

13.02.2012

# STAIRS

IBM search engine

How to rank?

Value of a term in

$$\text{a retrieved document} = f(a, b, c)$$

a → the freq. of the term in the doc.

b → the freq. of the term in the  
retrieved set

c → the no. of docs in the retrieved  
set in which the term appears

Value of a term = 
$$\frac{a \times b}{c}$$

### Score of a document

$\sum$  value of all query terms  
which appear in the document

### example

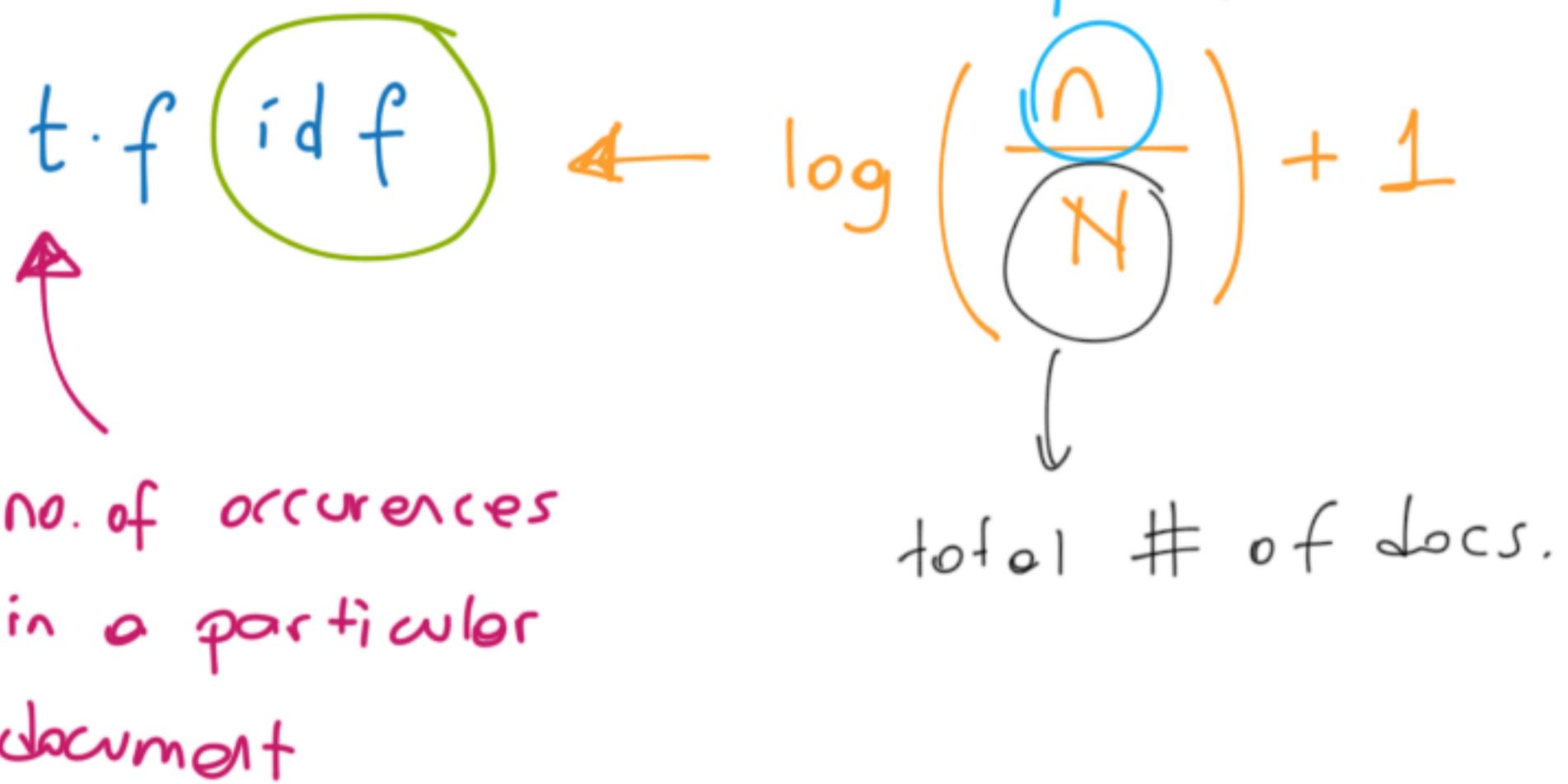
$a = 16$   $\Rightarrow$  appears this many no. of times  
in this particular document

$b = 12\#7$  appears this many no. of times  
in the retrieved documents

$c = 152$  it appears in this many no.  
of distinct documents

$$\text{Value} = \frac{16 \times 1247}{152} = 131,26$$

that



Van Rijsbergen 2 chapters  
 ↳ Information Retrieval

## Evaluation

Effectiveness

Efficiency

↳ space & time

## Search engines

ease of use

coverage of database

upto date

no broken links

response time

resource requirements

# Effectiveness

How to measure?

treceval

TREC = Text Retrieval Conference

Precision

Recall

↳ Assumes the availability of a test collection

- a set of documents
- a set of queries
- a set of relevant documents
- for each query

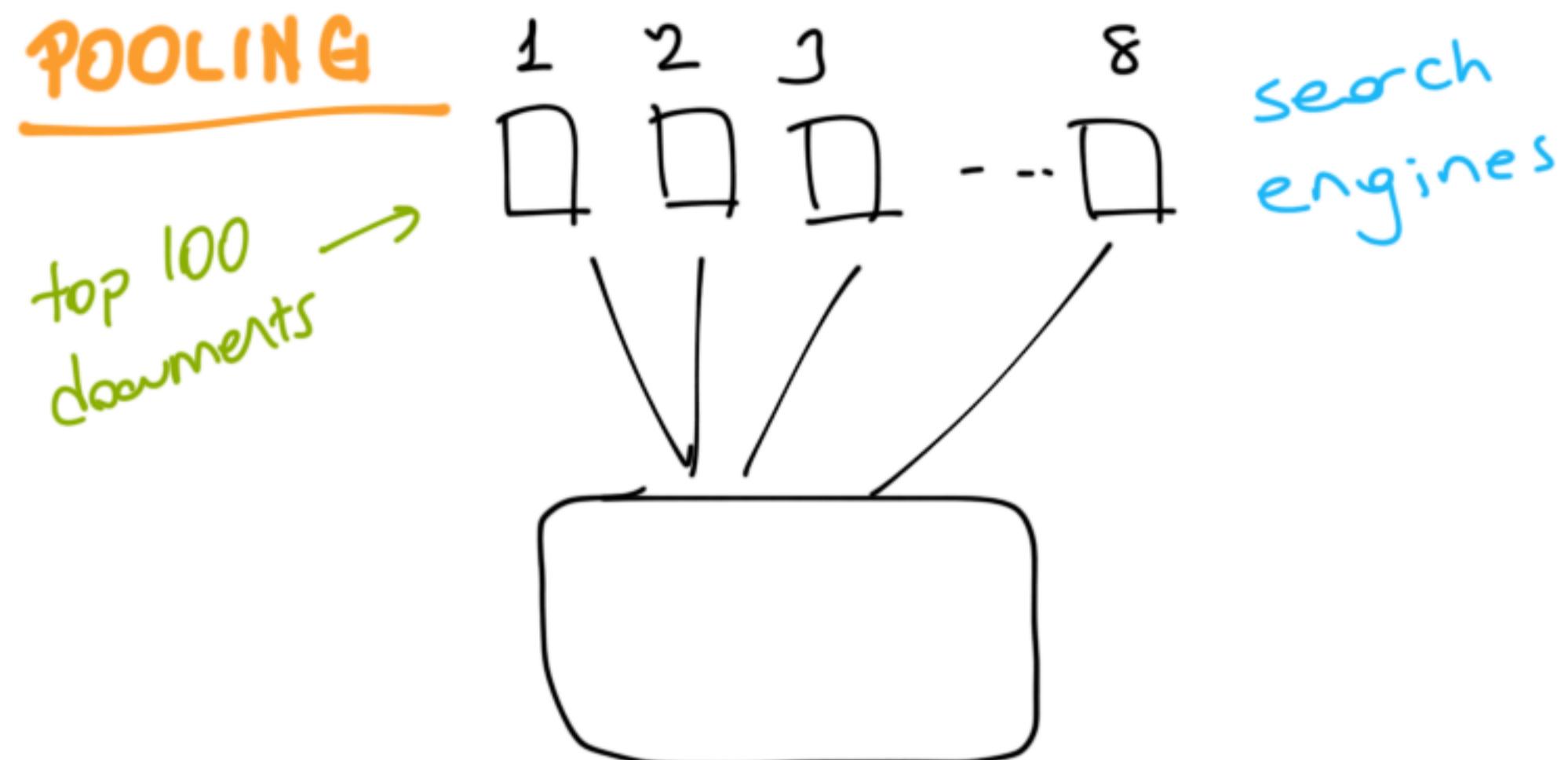
Recall =  $\frac{\text{# of retrieved and relevant docs.}}{\text{Total # of relevant docs in collection}}$

$$\frac{a}{a+c}$$

Precision =  $\frac{\text{# of retrieved and relevant docs.}}{\text{Total # of retrieved docs.}}$

$$\frac{a}{a+b}$$

	Relevant	Not relevant
Retrieved	a	b
Not Retrieved	c	d



pool

Is pooling reliable?

Problem of  
pooling

Pooling assumes that documents which are not evaluated are not relevant.

They may be relevant.

15.02.2012

Precision

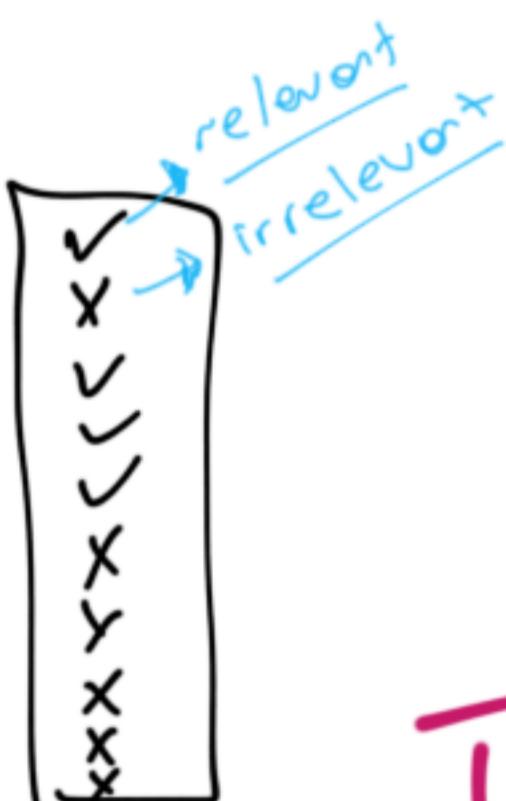
proportion of retrieved documents which are relevant.

Recall

proportion of relevant documents retrieved

example

For a given query let's assume that we have all together 20 relevant documents in the collection.



$$P@10 \rightarrow P = \frac{4}{10} = 0.4$$

$$R@20 \rightarrow R = \frac{4}{20} = 0.2$$

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

First 10 result

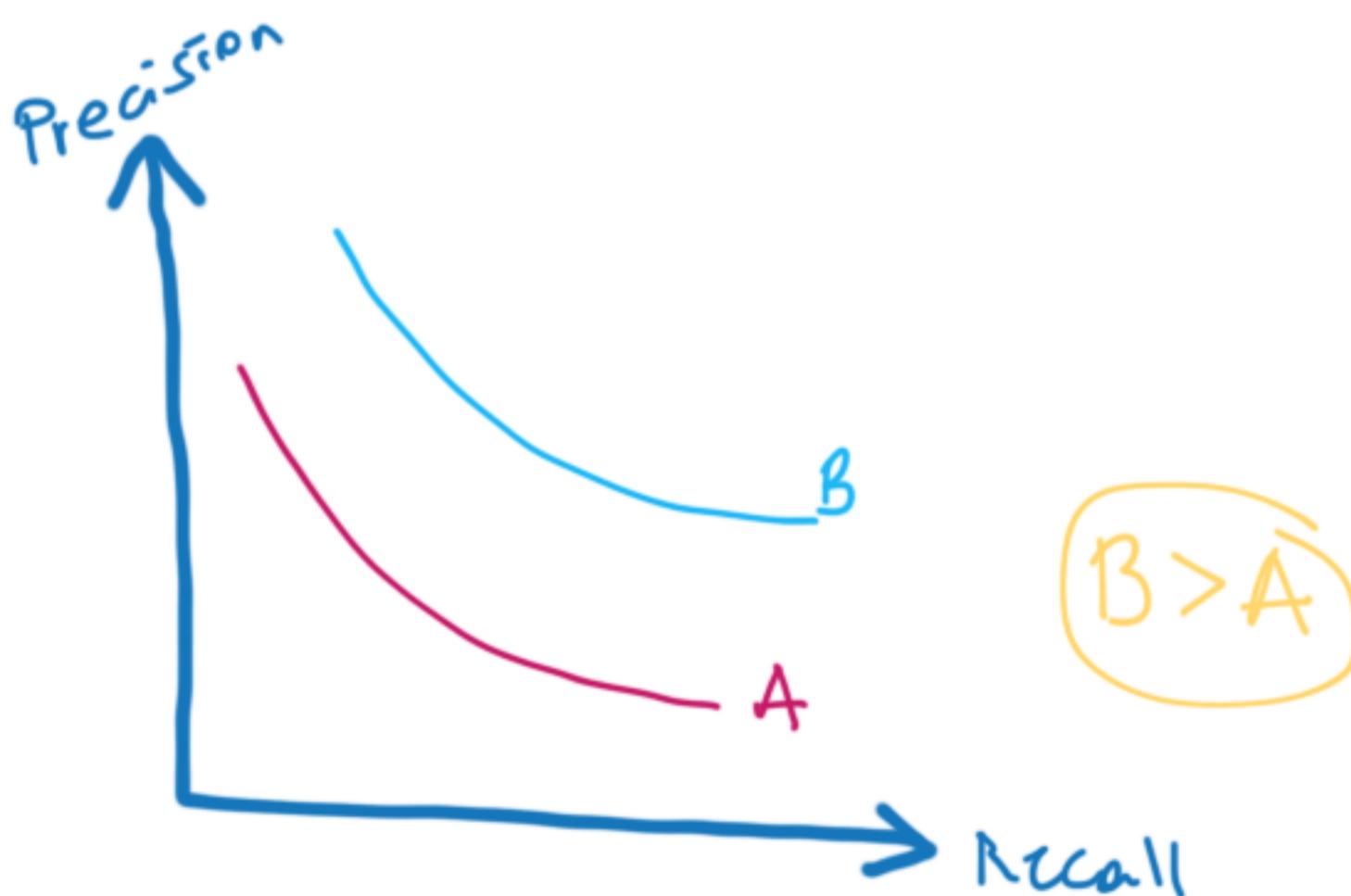
↳ harmonic mean

A document is decided as relevant or irrelevant by pooling.

### b-pref (Binary Preference)

It is an effectiveness measure which ignores documents which are not assessed by annotators.

annotators → pooling de query'nin  
gönderildiği search engine'ler.



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

Total relevance = 10,

## Similarity measures and Calculations

$Q : D \leftarrow$  query document matching  
 matching func.

$$\underline{s(d_i, d_j) = s(d_j, d_i)}$$

↳ similarity of  
document i to  
document j.

symmetric

not normalized

Inner product

Binary

Weighted

$$\sum_{i=1}^n x_i \cdot y_i$$

Dice coefficient

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

$$\frac{2 \sum_{i=1}^n x_i \cdot y_i}{\sum_{i=1}^n x_i^2 + \sum_{i=1}^n y_i^2}$$

	Binary	Weighted
Cosine coefficient	$\frac{ x \cap y }{ x ^{\frac{1}{2}}  y ^{\frac{1}{2}}}$	$\frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum x_i^2} \sqrt{\sum y_i^2}}$
Jaccard coefficient	$\frac{ x \cap y }{ x  +  y  -  x \cap y }$	$\frac{\sum x_i y_i}{\sum x_i^2 + \sum y_i^2 - \sum x_i y_i}$

example (binary)

$$x = (1 \ 0 \ 1 \ 1 \ 1) \quad |x|=4 = \sum x_i$$

$$y = (1 \ 1 \ 0 \ 1 \ 0) \quad |y|=3$$

Inner product = 2

$$\text{Dice coef.} = \frac{2 \times 2}{4 + 3}$$

$$\text{Cosine coef.} = \frac{2}{\sqrt{4} + \sqrt{3}}$$

monotonic w/  
respect to each  
other

example (weighted)

$$\mathbf{x} = (2 \ 0 \ 1 \ 3 \ 2)$$

$$\mathbf{y} = (1 \ 0 \ 2 \ 1 \ 5)$$

Dice coeff. =  $\frac{2 \times (2 \times 1 + 1 \times 2 + 3 \times 1 + 2 \times 5)}{(4 + 1 + 9 + 4) + (1 + 4 + 1 + 25)} \approx 0.69$

Cosine coeff. =  $\frac{(2 \times 1 + 1 \times 2 + 3 \times 1 + 2 \times 5)}{\sqrt{[(18) \cdot (31)]}} = 0.72$

Cluster hypothesis

documents similar to each other are relevant to the same query.

# Construction of Similarity (proximity) Matrices

	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	$t_6$
$d_1$	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

$$S = \begin{bmatrix} 1.0 & S_{12} & S_{13} & S_{14} & S_{15} \\ X & 1.0 & S_{23} & S_{24} & S_{25} \\ X & X & 1.0 & S_{34} & S_{35} \\ X & X & X & 1.0 & S_{45} \\ X & X & X & X & 1.0 \end{bmatrix}$$

$n$ : # of objects

$$(n-1) + (n-2) + \dots + 1 = \frac{n(n-1)}{2} = O(n^2)$$

Most documents don't have a common term.  
calculate similarity between docs that contain a common term.

↑ avoid unnecessary calculations

How to calculate  
similarities using  
the knowledge of  
term distribution  
in documents?

$$t_1 \rightarrow d_1, d_2$$

$$t_2 \rightarrow d_1, d_2$$

$$t_3 \rightarrow d_4, d_5$$

$$t_4 \rightarrow d_2, d_5$$

$$t_5 \rightarrow d_1, d_2$$

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



Consider  $d_1$

$t_1, t_2, t_5$

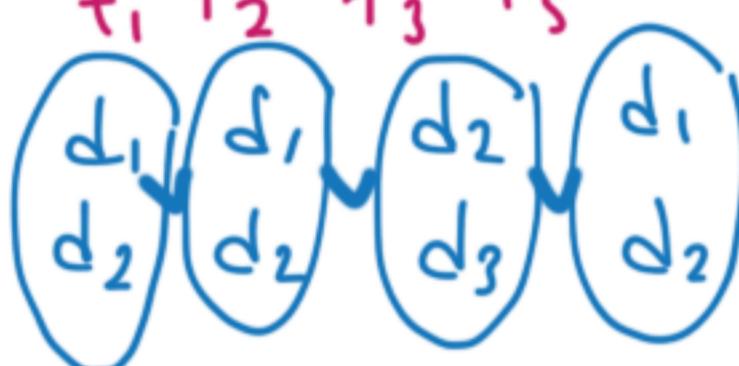


$$= (d_1, d_2)$$

just calculate  
S<sub>12</sub>

Consider  $d_2$

$t_1, t_2, t_3, t_5$



$$= (d_1, d_2, d_3)$$

just calculate  
S<sub>23</sub>

?

# How to calculate similarity btw 2 docs?

X: [5 | 7 | 8 | ...]

should be sorted!

Y: [1 | 3 | 10 | ...]

↑  
documents

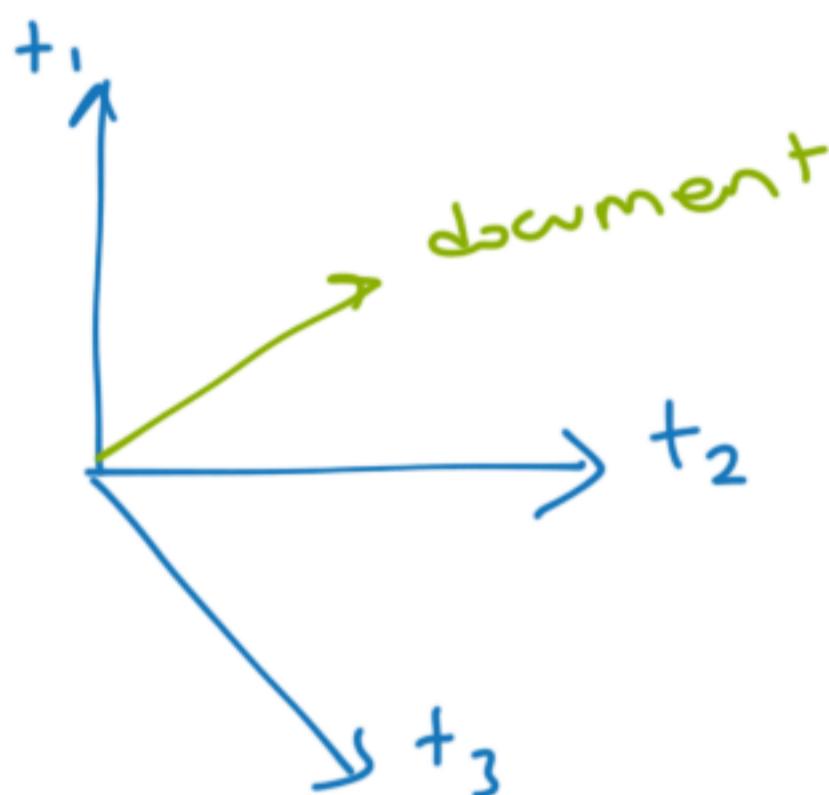
terms appear  
in that doc.

20.02.2012

---

## Similarity Matrix Calculation

Vector space



1. Straightforward approach

Compare all documents with each other.

2. Calculate similarity among documents  
that contain at least one common term.

3. Use inverted indexes



	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	$t_6$
$d_1$	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

posting  
list

- $t_1 \rightarrow <1, 1> <2, 1>$
- $t_2 \rightarrow <1, 1> <2, 1>$
- $t_3 \rightarrow <4, 1> <5, 1>$
- $t_4 \rightarrow <2, 1> <5, 1>$
- $t_5 \rightarrow <1, 1> <2, 1>$
- $t_6 \rightarrow <3, 1> <4, 1> <5, 1>$

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

Document length info.

$d_1$	$d_2$	$d_3$	$d_4$	$d_5$
3	4	1	2	3

Consider  $d_1$

contains  $t_1, t_2, t_3$

$s_{12}$	$s_{13}$	$s_{14}$	$s_{15}$
X	0	0	0

process  $t_1$

X	1	0	0	0
---	---	---	---	---

process  $t_2$

X	2	0	0	0
---	---	---	---	---

process  $t_3$

X	3	0	0	0
---	---	---	---	---

mail box



similarity of a document with the rest of the collection.

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

Consider d<sub>2</sub>

S<sub>21</sub>      S<sub>23</sub> S<sub>24</sub> S<sub>25</sub>

X	X	O	O	O
---	---	---	---	---

↳ already calculated in the prev. doc.

+<sub>1</sub>      +<sub>2</sub>      no change

+<sub>4</sub> →

X	X	O	O	I
---	---	---	---	---

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



• • •

$x_n$ : depth of indexing  
 $=$   
 avg. # of terms per document

$t_g$ : term generality  
 $=$   
 avg # of documents per term

$$\underbrace{x_d \cdot t_g}_{d_1} + \underbrace{x_d \cdot t_g}_{d_2} + \dots + \underbrace{x_d \cdot t_g}_{d_m}$$

m: # of docs  
 n: # of terms

⇒ If we store posting lists in the decreasing order of doc. numbers

≡

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

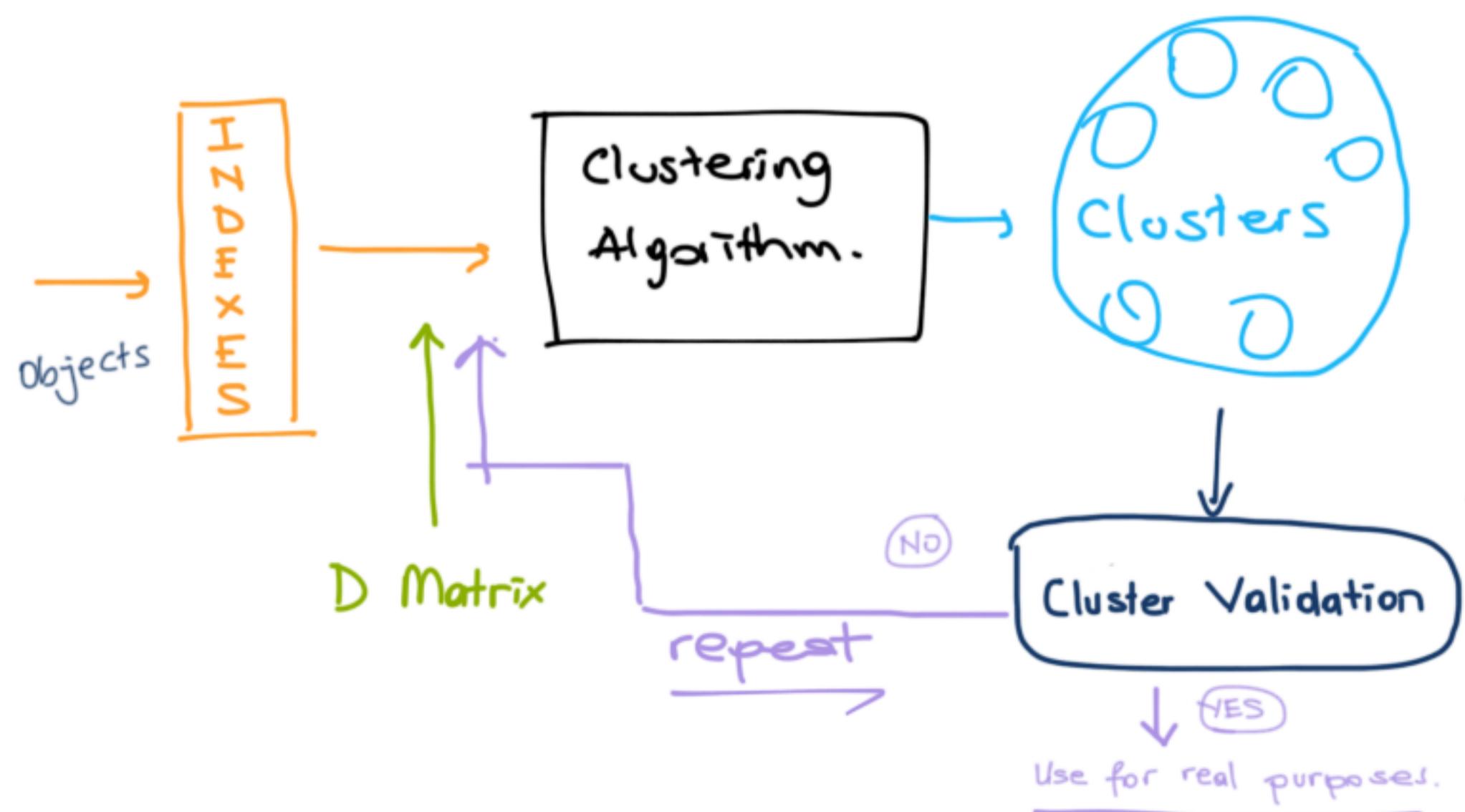
22.02.2012

---

## Document Clustering

Clustering is grouping objects into groups that contain similar objects.

↑ unsupervised  
 (without using previous knowledge)



## How to use clusters?

1. Choose the best cluster and present all of its contents to the user.
2. Cluster-based retrieval :  
Choose a number of best-matching clusters and rank their documents according to their similarity to the query.
3. Use clustering for browsing

What is our expectation?

- increase efficiency
- increase effectiveness

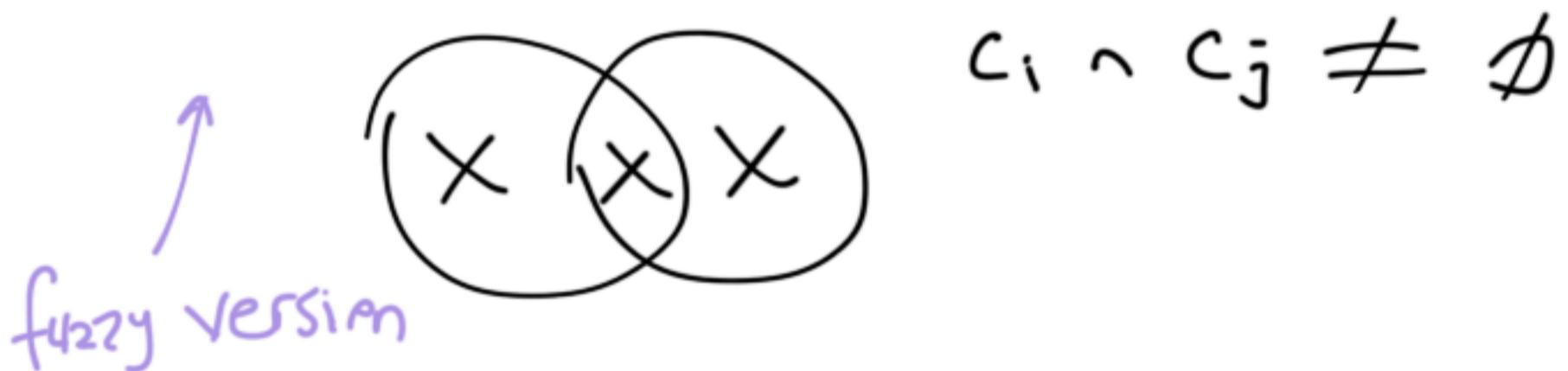
## Classification of Clustering Algorithms

A) According to structure

1. Partitioning  $C_i \cap C_j = \emptyset$   
 $i \neq j$

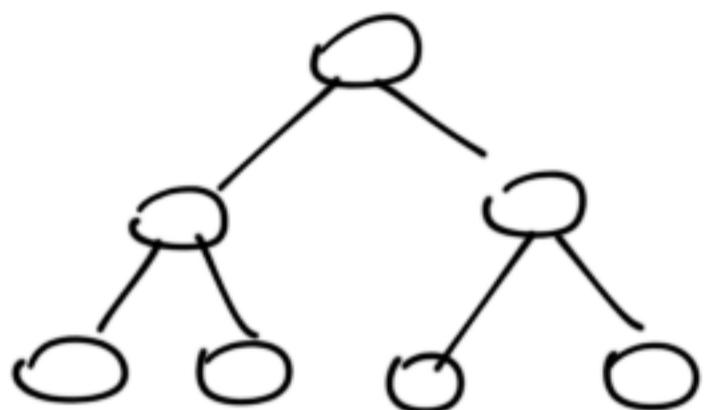
clusters do not have common members.

2. Overlapping



Objects can be a member of different clusters with a certain possibility.

### 3. Hierarchical



B. According to the working principles  
of Algorithms

1. Single pass
2. Multi pass
3. Graph theoretical  
(nodes and similarity or link among nodes)
4. Using user queries

## Single pass Algorithms

### a) seed oriented approach

Choose some objects as cluster initiators (cluster seeds)

Assign non-seed objects to cluster seeds.

How many clusters (seeds)?

$$n_c = \sqrt{m}$$

↑ # of objects  
  └ # of clusters

How to assign objects to seeds?

- choose a similarity measure
- in overlapping clusters choose a threshold based on this threshold, we may assign an object to more than one seed.

We may need to generate a rag bag cluster for mis-fits.

## How to choose seeds?

1. Choose objects that contain most # of keywords.

Clusters may not be separated from each other.

2. Seeds that contain most # of unique keywords.

Separates clusters from each other.

3. random

4. Use synthetic seeds

It brings supervision

5. Use inverted indexes

## b. Heuristic approaches

From van Rijsbergen

1. process objects one by one  
in the order they are numbered
2. the first one starts its own cluster.
3. 2nd or later objects are compared with already existing clusters and if they are different enough they start their own cluster.

### Questions

1. With what similarity value are we going to assign an object to a cluster?

2. How to define cluster representations  
(centroids)

first object

lost object

average object

order  
dependent

no problem if we are processing news articles or if we are doing new event detection.

on article which is different from all existing articles.

3. What are the characteristics of a good clustering algorithm?

1. order independence

2. cat-clustering algorithms  
=

non-parametric

effectiveness  
efficiency

3. robustness

errors in input do not affect the result.

4. stable

new additions do not significantly change the output.

5. easy to maintain

## Multi pass Algorithms

1. Generate an initial clustering structure using a quick & dirty method.
2. Try to make cluster better by iterative processing.

k-means

## Graph Theoretical Clustering

Agglomerative



obtain initial clusters  
using the objects then  
glue (join) existing clusters  
to obtain cluster of clusters.  
(super clusters)

single-link, complete-link, avg.-link

## Single-link

The similarity between a pair of clusters is taken to be the similarity between the most similar pair of objects, one of which appears in each cluster; thus each cluster member will be -- similar to at least one member that some cluster than to any members of another cluster.

$$S = \begin{bmatrix} 1.0 & 0.3 & 0.5 & 0.6 \\ - & 1.0 & 0.4 & 0.5 \\ - & - & 1.0 & 0.3 \\ - & - & - & 1.0 \end{bmatrix}$$

<u>Pair</u>	<u>Similarity</u>
AD	0.6
AC	0.5
BD	0.5
BC	0.4
AB	0.3
CD	0.3

## Complete-link

The similarity between the least similar pair of items from the two cluster is used as the cluster similarity. Each cluster member is more similar to the most dissimilar member of its own cluster than the most similar member of any other cluster.

5.3.2012

Seed based

**Method:** Calculate  $n_c$  (#of clusters)

Find  $n_c$  no. of cluster seeds

Assign non-seed documents to cluster seeds.

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

Inverse of row sums

$$\alpha_1 = 1/3, \alpha_2 = 1/4$$

$$\alpha_3 = 1/3, \alpha_4 = 1/2$$

$$\alpha_5 = 1/2$$

Columns

$$\beta_1 = 1/4, \beta_2 = 1/2$$

$$\beta_3 = 1/2, \beta_4 = 1/2$$

$$\beta_5 = 1/2, \beta_6 = 1/3$$

$$S = \begin{bmatrix} \frac{1}{3} & 0 & 0 & \frac{1}{3} & 0 & \frac{1}{3} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 & 0 \\ | & & | & & & \\ | & & | & & & \\ \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 & 0 \end{bmatrix}$$

$$S' = \begin{bmatrix} \frac{1}{4} & 0 & & \frac{1}{3} \\ \frac{1}{4} & 1 & & 0 \\ \frac{1}{4} & 0 & \dots & \frac{1}{3} \\ 0 & 0 & & \frac{1}{3} \\ \frac{1}{4} & 0 & & 0 \end{bmatrix}$$

$$S'^T = \begin{bmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 & \frac{1}{4} \\ 0 & 0 & \frac{1}{2} & 0 & \frac{1}{3} \end{bmatrix}$$

↓

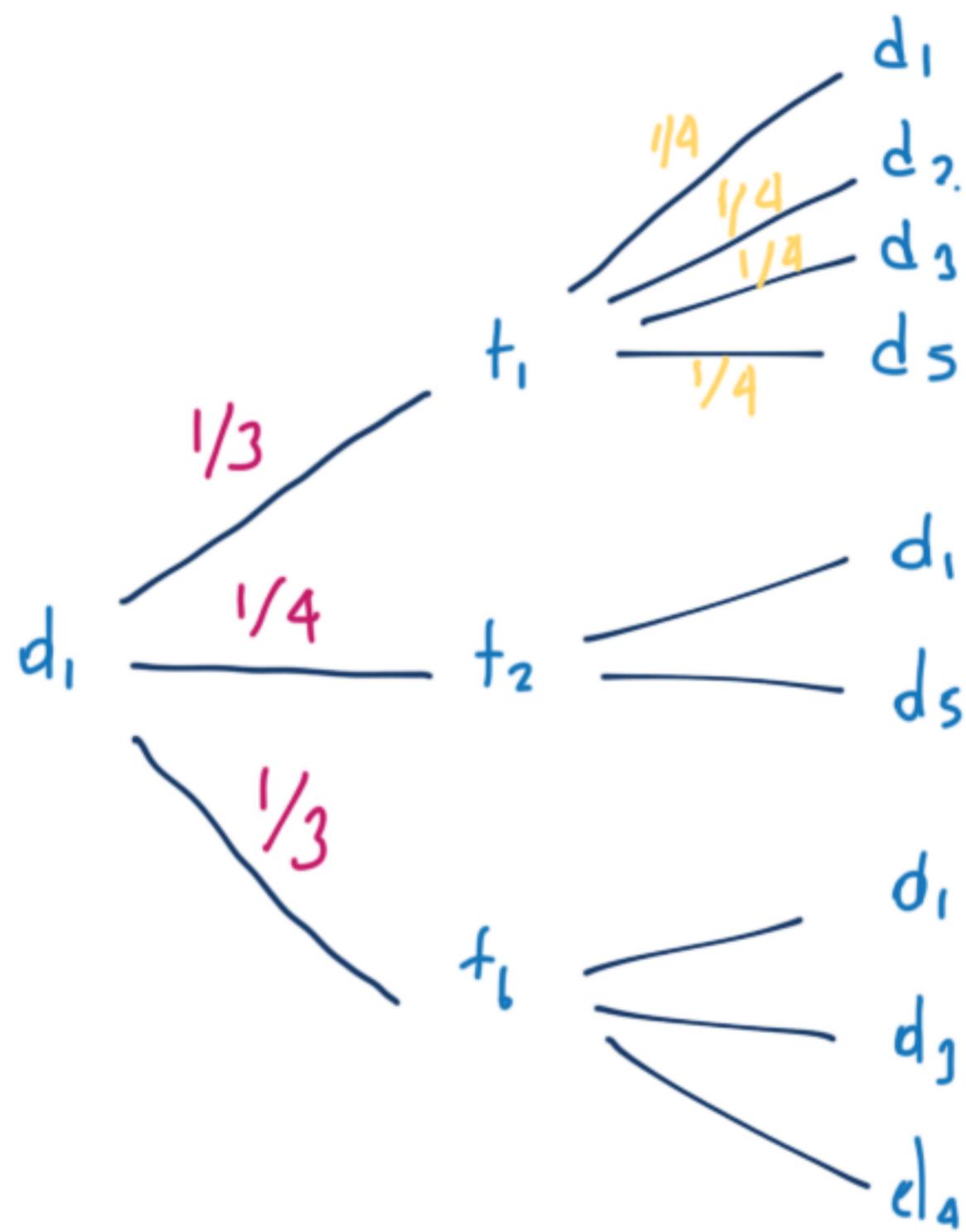
$$C = S \cdot S^T$$

*Cover Coefficient*

$C_{m \times m}$

$$\begin{bmatrix} C_{11} & \dots & C_{15} \\ C_{21} & \dots & C_{25} \\ \vdots & & \vdots \\ C_{51} & \dots & C_{55} \end{bmatrix}$$

$$C = \left[ \begin{array}{c} C_{11} = \frac{1}{3} \times \frac{1}{4} + \frac{1}{3} \times \frac{1}{2} + \frac{1}{3} \times \frac{1}{3} = \frac{13}{36} = 0.361 \\\vdots \end{array} \right]$$



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

$$C_{11} = \frac{1}{3} \times \frac{1}{4} + \frac{1}{3} \times \frac{1}{2} + \frac{1}{3} \cdot \frac{1}{3}$$

$C_{ij}$   $\Rightarrow$  probability of selecting  
any term of  $d_i$  from  
 $d_j$ .

$$C = \begin{bmatrix} 0.361 & 0.250 & 0.194 & 0.111 & 0.083 \\ \underline{0.188} & \underline{0.563} & \underline{0.063} & \underline{0} & \underline{0.188} \\ 0.194 & 0.083 & 0.361 & 0.277 & 0.083 \\ 0.167 & 0.000 & 0.417 & 0.417 & 0.000 \\ 0.125 & 0.375 & 0.125 & 0.000 & 0.375 \end{bmatrix}$$

$$0 < C_{ii} \leq 1, 0 \leq C_{ij} < 1. \quad i \neq j$$

$$C_{ii} \geq C_{ij}$$

