# Models for Retrieval

**1. HMM/N-gram-based**
**2. Latent Semantic Indexing (LSI)**
**3. Probabilistic Latent Semantic Analysis (PLSA)**

Berlin Chen 2003

References:

1. Berlin Chen et al., "An HMM/N-gram-based Linguistic Processing Approach for Mandarin Spoken Document Retrieval," EUROSPEECH 2001
2. M. W. Berry et al., "Using Linear Algebra for Intelligent Information Retrieval," technical report, 1994
3. Thomas Hofmann, "Unsupervised Learning by Probabilistic Latent Semantic Analysis," Machine Learning, 2001

# HMM/N-gram-based Model

- Model the query $Q$ as a sequence of input observations (index terms), $Q = q_1 q_2 .. q_n .. q_N$
- Model the doc $D$ as a discrete HMM composed of distribution of *N*-gram parameters
- The relevance measure, $P(Q|D \text{ is } R)$, can be estimated by the *N*-gram probabilities of the index term sequence for the query, $Q = q_1 q_2 .. q_n .. q_N$, predicted by the doc $D$
  - *A generative model for IR*

$$D^* = \arg \max_{D} P(D \text{ is } R | Q)$$

$$\approx \arg \max_{D} P(Q|D \text{ is } R) P(D \text{ is } R)$$

$$\approx \arg \max_{D} P(Q|D \text{ is } R)$$ <span style="color:red">with the assumption that ......</span>

# HMM/N-gram-based Model

$$P(W) \qquad \{W = w_1 w_2 .. w_n .. w_N\}$$
$$= P(w_1 w_2 .. w_n .. w_N)$$
$$= P(w_1)P(w_2|w_1)P(w_3|w_1 w_2).... P(w_N|w_1 w_2 .... w_{N-1})$$

- ***N*-gram approximation (Language Model)**
  - Unigram

$$P(W) = P(w_1)P(w_2)P(w_3).....P(w_N)$$

  - Bigram

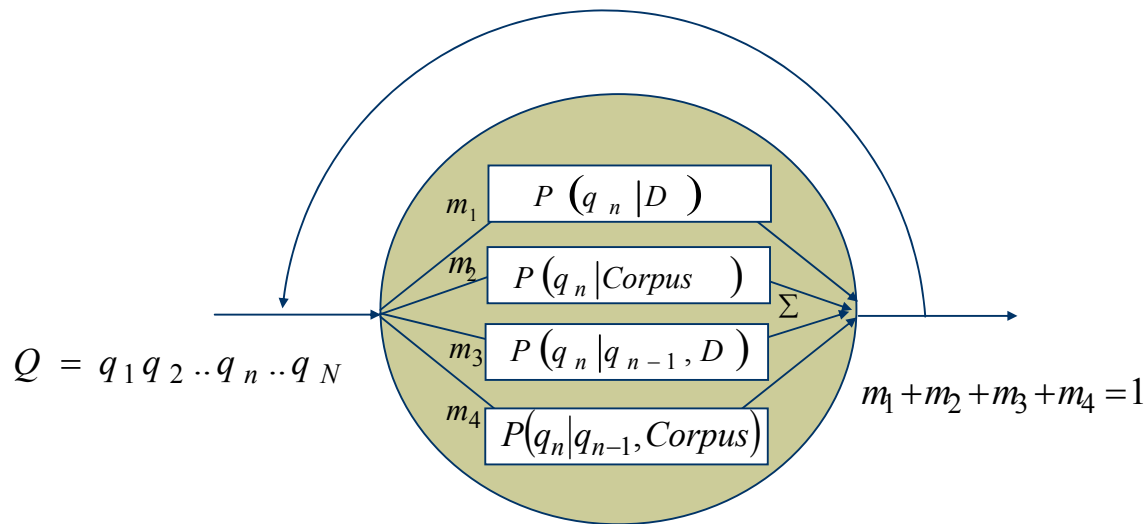$$P(W) = P(w_1)P(w_2|w_1)P(w_3|w_2).....P(w_N|w_{N-1})$$

  - Trigram

$$P(W) = P(w_1)P(w_2|w_1)P(w_3|w_1 w_2).....P(w_N|w_{N-2} w_{N-1})$$

  - ……..

# HMM/N-gram-based Model

- A discrete HMM composed of distribution of *N*-gram parameters

$$Q = q_1 q_2 .. q_n .. q_N$$

$$P(q_n|D)$$
$$P(q_n|Corpus)$$
$$P(q_n|q_{n-1},D)$$
$$P(q_n|q_{n-1},Corpus)$$

$m_1$, $m_2$, $m_3$, $m_4$

$\Sigma$

$$m_1 + m_2 + m_3 + m_4 = 1$$

$$P(Q|D \text{ is } R) = \left[ m_1 P(q_1|D) + m_2 P(q_1|Corpus) \right]$$

$$\cdot \prod_{n=2}^{N} \left[ m_1 P(q_n|D) + m_2 P(q_n|Corpus) + m_3 P(q_n|q_{n-1},D) + m_4 P(q_n|q_{n-1},Corpus) \right]$$

# HMM/N-gram-based Model

- Three Types of HMM Structures
  - Type I: Unigram-Based (Uni)

  $$P(Q|D \text{ is } R) = \prod_{n=1}^{N} \left[ m_1 P(q_n|D) + m_2 P(q_n|Corpus) \right]$$

  - Type II: Unigram/Bigram-Based (Uni+Bi)

  $$P(Q|D \text{ is } R) = \left[ m_1 P(q_1|D) + m_2 P(q_1|Corpus) \right] \\ \cdot \prod_{n=2}^{N} \left[ m_1 P(q_n|D) + m_2 P(q_n|Corpus) + m_3 P(q_n|q_{n-1}, D) \right]$$
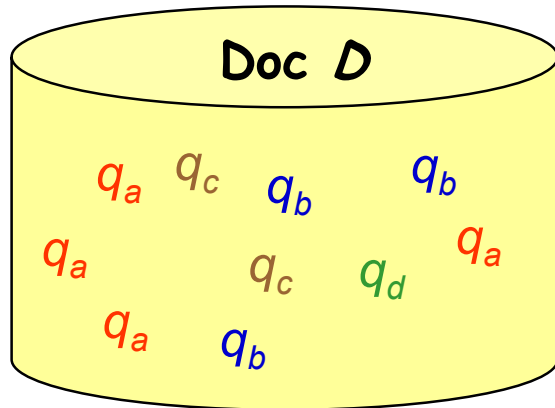
  - Type III: Unigram/Bigram/Corpus-Based (Uni+Bi*)

  $$P(Q|D \text{ is } R) = \left[ m_1 P(q_1|D) + m_2 P(q_1|Corpus) \right] \\ \cdot \prod_{n=2}^{N} \left[ m_1 P(q_n|D) + m_2 P(q_n|Corpus) + m_3 P(q_n|q_{n-1}, D) + m_4 P(q_n|q_{n-1}, Corpus) \right]$$

$P(陳水扁\ 總統\ 視察\ 阿里山\ 小火車|D)$
$=[m_1 P(陳水扁|D)+m_2 P(陳水扁|C)]$ x $[m_1 P(總統|D)+m_2 P(總統|C)+ m_3 P(總統|陳水扁, D)+m_4 P(總統|陳水扁, C)]$
x$[m_1 P(視察|D)+m_2 P(視察|C)+ m_3 P(視察|總統, D)+m_4 P(視察|總統, C)]$ x ..........

# HMM/N-gram-based Model

- The role of the corpus *N*-gram probabilities $\begin{array}{c} P(q_n|Corpus) \\ P(q_n|q_{n-1},Corpus) \end{array}$
  - Model the general distribution of the index terms
    - Help to solve zero-frequency problem $P(q_n|D) = 0$ !
    - Help to differentiate the contributions of different missing terms in a doc
  - The corpus *N*-gram probabilities were estimated using an outside corpus



Doc D

$q_a$ $q_c$ $q_b$ $q_b$
$q_a$ $q_c$ $q_d$ $q_a$
$q_a$ $q_b$

$P(q_a|D)=0.4$
$P(q_b|D)=0.3$
$P(q_c|D)=0.2$
$P(q_d|D)=0.1$
$P(q_e|D)=0.0$
$P(q_f|D)=0.0$

# HMM/N-gram-based Model

- Estimation of *N*-grams (Language Models)
  - Maximum likelihood estimation (MLE) for doc *N*-grams
    - Unigram

      Counts of term $q_i$ in the doc $D$

      $$P(q_i|D) = \frac{C_D(q_i)}{\sum_{q_j \in D} C_D(q_j)} = \frac{C_D(q_i)}{|D|}$$

      Length of the doc $D$

    - Bigram

      Counts of term pair $(q_j, q_i)$ in the doc $D$

      $$P(q_i|q_j, D) = \frac{C_D(q_j, q_i)}{C_D(q_j)}$$

      Counts of term $q_i$ in the doc $D$

  - Similar formulas for corpus N-grams

    $$P(q_i|Corpus) = \frac{C_{Corpus}(q_i)}{|Corpus|} \qquad P(q_i|q_j, D) = \frac{C_{Corpus}(q_j, q_i)}{C_{Corpus}(q_j)}$$

**Corpus**: an outside corpus or just the doc collection

# HMM/N-gram-based Model

- Basically, $m_1$, $m_2$, $m_3$, $m_4$, can be estimated by using the Expectation-Maximization (EM) algorithm
  because of the insufficiency of training data
  - All docs share the same weights here
  - The *N*-gram probability distributions also can be estimated using the EM algorithm instead of the maximum likelihood estimation
- For those docs with training queries, $m_1$, $m_2$, $m_3$, $m_4$, can be estimated by using the Minimum Classification Error (MCE) training algorithm
  - The docs can have different weights

# HMM/N-gram-based Model

- Expectation-Maximum Training
  - The weights are tied among the documents
  - E.g. **$m_1$ of Type I HMM** can be trained using the following equation:

the new weight

819 queries    $\leq$ 2265 docs     the old weight

$$m_1 = \frac{\displaystyle\sum_{Q \in [TrainSet]_Q} \sum_{D \in [Doc]_{R \text{ to } Q}} \sum_{q_n \in Q} \left[ \frac{\hat{m}_1 P(q_n|D)}{\hat{m}_1 P(q_n|D) + \hat{m}_2 P(q_n|Corpus)} \right]}{\displaystyle\sum_{Q \in [TrainSet]_Q} |Q| \cdot \left| [Doc]_{R \text{ to } Q} \right|}$$

- Where $[TrainSet]_Q$ is the set of training query exemplars, $[Doc]_{R \text{ to } Q}$ is the set of docs that are relevant to a specific training query exemplar $Q$, $|Q|$ is the length of the query , and $\left| [Doc]_{R \text{ to } Q} \right|$ is the total number of docs relevant to the query $Q$

# HMM/N-gram-based Model

- ## Expectation-Maximum Training

$$P(Q \mid D) > P(Q \mid \hat{D}) \ ?$$

The new model

The old model

**Empirical Derivation**

$$\log P(Q \mid D) - \log P(Q \mid \hat{D})$$

$$= \sum_{q_n \in Q} \left[ \sum_k P(k \mid q_n, \hat{D}) \log \left[ P(q_n \mid D) \frac{P(q_n, k \mid D)}{P(q_n, k \mid D)} \right] \right] - \sum_{q_n \in Q} \left[ \sum_k P(k \mid q_n, \hat{D}) \log \left[ P(q_n \mid \hat{D}) \frac{P(q_n, k \mid \hat{D})}{P(q_n, k \mid \hat{D})} \right] \right]$$

$$= \sum_{q_n \in Q} \sum_k P(k \mid q_n, \hat{D}) \log \frac{P(q_n, k \mid D)}{P(k \mid q_n, D)} - \sum_{q_n \in Q} \sum_k P(k \mid q_n, \hat{D}) \log \frac{P(q_n, k \mid \hat{D})}{P(k \mid q_n, \hat{D})}$$

$$= \sum_{q_n \in Q} \left[ \sum_k P(k \mid q_n, \hat{D}) \log P(q_n, k \mid D) - \sum_k P(k \mid q_n, \hat{D}) \log P(q_n, k \mid \hat{D}) \right]$$

$$\Phi(D, \hat{D}) - \Phi(\hat{D}, \hat{D})$$

$$+ \sum_{q_n \in Q} \left[ \sum_k P(k \mid q_n, \hat{D}) \log P(k \mid q_n, \hat{D}) - \sum_k P(k \mid q_n, \hat{D}) \log P(k \mid q_n, D) \right]$$
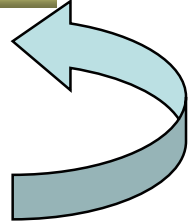
$$H(\hat{D}, \hat{D}) - H(D, \hat{D}) \geq 0$$

$$\geq \sum_{q_n \in Q} \left[ \sum_k P(k \mid q_n, \hat{D}) \log P(q_n, k \mid D) - \sum_k P(k \mid q_n, \hat{D}) \log P(q_n, k \mid \hat{D}) \right]$$

$$\left\{ \because \sum_i p_i \log q_i - \sum_i p_i \log p_i = \sum_i p_i \log \frac{q_i}{p_i} \leq \sum_i p_i \left( \frac{q_i}{p_i} - 1 \right) = 0 \right\}$$

*Jensen's inequality*

$$(\because \log x \leq x - 1)$$

$$\therefore \text{If} \sum_{q_n \in Q} \sum_k P(k \mid q_n, \hat{D}) \log P(q_n, k \mid D) \geq \sum_{q_n \in Q} \sum_k P(k \mid q_n, \hat{D}) \log P(q_n, k \mid \hat{D})$$

$$\text{then} \quad P(Q \mid D) > P(Q \mid \hat{D})$$

# HMM/N-gram-based Model

- ## Expectation-Maximum Training

Q function

$$\Phi(D,\hat{D}) = \sum_{q_n \in Q} \sum_{k} P\left(k|q_n,\hat{D}\right) \log P\left(q_n,k|D\right)$$

$$= \sum_{q_n \in Q} \sum_{k} \frac{P\left(q_n|k,\hat{D}\right)P\left(k|\hat{D}\right)}{P\left(q_n|\hat{D}\right)} \log \left[P\left(q_n|k,D\right)P\left(k|D\right)\right]$$

$$= \sum_{q_n \in Q} \sum_{k} \frac{P\left(q_n|k,\hat{D}\right)\hat{m}_k}{\sum_{j} P\left(q_n|j,\hat{D}\right)\hat{m}_j} \log \left[P\left(q_n|k,D\right)m_k\right]$$

$$\Phi'(D,\hat{D}) = \sum_{q_n \in Q} \sum_{k} \left\{ \frac{P\left(q_n|k,\hat{D}\right)\hat{m}_k}{\sum_{j} P\left(q_n|j,\hat{D}\right)\hat{m}_j} \log \left[P\left(q_n|k,D\right)m_k\right]\right\} + l\left(\sum_{i} m_i - 1\right)$$

normalization constraints using Lagrange multipliers

$$\frac{\partial \Phi'(D,\hat{D})}{\partial m_k} = \frac{1}{m_k}\left[ \sum_{q_n \in Q} \frac{P\left(q_n|k,\hat{D}\right)\hat{m}_k}{\sum_{j} P\left(q_n|j,\hat{D}\right)\hat{m}_j} \right] + l = 0$$

Assume $\quad G_k = \sum_{q_n \in Q} \frac{P\left(q_n|k,\hat{D}\right)\hat{m}_k}{\sum_{j} P\left(q_n|j,\hat{D}\right)\hat{m}_j} \quad \Rightarrow \quad \frac{G_1}{m_1} = \frac{G_2}{m_2} = \dots \frac{G_k}{m_k} = \dots = -l$

$G_k$

$$\therefore \quad l = -\sum_{s} G_s \qquad \therefore \quad m_k = \frac{\displaystyle\sum_{q_n \in Q} \frac{P\left(q_n|k,\hat{D}\right)\hat{m}_k}{\sum_{j} P\left(q_n|j,\hat{D}\right)\hat{m}_j}}{\displaystyle\sum_{s}\sum_{q_n \in Q} \frac{P\left(q_n|k,\hat{D}\right)\hat{m}_s}{\sum_{j} P\left(q_n|j,\hat{D}\right)\hat{m}_j}} = \frac{\displaystyle\sum_{q_n \in Q} \frac{P\left(q_n|k,\hat{D}\right)\hat{m}_k}{\sum_{j} P\left(q_n|j,\hat{D}\right)\hat{m}_j}}{|Q|}$$

$l$

# HMM/N-gram-based Model

- ## Experimental results with EM training
  - ### HMM/N-gram-based approach

| Average Precision | | Word-level | | | Syllable-level | | |
|---|---|---|---|---|---|---|---|
| | | Uni | Uni+Bi | Uni+Bi* | Uni | Uni+Bi | Uni+Bi* |
| TDT2 | TQ/TD | **0.6327** | 0.6069 | 0.5427 | 0.4698 | 0.5220 | 0.5718 |
| | TQ/SD | 0.5658 | **0.5702** | 0.4803 | 0.4411 | 0.5011 | 0.5307 |
| TDT3 | TQ/TD | **0.6569** | 0.6542 | 0.6141 | 0.5343 | 0.5970 | 0.6560 |
| | TQ/SD | 0.6308 | 0.6361 | 0.5808 | 0.5177 | 0.5678 | **0.6433** |

  - ### Vector space model

| Average Precision | | Word-level | | Syllable-level | |
|---|---|---|---|---|---|
| | | $S(N)$, $N=1$ | $S(N)$, $N=1\sim2$ | $S(N)$, $N=1$ | $S(N)$, $N=1\sim2$ |
| TDT2 | TQ/TD | 0.5548 | **0.5623** | 0.3412 | 0.5254 |
| | TQ/SD | 0.5122 | **0.5225** | 0.3306 | 0.5077 |
| TDT3 | TQ/TD | 0.6505 | **0.6531** | 0.3963 | 0.6502 |
| | TQ/SD | 0.6216 | 0.6233 | 0.3708 | **0.6353** |

  - ### HMM/N-gram-based approach is consistently better than vector space model

# Review: The EM Algorithm

- ## Introduction of EM (Expectation Maximization):
  - ### Why EM?
    - Simple optimization algorithms for likelihood function relies on the intermediate variables, called latent (隱藏的)data
      In our case here, **the state sequence** *is the latent data*
    - Direct access to the data necessary to estimate the parameters is impossible or difficult

  - ### Two Major Steps :
    - *E* : expectation with respect to the latent data using the current estimate of the parameters and conditioned on the observations
    - *M*: provides a new estimation of the parameters according to ML (or MAP)

# Review: The EM Algorithm

- The EM Algorithm is important to HMMs and other learning techniques
  - Discover new model parameters to maximize the log-likelihood of incomplete data $\log P(o|\lambda)$ by iteratively maximizing the expectation of log-likelihood from complete data $\log P(o, s|\lambda)$

- Example
  - The observable training data $o$
    - We want to maximize $P(o|\lambda)$, $\lambda$ is a parameter vector
  - The hidden (unobservable) data $s$
    - E.g. the component densities of observable data $o$, or the underlying state sequence in HMMs

$$\Theta\left(\lambda, \overline{\lambda}\right) = \sum_{o} E_{s|o,\lambda} \left[\log P(o, s|\lambda)\right]$$

# HMM/N-gram-based Model

- ## Minimum Classification Error (MCE) Training
  - Given a query $Q$ and a desired relevant doc $D^*$, define **the classification error function** as:

    $$E(Q, D^*) = \frac{1}{|Q|}\left[-\log P\left(Q\middle|D^* \text{ is } R\right) + \max_{D'} \log P\left(Q\middle|D' \text{ is not } R\right)\right]$$

    - \>0 means misclassified; <=0 means a correct decision
  - Transform the error function to the loss function

    $$L(Q, D^*) = \frac{1}{1 + \exp(-\alpha E(Q, D^*) + \beta)}$$

    - In the range between 0 and 1

# HMM/N-gram-based Model

- Minimum Classification Error (MCE) Training
  - Apply the loss function to the MCE procedure for iteratively updating the weighting parameters
    - Constraints:
    $$m_k \geq 0 \ , \quad \sum_k m_k = 1$$
    - Parameter Transforms, (e.g.,Type I HMM)
    $$m_1 = \frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} \quad \text{and} \quad m_2 = \frac{e^{\tilde{m}_2}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}}$$
  - Iteratively update $m_1$ (e.g., Type I HMM)
  $$\tilde{m}_1(i+1) = \tilde{m}_1(i) - \varepsilon(i) \cdot \left. \frac{\partial L(Q, D^*)}{\partial \tilde{m}_1} \right|_{D^* = D^*(i)}$$
    - Where,
    $$\nabla_{D^*, \tilde{m}_1} = \varepsilon(i) \cdot \frac{\partial L(Q, D^*)}{\partial \tilde{m}_1}$$
    $$= \varepsilon(i) \cdot \frac{\partial L(Q, D^*)}{\partial E(Q, D^*)} \cdot \frac{\partial E(Q, D^*)}{\partial \tilde{m}_1}, \quad \frac{\partial L(Q, D^*)}{\partial E(Q, D^*)} = \alpha \cdot L(Q, D^*) \cdot [1 - L(Q, D^*)]$$

# HMM/N-gram-based Model

- Minimum Classification Error (MCE) Training
  - Iteratively update $m_1$ (e.g., Type I HMM)

$$\frac{\partial E(Q, D^*)}{\partial \tilde{m}_1} = \frac{-1}{|Q|} \frac{\partial \left\{ \sum\limits_{q_n \in Q} \log \left[ \frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | D^*) + \frac{e^{\tilde{m}_2}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | Corpus) \right] \right\}}{\partial \tilde{m}_1}$$

$$= \frac{-1}{|Q|} \sum\limits_{q_n \in Q} \left\{ \frac{\frac{-e^{\tilde{m}_1}}{(e^{\tilde{m}_1} + e^{\tilde{m}_2})^2} \left[ e^{\tilde{m}_1} P(q_n | D^*) + e^{\tilde{m}_2} P(q_n | Corpus) \right] + \frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | D^*)}{\frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | D^*) + \frac{e^{\tilde{m}_2}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | Corpus)} \right\}$$

$$= \frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} - \frac{1}{|Q|} \sum\limits_{q_n \in Q} \left\{ \frac{\frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | D^*)}{\frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | D^*) + \frac{e^{\tilde{m}_2}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | Corpus)} \right\}$$

$$= -\left[ -m_1 + \frac{1}{|Q|} \sum\limits_{q_n \in Q} \frac{m_1 P(q_n | D^*)}{m_1 P(q_n | D^*) + m_2 P(q_n | Corpus)} \right],$$

# HMM/N-gram-based Model

- ## Minimum Classification Error (MCE) Training
  - Iteratively update $m_1$ (e.g., Type I HMM)

$$\nabla_{D^*,\tilde{m}_1}(i) = -\varepsilon(i) \cdot \alpha \cdot L(Q,D^*) \cdot [1 - L(Q,D^*)]$$

$$\cdot \left[ -m_1(i) + \frac{1}{|Q|} \sum_{q_n \in Q} \frac{m_1(i)P(q_n|D^*)}{m_1(i)P(q_n|D^*) + m_2(i)P(q_n|Corpus)} \right],$$

the new weight

$$m_1(i+1) = \frac{e^{\tilde{m}_1(i+1)}}{e^{\tilde{m}_1(i+1)} + e^{\tilde{m}_2(i+1)}}$$

$$\tilde{m}_1(i+1) = \tilde{m}_1(i) - \nabla_{D^*,\tilde{m}_1}(i)$$

$$= \frac{e^{\tilde{m}_1(i)}e^{-\nabla_{D^*,\tilde{m}_1}(i)}}{e^{\tilde{m}_1(i)}e^{-\nabla_{D^*,\tilde{m}_1}(i)} + e^{\tilde{m}_2(i)}e^{-\nabla_{D^*,\tilde{m}_2}(i)}}$$

$$= \frac{e^{\tilde{m}_1(i)}e^{-\nabla_{D^*,\tilde{m}_1}(i)} \big/ \left( e^{\tilde{m}_1(i)} + e^{\tilde{m}_2(i)} \right)}{\left[ e^{\tilde{m}_1(i)}e^{-\nabla_{D^*,\tilde{m}_1}(i)} \big/ \left( e^{\tilde{m}_1(i)} + e^{\tilde{m}_2(i)} \right) \right] + \left[ e^{\tilde{m}_2(i)}e^{-\nabla_{D^*,\tilde{m}_2}(i)} \big/ \left( e^{\tilde{m}_1(i)} + e^{\tilde{m}_2(i)} \right) \right]}$$

the old weight

$$= \frac{m_1(i) \cdot e^{-\nabla_{D^*,\tilde{m}_1}(i)}}{m_1(i) \cdot e^{-\nabla_{D^*,\tilde{m}_1}(i)} + m_2(i) \cdot e^{-\nabla_{D^*,\tilde{m}_2}(i)}},$$

# HMM/N-gram-based Model

- ## Minimum Classification Error (MCE) Training
  - ### Final Equations
    - Iteratively update $m_1$

$$\nabla_{D^*,\tilde{m}_1}(i) = -\varepsilon(i)\cdot\alpha\cdot L(Q,D^*)\cdot\left[1 - L(Q,D^*)\right]$$

$$\cdot\left[-m_1(i) + \frac{1}{|Q|}\sum_{q_n \in Q}\frac{m_1(i)P(q_n|D^*)}{m_1(i)P(q_n|D^*) + m_2(i)P(q_n|Corpus)}\right]$$

$$m_1(i+1) = \frac{m_1(i)\cdot e^{-\nabla_{D^*,\tilde{m}_1}(i)}}{m_1(i)\cdot e^{-\nabla_{D^*,\tilde{m}_1}(i)} + m_2(i)\cdot e^{-\nabla_{D^*,\tilde{m}_2}(i)}}$$
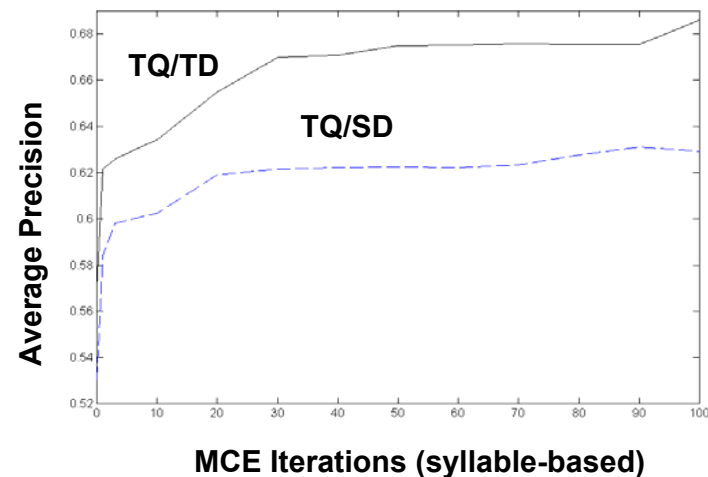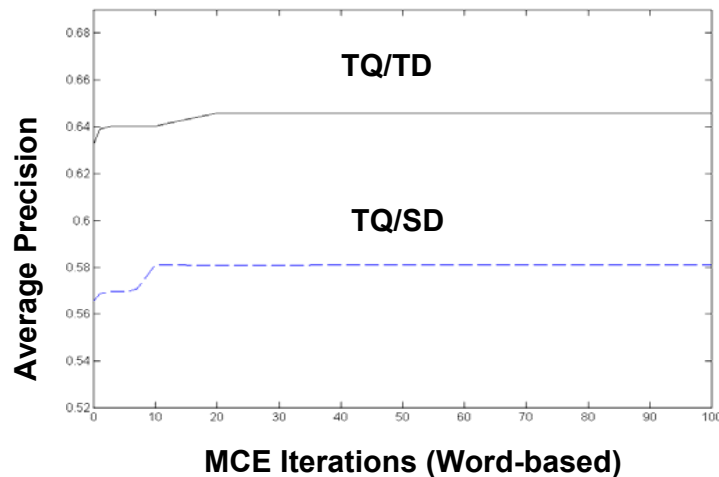
  - $m_2$ can be updated in the similar way

# HMM/N-gram-based Model

- ## Experimental results with MCE training

| Average Precision | | Word-level | Syllable-level | Fusion |
|---|---|---|---|---|
| | | Uni | Uni+Bi* | |
| TDT2 | TQ/TD | 0.6459 (0.6327) | 0.6858 (0.5718) | 0.7329 |
| | TQ/SD | 0.5810 (0.5658) | 0.6300 (0.5307) | 0.6914 |

**Before MCE Training**

Iterations=100



**MCE Iterations (Word-based)**



**MCE Iterations (syllable-based)**

– The results for the syllable-level index features were significantly improved

20

# HMM/N-gram-based Model

- Advantages
  - A formal mathematic framework
  - Use collection statistics but not heuristics
  - The retrieval system can be gradually improved through usage

- Disadvantages
  - Only literal term matching (or word overlap measure)
    - The issue of *relevance* or *aboutness* is not taken into consideration
  - The implementation relevance feedback or query expansion is not straightforward
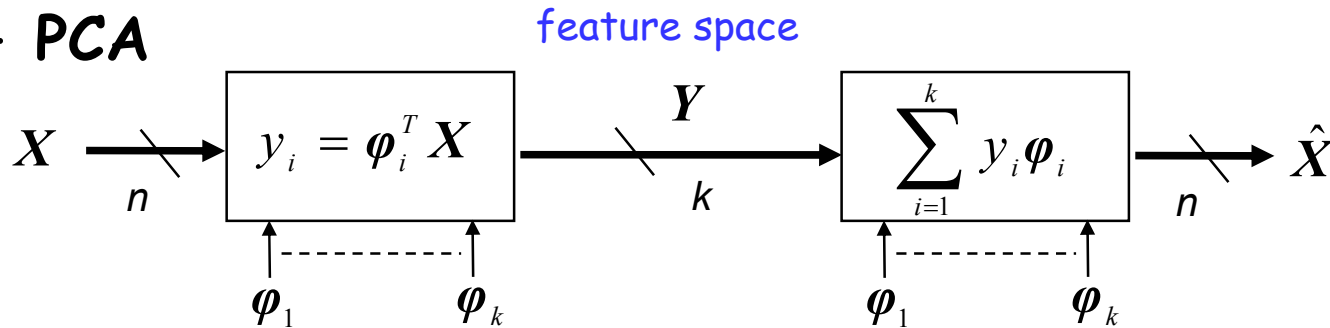
# Latent Semantic Indexing (LSI)

- LSI: a technique projects queries and docs into a space with "latent" semantic dimensions
  - Co-occurring terms are projected onto the same dimensions
  - In the latent semantic space (with fewer dimensions), a query and doc can have high cosine similarity even if they do not share any terms
  - Dimensions of the reduced space correspond to the axes of greatest variation
    - Closely related to Principal Component Analysis (PCA)

# Latent Semantic Indexing (LSI)
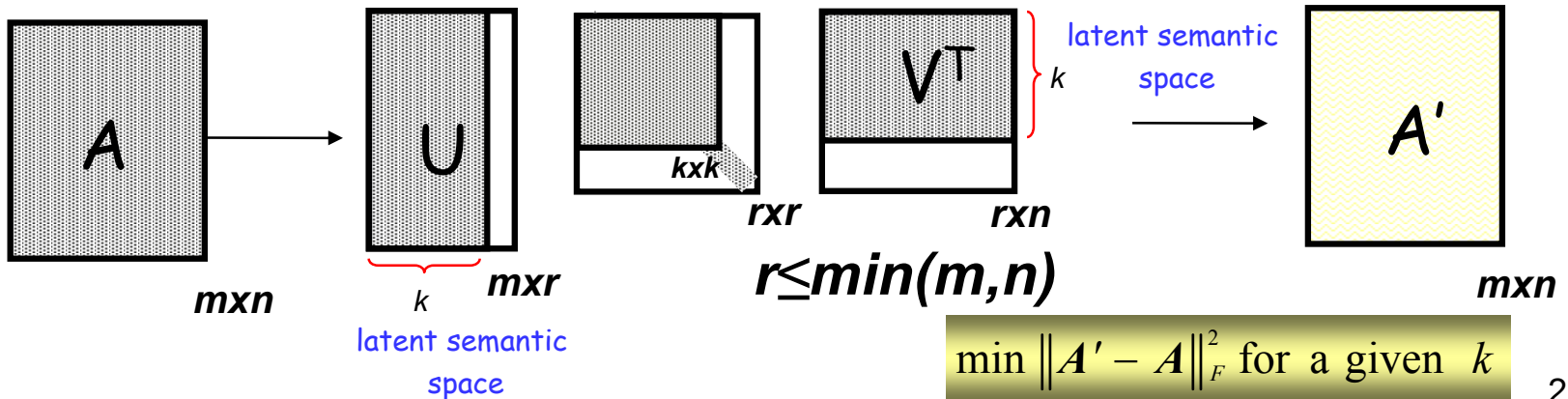
- Dimension Reduction and Feature Extraction
  - **PCA**

feature space

$$X \xrightarrow{n} \boxed{y_i = \boldsymbol{\varphi}_i^T X} \xrightarrow{Y \ \ k} \boxed{\sum_{i=1}^{k} y_i \boldsymbol{\varphi}_i} \xrightarrow{n} \hat{X}$$

$\boldsymbol{\varphi}_1 \qquad \boldsymbol{\varphi}_k \qquad\qquad \boldsymbol{\varphi}_1 \qquad \boldsymbol{\varphi}_k$

orthonormal basis

$$\min \left\| \hat{X} - X \right\|^2 \ \text{for a given} \ k$$

  - **SVD (in LSI)**



$A$ (mxn) → $U$ (mxr), $k$ latent semantic space | kxk rxr | $V^T$ (rxn) $k$ latent semantic space → $A'$ (mxn)

$$r \le min(m,n)$$

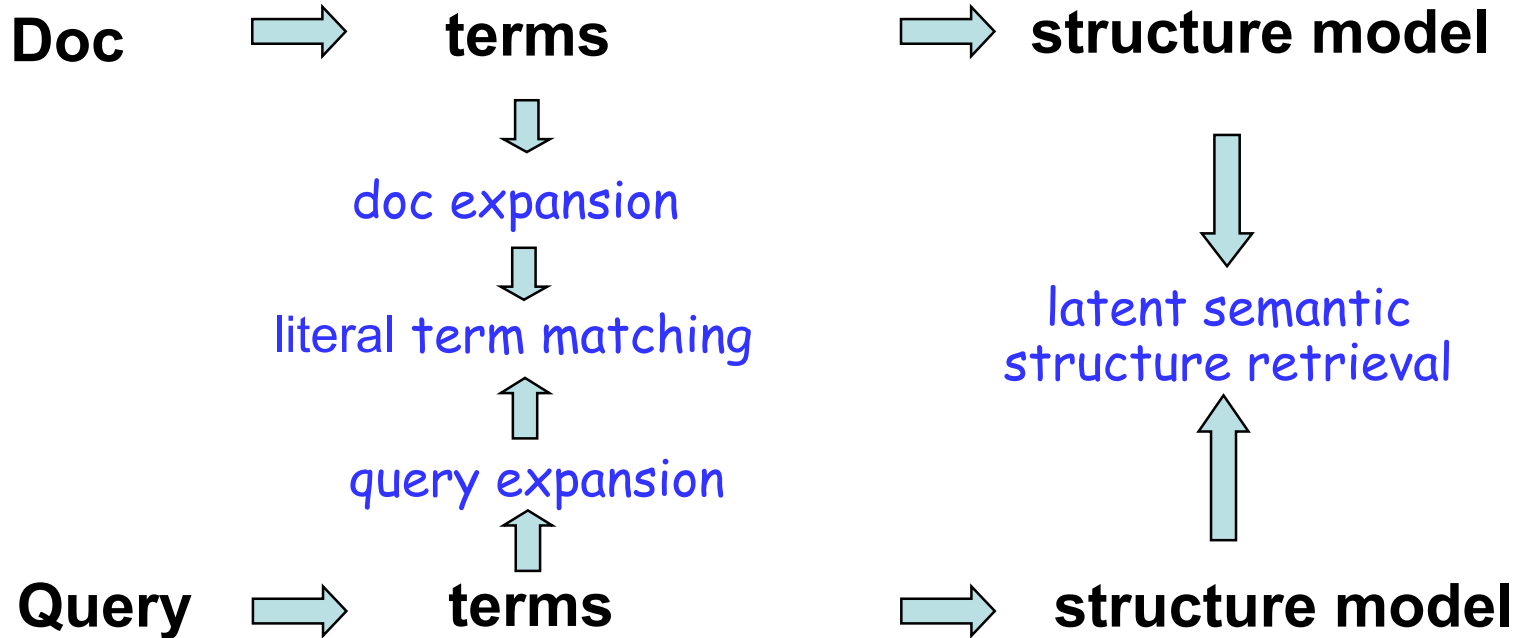$$\min \left\| A' - A \right\|_F^2 \ \text{for a given} \ k$$

23

# Latent Semantic Indexing (LSI)

– Singular Value Decomposition (SVD) used for the word-document matrix

  • A least-squares method for dimension reduction

|  | Term 1 | Term 2 | Term 3 | Term 4 |
|---|---|---|---|---|
| Query | user | interface | | |
| Document 1 | user | interface | HCI | interaction |
| Document 2 | | | HCI | interaction |

# Latent Semantic Indexing (LSI)

- Frameworks to circumvent vocabulary mismatch

**Doc** ⟹ **terms** ⟹ **structure model**

⟱

doc expansion

⟱

literal term matching       latent semantic structure retrieval

⟰

query expansion

⟰

**Query** ⟹ **terms**      **structure model**

# Latent Semantic Indexing (LSI)

## Titles

c1: Human machine interface for Lab ABC computer applications
c2: A survey of user opinion of computer system response time
c3: The EPS user interface management system
c4: System and human system engineering testing of EPS
c5: Relation of user-perceived response time to error measurement

m1: The generation of random, binary, unordered trees
m2: The intersection graph of paths in trees
m3: Graph minors IV: Widths of trees and well-quasi-ordering
m4: Graph minors: A survey

Terms | Documents

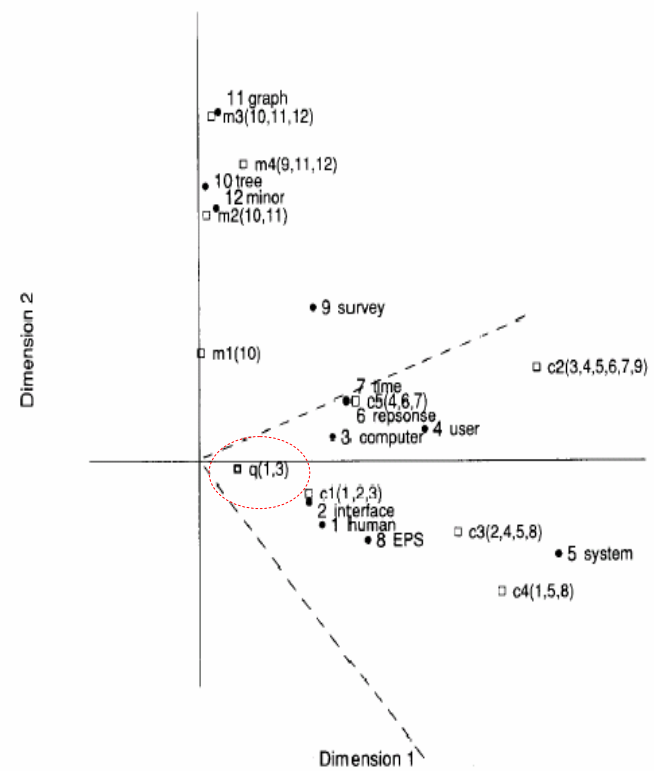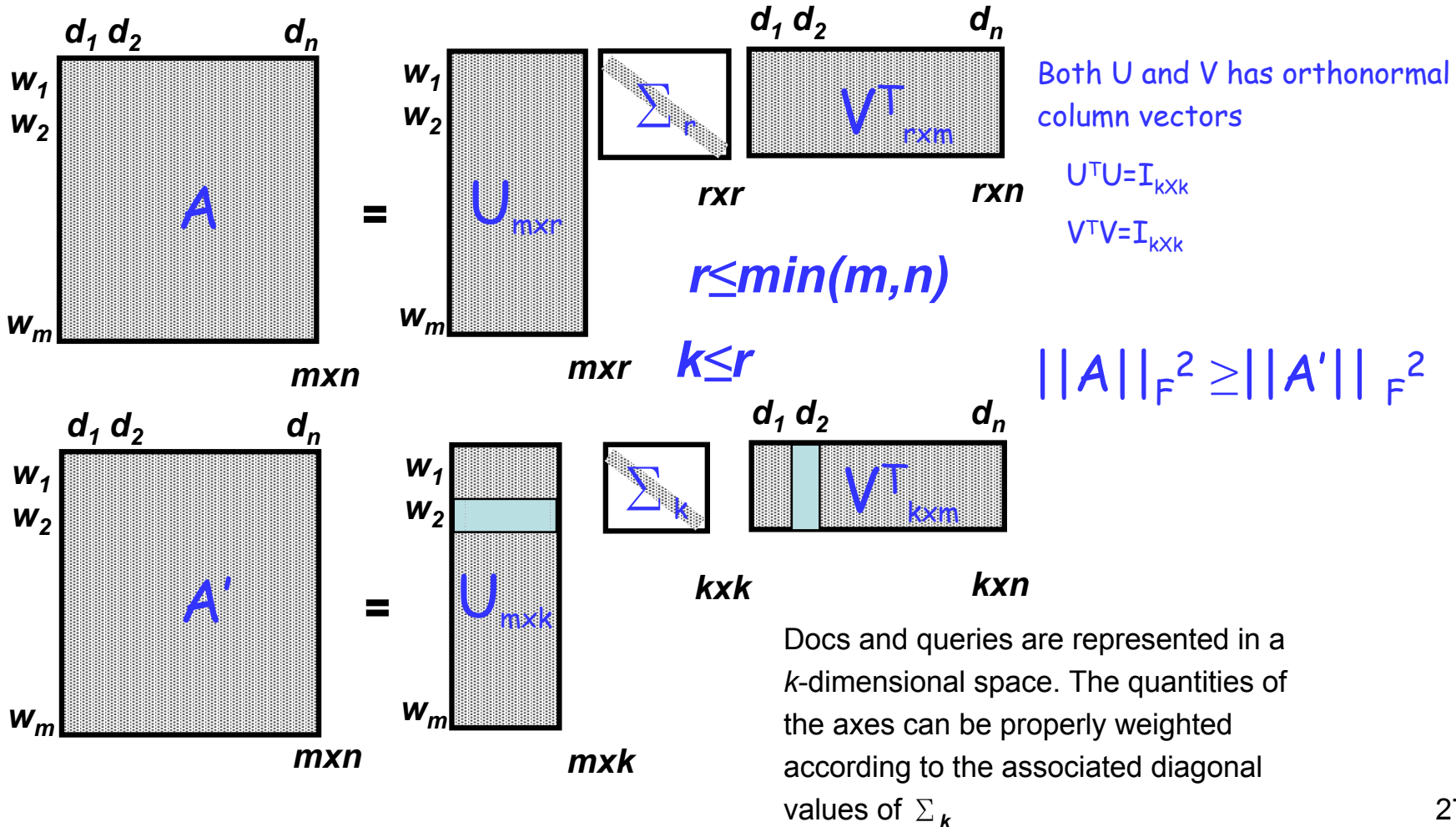| | c1 | c2 | c3 | c4 | c5 | m1 | m2 | m3 | m4 |
|---|---|---|---|---|---|---|---|---|---|
| human | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| interface | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| computer | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| user | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| system | 0 | 1 | 1 | 2 | 0 | 0 | 0 | 0 | 0 |
| response | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| time | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| EPS | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| survey | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| trees | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 |
| graph | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| minors | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |

2-D Plot of Terms and Docs from Example



FIG. 1. A two-dimensional plot of 12 Terms and 9 Documents from the sampe TM set. Terms are represented by filled circles. Documents are shown as open squares, and component terms are indicated parenthetically. The query ("human computer interaction") is represented as a pseudo-document at point q. Axes are scaled for Document-Document or Term-Term comparisons. The dotted cone represents the region whose points are within a cosine of .9 from the query q. All documents about human-computer (c1–c5) are "near" the query (i.e., within this cone), but none of the graph theory documents (m1=m4) are nearby. In this reduced space, even documents c3 and c5 which share no terms with the query are near it.

# Latent Semantic Indexing (LSI)

- ## Singular Value Decomposition (SVD)



$d_1\ d_2 \quad d_n$

$w_1$
$w_2$

$A$

$w_m$

$mxn$

=

$U_{mxr}$

$mxr$

$\Sigma_r$

$rxr$

$d_1\ d_2 \quad d_n$

$V^T_{rxm}$

$rxn$

Both U and V has orthonormal column vectors

$U^TU=I_{kXk}$

$V^TV=I_{kXk}$

$r \leq min(m,n)$

$k \leq r$

$||A||_F^2 \geq ||A'||_F^2$

$d_1\ d_2 \quad d_n$

$w_1$
$w_2$

$A'$

$w_m$

$mxn$

=

$U_{mxk}$

$mxk$

$\Sigma_k$

$kxk$

$d_1\ d_2 \quad d_n$

$V^T_{kxm}$

$kxn$

Docs and queries are represented in a $k$-dimensional space. The quantities of the axes can be properly weighted according to the associated diagonal values of $\Sigma_k$

27

# Latent Semantic Indexing (LSI)

- ## Singular Value Decomposition (SVD)
  - $A^T A$ is symmetric $n$x$n$ matrix
    - All eigenvalues $\lambda_j$ are nonnegative real numbers

      $$\lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_n \geq 0 \quad \Sigma^2 = diag(\lambda_1, \lambda_1, \ldots, \lambda_n)$$

    - All eigenvectors $v_j$ are orthonormal

      $$V = [v_1 v_2 \ldots v_n] \quad v_j^T v_j = 1 \quad (V^T V = I_{nxn})$$

    - Define **singular values** $\sigma_j = \sqrt{\lambda_j}, \ j = 1, \ldots, n$
      - As the square roots of the eigenvalues of $A^T A$
      - As the lengths of the vectors $Av_1, Av_2, \ldots, Av_n$

        *For $\lambda_i \neq 0, \ i=1,\ldots r,$*
        *{$Av_1, Av_2, \ldots, Av_r$} is an*
        *orthogonal basis of Col A*

        $$\sigma_1 = \|Av_1\|$$
        $$\sigma_2 = \|Av_2\|$$
        $$\ldots$$

# Latent Semantic Indexing (LSI)

- $\{Av_1, Av_2, \ldots, Av_r\}$ is an orthogonal basis of Col A

$$Av_i \bullet Av_j = (Av_i)^T Av_j = v_i^T A^T Av_j = \lambda_j v_i^T v_j = 0$$

- Suppose that $A$ (or $A^T A$) has rank $r \leq n$

$$\lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_r > 0, \quad \lambda_{r+1} = \lambda_{r+2} = \ldots = \lambda_r = 0$$

- Define an orthonormal basis $\{u_1, u_2, \ldots, u_r\}$ for Col A

$$u_i = \frac{1}{\|Av_i\|} Av_i = \frac{1}{\sigma_i} Av_i \Rightarrow \sigma_i u_i = Av_i$$

$$\Rightarrow [u_1 \, u_2 \ldots u_r] \Sigma_r = A[v_1 \, v_2 \quad v_r]$$

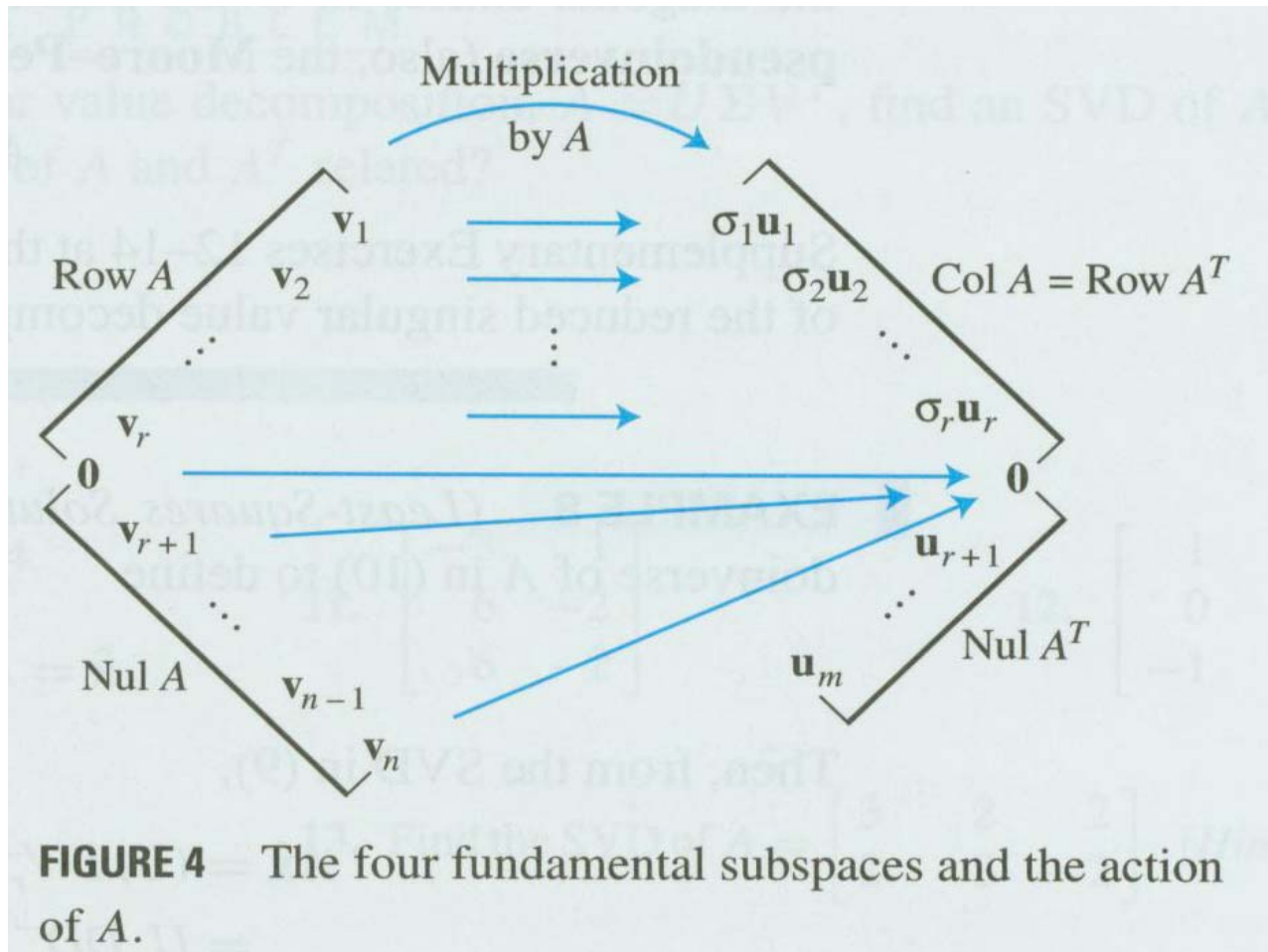  - Extend to an orthonormal basis $\{u_1, u_2, \ldots, u_m\}$ of $R^m$

$$\Rightarrow [u_1 \, u_2 \ldots u_r \ldots u_m] \Sigma = A[v_1 \, v_2 \ldots v_r \ldots v_m] \qquad \|A\|_F^2 = \sum_{i=1}^{m} \sum_{j=1}^{n} a_{ij}^2$$

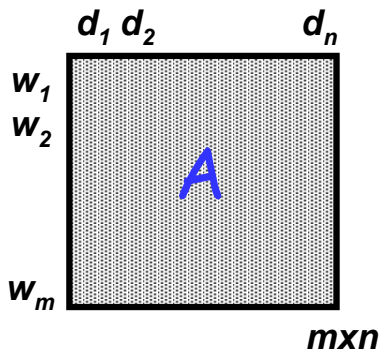$$\Rightarrow U\Sigma = AV$$

$$\Rightarrow A = U\Sigma V^T \qquad\qquad\qquad\qquad\qquad \|A\|_F^2 = \sigma_1^2 + \sigma_2^2 + \ldots + \sigma_r^2 \quad ?$$

# Latent Semantic Indexing (LSI)



**FIGURE 4**  The four fundamental subspaces and the action of $A$.

# Latent Semantic Indexing (LSI)

- ## Fundamental comparisons based on SVD

  - ### The original word-document matrix (A)

$d_1\ d_2 \qquad\qquad d_n$

$w_1$
$w_2$

$A$

$w_m$

$mxn$

  - compare two terms → dot product of two rows of A
    - or an entry in $AA^\top$
  - compare two docs → dot product of two columns of A
    - or an entry in $A^\top A$
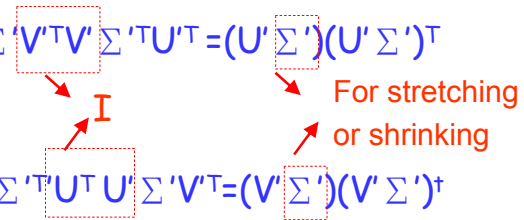  - compare a term and a doc → each individual entry of A

  - ### The new word-document matrix (A')

$U'=U_{mxk}$
$\Sigma'=\Sigma_k$
$V'=V_{nxk}$

  - compare two terms
    → dot product of two rows of $U'\Sigma'$

$A'A'^\top=(U'\Sigma'V'^\top)\ (U'\Sigma'V'^\top)^\top=U'\Sigma'V'^\top V'\Sigma'^\top U'^\top=(U'\Sigma')(U'\Sigma')^\top$

$I$

For stretching or shrinking

  - compare two docs
    → dot product of two rows of $V'\Sigma'$

$A'^\top A'=(U'\Sigma'V'^\top)^\top{}'(U'\Sigma'V'^\top)=V'\Sigma'^\top U'^\top U'\Sigma'V'^\top=(V'\Sigma')(V'\Sigma')^\dagger$

  - compare a query and a doc → each individual entry of A'

31

# Latent Semantic Indexing (LSI)

- **Fold-in**: find representations for pesudo-docs $q$
  - For objects (new queries or docs) that did not appear in the original analysis
    - Fold-in a new $m$x1 query (or doc) vector

$$\hat{q}_{1xk} = \left(q^T\right)_{1\,xm} U_{m \times k} \Sigma^{-1}_{k \times k}$$

Just like a row of V    Query represented by the weighted sum of it constituent term vectors    The separate dimensions are differentially weighted

  - Cosine measure between the query and doc vectors in the latent semantic space

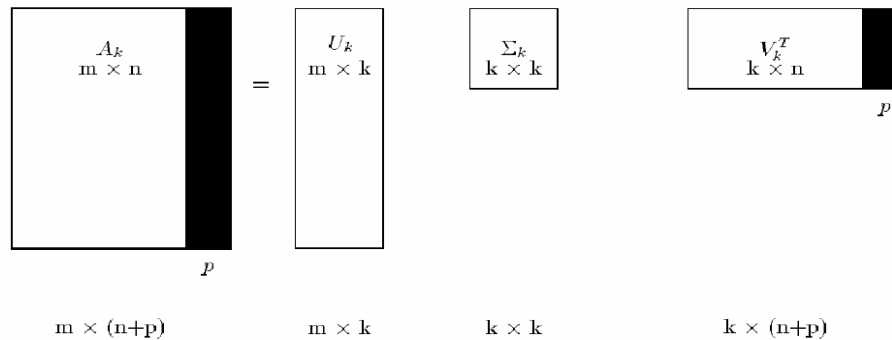$$sim\left(\hat{q}, \hat{d}\right) = coine\ (\hat{q}\Sigma, \hat{d}\Sigma) = \frac{\hat{q}\Sigma^2 \hat{d}^T}{|\hat{q}\Sigma| |\hat{d}\Sigma|}$$

row vectors

# Latent Semantic Indexing (LSI)

- Fold-in a new $1 \times n$ term vector

$$\hat{t}_{1\ xk} \ = \ t_{1\ xn} \ V_{\ n \times k} \ \Sigma_{\ k \times k}^{\ -1}$$



Mathematical representation of folding-in p documents.

Mathematical representation of folding-in q terms.

# Latent Semantic Indexing (LSI)

- Experimental results
  - HMM is consistently better than VSM at all recall levels
  - LSI is better than VSM at higher recall levels



Recall-Precision curve at 11 standard recall levels evaluated on
TDT-3 SD collection. (Using word-level indexing terms)

# Latent Semantic Indexing (LSI)

- Advantages
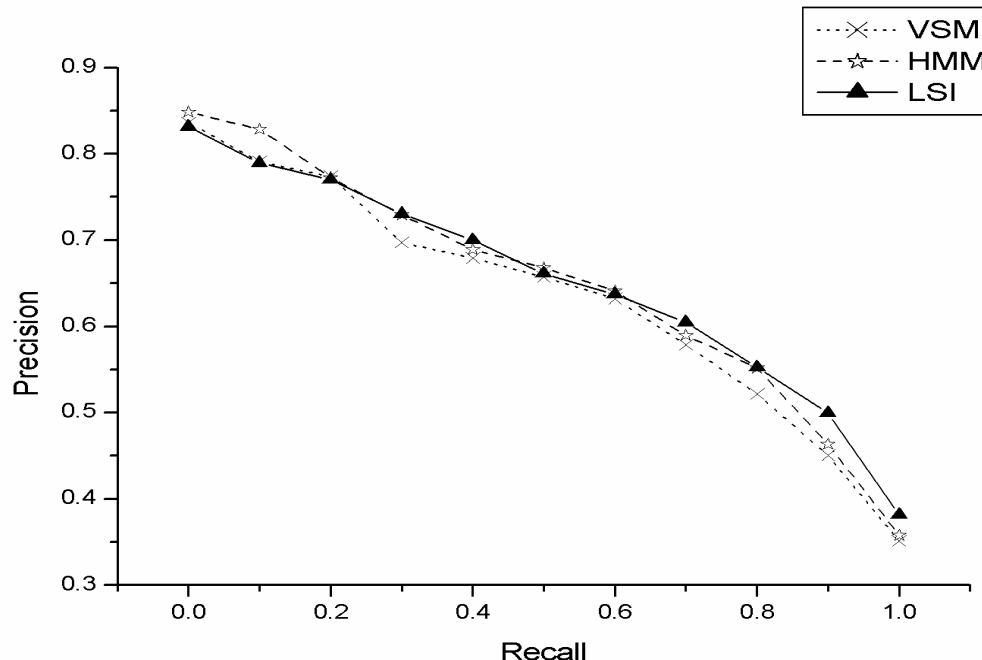  - A clean formal framework and a clearly defined optimization criterion (least-squares)
    - Conceptual simplicity and clarity
  - Handle synonymy problems ("heterogeneous vocabulary")
  - Good results for high-recall search
    - Take term co-occurrence into account
- Disadvantages
  - High computational complexity
  - LSI offers only a partial solution to polysemy
    - E.g. bank, bass,…

# Probabilistic Latent Semantic Analysis (PLSA)

- Also called The Aspect Model, Probabilistic Latent Semantic Indexing (PLSA)
  - Can be viewed as a complex HMM Model



The latent variables
=>The unobservable class variables $T_i$ (topics or domains)

$$Q = w_1 w_2 .. w_j .. w_J$$

$$P(T_1|D_i) \quad P(w_j|T_1)$$
$$P(T_2|D_i) \quad P(w_j|T_2)$$
$$P(T_K|D_i) \quad P(w_j|T_K)$$

$$\Sigma$$

$$sim(Q,D_i) = P(Q|D_i) = \prod_{w_j} P(w_j|D_i)$$

$$= \prod_{w_j} \left[ \sum_{k=1}^{K} P(w_j, T_k|D_i) \right] \quad ?$$

$$= \prod_{w_j} \left[ \sum_{k=1}^{K} P(w_j|T_k) P(T_k|D_i) \right]$$

$$sim(Q,D_i) = P(D_i|Q) = \frac{P(Q,D_i)}{P(Q)} \approx P(Q,D_i) = P(Q|D_i)P(D_i) \approx P(Q|D_i)$$

$$\Rightarrow sim(Q,D_i) \approx P(Q|D_i)$$

36

# Probabilistic Latent Semantic Analysis (PLSA)

- Definition
  - $P(D_i)$: the prob. when selecting a doc $D_i$
  - $P(T_k|D_i)$: the prob. when pick a latent class $T_k$ for the doc $D_i$
  - $P(w_j|T_k)$ : the prob. when generating a word $w_j$ from the class $T_k$

# Probabilistic Latent Semantic Analysis (PLSA)

- Assumptions
  - **Bag-of-words**: treat docs as *memoryless* source, words are generated independently

  - **Conditional independent**: the doc $D_i$ and word $w_j$ are independent conditioned on the state of the associated latent variable $T_k$

$$P\left(w_j, D_i \middle| T_k\right) \approx P\left(w_j \middle| T_k\right) P\left(D_i \middle| T_k\right)$$

$$P\left(w_j \middle| D_i\right) = \sum_{k=1}^{K} P\left(w_j, T_k \middle| D_i\right) = \sum_{k=1}^{K} \frac{P\left(w_j, D_i, T_k\right)}{P\left(D_i\right)} = \sum_{k=1}^{K} \frac{P\left(w_j, D_i \middle| T_k\right) P\left(T_k\right)}{P\left(D_i\right)}$$

$$= \sum_{k=1}^{K} \frac{P\left(w_j \middle| T_k\right) P\left(D_i \middle| T_k\right) P\left(T_k\right)}{P\left(D_i\right)} = \sum_{k=1}^{K} \frac{P\left(w_j \middle| T_k\right) P\left(T_k, D_i\right)}{P\left(D_i\right)}$$

$$= \sum_{k=1}^{K} P\left(w_j \middle| T_k\right) P\left(T_k \middle| D_i\right)$$

# Probabilistic Latent Semantic Analysis (PLSA)

- Probability estimation using EM (expectation-maximization) algorithm
  - **E** (expectation) step

take expectation

$$E\left[L^C\right] = \sum_{D_i} \sum_{w_j} n\left(w_j, D_i\right) E_{T_k \mid w_j, D_i}\left[\log P\left(w_j, T_k \mid D_i\right)\right]$$

complete data likelihood

$$= \sum_{D_i} \sum_{w_j} n\left(w_j, D_i\right) \sum_{T_k} \left[\hat{P}\left(T_k \mid w_j, D_i\right) \log P\left(w_j, T_k \mid D_i\right)\right]$$

empirical distribution          the model

$$\hat{P}\left(T_k \mid w_j, D_i\right) = \frac{\hat{P}\left(T_k, w_j \mid D_i\right)}{\hat{P}\left(w_j \mid D_i\right)} = \frac{\hat{P}\left(w_j \mid T_k\right)\hat{P}\left(T_k \mid D_i\right)}{\sum_{T_k} \hat{P}\left(w_j \mid T_k\right)\hat{P}\left(T_k \mid D_i\right)}$$

$$= \sum_{D_i} \sum_{w_j} n\left(w_j, D_i\right) \sum_{T_k} \left[\hat{P}\left(T_k \mid w_j, D_i\right) \log P\left(w_j \mid T_k\right)P\left(T_k \mid D_i\right)\right]$$

Kullback-Leibler divergence

$$= \sum_{D_i} \sum_{w_j} n\left(w_j, D_i\right) \sum_{T_k} \left[\frac{\hat{P}\left(w_j \mid T_k\right)\hat{P}\left(T_k \mid D_i\right)}{\sum_{T_k} \hat{P}\left(w_j \mid T_k\right)\hat{P}\left(T_k \mid D_i\right)} \log P\left(w_j \mid T_k\right)P\left(T_k \mid D_i\right)\right]$$

# Probabilistic Latent Semantic Analysis (PLSA)

- Probability estimation using EM
  - **M** (maximization) step

$$Q = E\left[L^C\right] + \sum_{T_k} \tau_k \left(1 - \sum_{w_j} P\!\left(w_j \middle| T_k\right)\right) + \sum_{D_i} \rho_i \left(1 - \sum_{T_k} P\!\left(T_k \middle| D_i\right)\right)$$

normalization constraints using Lagrange multipliers

$$Q_{P\left(w_j \middle| T_k\right)} = \sum_{D_i} \sum_{w_j} n\!\left(w_j, D_i\right) \hat{P}\!\left(T_k \middle| w_j, D_i\right) \log P\!\left(w_j \middle| T_k\right) + \tau_k \left(1 - \sum_{w_j} P\!\left(w_j \middle| T_k\right)\right)$$

$$Q_{P\left(T_k \middle| D_j\right)} = \sum_{w_j} n\!\left(w_j, D_i\right) \sum_{T_k} \hat{P}\!\left(T_k \middle| w_j, D_i\right) \log P\!\left(T_k \middle| D_i\right) + \rho_j \left(1 - \sum_{T_k} P\!\left(T_k \middle| D_i\right)\right)$$

# Probabilistic Latent Semantic Analysis (PLSA)

- Probability estimation using EM
  - **M** (maximization) step
    - Take differentiation

$$P(w_j|T_k) = \frac{\sum_{D_i} n(w_j, D_i)\hat{P}(T_k|w_j, D_i)}{\sum_{w_j}\sum_{D_i} n(w_j, D_i)\hat{P}(T_k|w_j, D_i)}$$

$$P(T_k|D_j) = \frac{\sum_{w_j} n(w_j, D_i)\hat{P}(T_k|w_j, D_i)}{\sum_{T_k}\sum_{w_j} n(w_j, D_i)\hat{P}(T_k|w_j, D_i)} = \frac{\sum_{w_j} n(w_j, D_i)\hat{P}(T_k|w_j, D_i)}{\sum_{w_j} n(w_j, D_i)}$$
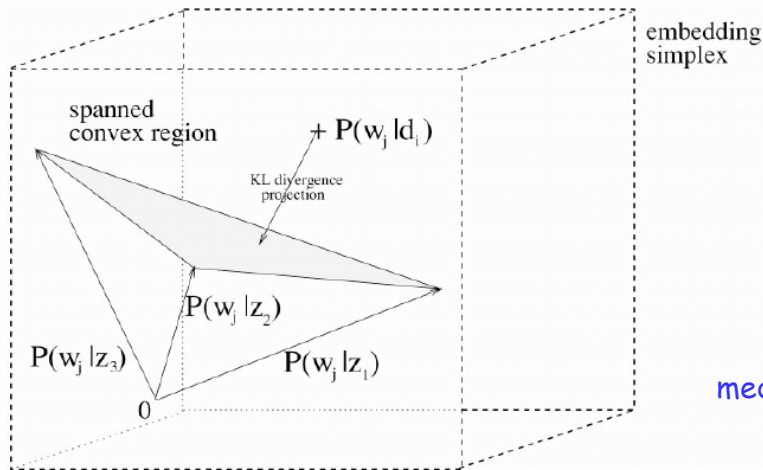
$$= \frac{\sum_{w_j} n(w_j, D_i)\hat{P}(T_k|w_j, D_i)}{n(D_i)}$$

The training formula

41

# Probabilistic Latent Semantic Analysis (PLSA)

- ## Latent Probability Spaces

<span style="color:blue">Dimensionality *K*=128 (latent classes)</span>



| Aspect 1 | Aspect 2 | Aspect 3 | Aspect 4 |
|----------|----------|----------|----------|
| imag | video | region | speaker |
| SEGMENT | sequenc | contour | speech |
| textur | motion | boundari | recogni |
| color | frame | descrip | signal |
| tissu | scene | imag | train |
| brain | SEGMENT | SEGMENT | hmm |
| slice | shot | precis | sourc |
| cluster | imag | estim | speakerindepend |
| mri | cluster | pixel | SEGMENT |
| algorithm | visual | paramet | sound |

<span style="color:blue">medical imaging    image sequence analysis    context of contour boundary detection    phonetic segmentation</span>

Sketch of the probability simplex and a convex region spanned by class-conditional probabilities in the aspect model.

$$P\left(w_j, D_i\right) = \sum_{T_k} P\left(w_j, T_k, D_i\right) = \sum_{T_k} P\left(w_j \middle| T_k, D_i\right) P\left(T_k, D_i\right)$$

$$= \sum_{T_k} P\left(w_j \middle| T_k\right) P\left(T_k\right) P\left(D_i \middle| T_k\right)$$

$$P\left(W, D\right) = \hat{U} : \left(P\left(w_j \middle| T_k\right)\right)_{j,k} \cdot \hat{\Sigma} : \mathrm{diag}\left(P\left(T_k\right)\right)_k \cdot \hat{V} : \left(P\left(D_i \middle| T_k\right)\right)_{i,k}$$

42

# Probabilistic Latent Semantic Analysis (PLSA)

- One more example on TDT1 dataset

| aviation | space missions | family love | Hollywood love |
| --- | --- | --- | --- |
| Aspect 1 | Aspect 2 | Aspect 3 | Aspect 4 |
| plane | space | home | film |
| airport | shuttle | family | movie |
| crash | mission | like | music |
| flight | astronauts | love | new |
| safety | launch | kids | best |
| aircraft | station | mother | hollywood |
| air | crew | life | love |
| passenger | nasa | happy | actor |
| board | satellite | friends | entertainment |
| airline | earth | cnn | star |

The 2 aspects to most likely generate the word 'flight' (left) and 'love' (right), derived from a $K = 128$ aspect model of the TDT1 document collection. The displayed terms are the most probable words in the class-conditional distribution $P(w_j | z_k)$, from top to bottom in descending order.

# Probabilistic Latent Semantic Analysis (PLSA)

- Comparison with LSI
  - Decomposition/Approximation
    - **LSI**: least-squares criterion measured on the L2- or Frobenius norms of the word-doc matrices
    - **PLSA**: maximization of the likelihoods functions based on the cross entropy or Kullback-Leibler divergence between the empirical distribution and the model
  - Computational complexity
    - LSI: SVD decomposition
    - PLSA: EM training, is time-consuming for iterations ?

# Probabilistic Latent Semantic Analysis (PLSA)

- Experimental Results
  **PLSI-U\***
  - Two ways to smoothen empirical distribution with PLSI
    - Combine the cosine score with that of the vector space model (so does LSI)
    - Combine the multinomials individually

$$P_{PLSI-Q*}(\omega_j \mid d_i) = \lambda P_{Empirical}(\omega_j \mid d_i) + (1-\lambda)P_{PLSA}(\omega_j \mid d_i)$$

$$P_{Empirical}(\omega_j \mid d_i) = \frac{n(w_j, d_i)}{n(d_i)}$$

  Both provide almost identical performance

  - It's not known if PLSA was used alone

# Probabilistic Latent Semantic Analysis (PLSA)

- Experimental Results

  **PLSI-Q\***

  - Use the low-dimensional representation $P(T_k | Q)$ and $P(T_k | D_i)$ (be viewed in a *k*-dimensional latent space) to evaluate relevance by means of cosine measure

  - Combine the cosine score with that of the vector space model

  - Use the ad hoc approach to reweight the different model components (dimensions) by

$$RW(T_k) = \sum_{w_j} \left[ P(w_j | T_k) \cdot idf(w_j) \right]$$

$$sim(Q, D) = \frac{\sum_{w_j \in Q} \left[ n(q, w_j) \sum_{T_k} RW^2(T_k) P(T_k | w_j) P(T_k | D_i) \right]}{\sqrt{\sum_{w_j \in Q} \left[ n(q, w_j) \sum_{T_k} RW^2(T_k) P^2(T_k | w_j) \right]} \sqrt{\sum_{T_k} RW^2(T_k) P^2(T_k | D_i)}}$$

# Probabilistic Latent Semantic Analysis (PLSA)

- Experimental Results