

# Text Retrieval

Sergej Sizov

Information Retrieval

Summer term 2009

# The document ranking problem

---

We have a collection of documents

User issues a query

A list of documents needs to be returned

## **Ranking method is core of an IR system:**

- ◆ In what order do we present documents to the user?
- ◆ We want the “best” document to be first, second best second, etc....

## **Idea: Rank by probability of relevance of the document w.r.t. information need**

- ◆  $P(\text{relevant}|\text{document})$

# Probabilistic IR: Principles

---

## Goal:

Ranking based on  $\text{sim}(\text{doc } d, \text{query } q) =$

$$P[R|d] = P [\text{doc } d \text{ is relevant for query } q \mid d \text{ has term vector } X_1, \dots, X_m ]$$

## Assumptions:

Document relevance does not depend of other documents

Relevant and irrelevant documents differ in their terms.

Binary Independence Retrieval (BIR) Model:

- ◆ Probabilities for term occurrence are pairwise
- ◆ independent for different terms.
- ◆ Term weights are binary:  $X_i \in \{0,1\}$ .

For terms that do not occur in query  $q$  the probabilities for such a term occurring are the same for relevant and irrelevant documents.

# Probabilistic IR with Term Independence

## Ranking Proportional to Relevance Odds

$$\text{sim}(d, q) = O(R | d) = \frac{P[R | d]}{P[\neg R | d]}$$

odds for relevance  
(ratio of relevant documents)

$$= \frac{P[d | R] \times P[R]}{P[d | \neg R] \times P[\neg R]}$$

Bayes' theorem

$$\sim \frac{P[d | R]}{P[d | \neg R]} = \prod_i \frac{P[X_i | R]}{P[X_i | \neg R]}$$

independence or  
linked dependence

$$\text{sim}(d, q)' = \log \prod_{i \in q} \frac{P[X_i | R]}{P[X_i | \neg R]}$$

$X_i = 1$  if  $d$  includes  
 $i$ -th term, 0 otherwise

$$= \sum_{i \in q} \log P[X_i | R] - \log P[X_i | \neg R]$$

# Probabilistic Retrieval

$$\text{sim}(d, q)' = \log \prod_{i \in q} \frac{P[X_i | R]}{P[X_i | \neg R]} = \sum_{i \in q} \log P[X_i | R] - \log P[X_i | \neg R]$$

$$= \sum_{i \in q} \log ( p_i^{X_i} (1 - p_i)^{1 - X_i} ) - \log ( q_i^{X_i} (1 - q_i)^{1 - X_i} )$$

with estimators  $p_i = P[X_i = 1 | R]$  and  $q_i = P[X_i = 1 | \neg R]$  (binary features)

$$= \sum_{i \in q} \log \left( \frac{p_i^{X_i} (1 - p_i)}{(1 - p_i)^{X_i}} \right) - \log \left( \frac{q_i^{X_i} (1 - q_i)}{(1 - q_i)^{X_i}} \right)$$

$$= \sum_{i \in q} X_i \log \frac{p_i}{1 - p_i} + \sum_{i \in q} X_i \log \frac{1 - q_i}{q_i} + \sum_{i \in q} \log \frac{1 - p_i}{1 - q_i}$$

$$\sim \sum_{i \in q} X_i \log \frac{p_i}{1 - p_i} + \sum_{i \in q} X_i \log \frac{1 - q_i}{q_i} = \text{sim}(d, q)''$$

# Probabilistic Retrieval

Robertson / Sparck Jones Formula

Estimate  $p_i$  and  $q_i$  based on training sample (query  $q$  on small sample of corpus) or based on intellectual assessment of first round's result (**relevance feedback**):

Let  $N$  be #docs in sample,  
 $R$  be # relevant docs in sample  
 $n_i$  #docs in sample that contain term  $i$ ,  
 $r_i$  # relevant docs in sample that contain term  $i$ .

$$\begin{aligned} \Rightarrow \text{Estimate: } p_i &= \frac{r_i}{R} & q_i &= \frac{n_i - r_i}{N - R} \\ \text{or: } p_i &= \frac{r_i + 0.5}{R + 1} & q_i &= \frac{n_i - r_i + 0.5}{N - R + 1} & \text{(Lidstone smoothing with } \lambda=0.5) \\ \Rightarrow \text{sim}(d, q)'' &= \sum_i X_i \log \frac{r_i + 0.5}{R - r_i + 0.5} + \sum_i X_i \log \frac{N - n_i - R + r_i + 0.5}{n_i - r_i + 0.5} \\ \Rightarrow \text{Weight of term } i \text{ in doc } d: & \log \frac{(r_i + 0.5)(N - n_i - R + r_i + 0.5)}{(R - r_i + 0.5)(n_i - r_i + 0.5)} \end{aligned}$$

# Probabilistic Retrieval: tf\*idf Formula

Assumptions (without training sample or relevance feedback):

- $p_i$  is the same for all  $i$ .
- Most documents are irrelevant.
- Each individual term  $i$  is infrequent.

This implies:

- $\sum_i X_i \log \frac{p_i}{1-p_i} = c \sum_i X_i$  with constant  $c$
- $q_i = P[X_i = 1 | \neg R] \approx \frac{df_i}{N}$
- $\frac{1-q_i}{q_i} = \frac{N-df_i}{df_i} \approx \frac{N}{df_i}$

$$\Rightarrow \text{sim}(d, q)'' = \sum_i X_i \log \frac{p_i}{1-p_i} + \sum_i X_i \log \frac{1-q_i}{q_i}$$

$$\approx c \sum_i X_i + \sum_i X_i \text{idf}_i$$

Scalar product over the product of tf and dampend idf values for query terms

# Probabilistic Retrieval: Example

Documents with relevance feedback:

	t1	t2	t3	t4	t5	t6	R
d1	1	0	1	1	0	0	1
d2	1	1	0	1	1	0	1
d3	0	0	0	1	1	0	0
d4	0	0	1	0	0	0	0
ni	2	1	2	3	2	0	
ri	2	1	1	2	1	0	
pi	5/6	1/2	1/2	5/6	1/2	1/6	
qi	1/6	1/6	1/2	1/2	1/2	1/6	

q: t1 t2 t3 t4 t5 t6

} R=2, N=4

Score of new document d5 (with Lidstone smoothing with  $\lambda=0.5$ ):

d5 $\cap$ q: <1 1 0 0 0 1>  $\rightarrow$  sim(d5, q) =  $\log 5 + \log 1 + \log 0.2$   
 $\log 5 + \log 5 + \log 5$

$$sim(d, q)'' = \sum_i X_i \log \frac{p_i}{1-p_i} + \sum_i X_i \log \frac{1-q_i}{q_i}$$



## BIR: Summery

---

- ◆ The most important classic prob. IR model
- ◆ Use only term presence/absence, thus also referred to as Binary Independence Model
- ◆ Most natural for relevance/pseudo feedback
- ◆ When without relevance judgments, the model parameters must be estimated in an ad hoc way
- ◆ Performance isn't as good as tuned VS model
- ◆ Extensions
  - ◆ Term correlations
  - ◆ Okapi
  - ◆ Smoothing

# Extensions of Probabilistic IR

Consider term correlations in documents (with binary  $X_i$ )

→ Problem of estimating m-dimensional prob.distribution

$$P[X_1=\dots \wedge X_2= \dots \wedge \dots \wedge X_m=\dots] =: f_X(X_1, \dots, X_m)$$

One possible approach: **Tree Dependence Model:**

a) Consider only 2-dimensional probabilities (for term pairs)

$$f_{ij}(X_i, X_j) = P[X_i=\dots \wedge X_j=\dots] = \sum_{X_1} \dots \sum_{X_{i-1}} \sum_{X_{i+1}} \dots \sum_{X_{j-1}} \sum_{X_{j+1}} \dots \sum_{X_m} P[X_1 = \dots \wedge \dots \wedge X_m = \dots]$$

b) For each term pair

estimate the error between independence and the actual correlation

c) Construct a tree with terms as nodes and the

m-1 highest error (or correlation) values as weighted edges

# Considering 2-dimensional Term Correlation

## Variant 1:

Error of approximating  $f$  by  $g$  (**Kullback-Leibler divergence**) with  $g$  assuming pairwise term independence:

$$\varepsilon(f, g) := \sum_{\vec{X} \in \{0,1\}^m} f(\vec{X}) \log \frac{f(\vec{X})}{g(\vec{X})} = \sum_{\vec{X} \in \{0,1\}^m} f(\vec{X}) \log \frac{f(\vec{X})}{\prod_{i=1}^m g_i(X_i)}$$

## Variant 2:

**Correlation coefficient** for term pairs:

$$\rho(X_i, X_j) := \frac{\text{Cov}(X_i, X_j)}{\sqrt{\text{Var}(X_i)} \sqrt{\text{Var}(X_j)}}$$

## Variant 3:

level- $\alpha$  values or p-values of **Chi-square independence test**

# Constructing the Term Dependence Tree

## Given:

complete graph  $(V, E)$  with  $m$  nodes  $X_i \in V$  and  $m^2$  undirected edges  $\in E$  with weights  $\varepsilon$  (or  $\rho$ )

## Wanted:

spanning tree  $(V, E')$  with maximal sum of weights

## Algorithm:

Sort the  $m^2$  edges of  $E$  in descending order of weight

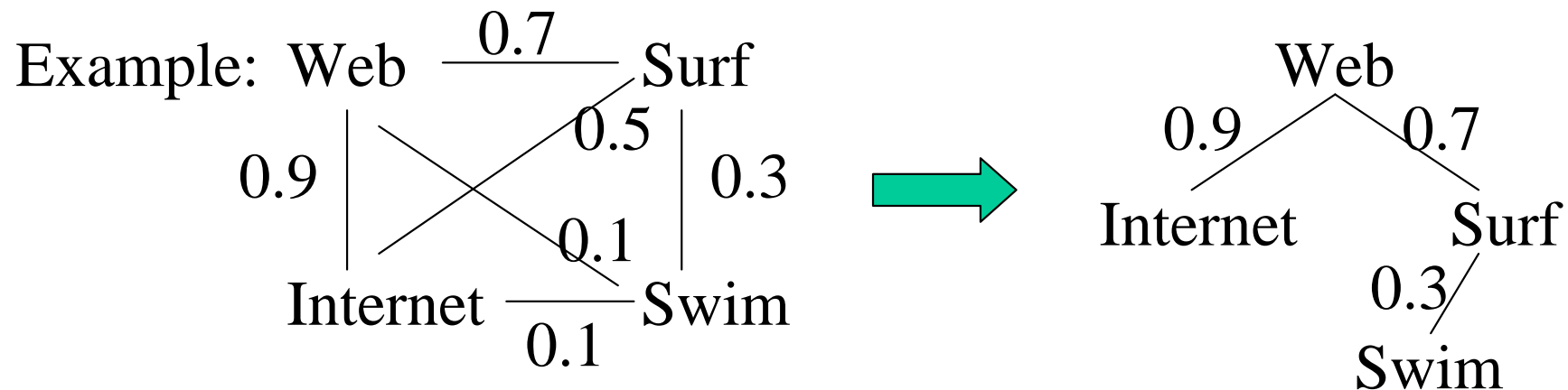
$E' := \emptyset$

Repeat until  $|E'| = m-1$

$E' := E' \cup \{(i,j) \in E \mid (i,j) \text{ has max. weight in } E\}$

provided that  $E'$  remains acyclic;

$E := E - \{(i,j) \in E \mid (i,j) \text{ has max. weight in } E\}$



# Estimation of Multidimensional Probabilities

with Term Dependence Tree

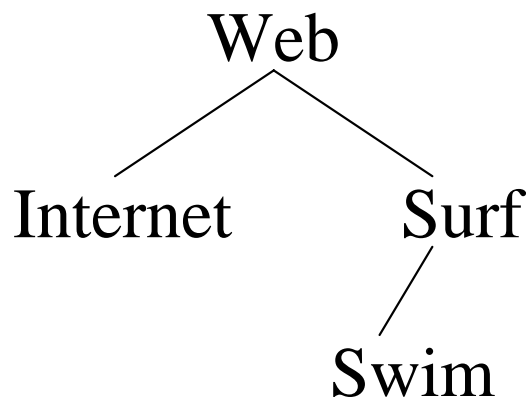
Given is a term dependence tree  $(V = \{X_1, \dots, X_m\}, E')$ .

Let  $X_1$  be the root, nodes are preorder-numbered, and assume that

$X_i$  and  $X_j$  are independent for  $(i, j) \notin E'$ . Then:

$$\begin{aligned} P[X_1 = \dots \wedge \dots \wedge X_m = \dots] &= P[X_1 = \dots] P[X_2 = \dots \wedge X_m = \dots | X_1 = \dots] \\ &= \prod_{i=1..m} P[X_i = \dots | X_1 = \dots \wedge X_{(i-1)} = \dots] \\ &= P[X_1] \cdot \prod_{(i,j) \in E'} P[X_j | X_i] \\ &= P[X_1] \cdot \prod_{(i,j) \in E'} \frac{P[X_i, X_j]}{P[X_i]} \end{aligned}$$

Example:



$P[\text{Web}, \text{Internet}, \text{Surf}, \text{Swim}] =$

$$P[\text{Web}] \frac{P[\text{Web}, \text{Internet}]}{P[\text{Web}]} \frac{P[\text{Web}, \text{Surf}]}{P[\text{Web}]} \frac{P[\text{Surf}, \text{Swim}]}{P[\text{Surf}]}$$

# Latent Semantic Analysis

---

- Lexical matching at term level inaccurate (claimed)
- Polysemy – words with number of ‘meanings’ – term matching returns irrelevant documents – impacts precision
- Synonymy – number of words with same ‘meaning’ – term matching misses relevant documents – impacts recall
- Fewer dimensions → dimension reduction
- Keep  $k$  strongest dimensions: remove noise

LSA assumes that there exists a LATENT structure in word usage – obscured by variability in word choice

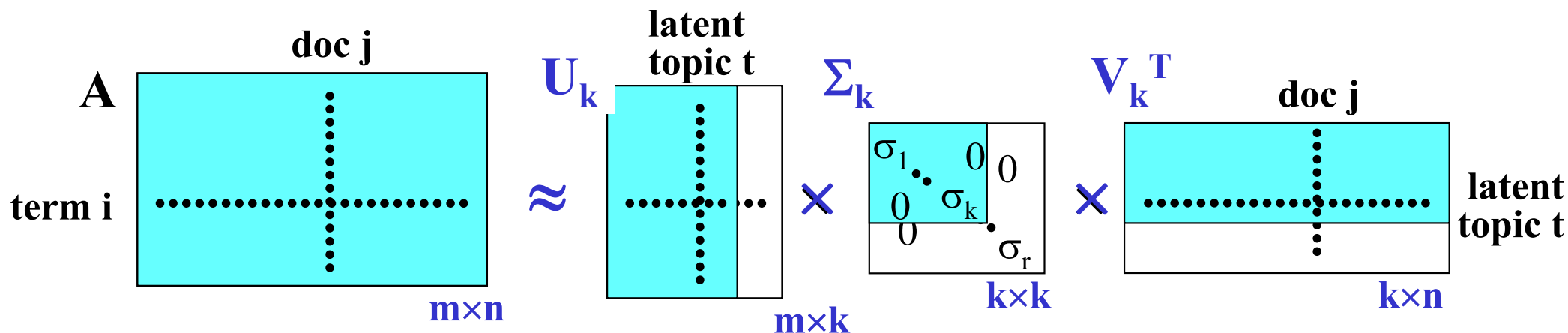
Word usage defined by term and document co-occurrence – matrix structure

# Latent Semantic Indexing (LSI)

[Deerwester et al. 1990]

A is the  $m \times n$  term-document matrix. Then:

- U and  $U_k$  are the  $m \times r$  and  $m \times k$  term-topic similarity matrices,
- V and  $V_k$  are the  $n \times r$  and  $n \times k$  document-topic similarity matrices,
- $A \times A^T$  and  $A_k \times A_k^T$  are the term-term similarity matrices,
- $A^T \times A$  and  $A_k^T \times A_k$  are the document-document similarity matrices



mapping of  $m \times 1$  vectors into latent-topic space:

$$d_j : \triangleright U_k^T \times d_j =: d_j'$$

$$q : \triangleright U_k^T \times q =: q'$$

scalar-product similarity in latent-topic space:

$$d_j'^T \times q' = ((\Delta_k V_k^T)_{*j})^T \times q'$$

# LSI: Indexing and Query Processing

- The matrix  $\Delta_k \mathbf{V}_k^T$  corresponds to a „**topic index**“ and is stored in a suitable data structure.  
Instead of  $\Delta_k \mathbf{V}_k^T$  the simpler **index**  $\mathbf{V}_k^T$  could be used.
- Additionally the **term-topic mapping**  $\mathbf{U}_k$  must be stored.
- A **query**  $\mathbf{q}$  (an  $m \times 1$  column vector) in the term vector space is transformed into query  $\mathbf{q}' = \mathbf{U}_k^T \times \mathbf{q}$  (a  $k \times 1$  column vector) and evaluated in the topic vector space (i.e.  $\mathbf{V}_k$ ) (e.g. by scalar-product similarity  $\mathbf{V}_k^T \times \mathbf{q}'$  or cosine similarity)
- A **new document**  $\mathbf{d}$  (an  $m \times 1$  column vector) is transformed into  $\mathbf{d}' = \mathbf{U}_k^T \times \mathbf{d}$  (a  $k \times 1$  column vector) and appended to the „index“  $\mathbf{V}_k^T$  as an additional column („**folding-in**“)



## LSI: Example

m=6 terms

t1: bak(e,ing)

t2: recipe(s)

t3: bread

t4: cake

t5: pastr(y,ies)

t6: pie

n=5 documents

d1: How to bake bread without bread recipes

d2: The classic art of Viennese Pastry

d3: Numerical recipes: the art of  
scientific computing

d4: Breads, pastries, pies and cakes:  
quantity baking recipes

d5: Pastry: a book of best French recipes

$$A = \begin{pmatrix} 0.4446 & 0.0000 & 0.0000 & 0.3422 & 0.0000 \\ 0.1083 & 0.0000 & 1.0000 & 0.0833 & 0.4002 \\ 0.8892 & 0.0000 & 0.0000 & 0.3422 & 0.0000 \\ 0.0000 & 0.0000 & 0.0000 & 0.6010 & 0.0000 \\ 0.0000 & 1.0000 & 0.0000 & 0.1908 & 0.9164 \\ 0.0000 & 0.0000 & 0.0000 & 0.6010 & 0.0000 \end{pmatrix}$$

# LSI: Example

$$A = \begin{pmatrix} -0.1337 & 0.4385 & -0.0916 & -0.0858 \\ -0.4039 & 0.0798 & 0.9089 & 0.06680 \\ -0.1909 & 0.7045 & -0.1092 & -0.5105 \\ -0.1336 & 0.3000 & -0.1298 & 0.5997 \\ -0.8642 & -0.3535 & -0.3464 & -0.0906 \\ -0.1336 & 0.3000 & -0.1298 & 0.5997 \end{pmatrix} \quad U$$

$$\times \begin{pmatrix} 1.4543 & 0.0000 & 0.0000 & 0.0000 \\ 0.0000 & 1.1764 & 0.0000 & 0.0000 \\ 0.0000 & 0.0000 & 0.9980 & 0.0000 \\ 0.0000 & 0.0000 & 0.0000 & 0.7115 \end{pmatrix} \quad \Delta$$

$$\times \begin{pmatrix} -0.1878 & -0.5942 & -0.2777 & -0.3233 & -0.65571 \\ 0.7061 & -0.3005 & 0.0678 & 0.5873 & -0.2482 \\ -0.0396 & -0.3471 & 0.9107 & -0.2155 & 0.0464 \\ -0.6817 & -0.1274 & 0.0940 & 0.7100 & -0.0792 \end{pmatrix} \quad V^T$$

## LSI: Example

---

$$A_2 = \begin{pmatrix} 0.4008 & -0.0395 & 0.0890 & 0.3659 & -0.0006 \\ 0.1766 & 0.3209 & 0.1695 & 0.2451 & 0.3619 \\ 0.6373 & -0.0841 & 0.1333 & 0.5766 & -0.0237 \\ 0.2857 & 0.0094 & 0.0779 & 0.2701 & 0.0398 \\ -0.0576 & 0.8718 & 0.3209 & 0.1621 & 0.9273 \\ 0.2857 & 0.0094 & 0.0779 & 0.2701 & 0.0398 \end{pmatrix} = U_2 \times \Delta_2 \times V_2^T$$

## LSI: Example

query q: baking bread

$$q = ( 1 \ 0 \ 1 \ 0 \ 0 \ 0 )^T$$

transformation into topic space with k=2

$$q' = U_k^T \times q = (-0.3246 \ 1.1430)^T$$

scalar product similarity in topic space with k=2:

$$\text{sim}(q, d1) = V_{k*1}^T \times q' \approx 0.87 \quad \text{sim}(q, d2) = V_{k*2}^T \times q' \approx -0.15$$

$$\text{sim}(q, d3) = V_{k*3}^T \times q' \approx 0.17 \quad \text{etc.}$$

Folding-in of a new document d6:

algorithmic recipes for the computation of pie

$$d6 = ( 0 \ 0.7071 \ 0 \ 0 \ 0 \ 0.7071 )^T$$

transformation into topic space with k=2

$$d6' = U_k^T \times d6 \approx (-0.3801 \ 0.2686)$$

d6' is appended to  $V_k^T$  as a new column

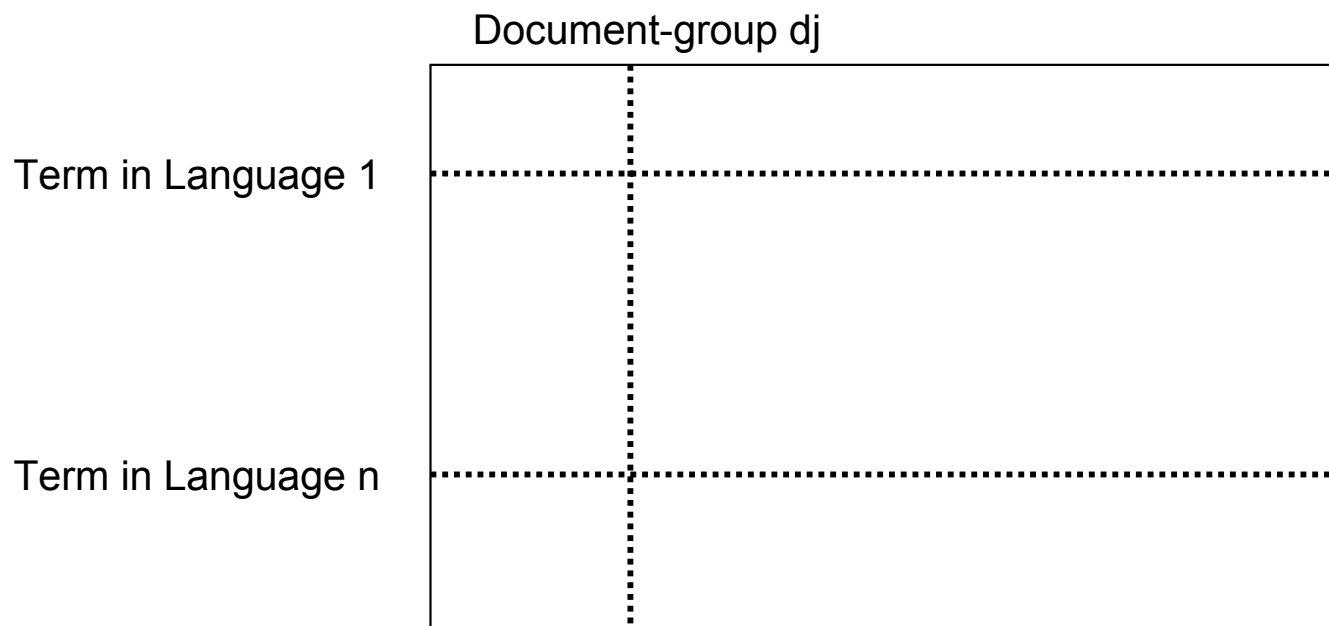
# Cross-Language LSI

---

- ◆ Build common feature-articles space
  - ◆ Don't differ between languages
  - ◆ One Document consists of several languages
- ◆ Decompose document-term-matrix by reduced SVD to build concept-document-space
- ◆ Map new documents to concept-document-space
- ◆ Mate retrieval – looking for corresponding document in a second language

# Cross-Language LSI

- Construct LSI model ( $U_k, \Delta_k, V_k^T$ ) from training documents that are available in multiple languages:
  - Consider all language variants of the same document as a single document and
  - Extract all terms or words for all languages



- SVD with regression brings terms in different language together
- Folding in documents in one language

# Datasets

---

## Wikipedia

- ◆ Online encyclopedia
- ◆ Using a subset of 160 000 articles which are present in English, German and French connected by language links

## Multext Corpus

- ◆ The Multext Corpus is an EU-Textcorpus
- ◆ Question and Answers of the European Commission in English, German, French, Spanish and Italian
- ◆ About 2780 Pairs of Question and Answer

## JRC-Arquis

- ◆ EU-Textcorpus
- ◆ Documents from beginning of the EU until today

# Test on Multext

---

Multext - Multext

60 % for training ca. 1670 documents

40 % for testing ca. 1113 documents

2 runs with different training and test sets on  
Question and Answer (as one document)

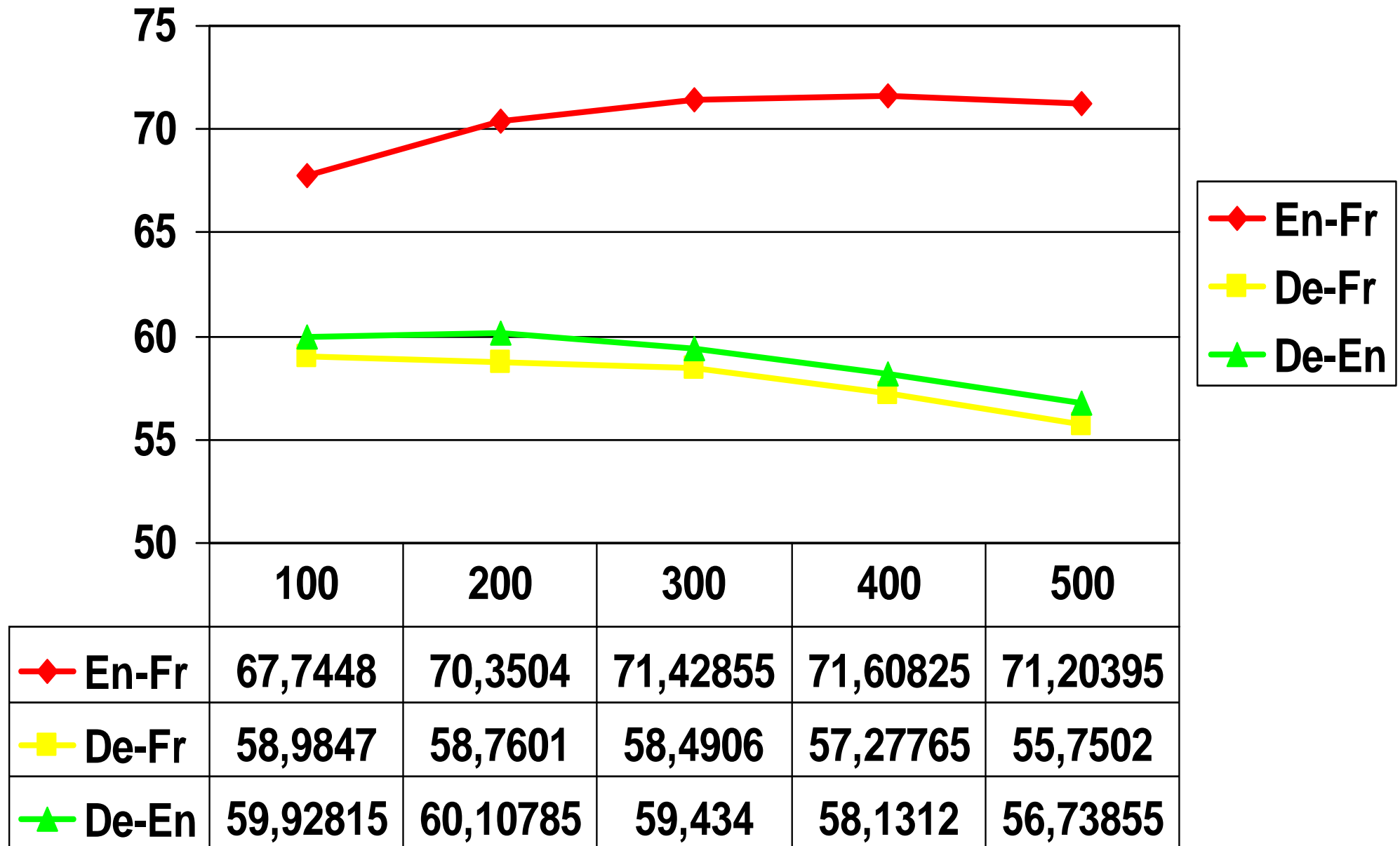
Wikipedia – Multext

Training on Wikipedia (ca. 50 000 documents)

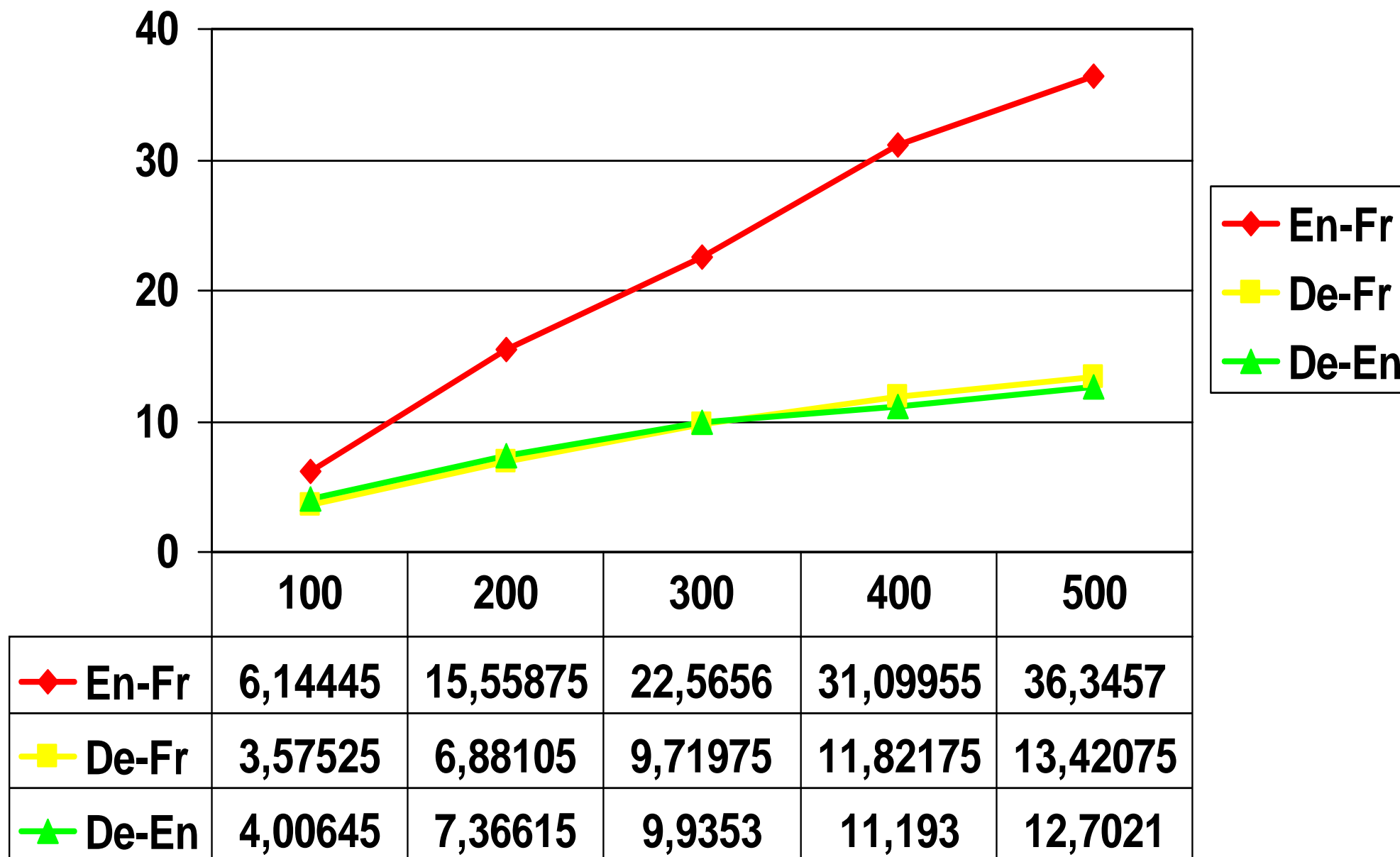
Mate Retrieval on Multext



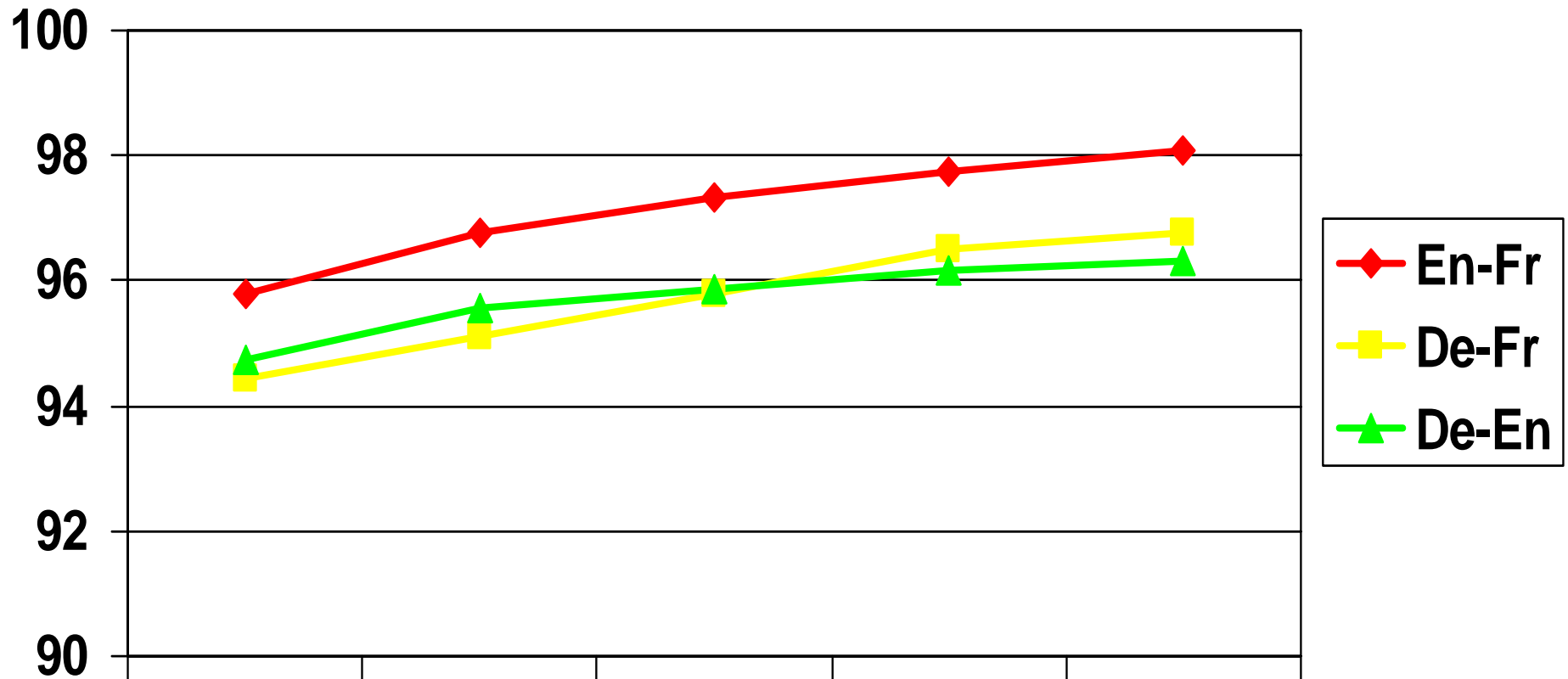
# MeanMate-Multex Top1



# MeanMate Multex-Wiki Top1

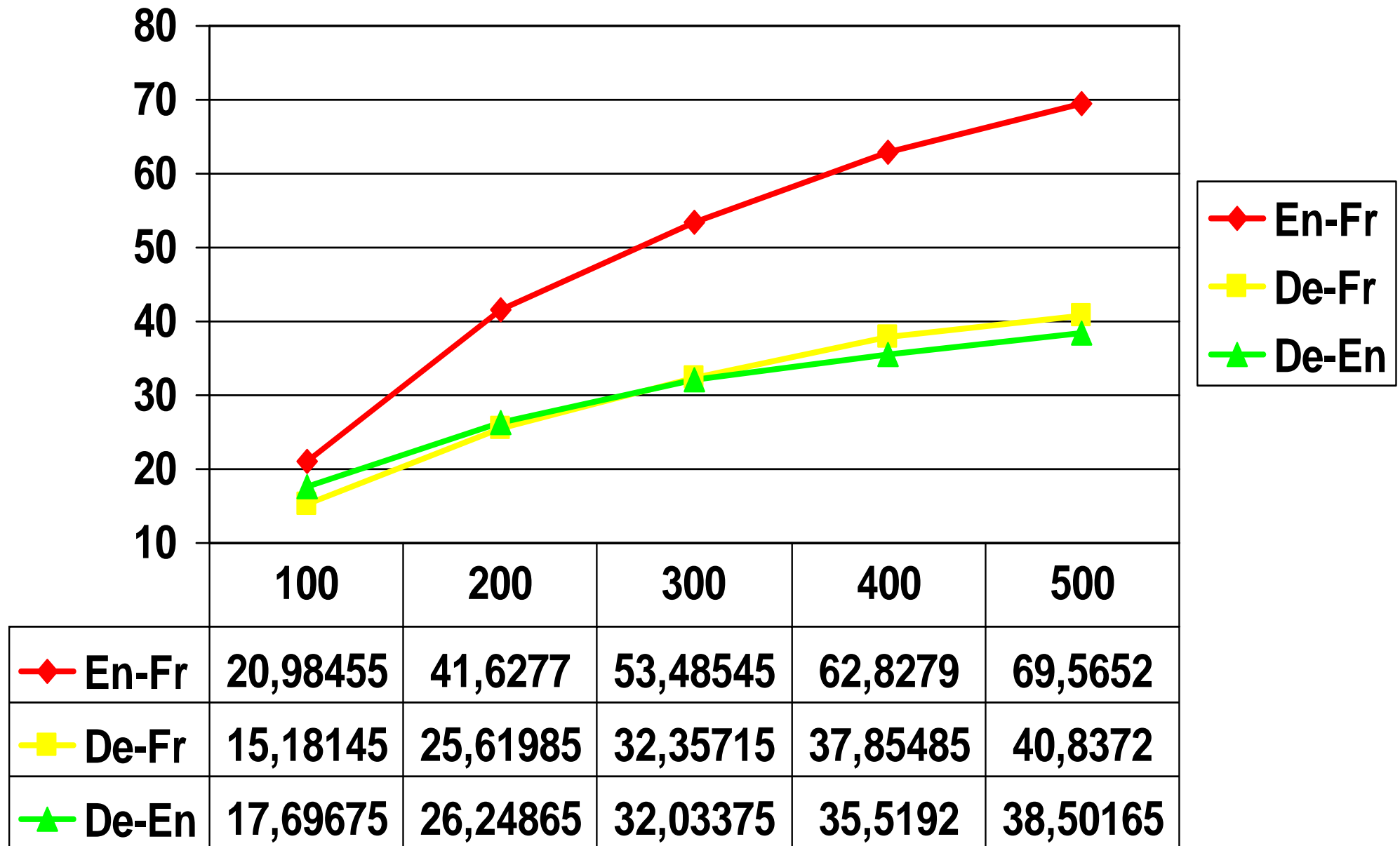


# MeanMate-Multex Top10



	100	200	300	400	400
<b>◆ En-Fr</b>	<b>95,7772</b>	<b>96,7655</b>	<b>97,3495</b>	<b>97,75385</b>	<b>98,0683</b>
<b>■ De-Fr</b>	<b>94,42945</b>	<b>95,1033</b>	<b>95,77715</b>	<b>96,496</b>	<b>96,7655</b>
<b>▲ De-En</b>	<b>94,74395</b>	<b>95,5526</b>	<b>95,86705</b>	<b>96,1815</b>	<b>96,31625</b>

# MeanMate Multex-Wiki Top10



# LSI: Summary

---

- + Elegant, mathematically well-founded model
- + „Automatic learning“ of term correlations (incl. morphological variants, multilingual corpus)
- + Implicit thesaurus (by correlations between synonyms)
- + Implicit discrimination of different meanings of polysems (by different term correlations)
- + Improved precision and recall on „closed“ corpora (e.g. TREC benchmark, financial news, patent databases, etc.) with empirically best  $k$  in the order of 100-200
- In general difficult choice of appropriate  $k$
- Computational and storage overhead for very large (sparse) matrices
- No convincing results for Web search engines (yet)

# Probabilistic Latent Semantic Analysis

---

Identification of Latent Classes (Aspects) using an Expectation Maximization (EM) Algorithm

Shown to solve

- ◆ Polysemy
  - Java could mean "coffee" and also the "PL Java"
  - Cricket is a "game" and also an "insect"
- ◆ Synonymy
  - "computer", "pc", "desktop" all could mean the same

Has a better statistical foundation than LSA

# What is a Statistical Language Model?

---

generative model for word sequence

(generates probability distribution of word sequences,  
or bag-of-words, or set-of-words, or structured doc, or ...)

Example:  $P[\text{„Today is Tuesday“}] = 0.001$

$P[\text{„Today Wednesday is“}] = 0.000000001$

$P[\text{„The Eigenvalue is positive“}] = 0.000001$

LM itself highly context- / application-dependent

Examples:

- **speech recognition**: given that we heard „Julia“ and „feels“, how likely will we next hear „happy“ or „habit“?
- **text classification**: given that we saw „soccer“ 3 times and „game“ 2 times, how likely is the news about sports?
- **information retrieval**: given that the user is interested in math, how likely would the user use „distribution“ in a query?

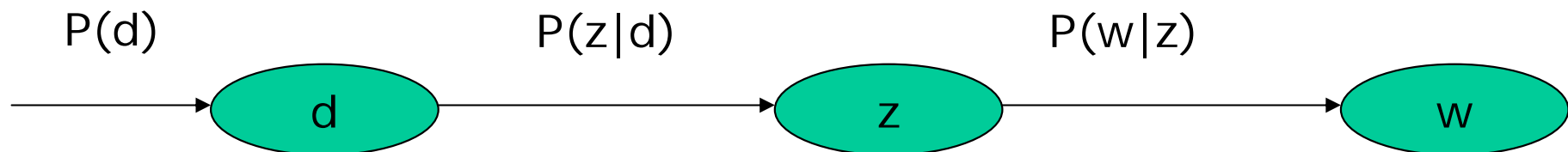
# Aspect Model

## Latent Variable model for general co-occurrence data

- ◆ Associate each observation  $(w,d)$  with a class variable  $z \in Z\{z_1, \dots, z_K\}$

## Generative Model

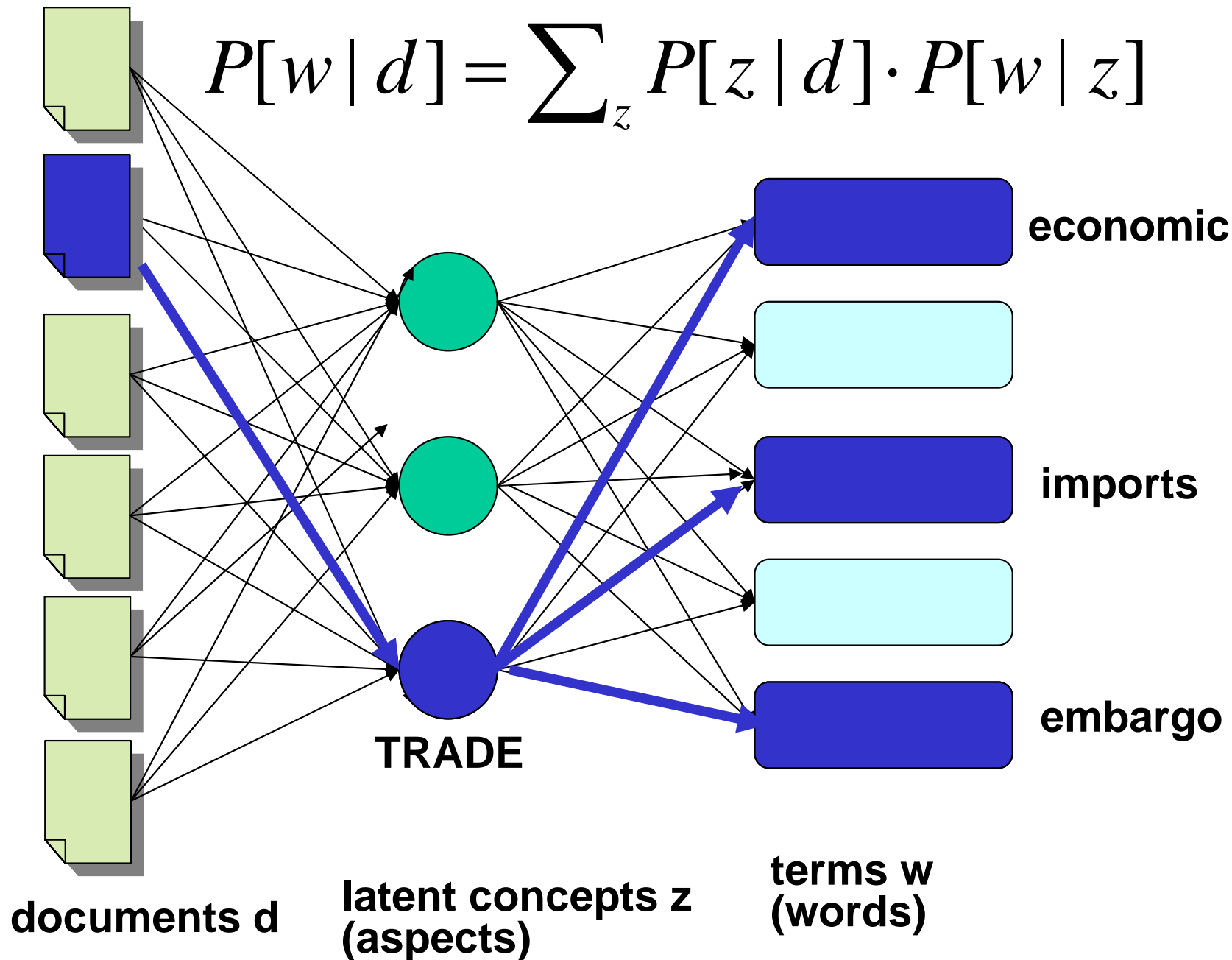
- ◆ Select a doc with probability  $P(d)$
- ◆ Pick a latent class  $z$  with probability  $P(z|d)$
- ◆ Generate a word  $w$  with probability  $P(w|z)$





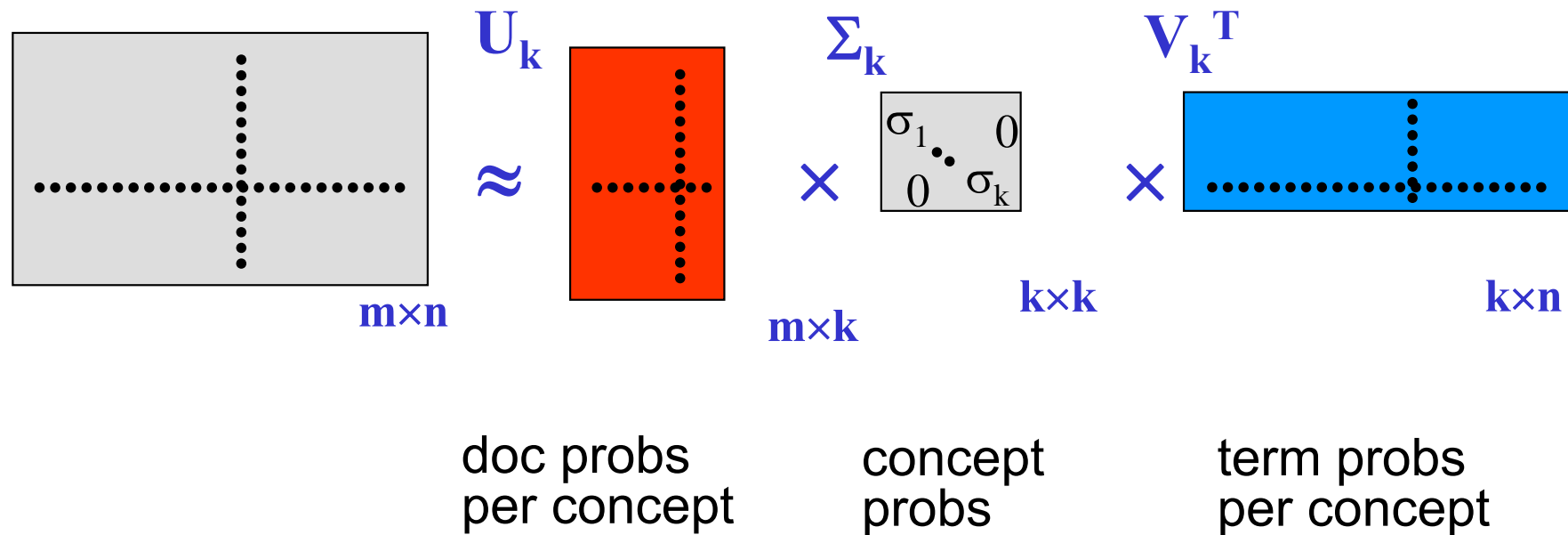
$$P[w | d] = \sum_z P[z | d] \cdot P[w | z]$$

**d and w  
conditionally  
independent  
given z**



# Relationship of pLSI to LSI

$$P[d, w] = \sum_z P[d/z] \cdot P[z] \cdot P[w/z]$$



## Key difference to LSI:

- non-negative matrix decomposition
- with L1 normalization

## Key difference to LMs:

- no generative model for docs
- tied to given corpus

# Aspect Model

---

## Aspect Model

- Document is a mixture of underlying (latent)  $K$  aspects
- Each aspect is represented by a distribution of words  $P(w|z)$

## Model fitting with Tempered EM

The joint distribution  $P(d, w)$  is the sum over all aspects  $z$

$$P(d, w) = \sum_z P(z) P(d, w | z)$$

Assuming conditional independence of  $d$  and  $w$  given  $z$ .

$$P(d, w) = \sum_z P(z) P(d | z) P(w | z)$$

# Expectation-Maximization Method (EM)

## Key idea:

when  $L(\theta, X_1, \dots, X_n)$  (where the  $X_i$  and  $\theta$  are possibly multivariate) is analytically intractable then

- introduce ***latent (hidden, invisible, missing) random variable(s) Z*** such that
  - the ***joint distribution  $J(X_1, \dots, X_n, Z, \theta)$  of the „complete“ data*** is tractable (often with  $Z$  actually being  $Z_1, \dots, Z_n$ )
- derive the incomplete-data likelihood  $L(\theta, X_1, \dots, X_n)$  by ***integrating (marginalization) J:***

$$\hat{\theta} = \arg \max_{\theta} \sum_z J[\theta, X_1, \dots, X_n, Z | Z = z] P[Z = z]$$

# EM Procedure

---

**Initialization:** choose start estimate for  $\theta^{(0)}$

**Iterate** ( $t=0, 1, \dots$ ) until convergence:

**E step (expectation):**

estimate posterior probability of  $Z$ :  $P[Z | X_1, \dots, X_n, \theta^{(t)}]$   
assuming  $\theta$  were known and equal to previous estimate  $\theta^{(t)}$ ,  
and compute  $E_{Z | X_1, \dots, X_n, \theta^{(t)}} [\log J(X_1, \dots, X_n, Z | \theta)]$   
by integrating over values for  $Z$

**M step (maximization, MLE step):**

Estimate  $\theta^{(t+1)}$  by maximizing  
 $E_{Z | X_1, \dots, X_n, \theta^{(t)}} [\log J(X_1, \dots, X_n, Z | \theta)]$

convergence is guaranteed

(because the E step computes a lower bound of the true L function,  
and the M step yields monotonically non-decreasing likelihood),  
but may result in local maximum of log-likelihood function

# EM at Indexing Time

(pLSI Model Fitting)

**observed data:**  $n(d, w)$  – absolute frequency of word  $w$  in doc  $d$

**model params:**  $P[z|d]$ ,  $P[w|z]$  for concepts  $z$ , words  $w$ , docs  $d$

**maximize log-likelihood**  $\sum_d \sum_w n(d, w) \cdot \log P[dw]$

**E step:** posterior probability of latent variables

$$P[z | d, w] = \frac{P[z | d]P[w | z]}{\sum_y P[y | d]P[w | y]}$$

prob. that occurrence of word  $w$  in doc  $d$  can be explained by concept  $z$

**M step:** MLE with completed data

$$P[w | z] \sim \sum_d n(d, w)P[z | d, w] \quad P[z | d] \sim \sum_w n(d, w)P[z | d, w]$$

freq. of  $w$  associated with  $z$

freq. of  $d$  associated with  $z$

actual procedure „perturbs“ EM for „smoothing“  
(avoidance of overfitting) → tempered annealing

## EM Details (PLSI Model Fitting)

$$P[z | d, w] = \frac{P[z | d] P[w | z]}{\sum_y P[y | d] P[w | y]} \quad \text{(E)}$$

$$P[w | z] = \frac{\sum_d n(d, w) P[z | d, w]}{\sum_{d, u} n(d, u) P[z | d, u]} \quad \text{(M1)}$$

$$P[z | d] = \frac{\sum_w n(d, w) P[z | d, w]}{\sum_{w, y} n(d, w) P[y | d, w]} \quad \text{(M2)}$$

or equivalently compute  $P[z]$ ,  $P[d|z]$ ,  $P[w|z]$  in M step  
(see S. Chakrabarti, pp. 110/111)

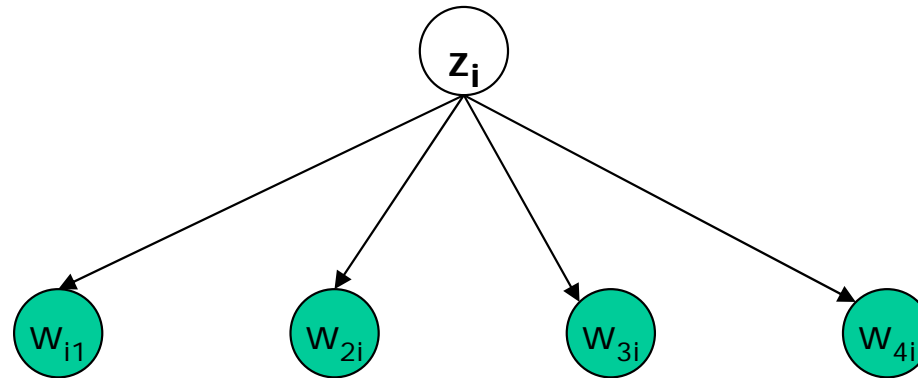
---

# Generalization: Language Models



# Mixture of Unigrams

---



Mixture of Unigrams Model (this is just Naïve Bayes)

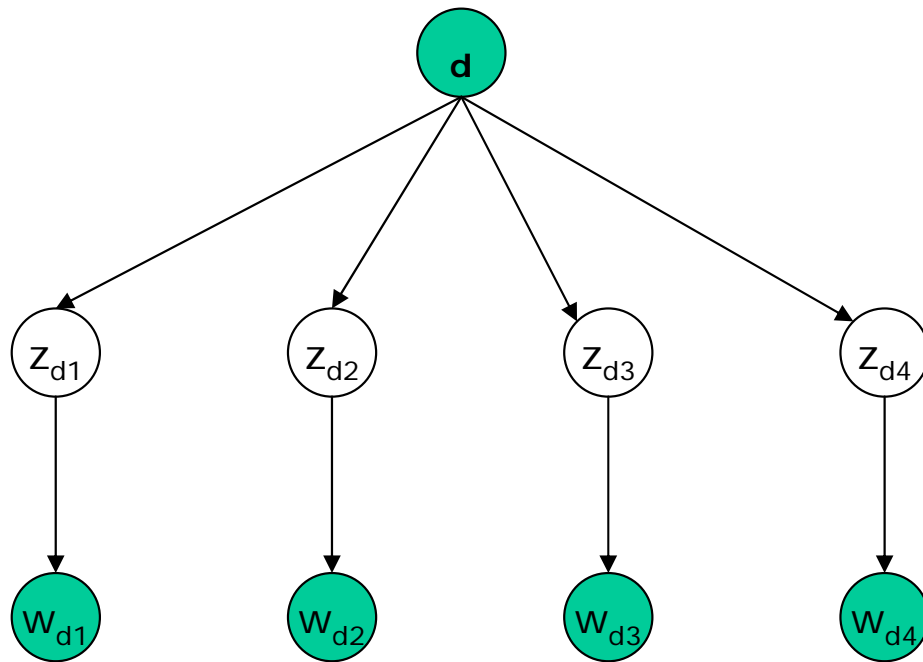
For each of  $M$  documents,

Choose a topic  $z$ .

Choose  $N$  words by drawing each one independently from a multinomial conditioned on  $z$ .

In the Mixture of Unigrams model, we can only have one topic per document!

# The pLSI Model



Probabilistic Latent Semantic Indexing (pLSI) Model

For each word of document  $d$  in the training set,

Choose a topic  $z$  according to a multinomial conditioned on the index  $d$ .

Generate the word by drawing from a multinomial conditioned on  $z$ .

In pLSI, documents can have multiple topics.

# Motivations for LDA

---

In pLSI, the observed variable  $d$  is an index into some training set. There is no natural way for the model to handle previously unseen documents.

The number of parameters for pLSI grows linearly with  $M$  (the number of documents in the training set).

We would like to be Bayesian about our topic mixture proportions.

# Dirichlet Distributions

---

In the LDA model, we would like to say that the *topic mixture proportions* for each document are drawn from some distribution.

So, we want to put a distribution on multinomials. That is,  $k$ -tuples of non-negative numbers that sum to one.

The space of all of these multinomials has a nice geometric interpretation as a  $(k-1)$ -*simplex*, which is just a generalization of a triangle to  $(k-1)$  dimensions.

Criteria for selecting our prior:

- ◆ It needs to be defined for a  $(k-1)$ -simplex.
- ◆ Algebraically speaking, we would like it to play nice with the multinomial distribution.

# Dirichlet Distributions

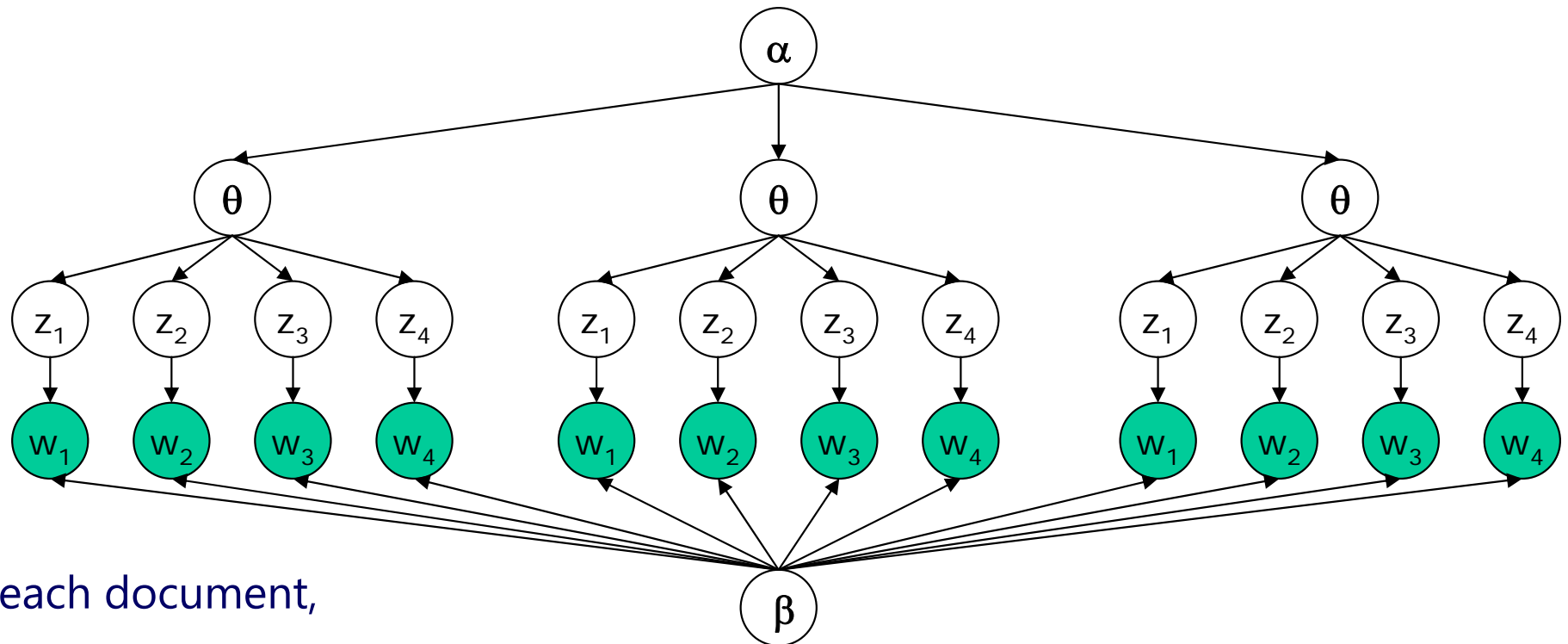
---

$$p(\theta|\alpha) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \prod_{i=1}^k \theta_i^{\alpha_i - 1}$$

## Useful Facts:

- ◆ This distribution is defined over a (k-1)-simplex. That is, it takes k non-negative arguments which sum to one. Consequently it is a natural distribution to use over multinomial distributions.
- ◆ In fact, the Dirichlet distribution is the conjugate prior to the multinomial distribution. (This means that if our likelihood is multinomial with a Dirichlet prior, then the posterior is also Dirichlet!)
- ◆ The Dirichlet parameter  $\alpha_i$  can be thought of as a prior count of the  $i^{\text{th}}$  class.

# The LDA Model



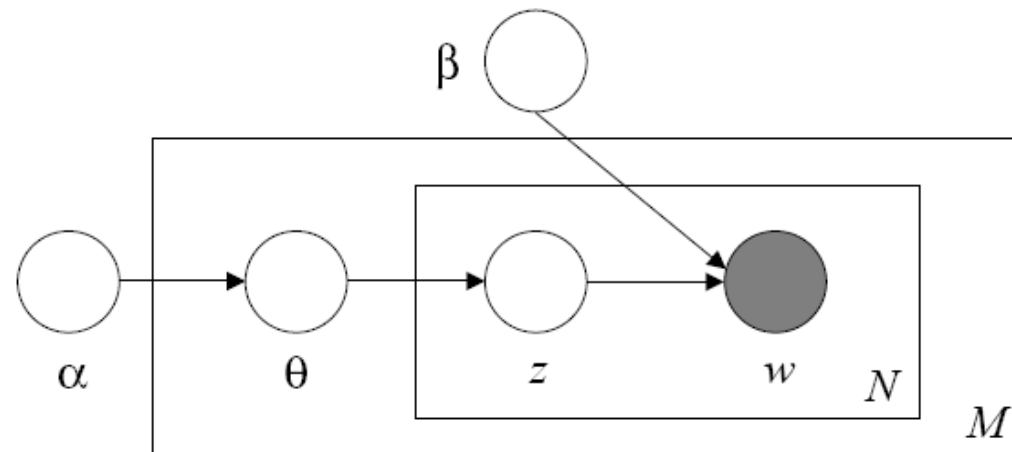
For each document,

Choose  $\theta \sim \text{Dirichlet}(\alpha)$

For each of the  $N$  words  $w_n$ :

- ◆ Choose a topic  $z_n \sim \text{Multinomial}(\theta)$
- ◆ Choose a word  $w_n$  from  $p(w_n | z_n, \beta)$ , a multinomial probability conditioned on the topic  $z_n$ .

# The LDA Model



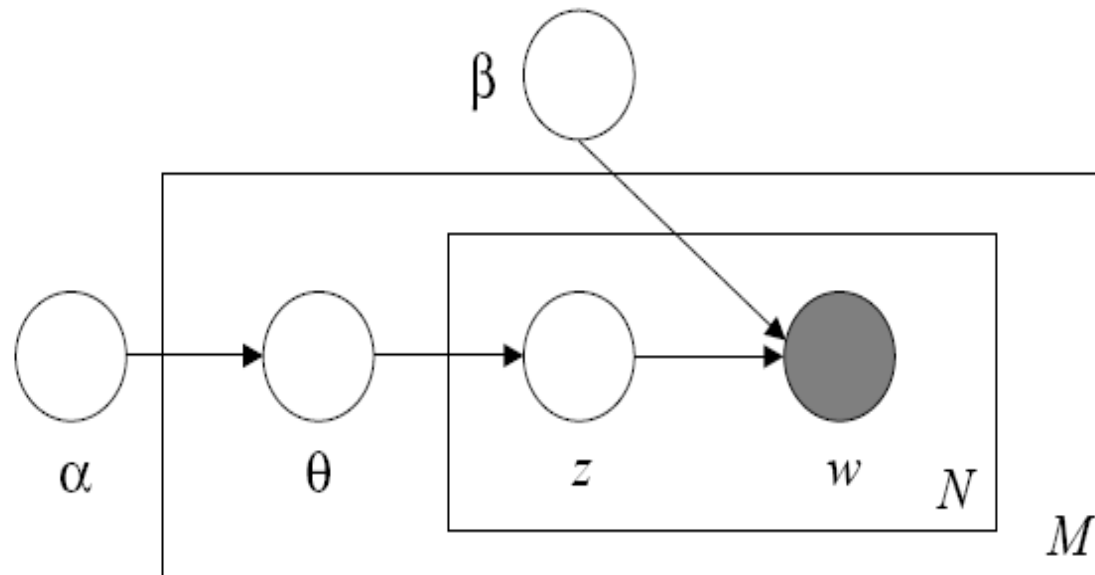
For each document,

Choose  $\theta \sim \text{Dirichlet}(\alpha)$

For each of the  $N$  words  $w_n$ :

- ◆ Choose a topic  $z_n \sim \text{Multinomial}(\theta)$
- ◆ Choose a word  $w_n$  from  $p(w_n|z_n, \beta)$ , a multinomial probability conditioned on the topic  $z_n$ .

# Inference



The inference problem in LDA is to compute the posterior of the hidden variables given a document and corpus parameters  $\alpha$  and  $\beta$ . That is, compute  $p(\theta, z | w, \alpha, \beta)$ .

Unfortunately, exact inference is intractable, so we turn to alternatives...



# Some Results

Given a topic, LDA can return the most probable words.

For the following results, LDA was trained on 10,000 text articles posted to 20 online newsgroups with 40 iterations of EM. The number of topics was set to 50.

"politics"	"sports"	"space"	"computers"	"christianity"
Political	Team	Space	Drive	God
Party	Game	NASA	Windows	Jesus
Business	Play	Research	Card	His
Convention	Year	Center	DOS	Bible
Institute	Games	Earth	SCSI	Christian
Committee	Win	Health	Disk	Christ
States	Hockey	Medical	System	Him
Rights	Season	Gov	Memory	Christians

---

Q: How can we exploit the topical model (pLSI, LDA) for ranked retrieval?!

Quick solution: roll out the vector space model with cosine similarity..

Better: exploit (dis)similarity measures known from Information Theory (how?!)

---

# Top-K query processing for text documents

# Simplified Query Processor

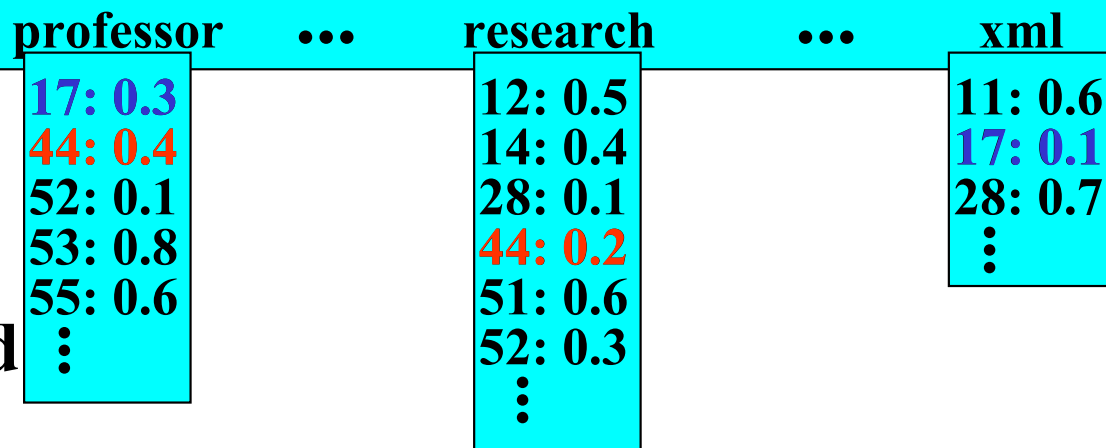
1. Remove stopwords from query and map all words onto terms (word stems)
2. Convert query  $q$  into normal form  
[ (t<sub>11</sub> AND t<sub>12</sub> AND ... AND t<sub>1k<sub>1</sub></sub>)  
OR ...  
(t<sub>r1</sub> AND t<sub>r2</sub> AND ... AND t<sub>rk<sub>r</sub></sub>) ]  
AND NOT t<sub>r+1</sub> AND NOT t<sub>r+2</sub> ...
3. For  $i=1$  to  $r$  do  
find DocId-Lists for  $t_{i_1}$  bis  $t_{i_{k_i}}$  and  
compute union  $V_i$  of these DocIds  
set  $rank(d)$  of  $d$  in  $V_i$  auf  $sim(d, q_i)$   
Compute union  $V$  of DocId-Sets  $V_1, \dots, V_r$
4. Sort all  $d$  in  $V$  in descending order of  $rank(d)$
5. Retrieve docs  $d$  in  $V$  in sorted order (possibly with heuristic termination criteria),  
eliminate docs  $d$  that contain one or more terms  $t_{r+1}, t_{r+2}, \dots$

# Top-k Query Processing with Scoring

Vector space model suggests *m*×*n* *term-document matrix*,  
but data is sparse and queries are even very sparse  
→ better use *inverted index lists* with terms as keys for B+ tree

q: professor  
research  
xml

B+ tree on terms



index lists with  
(DocId,  
s = tf\*idf)  
sorted by DocId

Google:  
> 10 mio. terms  
> 8 bio. docs  
> 4 TB index

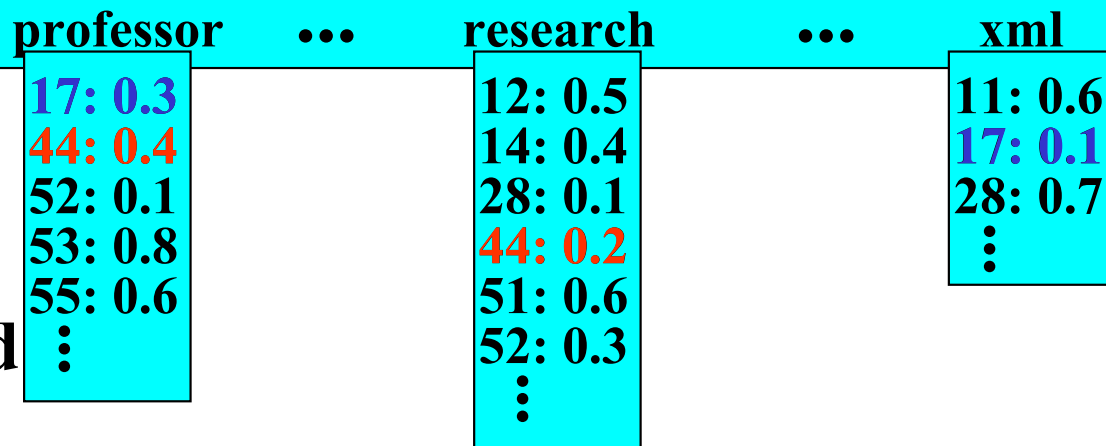
terms can be full words, word stems, word pairs,  
word substrings, etc.  
(whatever „dictionary terms“ we prefer for the application)

queries can be conjunctive or „andish“ (soft conjunction)

# DBS-Style Top-k Query Processing

q: professor  
research  
xml

B+ tree on terms



index lists with  
(DocId,  
s = tf\*idf)  
sorted by DocId

Google:  
> 10 mio. terms  
> 8 bio. docs  
> 4 TB index

**Given:** query  $q = t_1 t_2 \dots t_z$  with  $z$  (conjunctive) keywords  
similarity scoring function  $\text{score}(q,d)$  for docs  $d \in D$ , e.g.:  $\vec{q} \cdot \vec{d}$   
with precomputed scores (index weights)  $s_i(d)$  for which  $q_i \neq 0$

**Find:** top  $k$  results w.r.t.  $\text{score}(q,d) = \text{aggr}\{s_i(d)\}$  (e.g.:  $\sum_{i \in q} s_i(d)$ )

*Naive join&sort QP algorithm:*

top-k (

$\sigma[\text{term}=t_1]$ (index)	×	DocId
$\sigma[\text{term}=t_2]$ (index)	×	DocId
...	×	DocId
$\sigma[\text{term}=t_z]$ (index)	×	DocId

order by s desc)

## Index List Processing by Merge Join

Keep  $L(i)$  in **ascending order of doc ids**

Compress  $L(i)$  by actually storing the gaps between successive doc ids  
(or using some more sophisticated prefix-free code)

QP may start with those  $L(i)$  lists that are short and have high idf

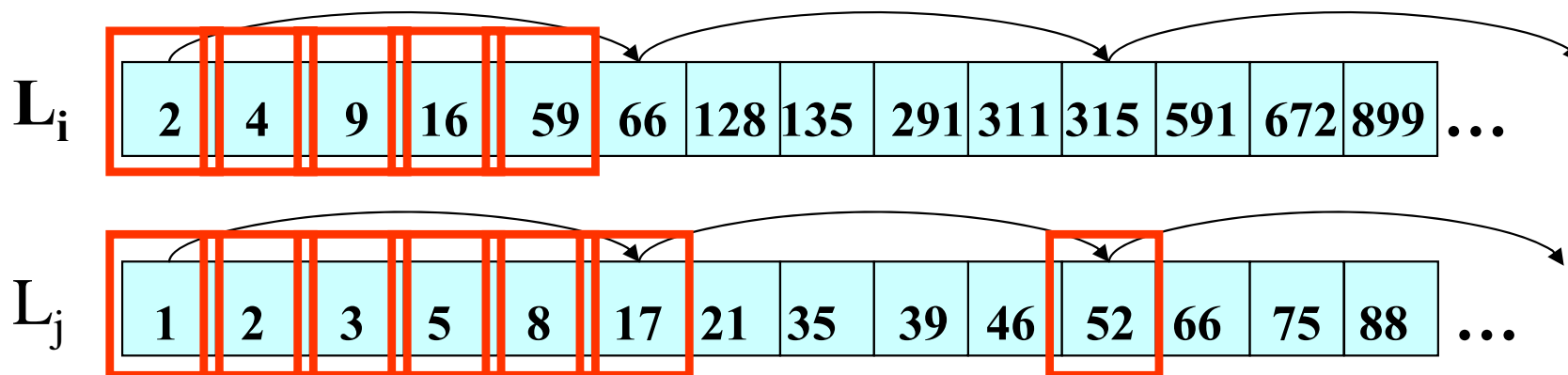
Candidate results need to be looked up in other lists  $L(j)$

To avoid having to uncompress the entire list  $L(j)$ ,

$L(j)$  is encoded into groups of entries

with a **skip pointer** at the start of each group

→  $\sqrt{n}$  evenly spaced skip pointers for list of length  $n$



## Computational Model for Top-k Queries

Assume *local scores*  $s_i$  for query  $q$ , data item  $d$ , and dimension  $i$ , and *global scores*  $s$  of the form  $s(q, d) = \text{aggr}\{s_i(q, d) / i = 1..m\}$  with a *monotonic* aggregation function  $\text{aggr} : [0,1]^m \rightarrow [0,1]$

Examples:  $s(q, d) = \sum_{i=1}^m s_i(q, d)$        $s(q, d) = \max\{s_i(q, d) / i = 1..m\}$

**Find top-k data items with regard to global scores:**

- process  $m$  *index lists*  $L_i$  with *sorted access (SA)* to entries  $(d, s_i(q, d))$  in *ascending order of doc ids* or *descending order of  $s_i(q, d)$*
- maintain for each candidate  $d$  a set  $E(d)$  of evaluated dimensions and a *partial score „accumulator“*
- for candidate  $d$  with incomplete  $E(d)$  consider looking up  $d$  in  $L_i$  for all  $i \in R(d)$  by *random access (RA)*
- terminate index list scans when enough candidates have been seen
- if necessary sort final candidate list by global score

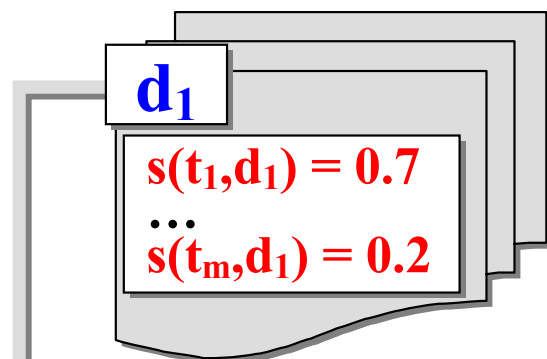


# Efficient Top-k Search

[Buckley85, Güntzer/Balke/Kießling 00, Fagin01]

*threshold algorithms: efficient & principled top-k query processing with monotonic score aggr.*

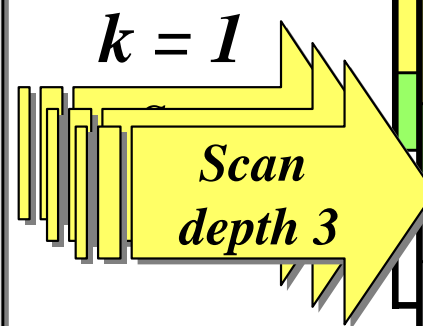
Data items:  $d_1, \dots, d_n$



Query:  $q = (t_1, t_2, t_3)$

	Index lists					
$t_1$	d78 0.9	d23 0.8	d10 0.8	d1 0.7	d88 0.2	...
$t_2$	d64 0.8	d23 0.6	d10 0.6	d10 0.2	d78 0.1	...
$t_3$	d10 0.7	d78 0.5	d64 0.4	d99 0.2	d34 0.1	...

**TA with sorted access only (NRA):**  
 can index lists; consider  $d$  at  $\text{pos}_i$  in  $L_i$ ;  
 $E(d) := E(d) \cup \{i\}$ ;  $\text{high}_i := s(t_i, d)$ ;  
 $\text{worstscore}(d) := \text{aggr}\{s(t_v, d) \mid v \in E(d)\}$ ;  
 $\text{bestscore}(d) := \text{aggr}\{\text{worstscore}(d), \text{aggr}\{\text{high}_v \mid v \notin E(d)\}\}$ ;  
 if  $\text{worstscore}(d) > \text{min-k}$  then add  $d$  to top-k  
 $\text{min-k} := \min\{\text{worstscore}(d') \mid d' \in \text{top-k}\}$ ;  
 else if  $\text{bestscore}(d) > \text{min-k}$  then  
 $\text{cand} := \text{cand} \cup \{d\}$ ;  $s$   
 $\text{threshold} := \max\{\text{bestscore}(d') \mid d' \in \text{cand}\}$ ;  
 if  $\text{threshold} \leq \text{min-k}$  then exit;



Rank	Doc	Worst-score	Best-score
1	d10	2.1	2.1
2	d78	1.4	2.0
3			1.8
4		1.2	2.0

A red starburst with the word "STOP!" is overlaid on the table, indicating that the search process has terminated.

*keep  $L(i)$  in descending order of scores*

# Threshold Algorithm (TA, Quick-Combine, MinPro)

(Fagin'01; Güntzer/Balke/Kießling; Nepal/Ramakrishna)

```
scan all lists  $L_i$  ( $i=1..m$ ) in parallel:
  consider  $d_j$  at position  $pos_i$  in  $L_i$ ;
   $high_i := s_i(d_j)$ ;
  if  $d_j \notin \text{top-k}$  then {
    look up  $s_v(d_j)$  in all lists  $L_v$  with  $v \neq i$ ; // random access
    compute  $s(d_j) := \text{aggr } \{s_v(d_j) \mid v=1..m\}$ ;
    if  $s(d_j) > \text{min score among top-k}$  then
      add  $d_j$  to top-k and remove min-score  $d$  from top-k; }
   $\text{threshold} := \text{aggr } \{high_v \mid v=1..m\}$ ;
  if  $\text{min score among top-k} \geq \text{threshold}$  then exit;
```

*but random accesses  
are expensive !*

$m=3$

aggr: sum

$k=2$

f: 0.5
b: 0.4
c: 0.35
a: 0.3
h: 0.1
d: 0.1

a: 0.55
b: 0.2
f: 0.2
g: 0.2
c: 0.1

h: 0.35
d: 0.35
b: 0.2
a: 0.1
c: 0.05
f: 0.05

top-k:

~~f: 0.75~~

a: 0.95

b: 0.8

# No-Random-Access Algorithm

(NRA, Stream-Combine, TA-Sorted)

scan index lists in parallel:

consider  $d_j$  at position  $pos_i$  in  $L_i$ ;

$E(d_j) := E(d_j) \cup \{i\}$ ;  $high_i := si(q, d_j)$ ;

$bestscore(d_j) := aggr\{x_1, \dots, x_m\}$

with  $x_i := si(q, d_j)$  for  $i \in E(d_j)$ ,  $high_i$  for  $i \notin E(d_j)$ ;

$worstscore(d_j) := aggr\{x_1, \dots, x_m\}$

with  $x_i := si(q, d_j)$  for  $i \in E(d_j)$ , 0 for  $i \notin E(d_j)$ ;

$top-k := k$  docs with largest  $worstscore$ ;

$threshold := bestscore\{d \mid d \text{ not in } top-k\}$ ;

if  $\min worstscore$  among  $top-k \geq threshold$  then exit;

$m=3$

aggr: sum

$k=2$

f: 0.5  
b: 0.4  
c: 0.35  
a: 0.3  
h: 0.1  
d: 0.1

a: 0.55  
b: 0.2  
f: 0.2  
g: 0.2  
c: 0.1

h: 0.35  
d: 0.35  
b: 0.2  
a: 0.1  
c: 0.05  
f: 0.05

top-k:

a: 0.95

b: 0.8

candidates:

f:  $0.7 + ? \leq 0.7 + 0.1$

h:  $0.35 + ? \leq 0.35 + 0.5$

c:  $0.35 + ? \leq 0.35 + 0.3$

d:  $0.35 + ? \leq 0.35 + 0.5$

g:  $0.2 + ? \leq 0.2 + 0.4$

## Optimality of TA

### Definition:

For a class  $\mathcal{A}$  of algorithms and a class  $\mathcal{D}$  of datasets, let  $\text{cost}(A, D)$  be the execution cost of  $A \in \mathcal{A}$  on  $D \in \mathcal{D}$ .

Algorithm B is *instance optimal* over  $\mathcal{A}$  and  $\mathcal{D}$  if

for every  $A \in \mathcal{A}$  on  $D \in \mathcal{D}$ :  $\text{cost}(B, D) = O(\text{cost}(A, D))$ ,

that is:  $\text{cost}(B, D) \leq c \cdot \text{cost}(A, D) + c'$

with optimality ratio (competitiveness)  $c$ .

### Theorem:

- TA is instance optimal over all algorithms that are based on sorted and random access to (index) lists (no „wild guesses“).

TA has optimality ratio  $m + m(m-1) C_{RA}/C_{SA}$

with random-access cost  $C_{RA}$  and sorted-access cost  $C_{SA}$

- NRA is instance-optimal over all algorithms with SA only.

*if „wild guesses“ are allowed,*

*then no deterministic algorithm is instance-optimal*

## Execution Cost of TA Family

---

**Run-time cost** is  $O\left(n^{\frac{m-1}{m}} \cdot k^{\frac{1}{m}}\right)$  with arbitrarily high probability

(for independently distributed Li lists)

**Memory cost** is  $O(k)$  for TA

and  $O(n^{(m-1)/m})$  for NRA (priority queue of candidates)

## Approximation TA:

A  *$\theta$ -approximation*  $T'$  for top-k query  $q$  with  $\theta > 1$  is a set  $T'$  of docs with:

- $|T'|=k$  and
- for each  $d' \in T'$  and each  $d'' \notin T'$ :  $\theta * \text{score}(q, d') \geq \text{score}(q, d'')$

## Modified TA:

...

Stop when  $\min_k \geq \text{aggr}(\text{high}_1, \dots, \text{high}_m) / \theta$

## General heuristics:

- **disregard index lists with idf below threshold**
- **for index scans give priority to index lists that are short and have high idf**

---

# String similarity



**Hamming-Distanz** von Strings  $s_1, s_2 \in \Sigma^*$  mit  $|s_1|=|s_2|$ :

Anzahl verschiedener Zeichen (Kardinalität von  $\{i: s_{1i} \neq s_{2i}\}$ )

**Levenshtein-Distanz (Editierdistanz)** von Strings  $s_1, s_2 \in \Sigma^*$ :

minimale Anzahl der Editieroperationen auf  $s_1$

(Ersetzen, Löschen oder Einfügen eines Zeichens),

um  $s_1$  in  $s_2$  zu ändern

Für

$\text{edit}(i, j)$ : Levenshtein-Distanz von  $s_1[1..i]$  und  $s_2[1..j]$

gilt:  $\text{edit}(0, 0) = 0$ ,  $\text{edit}(i, 0) = i$ ,  $\text{edit}(0, j) = j$

$\text{edit}(i, j) = \min \{ \text{edit}(i-1, j) + 1,$

$\text{edit}(i, j-1) + 1,$

$\text{edit}(i-1, j-1) + \text{diff}(i, j) \}$

mit  $\text{diff}(i, j) = 1$  falls  $s_{1i} \neq s_{2j}$ , 0 sonst

→ Berechnung durch dynamische Programmierung

## String similarity : Damerau-Levenshtein

**Damerau-Levenshtein-Distanz** von Strings  $s_1, s_2 \in \Sigma^*$ :

minimale Anzahl an Ersetzungs-, Einfüge-, Lösch- oder Transpositionsoperationen (Vertauschen benachbarter Zeichen), um  $s_1$  in  $s_2$  zu ändern

Für  $\text{edit}(i, j)$ : Damerau-Levenshtein-Distanz von  $s_1[1..i]$  und  $s_2[1..j]$

gilt:  $\text{edit}(0, 0) = 0$ ,  $\text{edit}(i, 0) = i$ ,  $\text{edit}(0, j) = j$

$$\text{edit}(i, j) = \min \left\{ \begin{array}{l} \text{edit}(i-1, j) + 1, \\ \text{edit}(i, j-1) + 1, \\ \text{edit}(i-1, j-1) + \text{diff}(i, j), \\ \text{edit}(i-2, j-2) + \text{diff}(i-1, j) + \text{diff}(i, j-1) + 1 \end{array} \right\}$$

mit  $\text{diff}(i, j) = 1$  falls  $s_1[i] \neq s_2[j]$ , 0 sonst

# String Similarity: N-Grams

---

Bestimme für String  $s$  die Menge seiner  $N$ -Gramme:

$G(s) = \{\text{Substrings von } s \text{ mit der Länge } N\}$

(häufig werden Trigramme betrachtet mit  $N=3$ )

Ähnlichkeit von Strings  $s_1$  und  $s_2$ :

$$|G(s_1)| + |G(s_2)| - 2|G(s_1) \cap G(s_2)|$$

Beispiel:

$G(\text{rodney}) = \{\text{rod, odn, dne, ney}\}$

$G(\text{rhodnee}) = \{\text{rho, hod, odn, dne, nee}\}$

$$\text{Ähnlichkeit}(\text{rodney}, \text{rhodnee}) = 4 + 5 - 2 \cdot 2 = 5$$

# Phonetic Similarity (1)

---

## Soundex-Code:

Abbildung von Wörtern (speziell: Nachnamen) auf 4-Zeichen-Codes,  
so daß ähnlich ausgesprochene Wörter denselben Code haben

- Erstes Codezeichen = erstes Zeichen des Wortes
- Codezeichen zwei bis vier (a, e, i, o, u, y, h, w werden ignoriert):

b, p, f, v → 1 c, s, g, j, k, q, x, z → 2

d, t → 3 l → 4

m, n → 5 r → 6

- Aufeinanderfolgende identische Codes werden zusammengefasst  
(es sei denn, sie sind durch h getrennt)

Beispiele:

Powers → P620 , Perez → P620

Penny → P500, Pennee → P500

Tymczak → T522, Tanshik → T522

# Phonetic Similarity: Editex

**Idea:** Editierdistanz mit Berücksichtigung phonetischer Codes

Für  $\text{editex}(i, j)$ : Editex-Distanz von  $s1[1..i]$  und  $s2[1..j]$

gilt:  $\text{editex}(0, 0) = 0$ ,

$\text{editex}(i, 0) = \text{editex}(i-1, 0) + d(s1[i-1], s1[i])$ ,

$\text{editex}(0, j) = \text{editex}(0, j-1) + d(s2[j-1], s2[j])$ ,

$\text{editex}(i, j) = \min \{ \text{editex}(i-1, j) + d(s1[i-1], s1[i]),$   
 $\text{editex}(i, j-1) + d(s2[j-1], s2[j]),$   
 $\text{edit}(i-1, j-1) + \text{diffcode}(i, j) \}$

mit  $\text{diffcode}(i, j) = 0$  falls  $s1i = s2j$ ,

1 falls  $\text{group}(s1i) = \text{group}(s2j)$ , 2 sonst

und  $d(X, Y) = 1$  falls  $X \neq Y$  und  $X$  ist h oder w,

$\text{diffcode}(X, Y)$  sonst

mit  $\text{group}$ :

$\{a e i o u y\}$ ,  $\{b p\}$ ,  $\{c k q\}$ ,  $\{d t\}$ ,  $\{l r\}$ ,

$\{m n\}$ ,  $\{g j\}$ ,  $\{f p v\}$ ,  $\{s x z\}$ ,  $\{c s z\}$

---

## Demos: advanced models for text IR