choosing which IRS to fuse.

### 1. Reciprocal Rank

$$r(d_i) = \frac{1}{\sum_{j=1}^n \frac{1}{\text{position } d_j}}$$

↑ IRS

} After finding all r values, rank acc. to r values (in the increasing order)  
(i.e. lower r(i) = better rank)

Example.

IRS: A, B, C, D      Documents: a, b, c, d, e, f, g

A: (a, b, c, d)

B: (a, d, b, e)

C: (c, a, f, e)

D: (b, g, e, f)

$$r(a) = \frac{1}{\frac{1}{1} + \frac{1}{1} + \frac{1}{2}} = \frac{1}{2.5} = 0.4$$

$$r(b) = \frac{1}{\frac{1}{2} + \frac{1}{3} + \frac{1}{1}} = \frac{6}{3+2+6} = \frac{6}{11} \approx 0.54$$

$$r(c) = \frac{1}{\frac{1}{3} + \frac{1}{1}} = \frac{3}{1+3} = 0.75$$

$$r(d) = \frac{1}{\frac{1}{4} + \frac{1}{2}} = \frac{4}{1+2} \approx 1.33$$

$$r(e) = \frac{1}{\frac{1}{4} + \frac{1}{4} + \frac{1}{3}} = \frac{12}{3+3+4} = \frac{12}{10} = 1.2$$

$$r(f) = \frac{1}{\frac{1}{3} + \frac{1}{4}} = \frac{12}{4+3} = \frac{12}{7} \approx 1.71$$

$$r(g) = \frac{1}{\frac{1}{2}} = 2$$

- Sort in increasing order:

$$0.4 < 0.54 < 0.75 < 1.2 < 1.33 < 1.71 < 2$$

a      b      c      e      d      f      g

Ranks: a > b > c > e > d > f > g.

## 2. Borda Count (Jean Charles, Cavalier de Borda)

- The Highest ranked individual (in an n-votes) gets n votes, and each subsequent indv. gets one vote less; e.g. no 2 gets  $n-1$  votes.
- If there are candidates left unranked by the voter, the remaining points are divided evenly among the unranked candidates.
- The choice with the Highest number of votes wins

Example.  $A = (a, c, b, d)$        $C = (c, a, b, e)$   
 $B = (b, c, a, e)$       Docs/candidates = {a, b, c, d, e}

$$BC(i) = BC_A(i) + BC_B(i) + BC_C(i)$$

$$\left. \begin{array}{l} \# \text{docs} = n = 5 \Rightarrow \text{max vote} \\ BC(a) = 5 + 3 + 4 = 12 \\ BC(b) = 3 + 5 + 3 = 11 \\ BC(c) = 4 + 4 + 5 = 13 \\ BC(d) = 2 + 0 + 0 = 2 \\ BC(e) = 0 + 2 + 2 = 4 \end{array} \right\} \text{rankings: } c > a > b > e > d.$$

### 3. Condorcet Method (Marie Jean Antoine Nicholas de Caritat Condorcet)

The voting process takes "each preference of each voter for one candidate over another" into account.

Example. candidates/doc.s = a, b, c

Five voters = A, B, C, D, E

A: a > b > c

C: a > b = c

E: c > a

B: a > c > b

D: b > a

Here, a's ranking is better than b, since b doesn't exist in the comparison

Each non diagonal entry  $(i,j)$  of the matrix below shows # votes i over j. (e.g. cell [a,b] shows # wins, # loses and # ties of doc. a over doc. b)  
(win, lose, tie)

	a	b	c
a	-	(4, 1, 0)	(2, 1, 0)
b	(1, 4, 0)	-	(2, 2, 0)
c	(1, 2, 0)	(2, 2, 0)	-

- At the end, document having the most wins get the first rank.
- If two docs have the same # wins, then the one w/ smaller # loses wins.
- else, tied.

Won over b,c

	Win	Lose	Tie
a	2	0	0
b	0	1	1
c	0	1	1

ranks:  $a > b = c$ .

tie with c

lost to a

- We don't want to fuse all of the data coming from each IDS. We need to choose systems for data fusion.
- But how to choose?
  - + Select systems far away from the avg behavior.

## Choosing systems to fuse

Docs : a,b,c,d,e,f,g

$$A = \begin{bmatrix} a & b & c & d \\ b & a & c & d \\ a & b & c & e \end{bmatrix} \xrightarrow{\text{results of query } Q_i} \begin{bmatrix} b & f & c & e \\ b & c & f & g \\ c & f & g & e \end{bmatrix} = B$$

$\xrightarrow{\quad \quad \quad \quad \quad \quad}$

$\rightarrow$  a is retrieved 3 times

$$X_A = (3, 3, 3, 2, 1, 0, 0)$$

$$X_B = (0, 2, 3, 0, 2, 3, 2)$$

*a b c d e f g*

$X = (3, 5, 6, 2, 3, 3, 2) \leftarrow$  the norm

$$S(V, W) = \frac{\sum V_i * W_i}{\sqrt{\sum V_i^2 * \sum W_i^2}}$$

$$\text{Sim}(X_A, X) = \frac{9+15+18+4+3+0+0}{\sqrt{(9+9+9+4+1) \cdot (9+25+36+4+9+9+2)}} = 0.8841$$

$$\text{Sim}(X_B, X) = 0.8758$$

$$\text{Bias}(A) = 1 - \text{Sim}(X_A, X) = 0.1159$$

$$\text{Bias}(B) = 1 - \text{Sim}(X_B, X) = 0.1242$$

→ B has higher bias

higher bias ⇒ being farther away from the norm.

Therefore we choose B during ranking.

14.11.2023

### Maximal Margin Relevance (MMR)

- IR Query result diversification.

Query: Jaguar → Cat  
→ Drink  
→ Car  
→ OS } Which one?

ACM SIGIR, 1998

$$\text{MMR} = \underset{d_i \in R/S}{\operatorname{argmax}} \left[ \lambda \text{Sim}_1(d_i, Q) - (1-\lambda) \max_{d_j \in S} \text{Sim}_2(d_i, d_j) \right]$$

R: relevant doc.s in collection C

S: current result set

$\lambda$ : hyperparameter

- High  $\lambda \uparrow$  higher accuracy
- Low  $\lambda \downarrow$  " diversity
- The motivation create a set of docs that not only contains highly relevant info but also provides a diverse & comprehensive coverage of the topic.
- Showing multiple similar doc.s , even if they are relevant, introduces redundancy to the query results. MMR provides a balance between relevance & redundancy.

Example.

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$q$	<i>the query</i>
$d_1$	1	0.11	0.23	0.76	0.25	0.91	
$d_2$		1	0.29	0.57	0.51	0.90	
$d_3$			1	0.02	0.20	0.50	
$d_4$				1	0.33	0.06	
$d_5$					1	0.63	
$q$						1	

1<sup>st</sup> iteration:  $S = \emptyset$ ,  $R/S = \{d_1, d_2, d_3, d_4, d_5\}$

$$MMR = \operatorname{argmax}_{d_i \in R/S} \left[ \lambda \text{Sim}_1(d_i, Q) - (1-\lambda) \max_{d_j \in S} \text{Sim}_2(d_i, d_j) \right]$$

$\hookrightarrow S = \emptyset$

$$MMR = \operatorname{argmax}_{d_i \in R/S} (\lambda \text{Sim}_1(d_i, q))$$

↓ doesn't change argmax here

$$= d_1 \rightarrow \text{has max sim w/ } q.$$

$S = \{d_1\}$ ,  $R/S = \{d_2, d_3, d_4, d_5\}$ .

2<sup>nd</sup> iteration:  $\lambda = 0.5$

$$MMR = \operatorname{argmax}_{d_i \in R/S} \left[ 0.5 \left( \text{Sim}_1(d_i, Q) - \max_{d_j \in S} \text{Sim}_2(d_i, d_j) \right) \right]$$

λ doesn't affect results again.

$$MMR(d_i) = 0.5 \left[ \text{Sim}_1(d_i, Q) - \text{Sim}_2(d_i, d_1) \right]$$

$$MMR(d_2) = 0.5 (0.90 - 0.11) = 0.395$$

$$MMR(d_3) = 0.5 (0.50 - 0.23) = 0.135$$

$$\text{MMR}(d_4) = 0.5(0.06 - 0.76) = -0.35$$

$$\text{MMR}(d_5) = 0.5(0.63 - 0.25) = 0.19$$

$$\text{MMR} = \underset{d_i \in R/S}{\operatorname{argmax}} (\text{MMR}(d_i)) = d_2.$$

Update:

$$S = \{d_1, d_2\} \quad R/S = \{d_3, d_4, d_5\}$$

3<sup>rd</sup> iteration:

$$\text{MMR} = \underset{d_i \in R/S}{\operatorname{argmax}} \left[ 0.5 \left( \overline{\text{Sim}_1(d_i, Q)} - \max_{d_j \in S} \text{Sim}_2(d_i, d_j) \right) \right]$$

$$\begin{aligned} \text{MMR}(d_3) &= 0.5(0.50 - \max \{ \text{Sim}_2(d_3, d_2), \text{Sim}_2(d_3, d_1) \}) \\ &= 0.5(0.50 - \max \{ 0.23, 0.29 \}) \\ &= 0.50(0.50 - 0.29) = 0.105 \end{aligned}$$

$$\text{MMR}(d_4) = 0.5(0.06 - \max \{ 0.76, 0.57 \}) = -0.35$$

$$\text{MMR}(d_5) = 0.5(0.63 - \max \{ 0.25, 0.51 \}) = 0.065$$

$$\text{MMR} = \underset{d_i \in R/S}{\operatorname{argmax}} \text{ MMR}(d_i) = d_3$$

Update :

$$S = \{d_1, d_2, d_3\} \quad R/S = \{d_4, d_5\}$$

After 3<sup>rd</sup> iteration for  $\lambda=5$ ,  $S = \text{result set} = \{d_1, d_2, d_3\}$ .

Total pairwise similarity for S:

$$\text{sim}(d_1, d_2) + \text{sim}(d_2, d_3) + \text{sim}(d_1, d_3) = 0.63$$

For  $\lambda=1$ ,  $S'$  would be  $\{d_1, d_2, d_5\}$ .

Total pairwise similarity for  $S'$ :

$$\text{sim}(d_1, d_2) + \text{sim}(d_2, d_5) + \text{sim}(d_1, d_5) = 0.97.$$

$\lambda \downarrow$ , diversity  $\uparrow$

21.11.2023

## PAT trees & PAT Arrays (used for string search)

Info ret.: algs.s & data structures

Example binary alphabet

1 2 3 4 5 6 7 8 9 10 11 12  
0 1 1 0 0 1 0 0 0 1 0 1

(semi infinite string)

→ sis

sis1: 011001000101...

sis2: 11001000101...

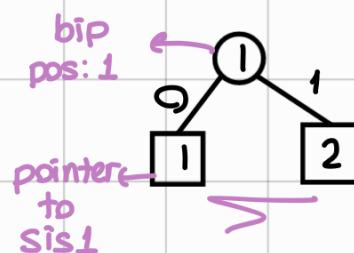
sis3: 1001000101...

sis4: 001000101...

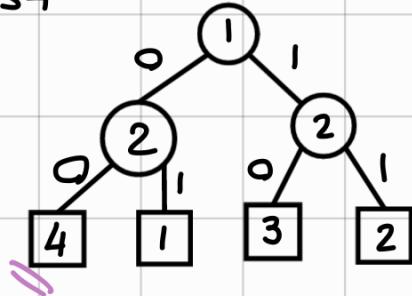
sis5: 01000101...

sis6: 1000101...

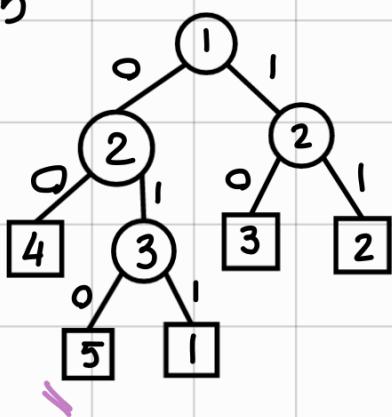
sis7: 000101...



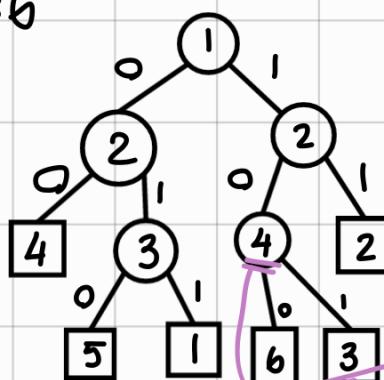
sis4



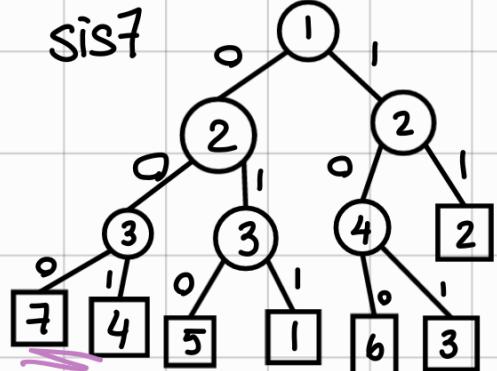
sis5



sis6

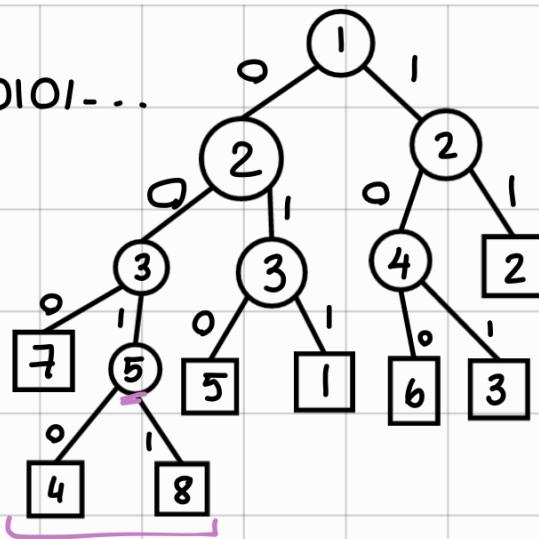


sis7



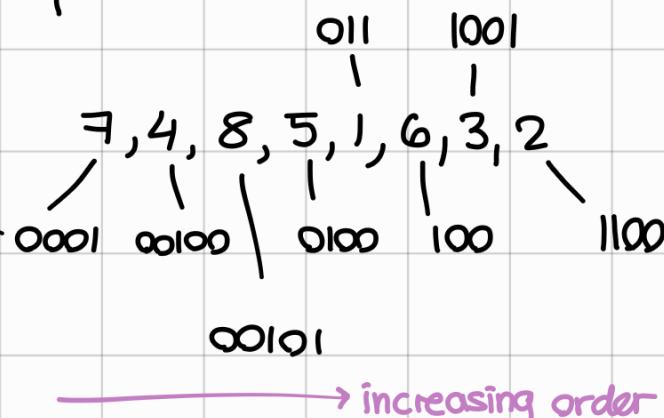
diff. at the 80<sup>th</sup> pos, we could dif at 4  
because we couldn't

sis8: 00101...



PAT Array

(bit patterns)



Gonnet et. al

- Perform a binary search in the PAT array to see if a string appears for a string query.

(You can keep the PAT tree as well,

⚠️ but you don't need to keep the tree structure, the array is enough for binary search.)

- Look at the middle item, if query > item : search right subtree ; else, left subtree.
- This method is for binary alphabet.

PAT tree of a text:

- same principle, we should be able to decide  $>$ ,  $<$ ,  $=$

sis1: that is that that is not that

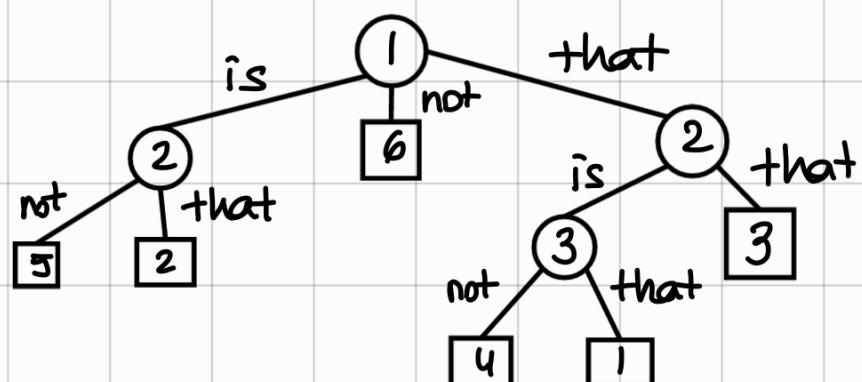
sis2: is that that is not that

sis3: that that is not that

sis4: that is not that

sis5: is not that

sis6: not that



## Proximity Search

string1  $\leftrightarrow$  string2  
dist

A string  $s_1$  is at most a fixed # characters away from another string  $s_2$ .

Distance between 01 and 10

5,1

6,3

ASK

what does this mean?

28.11.2023

## Evaluation Under Diversity

Query: James Bond

Actors, Music, Director, Places, Singer

→ Ambiguous query

Remark: we used MMR w/ low  $\lambda$  to increase diversity of query results.



S-Recall (Subtopic recall) =

$$\frac{(\text{Unique \# topics covered until } n^{\text{th}} \text{ rank})}{(\text{Total \# topics})}$$

Example:

<u>Rank</u>	<u>Document</u>	<u>Subtopic</u> ↪	
1	$d_1$	$m_3$	Assume $\exists 6$ subtopics
2	$d_2$	$m_4$	$\rightarrow S\text{-Recall @ 2} = \frac{ \{m_3, m_4\} }{6}$
3	$d_3$	$m_1, m_2$	$= 2/6 \approx 0.33$
4	$d_4$	$m_5, m_6$	
5	$d_5$	$m_6$	$\rightarrow S\text{-Recall @ 6} = \frac{ \{m_1, m_2, \dots, m_6\} }{6}$
6	$d_6$	$m_5$	$= 6/6 = 1.$
7	$d_7$	$m_4$	
8	$d_8$	$m_3$	
9	$d_9$	$m_2$	
10	$d_{10}$	$m_1$	

Rank	$m_1$	$m_2$	$m_3$	$m_4$	$m_5$	$m_6$
1			1			
2				1		
3	1	1				
4				1	1	
5						1
6					1	
7			1			
8				1		
9		1				
10	1					
P@5	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{2}{5}$

IA - Precision

Precision - IA

Interest Aware

$$\begin{aligned}
 \text{Precision-IA}@5 &= \frac{1}{|m_i|} \sum_{m_i} P@5 \\
 &= \frac{1}{6} \cdot \left( \frac{1}{5} + \frac{1}{5} + \frac{1}{5} + \frac{1}{5} + \frac{1}{5} + \frac{2}{5} \right) \\
 &\approx 0.23
 \end{aligned}$$

- Recommendations should cover diff. types of mentions
- NOTE: this topic "eval. under div." belongs w/ MM2.

# Signature Files

Document → Document Signature



Query Signature →



Matching

Document

True Match



False Match (eliminate them!)



Term<sub>1</sub> → Term Sig<sub>1</sub>

Term<sub>2</sub> → Term Sig<sub>2</sub>

Superimpose →  
(or, v)

Doc

Sig

F = signature size = # bits  $\approx 500$  bits

- Now let's assume sig. size F = 8.

<u>Term</u>	<u>Term Signature</u>	
Object	1000	1000
Signature	0010	0100
Generation	1000	0100

1010 1100 ← doc sig. (DS)  
85

$m = \#$  bits to be set  $\forall$  term  $i$ .  
(for each term)

Object  $\rightarrow$  obtain its doc equivalent



Consider conjunctive queries  $= q_1 \& q_2$

Query  $\rightarrow$  Generation

Query Sig. (QS) : 1000 0100

! if  $Q_s \& D_s == Q_s$ , then retrieve doc.  
else, skip.

Example.

$$D_s = 1010 1100$$

$$\wedge Q_s = 1000 0100$$

$$\underline{1000 \ 0100} = Q_s \rightarrow \text{doc contains the term,}$$

true match

- For matchers, obtain the doc. Make sure it contains the term.

Query : "Operation"

$$Q_s = 1000 \ 0100$$

(why is it the same  
as query "generation")

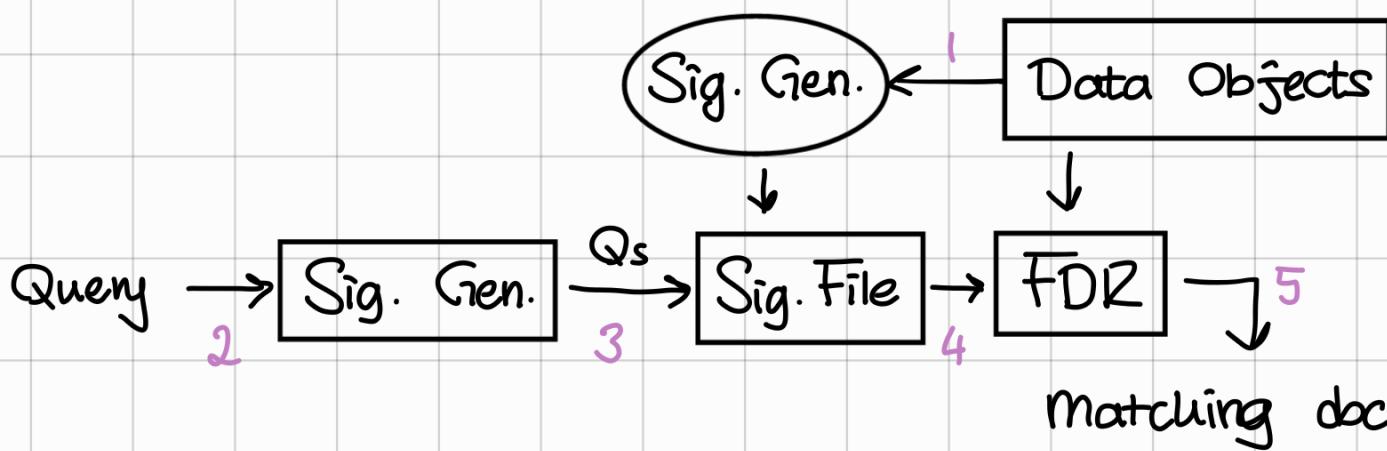
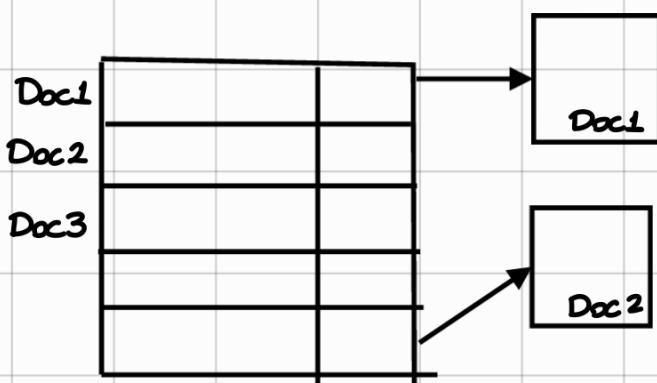
$$D_s = 1010 \ 1100$$

$$\begin{array}{r} \wedge Q_s = 1000 \ 0100 \\ \hline \end{array}$$

$$\begin{array}{r} 1000 \ 0100 \\ \downarrow \\ \hline \end{array} = Q_s$$

This is false match because "operation" term is not in our real doc.

False Drop Elimination = False Drop Resolution (FDR)



- MIT ms 1991  
UNN of Toronto  
(Christodolay?  
Faloutsos , Astor Zobel)

- Docs are divided into blocks w/ the same  
# unique term.



$F$  = signature size (fixed size)

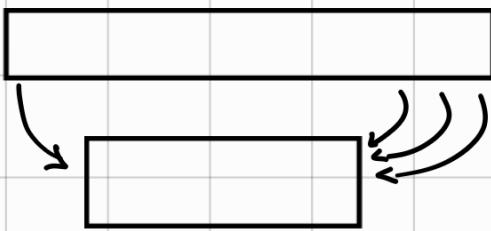
$D$  = avg # unique term / block

$m$  = # bits to set per term

$$m = \frac{F \cdot \ell_1 2}{D} \rightarrow \text{Prob. of having a '0' or a '1' is the same } \forall \text{ bits positions of}$$

- Lee & Leng 1993

They claimed to generate block signature w/ the same # '1's.



update sig.

- In a doc.sig.

$$\#\text{"1"}\text{ s} = \#\text{"0"}\text{ s}$$

- We will see different signature file organizations:

- extendible hashing
- linear hashing

## Storage Structures for Signature Files

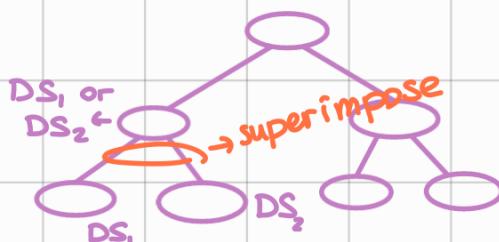
- Aims to avoid query-doc signature file comparison.

### 1. Single Level

- Sequential (problem here is we would have to compare all)

### 2 Multi Level

- Signature tree →



### 3. Vertical Partitioning

- bit sliced

- frame sliced

### 4. Horizontal Partitioning

- fixed prefix

- extendible hashing

- linear hashing

## Bit Sliced Signature Files

$S_1: 0001\ 1110$

$S_2: 1101\ 0001$  columns

$S_3: 0011\ 1100$   $\xrightarrow{\text{of}}$

$S_4: 1100\ 0011$  these  
sigs

$S_5: 0011\ 0110$

$S_6: 1100\ 1001$

$C1: 0\ 1\ 0\ 1\ 0\ 1$

$C2: 0\ 1\ 0\ 1\ 0\ 1$

$C3: 0\ 0\ 1\ 0\ 1\ 0$

$C4: 1\ 1\ 1\ 0\ 1\ 0$

$C5: 1\ 0\ 1\ 0\ 0\ 1$

$C6: 1\ 0\ 1\ 0\ 1\ 0$

$C7: 1\ 0\ 0\ 1\ 0\ 0$

$C8: 0\ 1\ 0\ 1\ 0\ 1$

$C1$	$C2$	$\dots$	$C8$
------	------	---------	------

↑ extra new space for possible new columns ?

- Insertion of a new signature is difficult,  
since col.s are stored one after another. That's  
why we store an extra space for possible new  
comes.

Ex.  $Q_s = 1011 \ 0000$

↓      ↘  
column1    column3

C1 & C3

$$\begin{array}{cc} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{array} \wedge \begin{array}{c} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} = \left\{ \begin{array}{l} \text{None of the 6 sigs have the terms} \\ \text{in } Q_s. \text{ Don't proceed further.} \end{array} \right.$$

↳ intermediate result

We compared 6 bits instead of 8. ↶

(works when  $|sigs| < |columns in a sig|$ )

6                8

If not all sigs would have 0, we would proceed with column 4, by (intermediate result)  $\wedge$  C4.

- Here bit density is 50% (half '0's, half '1's)

$\Rightarrow$  At step (of this col. comparison), we eliminate

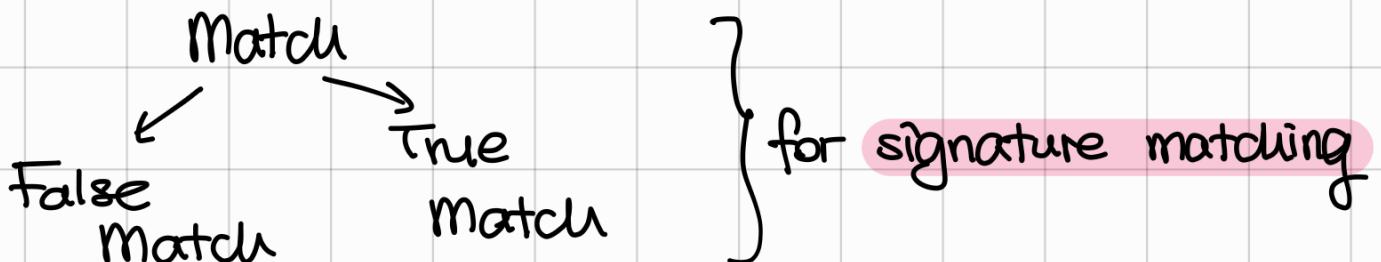
half of them. → ASK: half of cols or sigs?

- If remaining '1's are so small, we switch to FDR.

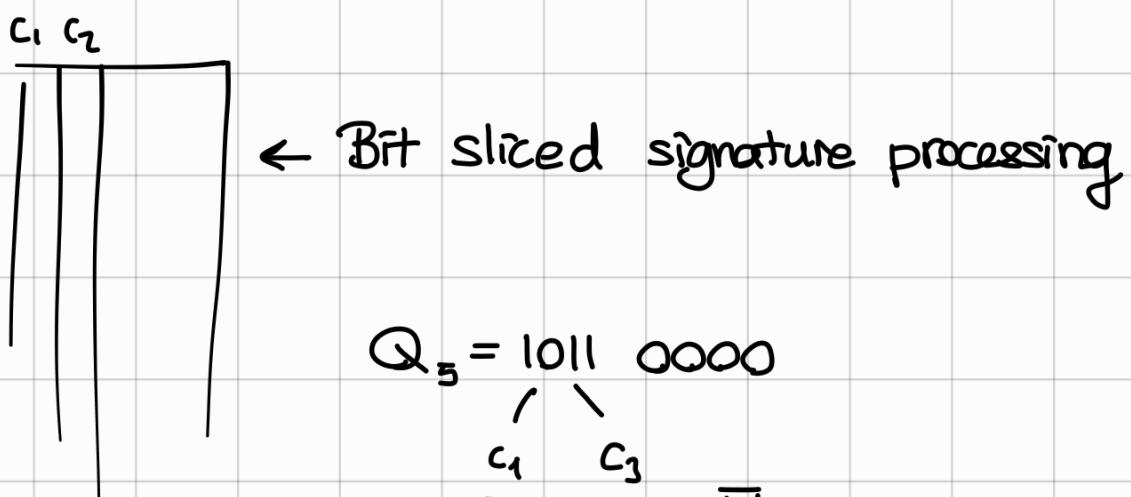
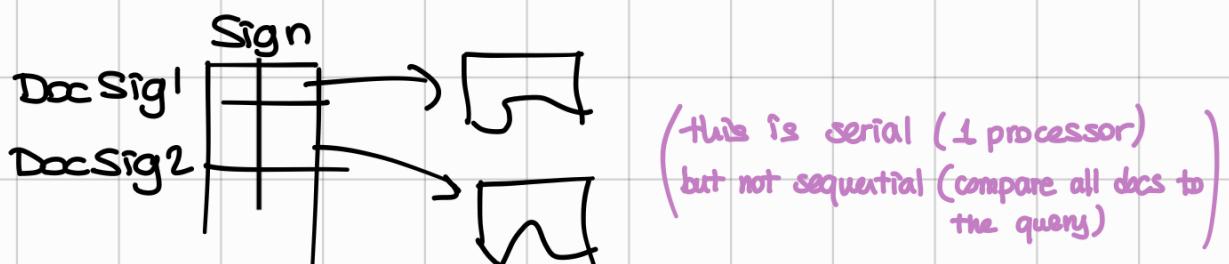
Since now FDR is less expensive than ANDing.

- Karaburk & Can: Partial Eval. of Queries for Bit-Sliced Sig. Files

05.12.2023



Falsedrop resolution: elimination of false drops.  
(false matches)



$$Q_5 = 1011 \quad 0000$$

$$\boxed{\phantom{0}} \& \boxed{\phantom{0}} = \boxed{\phantom{0}}$$

50% "0"

50% "1"

Example. (we are repeating from a prev. lecture)

$S_1$ : 0001 1110

$S_2$ : 1101 0001

$S_3$ : 0011 1100

$S_4$ : 1100 0011

$S_5$ : 0011 0110

$S_6$ : 1100 1001

( $\neg f 1001$ )  $\rightarrow$  1<sup>st</sup> and 4<sup>th</sup> cols.

$Q_S$ : 1011 0000

$S_1$ : 0001 1110

$\wedge Q_S: 1011 0000$

0001

0000  $\rightarrow$  not a match

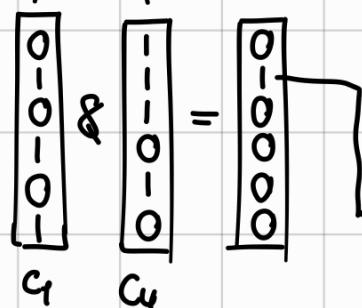
$\neq Q_S$

$C_1$ : 010 101

$C_2$ : 010 100  $\rightarrow Q_S: \begin{array}{|c} 1001 \\ \hline 0000 \end{array}$

$C_3$ :

$C_4$ : 111 010



$s_2$  satisfies

Query weight:

# 1's in a  
query's  $Q_S$

$\hookrightarrow w(Q)$

Ex.  $Q_S: 1001 0000$

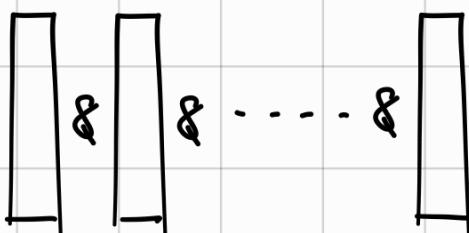
$w(Q_S) = 2$ .

Paper:

# Partial Evaluation of Queries for Bit Sliced Signature Files (Information Processing Letters, Kochberber)

$$OP = \frac{\# 1's \text{ in a signature}}{\text{signature length}} = \text{bit density}$$

$$\text{false drop probability} = \frac{\# \text{ matches}}{\underbrace{\# \text{ signatures}}_N}$$



After processing  $i$  many bits of the query:

$$f_{d_i} = OP$$

freq. of  $d_i = s_i$

on bit density of  $s_i$

(this is I think for bits  
not the same)  $i$  in  $d_i$ .

ASK

so maybe revise  
the formula as

$$f_{d_i} = OP_i ?$$

- We assume all matches are false matches.

$$RT(i) = \text{response time after processing } i \text{ many bits}$$

$$= \underbrace{i \cdot T_{\text{slice}}}_{\substack{\downarrow \\ \text{time needed} \\ \text{to process a} \\ \text{signature}}} + \underbrace{N \cdot op^i \cdot T_{\text{resolve}}}_{\substack{\downarrow \\ \# \text{ matching} \\ \text{signatures}}} \xrightarrow{\text{assume false} \\ \text{match}}$$

$N = \# \text{sigs}$

To find the  $i$  value that minimizes the  $RT(i)$ :  
 (don't need to memorize the derivation)

↳ We take the derivative of  $RT(i)$  wrt  $i$ :

$$\frac{d RT(i)}{di} = T_{\text{slice}} + N \cdot T_{\text{resolve}} \cdot op^i \cdot \ln op \quad (1)$$

To find the optimum  $i$ :

↳ Let the eq. (1) be 0 and solve for  $i$ :

$$op^i = \frac{T_{\text{slice}}}{N \cdot T_{\text{resolve}} \cdot (-\ln op)}$$

$$\ln op^i = \ln \left( \frac{T_{\text{slice}}}{N \cdot T_{\text{resolve}} \cdot (-\ln op)} \right)$$

$$\Rightarrow i = \lfloor \left( \frac{T_{\text{slice}}}{N \cdot T_{\text{resolve}} \cdot (-\ln p)} \right) / \ell_i(\text{op}) \rfloor$$

$w(Q)$  = # 1's in a query

If  $i > w(Q)$  then  $i = w(Q)$

[i.e. if  $i = \arg \min R_T(i) > w(Q)$ , then stop preprocessing at  $w(Q)$ ]

else, preprocess until  $i$  (don't go through all  $w(Q)$ ) ]

## Horizontal Partitioning

(Lee & Leng, ACM TOIS, 1993)

Signature density = 0.5 (half of the bits are 1)  
50%

$s_1: 0111 1000$

$s_5: 0110 1100$

$s_2: 1000 1011$

$s_6: 1001 0011$

$s_3: 0011 1100$

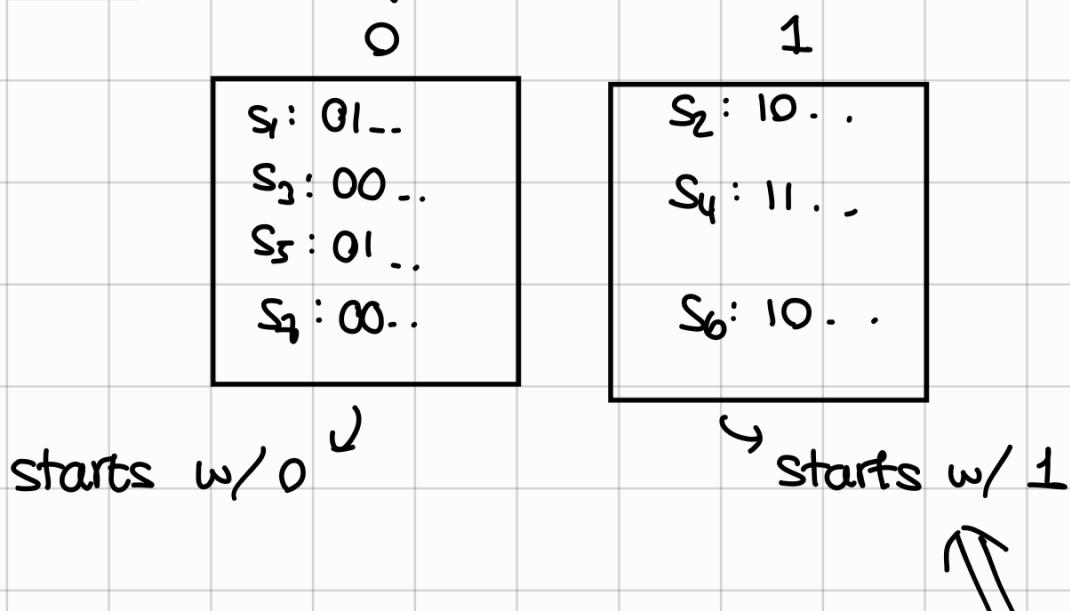
$s_7: 0000 1111$

$s_4: 1100 0011$

$Q = 1001 0001$   
||

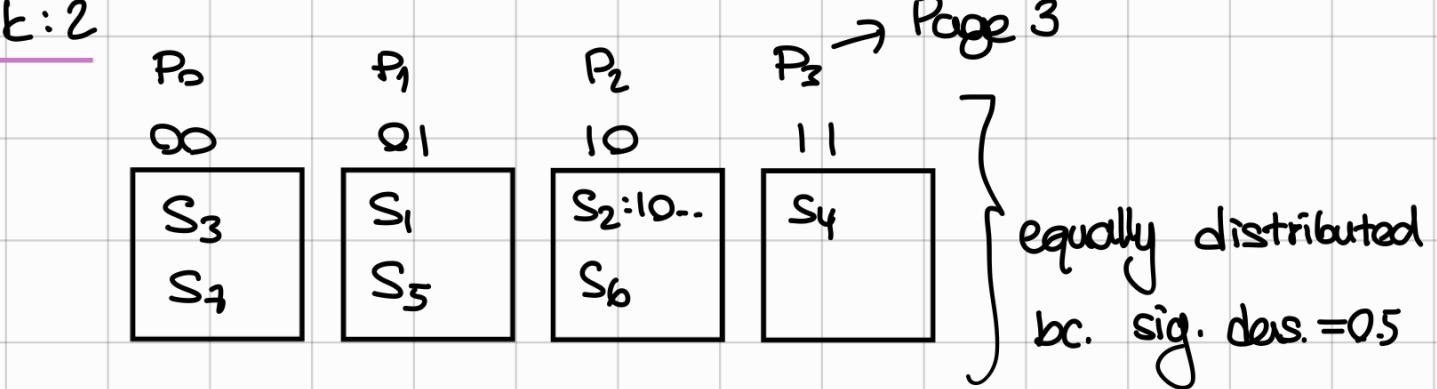
## Fixed prefix partitioning

k=1 (use 1-prefix):



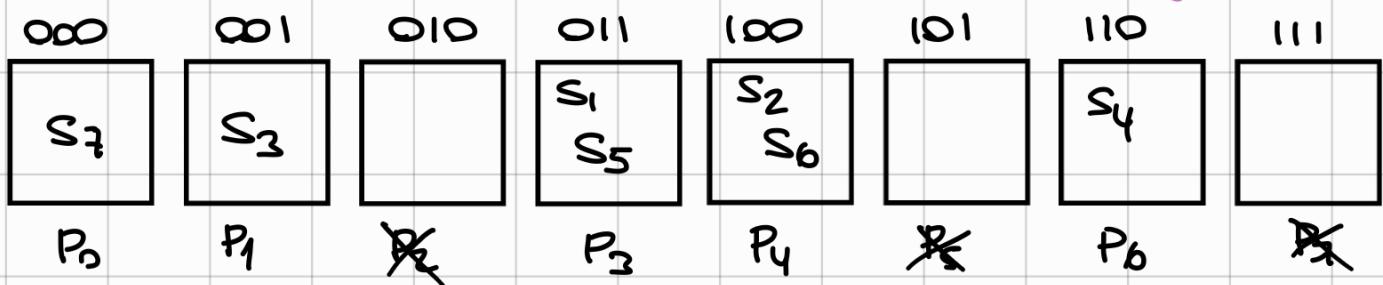
Query  $\rightarrow Q = 1001\ 0001 \Rightarrow$  start's w/ 1

k=2



k=3

Why not equally dist. now?



$(100, 0001)$

Q matching partitions:

100

$\Rightarrow 1XX$

97

(X: don't care)

$P_{100}, P_{101}, P_{110}, P_{111}$ .

Access a page if  $Q_{sk} \& P_s = Q_s$ , then access that page for false drop elimination.

If  $Q = 0000\ldots$  we would have to access all the pages.

" " = 1110... " " = " = 1 page.

$\Rightarrow$  # partitions to access:  $\begin{cases} 2 & (\text{# } 0\text{'s in the first } k \text{ bits}) \\ & \text{of } Q_{sk} \end{cases}$

<u>Query</u>	<u><math>k=1</math></u>	<u><math>k=2</math></u>	<u><math>k=3</math></u>
$Q_1 = 1000\ldots$	1(1)	2(10, 11)	4 (100, 101, 110, 111)
$Q_2 = 111\ldots$	1(1)	1 (11)	1 (111)
$Q_3 = 110\ldots$	1(1)	1 (11)	2 (110, 111)

↑

# pages to be accessed (page no)

# Connection Machine, paper, 1987

Parallel Environment

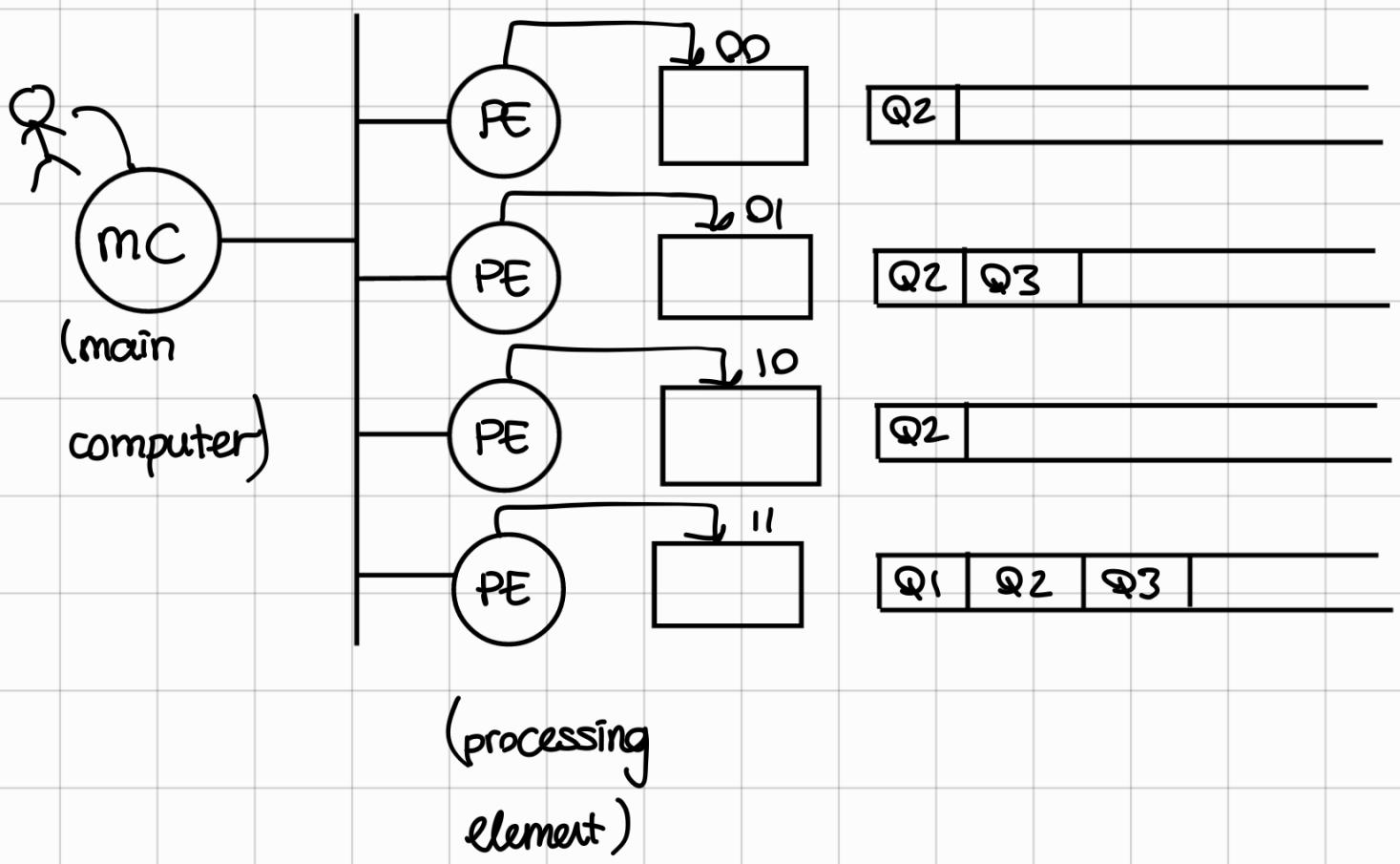
for  $t=2$ :

$t(2)$

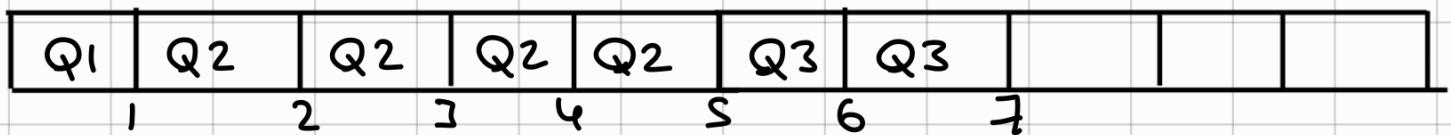
$Q_1 : 1110 \dots 1(11)$

$Q_2 : 0000 \dots 4(00,01,10,11)$

$Q_3 : 0110 \dots 2(01,11)$



Sequential processing:

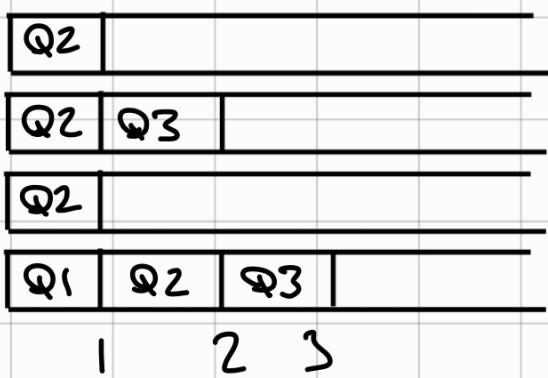


$Q_1 : 1 \text{ tu}$

$Q_2 : 5 \text{ tu}$

$Q_3 : 7 \text{ tu}$

## Parallel processing:



Q1: 1tu

tu: time unit

Q2: 2tu

Q3: 3tu

Avg turnaround time = (completion time - arrival time)

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

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

Throughput: # jobs completed with time

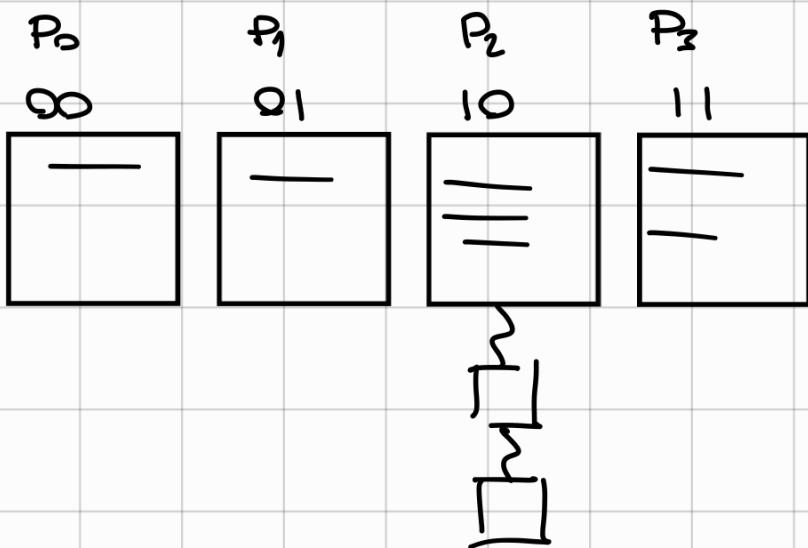
$$\text{Seq: } \frac{3 \text{ jobs}}{7 \text{ tu}}$$

$$\text{Parallel: } \frac{3 \text{ jobs}}{3 \text{ tu}}$$

Speed up ratio: total seq time / total par time  
7/3.

## Fixed Prefix

$k=2$



How many "# Qs vs sig. comparisons"?

$$Q = \underline{1}001 \ 0011 \quad k=2$$

IX:  $P_{10}$  and  $P_{11}$  will be accessed.

- There is an **overflow** problem:  
if  $\nexists$  place in the page then we have overflow pages. It slows the response time down.
- If we have too many overflow pages, we have to maintain the file.  
we want data structures that maintain themselves.

## Extendable Hashing

- Extendable Hashing paper } uses prefix  
Fagin et al, ACM TODS 19 }
  - Linear Hashing paper } uses suffix  
Litwin et al, 1980 VLDB }
- } both horizontal partitioning

Example. S1: 0011 1000      Blocking factor = BF = 2.

S2 : 0100 0011

||

S3: 0100 1011

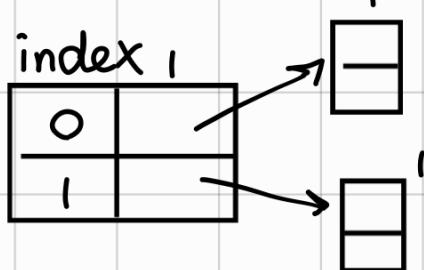
2 signatures

per block.

S4: 0100 1111

S5: 0000 1111

Initial Condition.

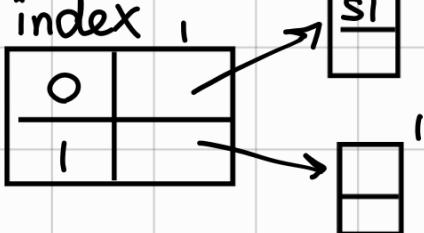


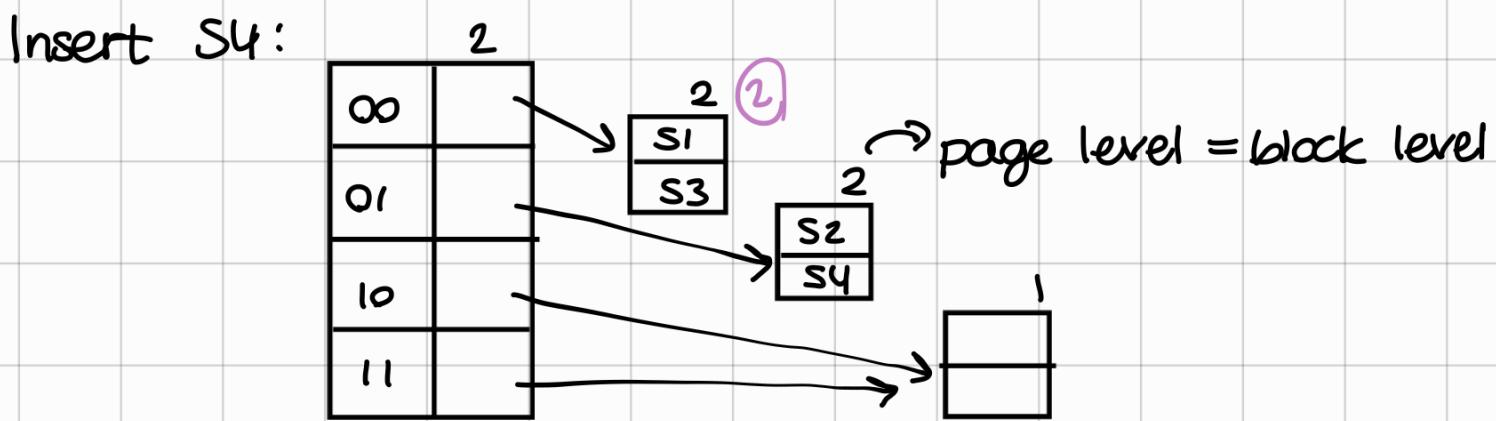
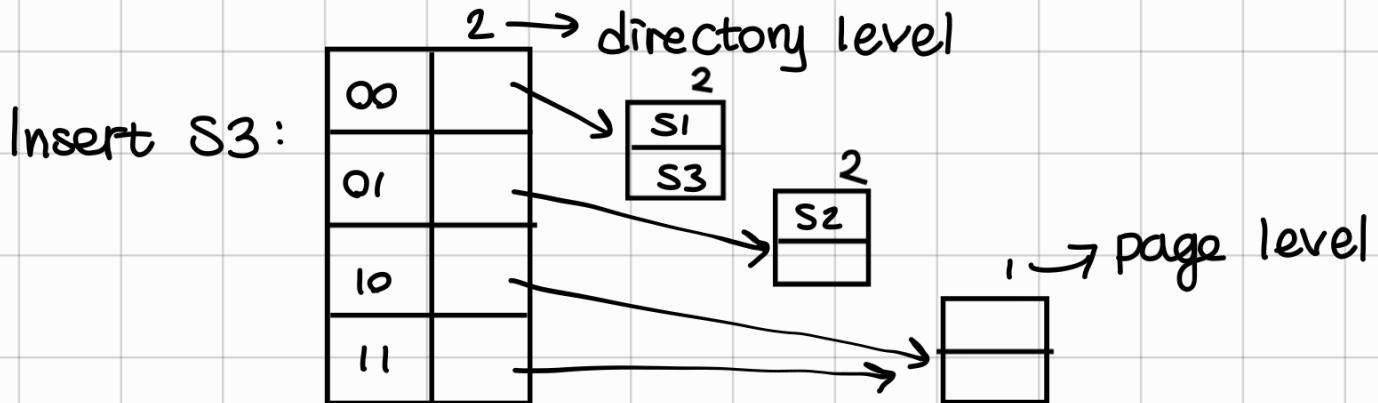
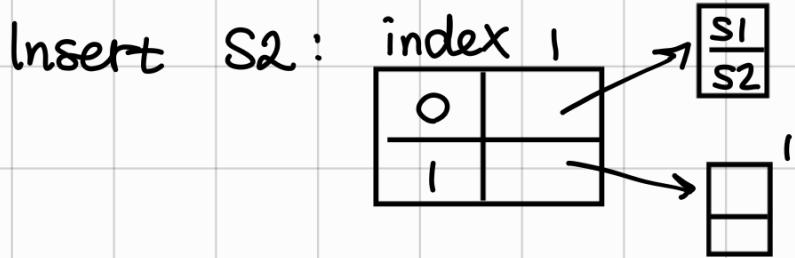
index = directory

directory level = 1

page level = 1

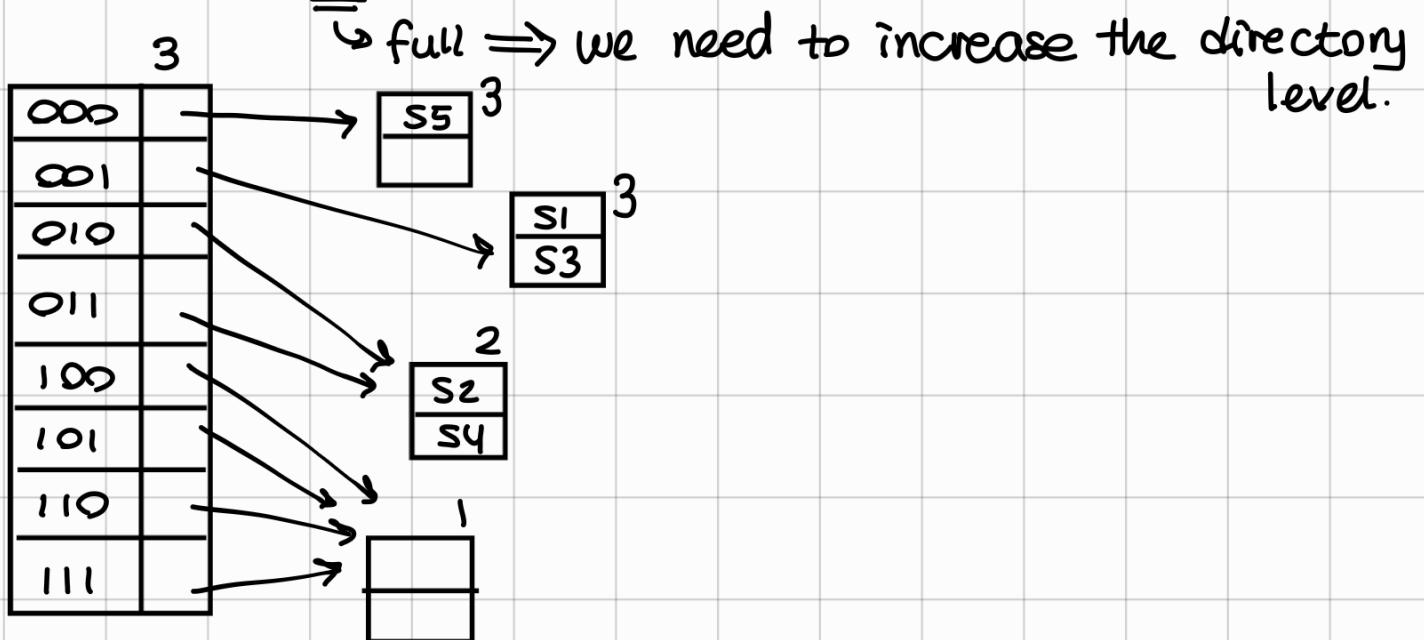
Insert S1: index 1





- If the storage is critically low, you don't need to store the first 2 bits for the above structure e.g. we know that the first 2 bits of S1 and S3 are "00".

Insert S5 : 0000 1111



$1 \leq$  page (block) level  $\leq$  directory (index) level

Let  $Q = 0110\ 0000$ .

Directory level = 3  $\Rightarrow k=3 \Rightarrow Q = \underline{0110}\ 0000$ .

$\Rightarrow$  Access Page<sub>011</sub> or Page<sub>111</sub>.

+ Extendible hashing prevents overflow pages!

vs

- There can be overflow in Linear Hashing  $\Downarrow$
- + We don't need to keep indices like directory level in Linear Hashing.  $\Downarrow$

## Linear Hashing vs Extendible Hashing

both horizontal partitioning

x		no overflow
no index keeping		x
suffix based		prefix based

12.12.2023

- Fixed prefix : dbg , needs cleaning
- Linear Hashing      ] : cat , cleans
- Extendible Hashing ] : itself

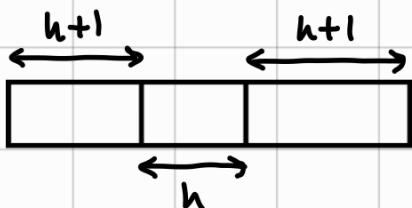
Ex. 0101 0011      2-suffix 11  
                            3-suffix 011

1993 QuickFilter

# Linear Hashing

## Insertion Algorithm

1. Use the  $h$ -bit suffix or  $h+1$  suffix to find the page to insert.



$$LF = \text{Load Factor} = \frac{\# \text{ recs in file}}{\# \text{ recs in prime area}}$$

Bkfr = blocking factor = # sig.s that can be stored in a page

2. When LF is at the desired level after adding  $LF \times Bkfr$  sig.s , add 1 more block to the file; and distribute the records of the block (pointed by bv) between this new block and the bie block.

3. Update bie (boundary value):

$$bv = bv + 1$$

if  $bv = 2^h$  then

$$h = h + 1$$

$$bv = 0$$

Example.

$S_1: 0011 \underline{1000}$	$S_5: 0000 \quad 1111$	$S_9: 0001 \quad 1100$
$S_2: 0100 \quad 0011$	$S_6: 1011 \quad 0001$	$S_{10}: 1100 \quad 0100$
$S_3: 0010 \quad 1011$	$S_7: 0001 \quad 0101$	$S_{11}: 1101 \quad 0100$
$S_4: 0100 \quad 1111$	$S_8: 0100 \quad 0000$	$S_{12}: 1111 \quad 0010$

Initially:  $br=0$ ,  $Bkfr=3$ ,  $h=1$ , desired  $LF = \frac{2}{3}$ .

desired hashing  $\rightarrow$   
level

be $\rightarrow$	0	
	1	

Add  $LF \times Bkfr = \frac{2}{3} \cdot 3 = 2$  blocks at a time.

Insert  $S_1, S_2$ :

be $\rightarrow$	0	$S_1$		
	1	$S_2$		

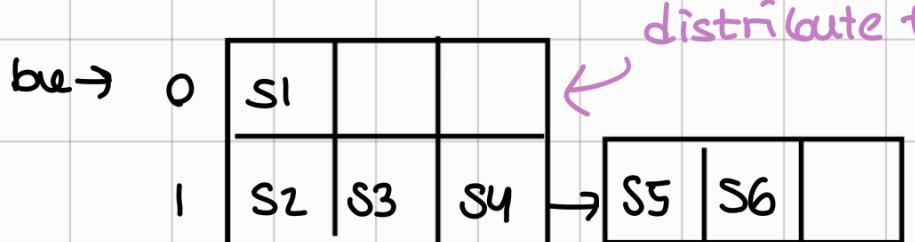
Insert  $S_3, S_4$ :

be $\rightarrow$	0	$S_1$		
	1	$S_2$	$S_3$	$S_4$

$LF = \frac{4}{6} = \frac{2}{3} \Rightarrow$  add 1 block & update by distributing

Insert S5, S6 :

→ added 2 sig.s,

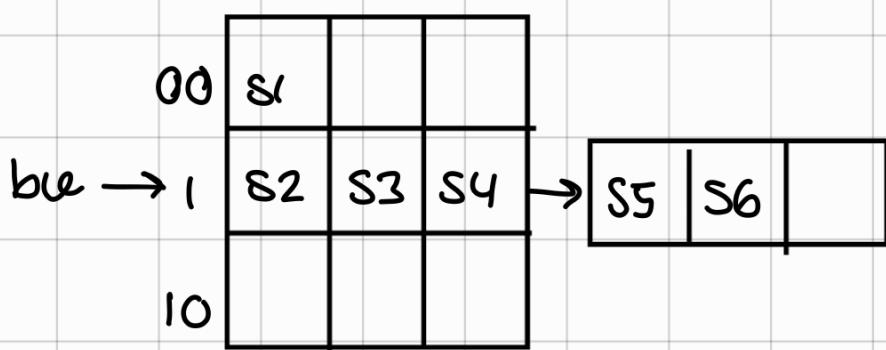


distribute these

now update

(the file & bu)

$$bue = bu + 1 \\ = 1.$$

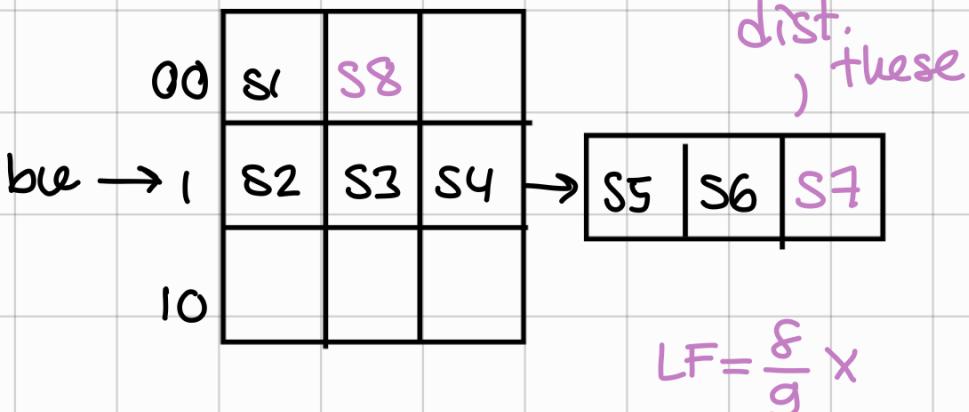


$$LF = \frac{6}{9} = \frac{2}{3}.$$

↓

add 2 more

Insert S7, S8:



dist.  
these

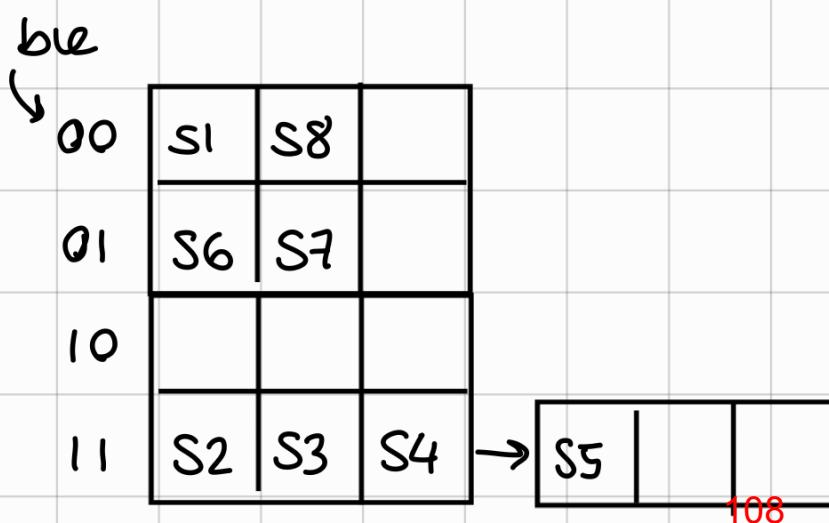
$$bue = bu + 1 = 2$$

$$= 2^h = 2^1$$

$$\hookrightarrow bu = 0$$

$$h = h+1 = 2.$$

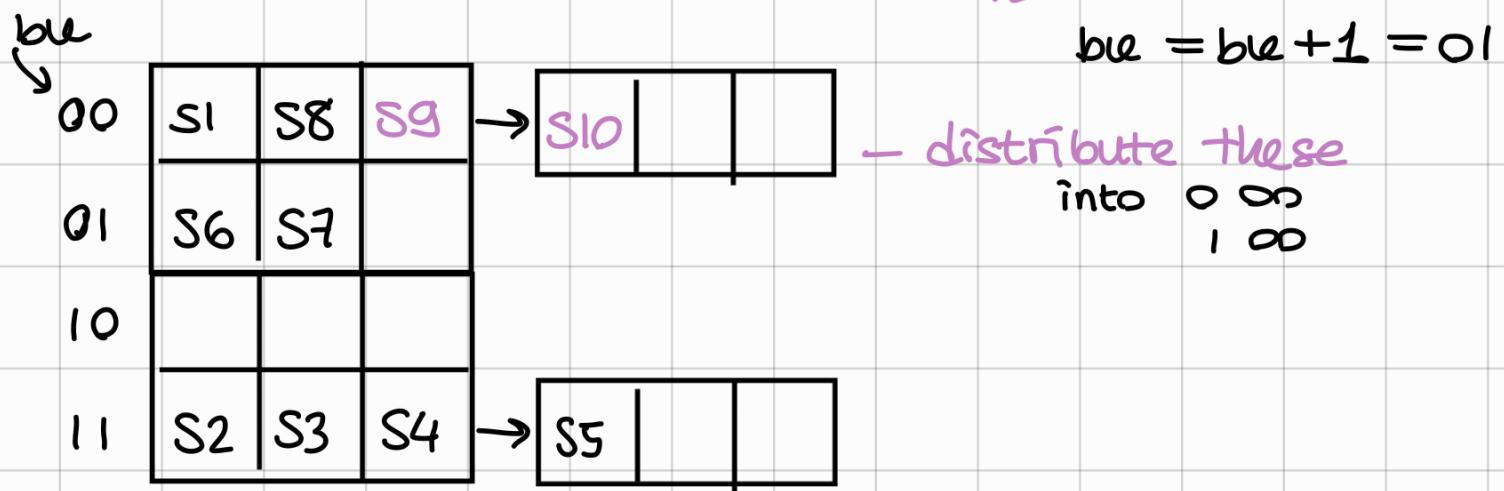
$$LF = \frac{8}{12} \times$$



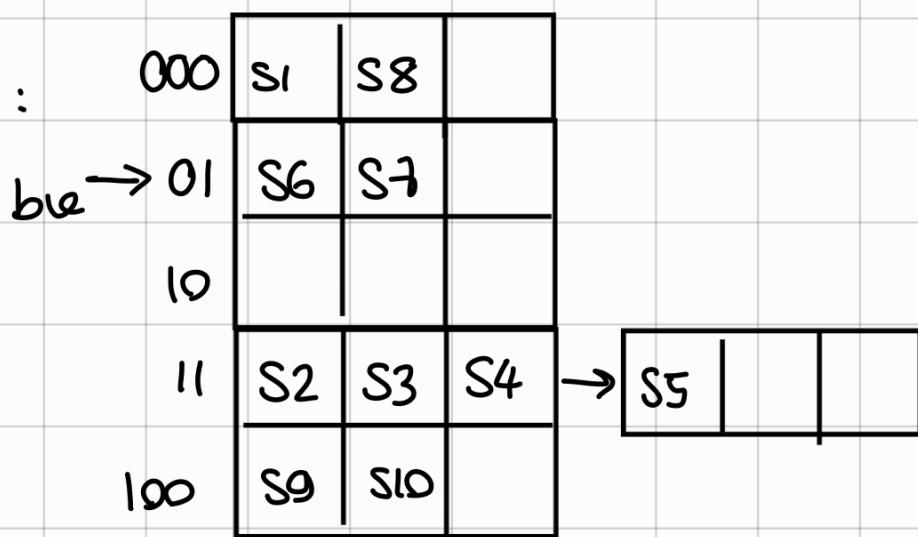
$$LF = \frac{8}{12} = \frac{2}{3} \checkmark$$

Insert S9, S10:

$$LF = \frac{10}{12} X$$

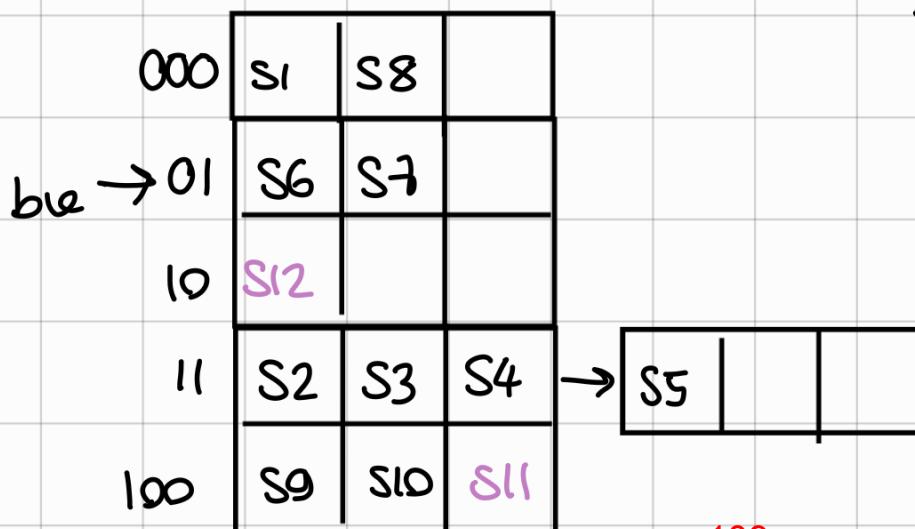


Update file:



$$LF = \frac{10}{15} \checkmark$$

Insert S11, S12:



$$bue = bu + 1 = 10.$$

$$h = 2.$$

$$LF = \frac{12}{15} X$$

Distribute 01 into 001 & 101 to maintain LF.

000	S1	S8	
001	S6		
10	S12		
11	S2	S3	S4
100	S9	S10	S11
101	S7		

$$LF = \frac{12}{18} \checkmark$$

$$ble = ble + 1 -$$

$$= 01 + 1 = 10.$$

Insert 2 sigs, if LF ✓

if LF > desired

dist. & update  
ble.

000		$\} h+1$
00		
10		$\} h$
11		
000		$\} h+1$
101		

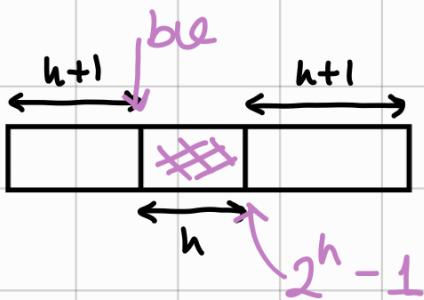
If we add a new sig, the structure would change

- $n = \# \text{ pages}$

$$\rightarrow n = 6$$

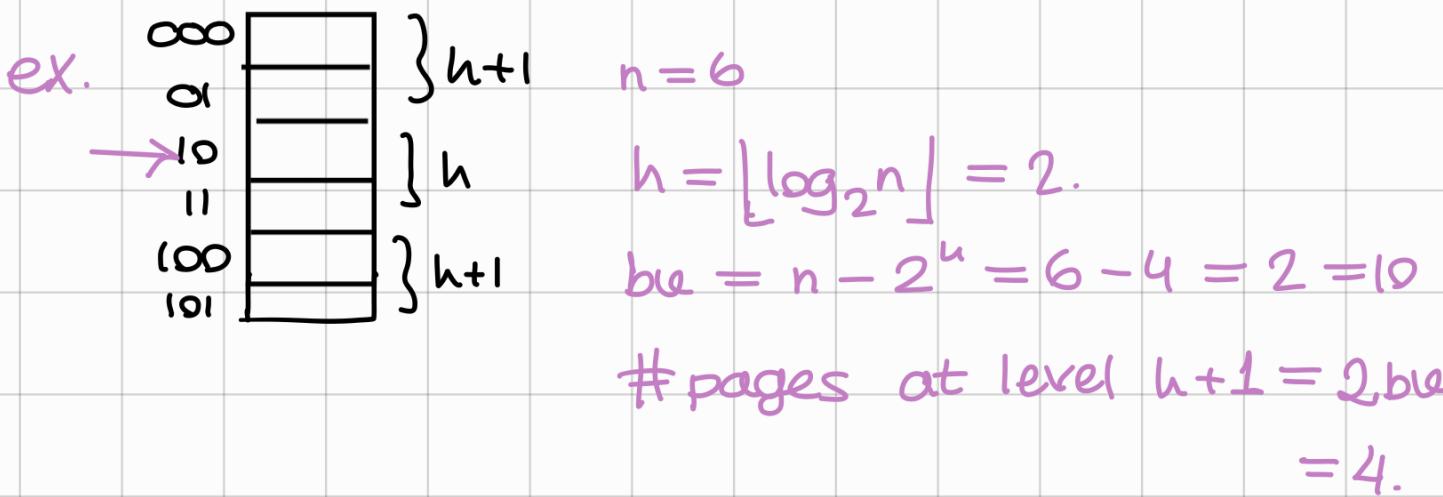
- $h = \lfloor \log_2 n \rfloor$

$$h = \lfloor \log_2 n \rfloor = 2 \checkmark$$



Pages  $[ble, 2^h - 1]$  are at level  $h$ .

- # pages at level  $h+1 = 2 \cdot ble$
- # " " " =  $h = n - 2^{ble}$
- $ble = f(n, h) = n - 2^h$



## Query Processing

Q: 0010 0110  
            
 if  $Q_k \wedge P_k = Q_k$   
 access page

$P_k$	$Q_k$	$P_k \wedge Q_k$
000	110	000 X
01	10	00 X
10	10	10 ✓
11	10	10 ✓
100	110	100 X
101	110	100 X

Info. ret. & info. filtering : Two Sides of the same Coin

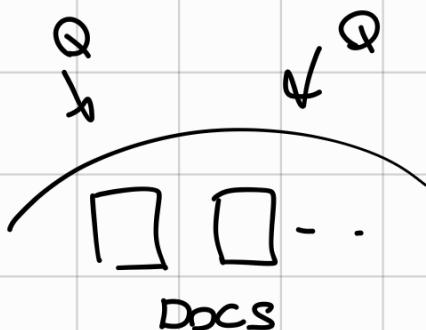
Bruce Croft UMass

1993

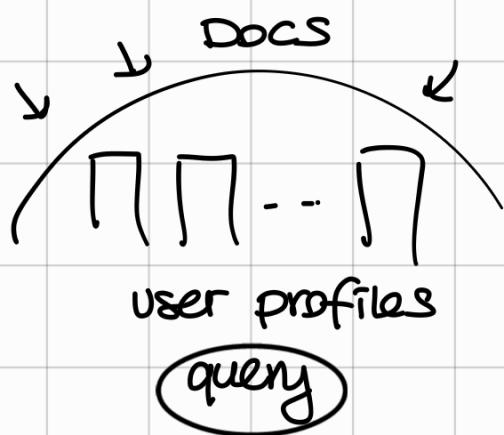
Berkin Rutgers

ACM

Info retrieval



Info. Filtering



One last approach in information filtering systems

smth like: new event detection tracking

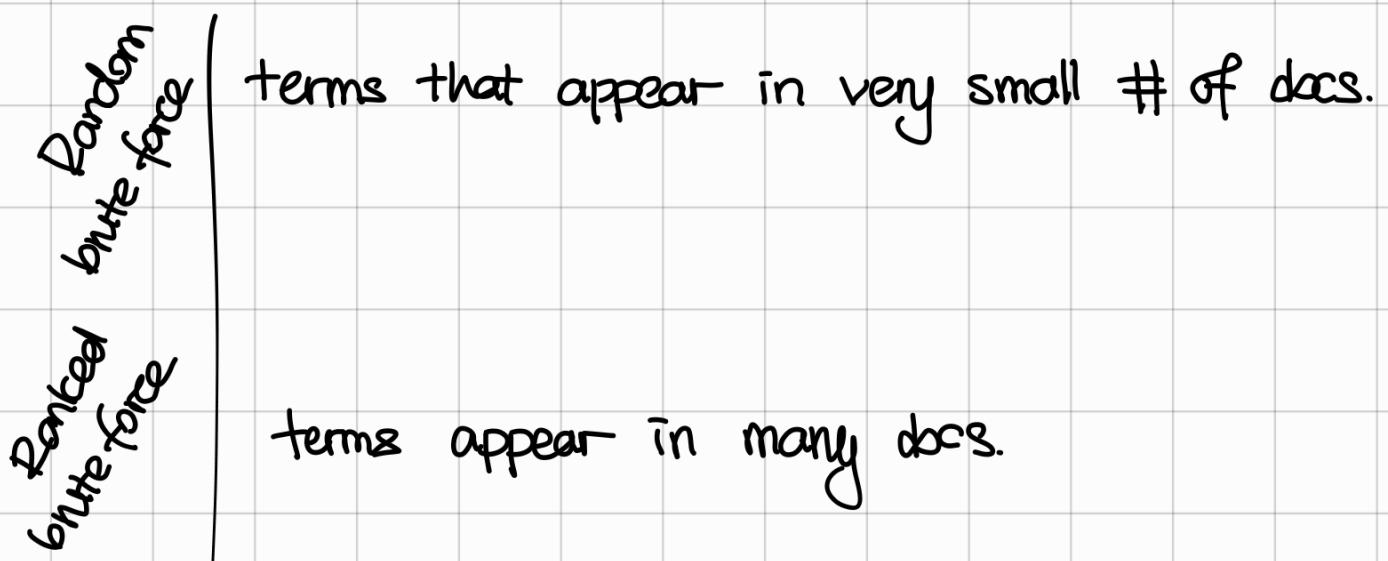
The filter bubble : how you are trapped in an echo chamber of your own interests, and not challenged

ACM TODS , TW Yang, Hector Garcia Morina:

- conjunctive user profiles: a doc satisfies the user profile if it has all of the terms in the profile.

## 1. Brute Force Approach

- a) Compare the incoming doc. with each profile one by one.  
(Random Brute Force)
- b) Use an occurrence table  
(Ranked Brute Force)



## 2. The Count Method

Sample profiles

$P_1 = (a, b)$  → user  $P_1$  is interested in terms a and b.

$P_2 = (a, d)$

$P_3 = (a, d, e)$

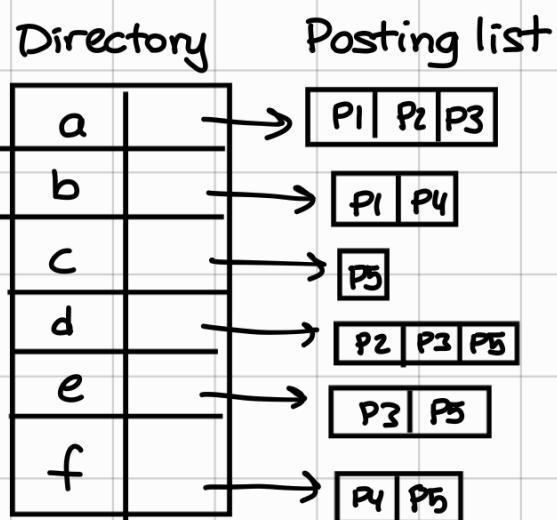
$P_4 = (b, f)$

$P_5 = (c, d, e, f)$

Unique doc terms = {a, b, c, d, e, f}

sample doc = a, c, a, f, b, c

inverted index →  
structure for  
profiles



size of each profile in terms of the # terms it has

Count Array

	TOTAL	a	b	c	f
P <sub>1</sub>	2	0	1	2	2
P <sub>2</sub>	2	0	1	1	1
P <sub>3</sub>	3	0	1	1	1
P <sub>4</sub>	2	0	0	1	1
P <sub>5</sub>	4	0	0	0	1

sample doc's  
terms = {a, b, c, f}

$$\left. \begin{array}{l} \text{Total (P1)} = \text{Count Array (P1)} \\ \text{Total (P4)} = \text{Count Array (P4)} \end{array} \right\} \begin{array}{l} \text{This doc satisfies P1} \\ \text{and P4.} \end{array}$$

$\text{Total (P5)} > \text{Count Array (P5)}$   $\rightarrow$  P5 is not satisfied  
w/ doc: a,c,a,f,b,c

### 3. The key method

In the key method a profile only appears in one of the posting lists. That term is called the key.  
( $\exists!$  term, i.e. key, whose posting list has the profile.)

#### a. Random key

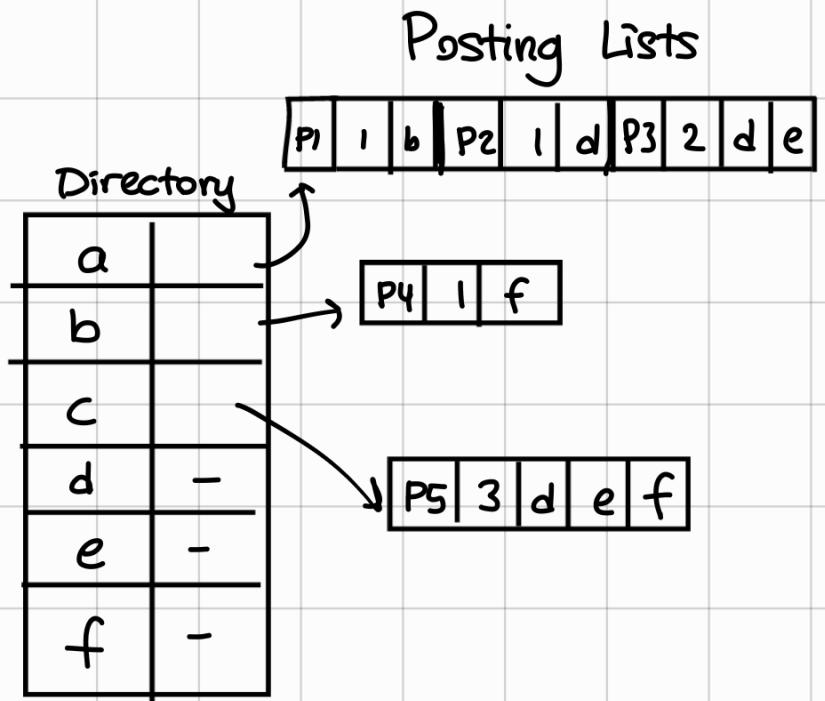
Any word is the key

#### b. Ranked key

Store the profile in the list of the word w/ lowest rank

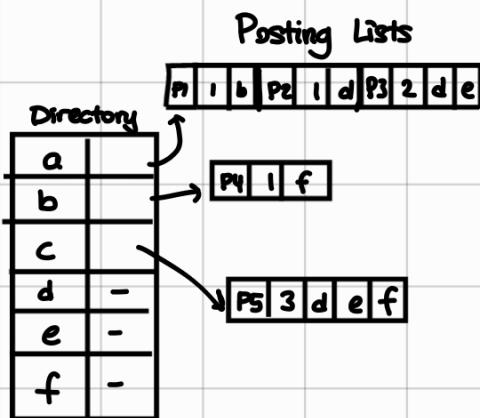
		lowest rank
a	1	→ least freq in docs - better discriminator
b	2	
c	1	
d	1	
e	1	
f	6	f is the most freq. in the docs
		highest rank

- A user profile can appear in only one posting list.



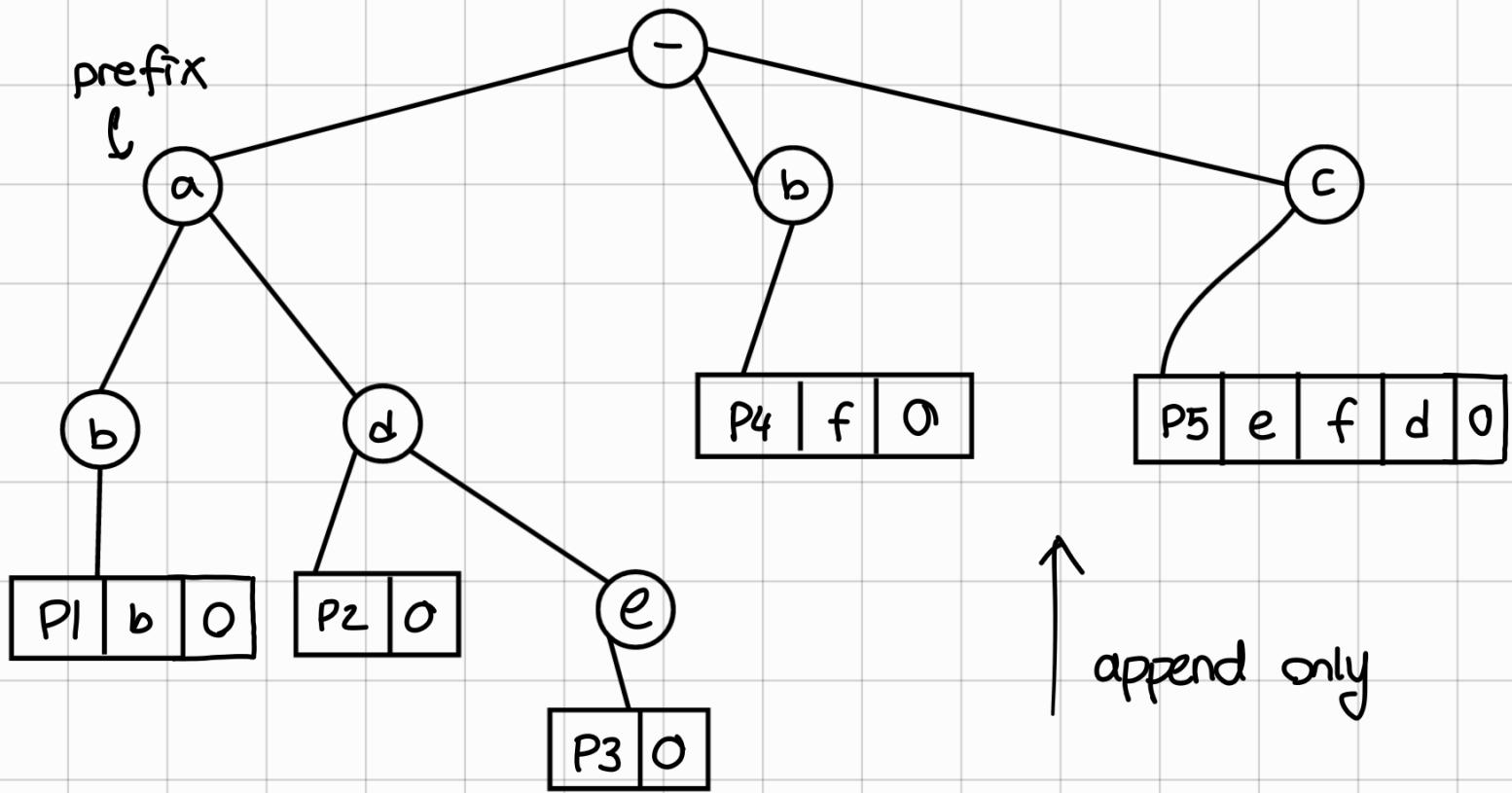
Doc : a,b,c,f

Consider posting list of a: P1(b) P2(d) P3(d) ]  
 ✓ b: P4(f)  
 ✓ c: P5(d,e,f) } doc satisfies  
 P1 & P4



graph for # elems in  
 posting list. It polynomially  
 decreases as the freq of the  
 term increase.

#### 4. The Tree Method



- Tree structure saves space when  $\exists$  many common prefixes among the profiles.

! These methods are only for conjunctive queries

Notes are by Ecem İlgün  
From Fall 2023