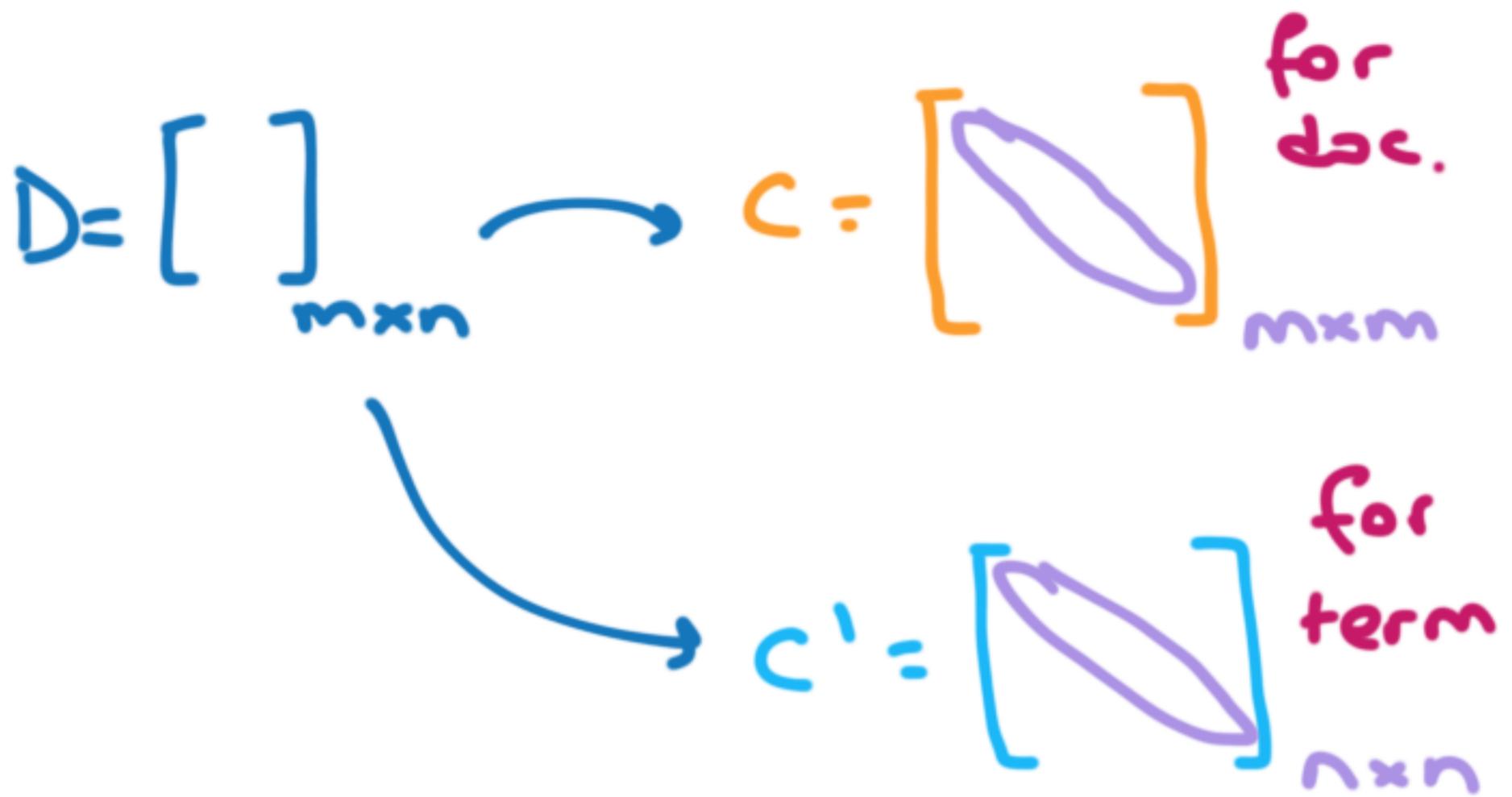


12.03.2012



$$\gamma_C = \sum_{i=1}^m c_{ii} = \sum \delta_i$$

decoupling
coef.

$$\gamma'_C = \gamma_i = \sum_{i=r}^m c_{ii}$$

D contains unique documents

documents
containing no
common terms

$$C = \begin{bmatrix} 1, 1, 0 \\ 0, 1, 1 \end{bmatrix} n_c = m$$

D contains identical documents

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

C³M

1. Find n_c
2. Select cluster seeds (initiators)
3. Assign no seeds to cluster seeds.

terms

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

documents

$$C = \begin{bmatrix} 0.361 & 0.250 & 0.194 & 0.111 & 0.083 \\ 0.188 & 0.563 & 0.063 & 0.000 & 0.188 \\ 0.194 & 0.053 & 0.361 & 0.277 & 0.083 \\ 0.167 & 0.000 & 0.417 & 0.417 & 0.083 \\ 0.125 & 0.375 & 0.125 & 0.000 & 0.175 \end{bmatrix}$$

m x m

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

of seeds

of non-seeds

$$n_c \ll m$$

calculate this
many elements
of C. matrix

C³M (revised)

1. Find $n_c = \sum \delta_i = \sum c_i \approx 2$

2. Select cluster seeds

$$P_i = \delta_i \times \Psi_i \times X_{d_i}$$

$$= C_{ii} \times (1 - C_{ii})$$

depth of indexing
for d_i .
of elem.
in d_i .

P_i = seed power of d_i

Calculate cluster seed power of each document

Choose n_c of them as cluster seeds.

3. Assign non-seeds to cluster seeds

Assign a non-seed to the seed documents which provides the max. coverage.

$\max(c_{ij})$ where $d_j \in$ cluster seed



How to identify identical seeds

$$P_i \approx P_j;$$

$$P_i = \delta_i * \psi_i * x_{d_i}$$

$$= C_{11} * (1 - C_{11}) * (\text{row sum})$$

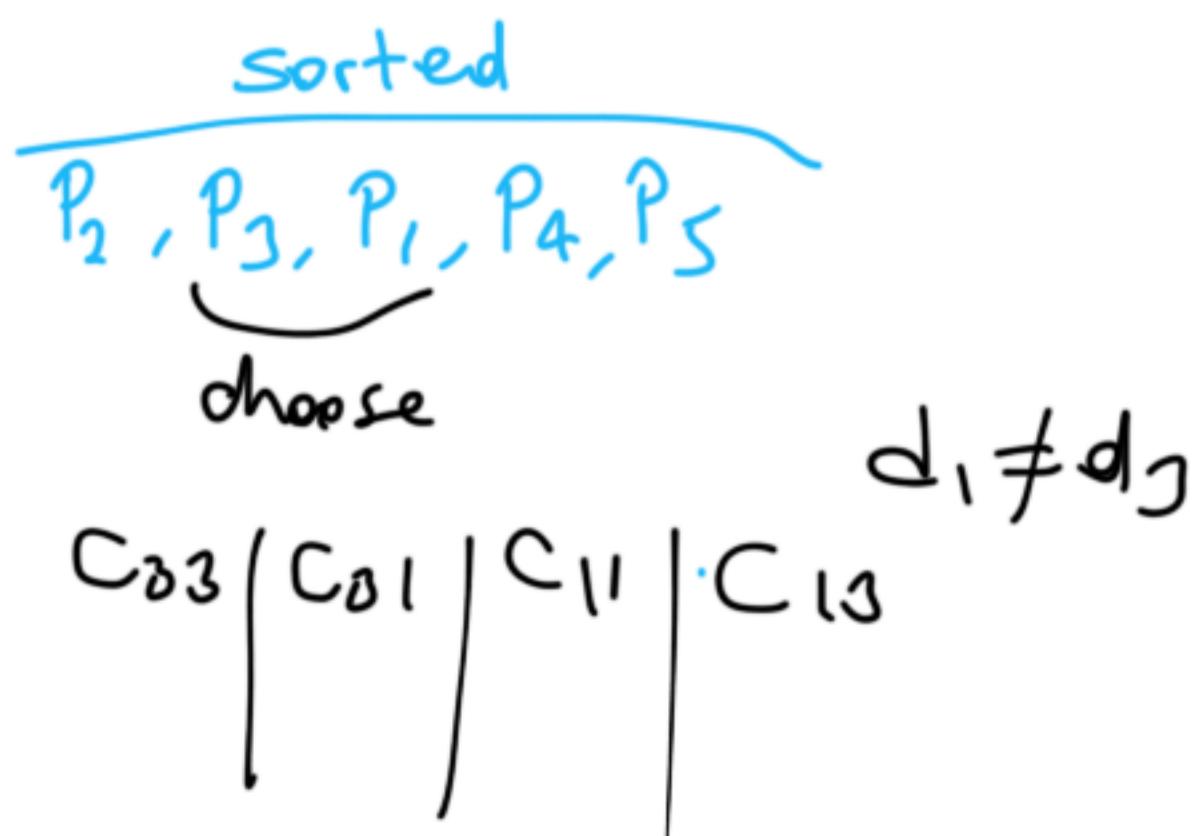
$$= 0.361 * (0.361) * 3 = \underline{\underline{0.692}}$$

$$P_2 = 0.984$$

$$P_3 = 0.692$$

$$P_4 = 0.484$$

$$P_5 = 0.469$$



d_2 is a cluster seed

\rightarrow contains t_1, t_2, t_3, t_4

$d_3 \rightarrow +_1, +_5, +_6$

$d_1 \rightarrow +_1, +_4, +_6$

it contains more
 # of terms which
 are not used in
 previously selected
 seeds. So select
 d_3 as the 2nd
 seed.

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

$d_1 \rightarrow \langle t_1, 1 \rangle \langle t_2, 2 \rangle \langle t_6, 1 \rangle$

$d_2 \rightarrow \langle t_1, 1 \rangle \langle t_2, 1 \rangle \langle t_3, 1 \rangle \langle t_4, 1 \rangle$

Inverted Index for Seed Documents

$t_1 \rightarrow \langle d_2, 1 \rangle \langle d_3, 1 \rangle$

$t_2 \rightarrow \langle d_2, 1 \rangle$

$t_3 \rightarrow \langle d_2, 1 \rangle$

$t_4 \rightarrow \langle d_2, 1 \rangle$

$t_5 \rightarrow \langle d_3, 1 \rangle$

$t_6 \rightarrow \langle d_3, 1 \rangle$

Calculate

$$C_{I2} = 0 \quad C_{I3} = 0$$

~~t₁~~

$$C'_{I2} = C_{r2} + \alpha_1 \times (d_{11} \times \beta_1 \times d_{21}) = 1/12$$

$$C'_{I3} = C_{r3} + \alpha_1 \times (d_{11} \times \beta_1 \times d_{31}) = 1/12$$

14.03.2012

$$n_c = \sum c_{ii}$$

$$n_c = \frac{m \times n}{t} \rightarrow \begin{matrix} \text{\# of non-zero} \\ \text{elements in the} \\ D\text{-matrix} \end{matrix}$$

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

$$n_c = \frac{2 \times 3}{6} = 1$$

$$n_c = \frac{m \times n}{t} = \frac{n}{x_d} \rightarrow \begin{matrix} \text{depth of} \\ \text{indexing.} \\ (\text{avg. number} \\ \text{of terms/doc}) \\ = t/m \end{matrix}$$

$$n_c = \frac{m \times n}{t} = \frac{m}{\frac{t}{n}} \rightarrow \begin{array}{l} \text{term} \\ \text{generality} \\ (\# \text{ of docs/} \\ \text{term}) \end{array}$$

$$t_g = t/n$$

$$\max(t) = m \times n$$

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

$$\max(t_g) = \max(t)/n = m \times n / n = m$$

$$\min(t_g) = \min(t)/n = \max(m, n)/n$$

$$= \max(m/n, 1)$$

$$n_c = m/t_g$$

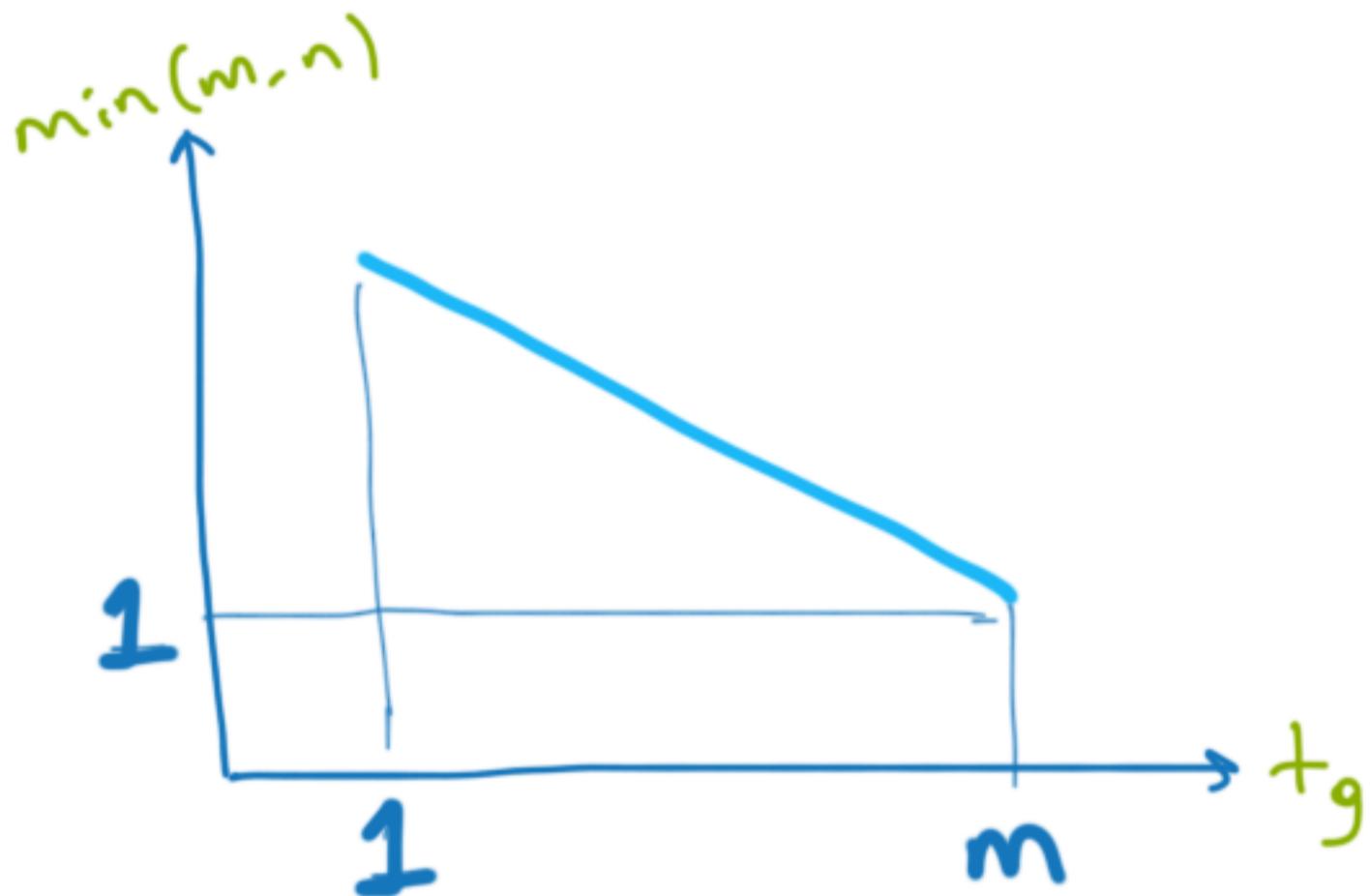
$$\max(n_c) = m/\min(t_g) = m/\max(m/n, 1)$$

$$= 1/\max(1/n, 1/m)$$

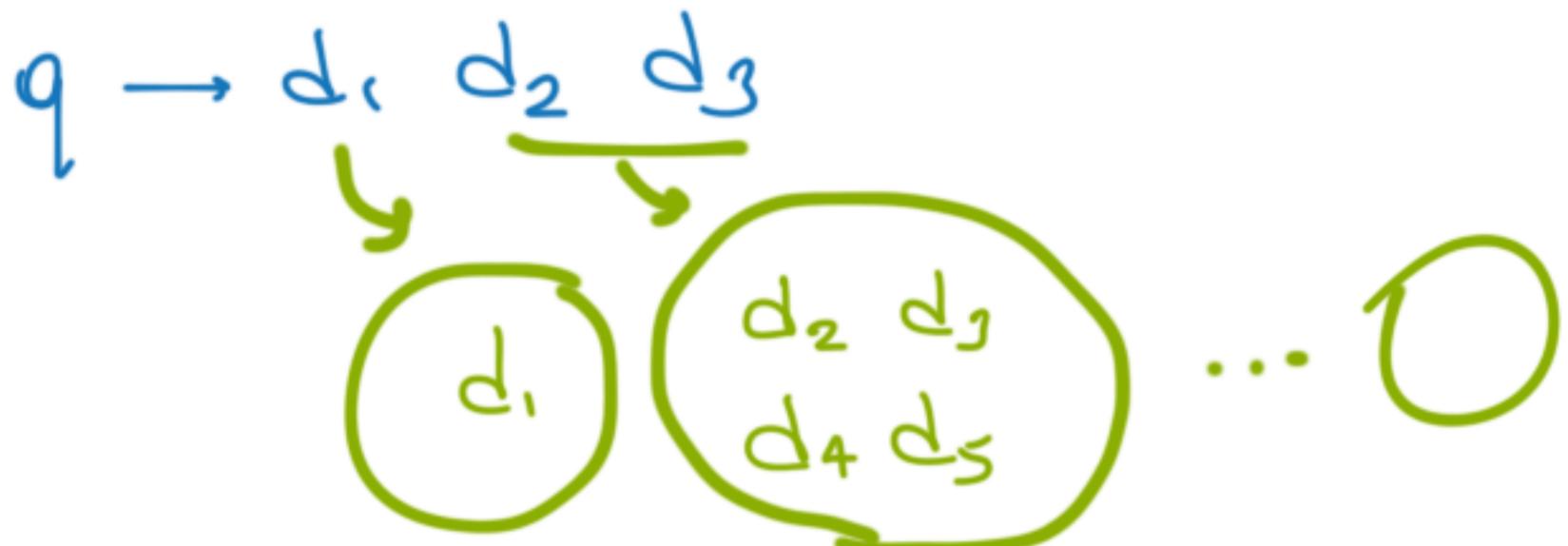
$$= \frac{1}{[\min(m, n)]^{-1}} = \min(m, n)$$

$$\min(n_c) = 1$$

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



Yao's Formula



for q we have 2 tar. clus.

A cluster that contains a relevant document is called

target cluster.

n_t = # of target clusters using
the clustering structure obtained
by our clustering algorithm.

n_{tr} = # of target clusters when
documents are randomly
distributed among the clusters.

If we have a meaningful
clustering structure $n_t < n_{tr}$



significant

n : # of records

m : # of blocks (n_c)

block size: n/m

k : # of records to be accessed

$p = n/m$ # of records in the
jth block

$n-p$: # of records in the
other blocks

C_k^n : # of combinations that we can have if
we try to select k documents out of n documents.

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

$$C_3^3 = \frac{3!}{2! (3-2)!} = 3$$

a, b, c
ab, ac, bc

C_k^{n-p} : different ways of selecting k documents from $(n-p)$ documents

The probability that no records are selected from the jth block.

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

Let $d = 1 - \frac{1}{m}$

$$n-p = n - \gamma/m = \gamma \left(1 - \frac{1}{m}\right) = \gamma \cdot d$$

$E(I_j)$: probability of selecting at least a record from the j th block

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

Expected # of blocks to be accessed

↳ $\sum_{j=1}^m E(I_j) = m \times (1 - C_k^{nd} / C_k^n)$



$$\begin{aligned}
 &= m \times \left[1 - \frac{\frac{nd!}{k!(nd-k)!}}{\frac{n!}{k!(n-k)!}} \right] = m \times \left[1 - \frac{nd!}{k!(nd-k)!} \cdot \frac{k!(n-k)!}{n!} \right] \\
 &= m \times \left[1 - \frac{\frac{1 \cdot 2 \cdot \dots \cdot nd}{1 \cdot 2 \cdot \dots \cdot (nd-k)}}{\frac{1 \cdot 2 \cdot \dots \cdot n}{1 \cdot 2 \cdot \dots \cdot (n-k)}} \right] \Rightarrow m \cdot \left[1 - \prod_{i=1}^k \frac{\frac{nd-i+1}{n-i+1}}{\frac{(nd-k+1)(nd-k+2)\dots nd}{(n-k+1)(n-k+2)\dots n}} \right]
 \end{aligned}$$

Assume that we have clusters with different sizes

cluster size $|C_j| \quad 1 \leq j \leq n_c$

$$m_j = m \cdot |C_j|$$

$$|C_j|, m_j = 1 - |C_j|$$

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

of documents

Example

$$k = 3$$

$$m = 100 \text{ (# of docs)}$$

$$|C_1| = 5$$

$$m_1 = 100 \cdot 5 = 95$$

$$P_1 = 0.14$$

19.03.2012

Cluster hypothesis

$n_t < n_{tr}$ should be significant

↳ random clustering
with clustering algorithm

m : # of docs.

$$m_j = m - |C_j| \quad 1 \leq j \leq n_c$$

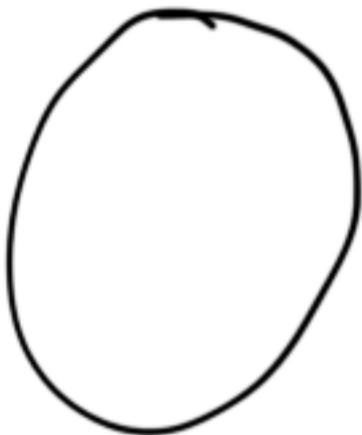
k : # of relevant docs for the query

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

$$n_{tr} = \sum_{j=1}^{n_c} P_j$$

Cluster Validation

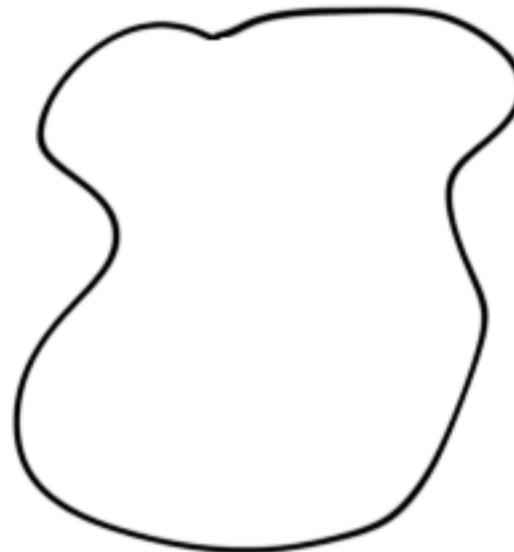
Rand Index



actual
clustering
structure

$\{a, b, c\}$

$\{d', e', f'\}$



hypothesized
clustering

↑ structure

generated by our
algorithm



$\{a, b\}$

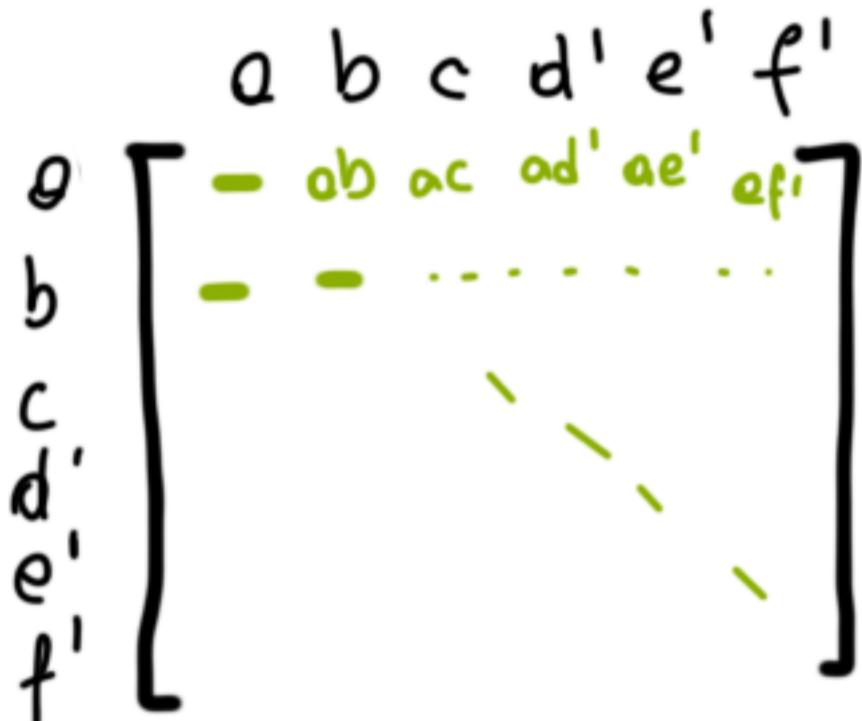


$\{c, d'\}$



$\{e', f'\}$

$$\binom{6}{2} = \frac{6!}{4! 2!} = 15$$



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

$$RI: \frac{2+8}{15} = 0.66$$

$$\text{Precision} = \frac{2}{3} : 0.66$$

$$\text{Recall} = \frac{2}{6} = 0.33$$

$$F = \frac{2 \cdot P \cdot R}{P + R} = 0.44$$

TP (true positive) two similar documents are assigned to the same h-cluster

TN (true negative) two dissimilar docs. are assigned to two different h-clusters

FP (false positive) two dissimilar docs. are assigned to the same h-cluster

FN (false negative) two similar docs. are assigned to different h-clusters

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

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

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

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

↓
f-measure

Corrected Rand Index: subtract the positive effect of random clusters

Internal Cluster

Validation Criterion

Purity: Each cluster is assigned to the class which is most frequent in the cluster, count the number of correct assignments and divide it by total # of elem.

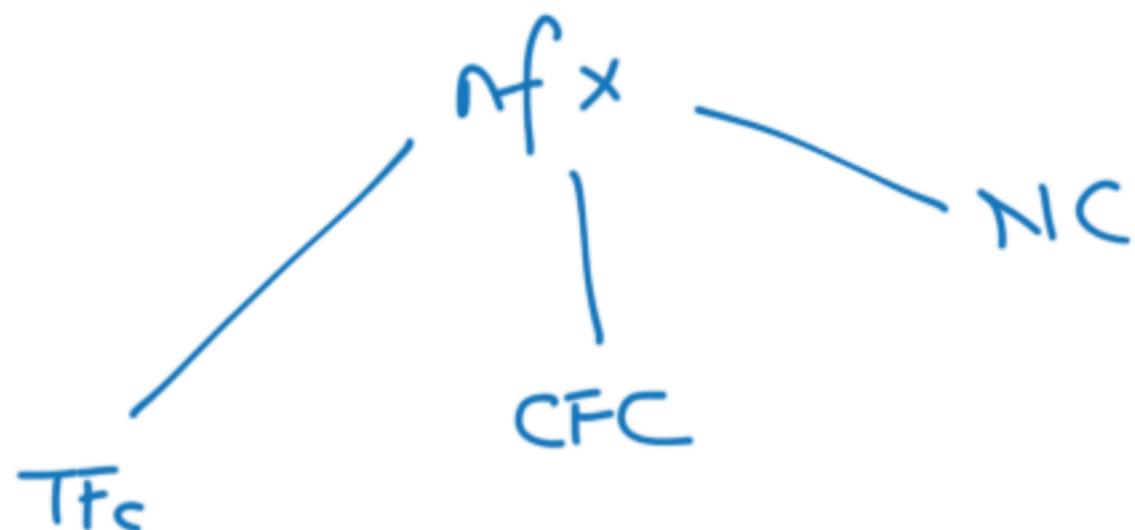
Example:



26.03.2012

Term weighting (cont.)

Query Weight Normalization



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

$$t_g \ 4 \ 2 \ 2 \ 4 \ 3$$

$$TFC = \cap \left(0.5 + 0.5 \frac{1}{2} \ 0 \ 0 \ 0.5 \cdot \frac{2}{2} \ 0 \right)$$

$$CFC = f \ \ln \frac{m}{t_g} + 1 = \left(1.22 \ 1.92 \ 1.92 \ 1.22 \ 1.51 \right)$$

nf^x:

$$\begin{aligned} & (0.75 \cdot 1.22 \ 0 \ 0 \ 1 \cdot 1.22 \ 0) \\ & (0.92 \ 0 \ 0 \ 1.22 \ 0) \end{aligned}$$

do nothing as normalization

$$D = \begin{bmatrix} 0.62 & 0 & 0.49 & 0.62 & 0 \end{bmatrix}$$

$$\text{Sim}(Q, d_1) = 0.62 \times 0.92 + 0 \times 0 + 0.62 \times 1.22 + 0 \\ = 1.33_{//}$$

$$\text{Sim}(Q, d_2) = 0.77$$

$$d_3 = 1.20$$

$$d_4 = 0.58$$

$$d_5 = 1.12$$

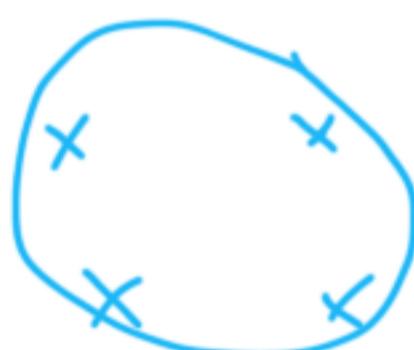
ranking

$$d_1 > d_3 > d_5 > d_2 > d_4$$

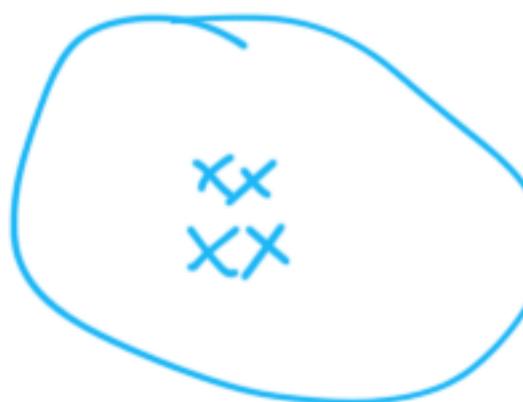
Term Discrimination Value

This concept assumes that terms which make documents more distinguishable from each other should have higher importance.

Example :



After assigning a
good discriminator



After assigning a
bad term (a term
which is not a
good indicator)
documents become more similar
to each other

How to calculate TDVs

1. Using similarity values
2. " cover coefficient concept"