

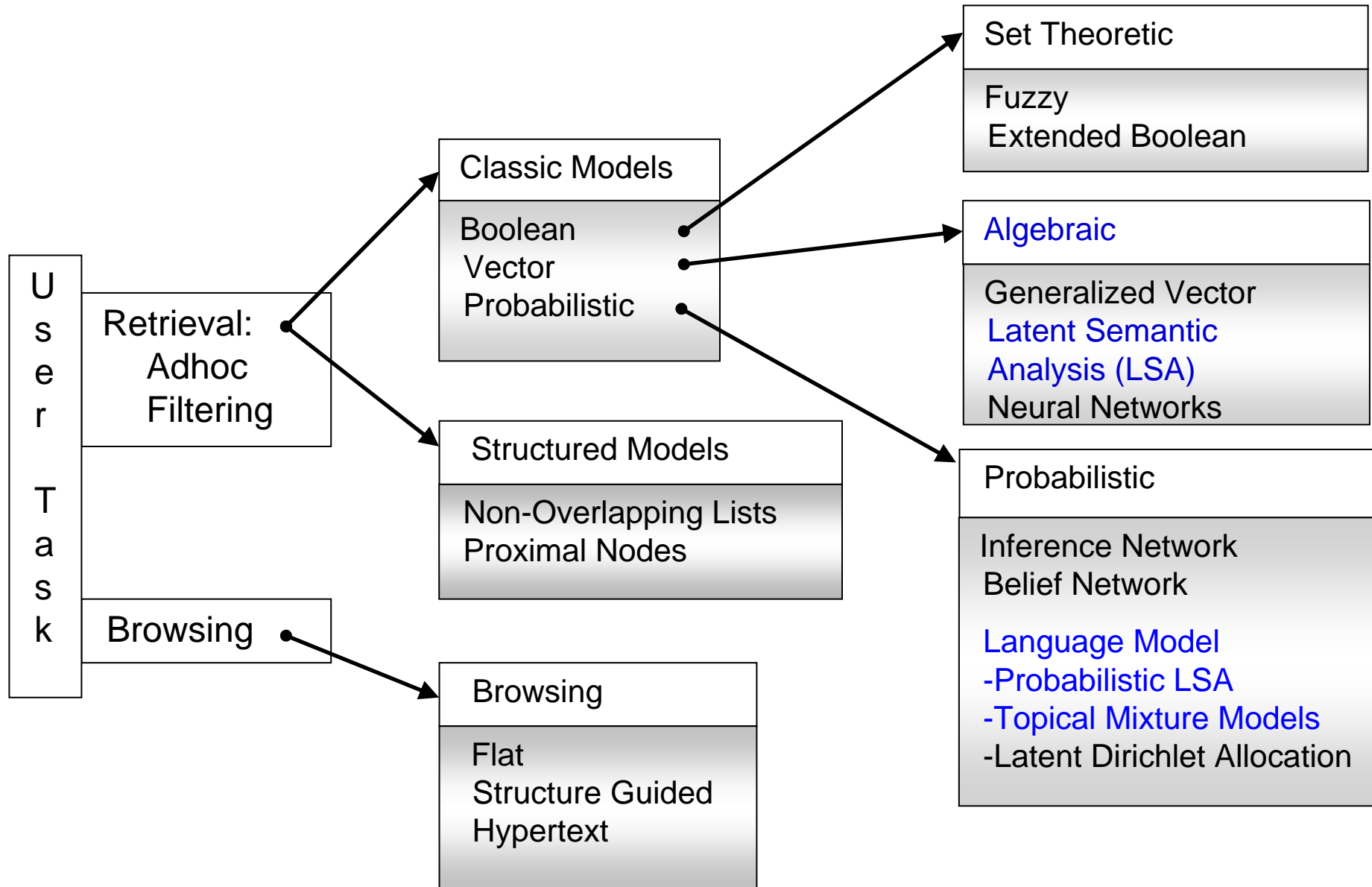
Latent Semantic Approaches for Information Retrieval and Language Modeling



Berlin Chen
Department of Computer Science & Information Engineering
National Taiwan Normal University



Taxonomy of Classic IR Models



Classification of IR Models Along Two Axes

- Matching Strategy

- Literal term matching

- E.g., Vector Space Model (VSM), Hidden Markov Model (HMM), Language Model (LM)

- Concept matching

- E.g., Latent Semantic Analysis (LSA), Probabilistic Latent Semantic Analysis (PLSA), Topical Mixture Model (TMM)

- Learning Capability

- Heuristic approaches for term weighting, query expansion, document expansion, etc.

- E.g., Vector Space Model, Latent Semantic Analysis
- Most approaches are based on linear algebra operations

- Solid statistical foundations (optimization algorithms)

- E.g., Hidden Markov Model (HMM), Probabilistic Latent Semantic Analysis, Latent Dirichlet Allocation (LDA)
- Most models belong to the language modeling approach

Two Perspectives for IR Models (cont.)

- Literal Term Matching vs. Concept Matching



香港星島日報篇報導引述軍事觀察家的話表示，到二零零五年台灣將完全喪失空中優勢，原因是中國大陸戰機不論是數量或是性能上都將超越台灣，報導指出中國在大量引進俄羅斯先進武器的同時也得加快研發自製武器系統，目前西安飛機製造廠任職的改進型飛豹戰機即將部署尚未與蘇愷三十通道地對地攻擊住宅飛機，以督促遇到挫折的監控其戰機目前也已經取得了重大階段性的認知成果。根據日本媒體報導在台海戰爭隨時可能爆發情況之下北京方面的基本方針，使用高科技答應局部戰爭。因此，解放軍打算在二零零四年前又有包括蘇愷三十二期在內的兩百架蘇霍伊戰鬥機。

- There are usually many ways to express a given concept (an information need), so literal terms in a user's query may not match those of a relevant document

Latent Semantic Analysis (LSA)

- Also called Latent Semantic Indexing (LSI), Latent Semantic Mapping (LSM), or Two-Mode Factor Analysis
 - Original formulated in the context of information retrieval
 - Users tend to retrieve documents on the basis of conceptual content
 - Individual terms (**units**) provide unreliable evidence about the conceptual topic or meaning of a document (**composition**)
 - There are many ways to express a given concept
 - LSA attempts to explore some underlying latent semantic structure in the data (documents) which is partially obscured by the randomness of word choices
 - LSA results in a parsimonious description of terms and documents
 - Contextual or positional information for words in documents is discarded (the so-called **bag-of-words** assumption)

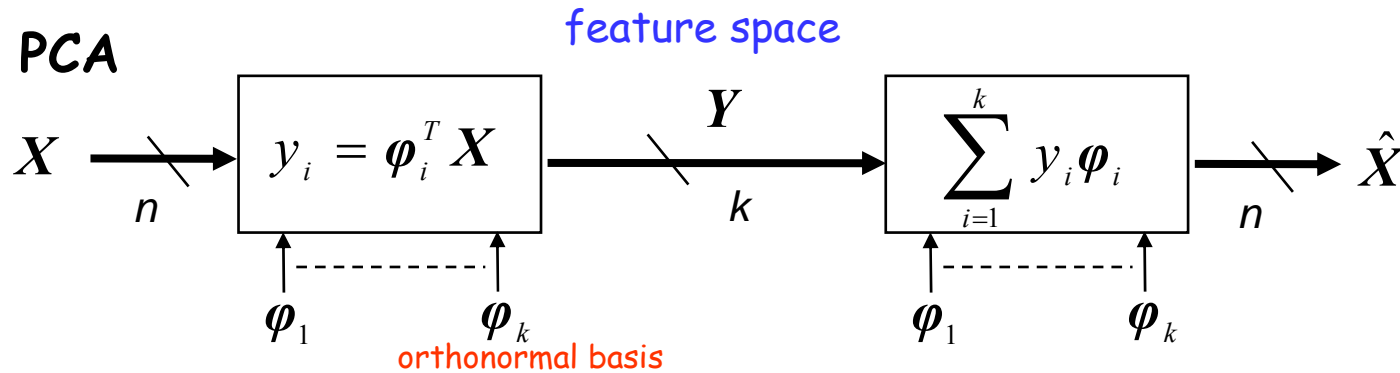
Applications of LSA

- Information Retrieval
- Word/document/Topic Clustering
- Language Modeling
- Automatic Call Routing
- Language Identification
- Pronunciation Modeling
- Speaker Verification (Prosody Analysis)
- Utterance Verification
- Text/Speech Summarization
- Automatic Image Annotation
-

LSA : Schematic Depiction

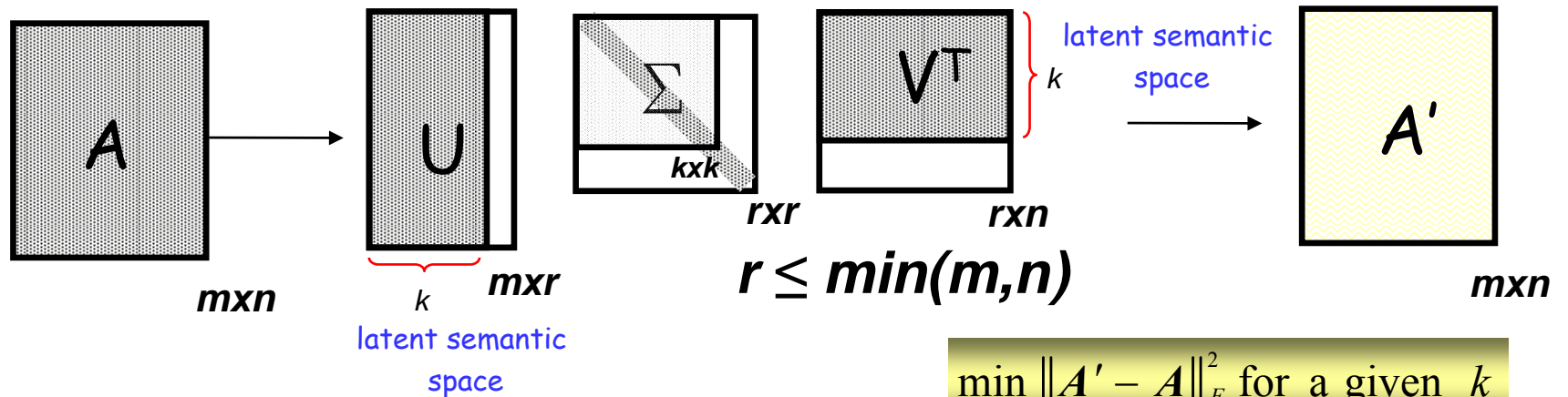
- Dimension Reduction and Feature Extraction

- PCA



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

- SVD (in LSA)



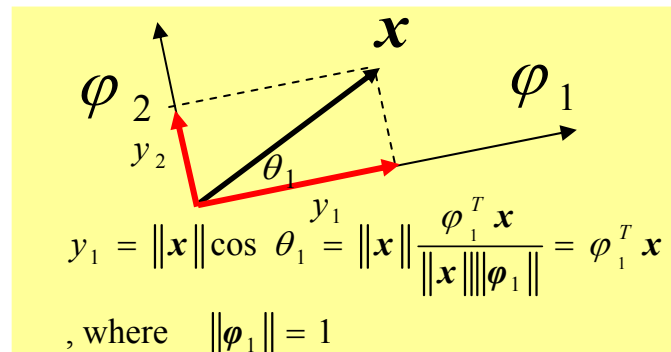
$$\min \left\| A' - A \right\|_F^2 \text{ for a given } k$$

LSA: An Example

- Singular Value Decomposition (SVD) used for the word-document matrix
 - A least-squares method for dimension reduction

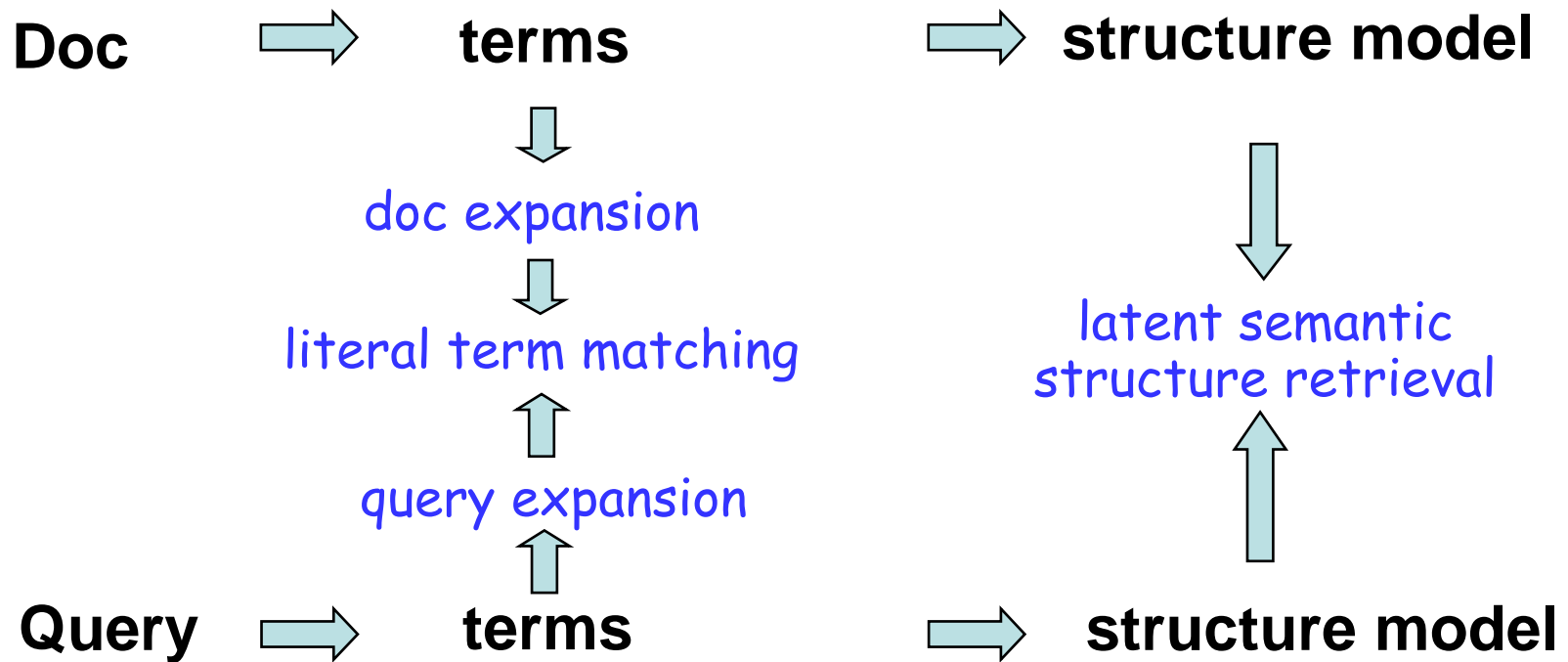
	Term 1	Term 2	Term 3	Term 4
Query	user	interface		
Document 1	user	interface	HCI	interaction
Document 2			HCI	interaction

Projection of a Vector \mathbf{x} :



LSA: Latent Structure Space

- Two alternative frameworks to circumvent vocabulary mismatch



LSA: Another Example (1/2)

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
1.	<i>human</i>	1	0	0	1	0	0	0	0	0
2.	<i>interface</i>	1	0	1	0	0	0	0	0	0
3.	<i>computer</i>	1	1	0	0	0	0	0	0	0
4.	<i>user</i>	0	1	1	0	1	0	0	0	0
5.	<i>system</i>	0	1	1	2	0	0	0	0	0
6.	<i>response</i>	0	1	0	0	1	0	0	0	0
7.	<i>time</i>	0	1	0	0	1	0	0	0	0
8.	<i>EPS</i>	0	0	1	1	0	0	0	0	0
9.	<i>survey</i>	0	1	0	0	0	0	0	0	1
10.	<i>trees</i>	0	0	0	0	0	1	1	1	0
11.	<i>graph</i>	0	0	0	0	0	0	1	1	1
12.	<i>minors</i>	0	0	0	0	0	0	0	1	1

LSA: Another Example (2/2)

2-D Plot of Terms and Docs from Example

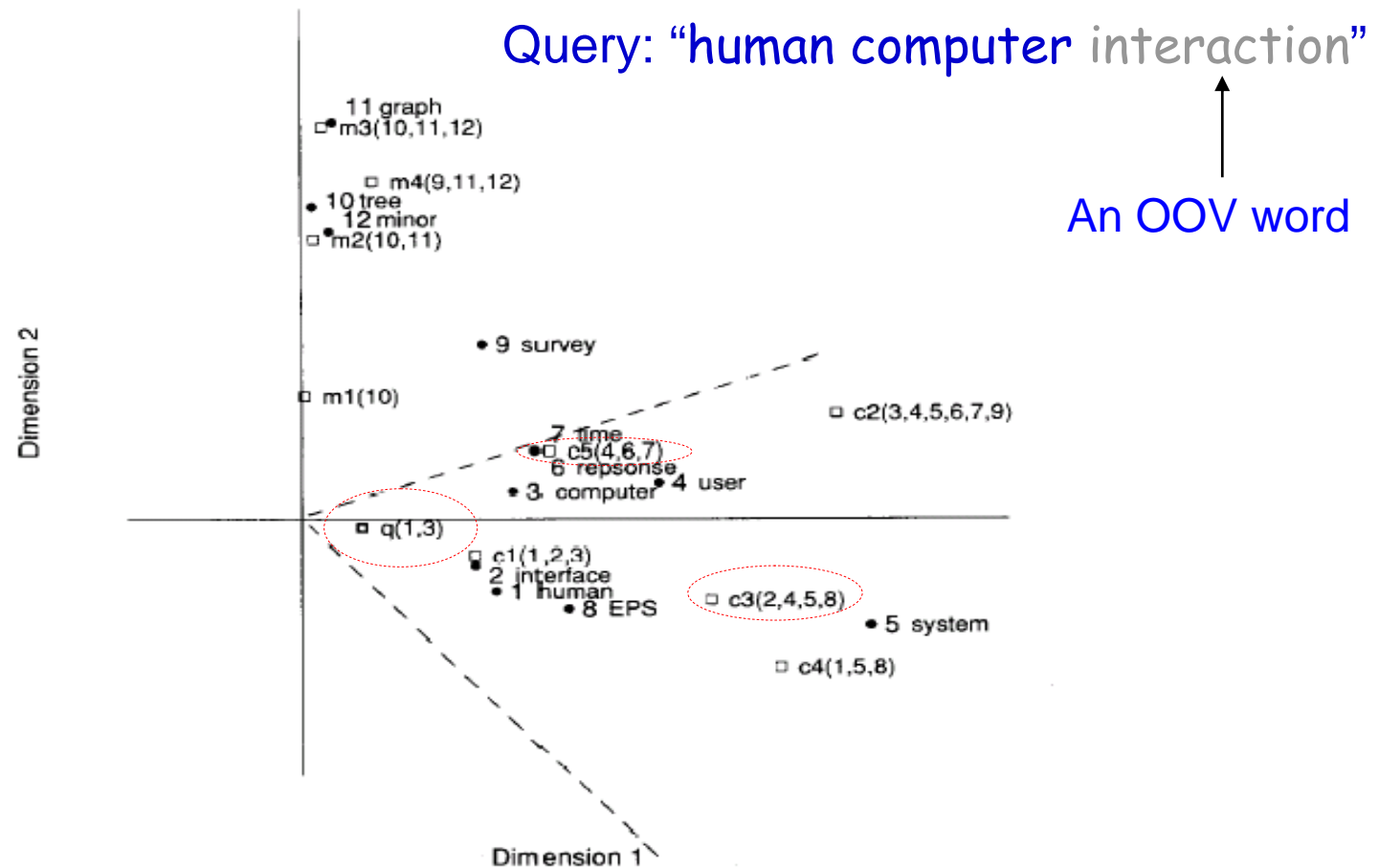
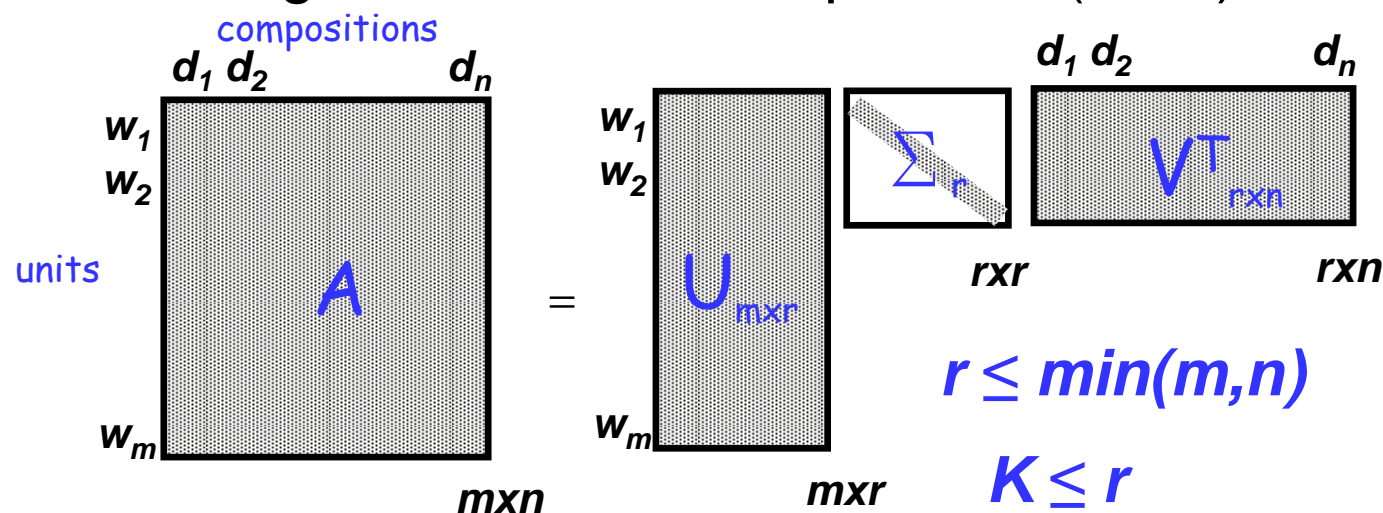


FIG. 1. A two-dimensional plot of 12 Terms and 9 Documents from the same TM set. Terms are represented by filled circles. Documents are shown as open squares, and component terms are indicated parenthetically. The query ("human computer interaction") is represented as a pseudo-document at point q . Axes 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.

LSA: Theoretical Foundation

- Singular Value Decomposition (SVD)

Row $A \in R^n$
 Col $A \in R^m$

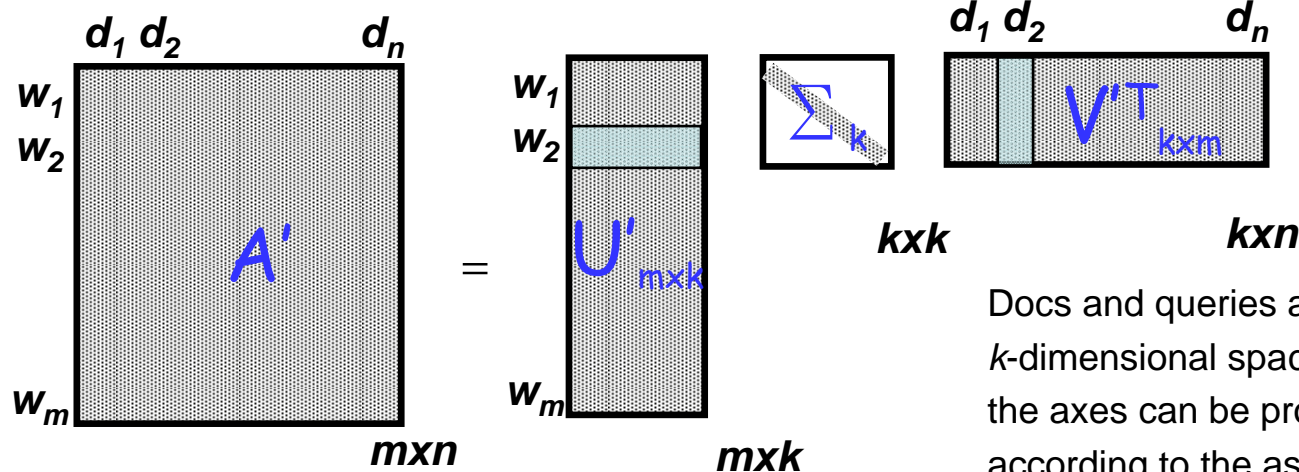


Both U and V has orthonormal column vectors

$$U^T U = I_{r \times r}$$

$$V^T V = I_{r \times r}$$

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



$$\|A\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n a_{ij}^2$$

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 Σ_k

LSA: Theoretical Foundation

- “term-document” matrix A has to do with the co-occurrences between terms (**units**) and documents (**compositions**)
 - Contextual or positional information for words in documents is discarded
 - “**bag-of-words**” modeling

- **Feature extraction** for the entities $a_{i,j}$ of matrix A

1. Conventional *tf-idf* statistics

2. Or, $a_{i,j}$: occurrence frequency weighted by negative entropy

occurrence count of term i in document j

$$a_{i,j} = \frac{f_{i,j}}{|d_j|} \times (1 - \varepsilon_i), \quad |d_j| = \sum_{i=1}^m f_{i,j}$$

negative normalized entropy

document length

normalized entropy of term i

$$\varepsilon_i = -\frac{1}{\log n} \sum_{j=1}^n \left(\frac{f_{i,j}}{\tau_i} \log \frac{f_{i,j}}{\tau_i} \right), \quad \tau_i = \sum_{j=1}^n f_{i,j}$$

occurrence count of term i in the collection

$$0 \leq \varepsilon_i \leq 1$$

LSA: Theoretical Foundation

- Singular Value Decomposition (SVD)

- $A^T A$ is symmetric $n \times n$ matrix

- All eigenvalues λ_j are nonnegative real numbers

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0 \quad \Sigma^2 = \text{diag}(\lambda_1, \lambda_1, \dots, \lambda_n)$$

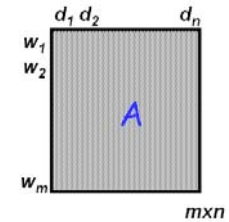
- All eigenvectors v_j are orthonormal ($\in \mathbb{R}^n$)

$$V = [v_1 \ v_2 \ \dots \ v_n] \quad v_j^T v_j = 1 \quad (V^T V = I_{n \times n})$$

- Define **singular values**: sigma $\sigma_j = \sqrt{\lambda_j}$, $j = 1, \dots, n$

- As the square roots of the eigenvalues of $A^T A$

- As the lengths of the vectors Av_1, Av_2, \dots, Av_n



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

$$\sigma_1 = \|Av_1\|$$

$$\sigma_2 = \|Av_2\|$$

.....

$$\|Av_i\|^2 = v_i^T A^T A v_i = v_i^T \lambda_i v_i = \lambda_i$$

$$\Rightarrow \|Av_i\| = \sigma_i$$

LSA: Theoretical Foundation

- $\{Av_1, Av_2, \dots, Av_r\}$ is an **orthogonal** basis of **Col A** ($\in R^m$)

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

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

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

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

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

*U is also an
orthonormal matrix
(m x r)*

$$\Rightarrow [u_1 \ u_2 \ \dots \ u_r] \Sigma_{r \times r} = A [v_1 \ v_2 \ \dots \ v_r]$$

V: an orthonormal matrix

Known in advance

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

$$\Rightarrow [u_1 \ u_2 \ \dots \ u_r \ \dots \ u_m] \Sigma_{m \times n} = A [v_1 \ v_2 \ \dots \ v_r \ \dots \ v_n]$$

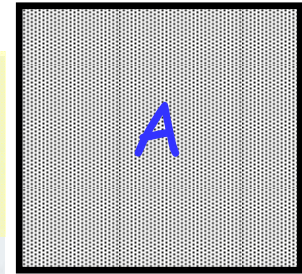
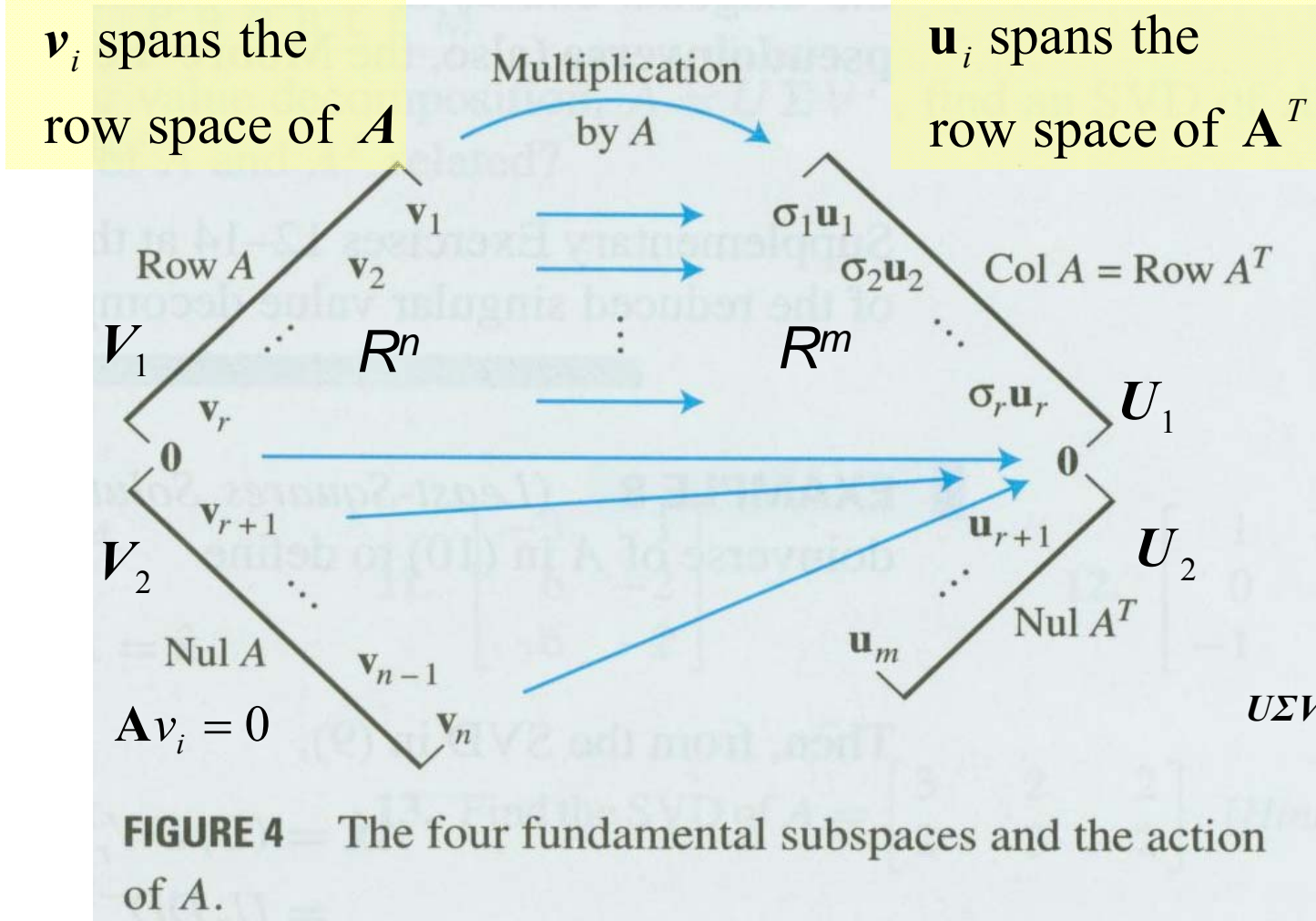
$$\|A\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n a_{ij}^2$$

$$\Rightarrow U \Sigma = AV \Rightarrow U \Sigma V^T = A$$

$$\|A\|_F^2 = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_r^2 \quad ?$$

$$\Rightarrow A = U \Sigma V^T \quad \Sigma_{m \times n} = \begin{pmatrix} \Sigma_r & \mathbf{0}_{r \times (n-r)} \\ \mathbf{0}_{(m-r) \times r} & \mathbf{0}_{(m-r) \times (n-r)} \end{pmatrix} I_{n \times n} \quad ?$$

LSA: Theoretical Foundation



$m \times n$

$$\begin{aligned}
 \mathbf{U} & \begin{matrix} \text{U} \\ \text{U}_1 \\ \text{U}_2 \end{matrix} & \mathbf{V}^T & \begin{matrix} \text{V}^T \\ \text{V}_1^T \\ \text{V}_2^T \end{matrix} \\
 \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T &= (\mathbf{U}_1 \quad \mathbf{U}_2) \begin{pmatrix} \mathbf{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{V}_1^T \\ \mathbf{V}_2^T \end{pmatrix} \\
 &= \mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^T \\
 &= \mathbf{A} \mathbf{V}_1 \mathbf{V}_1^T \\
 &= \mathbf{A}
 \end{aligned}$$

$\mathbf{U} \mathbf{\Sigma} = \mathbf{A} \mathbf{V}$

LSA: Theoretical Foundation

- Additional Explanations

- Each row of U is related to the projection of a corresponding row of A onto the basis formed by columns of V

$$A = U\Sigma V^T$$

$$\Rightarrow AV = U\Sigma V^T V = U\Sigma \Rightarrow U\Sigma = AV$$

- the i -th entry of a row of U is related to the projection of a corresponding row of A onto the i -th column of V

- Each row of V is related to the projection of a corresponding row of A^T onto the basis formed by U

$$A = U\Sigma V^T$$

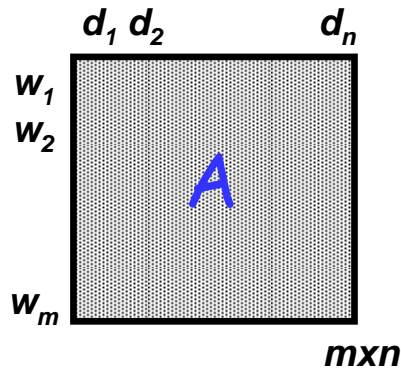
$$\Rightarrow A^T U = (U\Sigma V^T)^T U = V\Sigma U^T U = V\Sigma$$

$$\Rightarrow V\Sigma = A^T U$$

- the i -th entry of a row of V is related to the projection of a corresponding row of A^T onto the i -th column of U

LSA: Theoretical Foundation

- Fundamental comparisons based on SVD
 - The original word-document matrix (A)



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

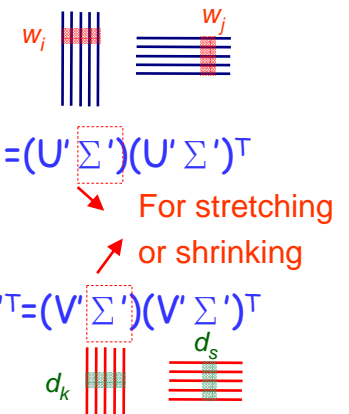
- The new word-document matrix (A')

$$U' = U_{m \times k}$$

$$\Sigma' = \Sigma_k$$

$$V' = V_{n \times k}$$

- compare two terms $A'A^T = (U' \Sigma' V'^T)(U' \Sigma' V'^T)^T = U' \Sigma' V'^T V' \Sigma'^T U'^T = (U' \Sigma') (U' \Sigma')^T$
 - \rightarrow dot product of two rows of $U' \Sigma'$
- compare two docs $A^T A = (U' \Sigma' V'^T)^T (U' \Sigma' V'^T) = V' \Sigma'^T U'^T U' \Sigma' V'^T = (V' \Sigma') (V' \Sigma')^T$
 - \rightarrow dot product of two rows of $V' \Sigma'$
- compare a query word and a doc \rightarrow each individual entry of A'



LSA: Fold-in

- Find representations for pseudo-docs
 - For objects (new queries or docs) that did not appear in the original analysis
 - Fold-in a new $m \times 1$ query (or doc) vector

See Figure A in next page

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

Just like a row of V

Query represented by the weighted sum of its constituent term vectors

The separate dimensions are differentially weighted

- Represented as the weighted sum of its component word (or term) vectors
- Cosine measure between the query and doc vectors in the latent semantic space

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

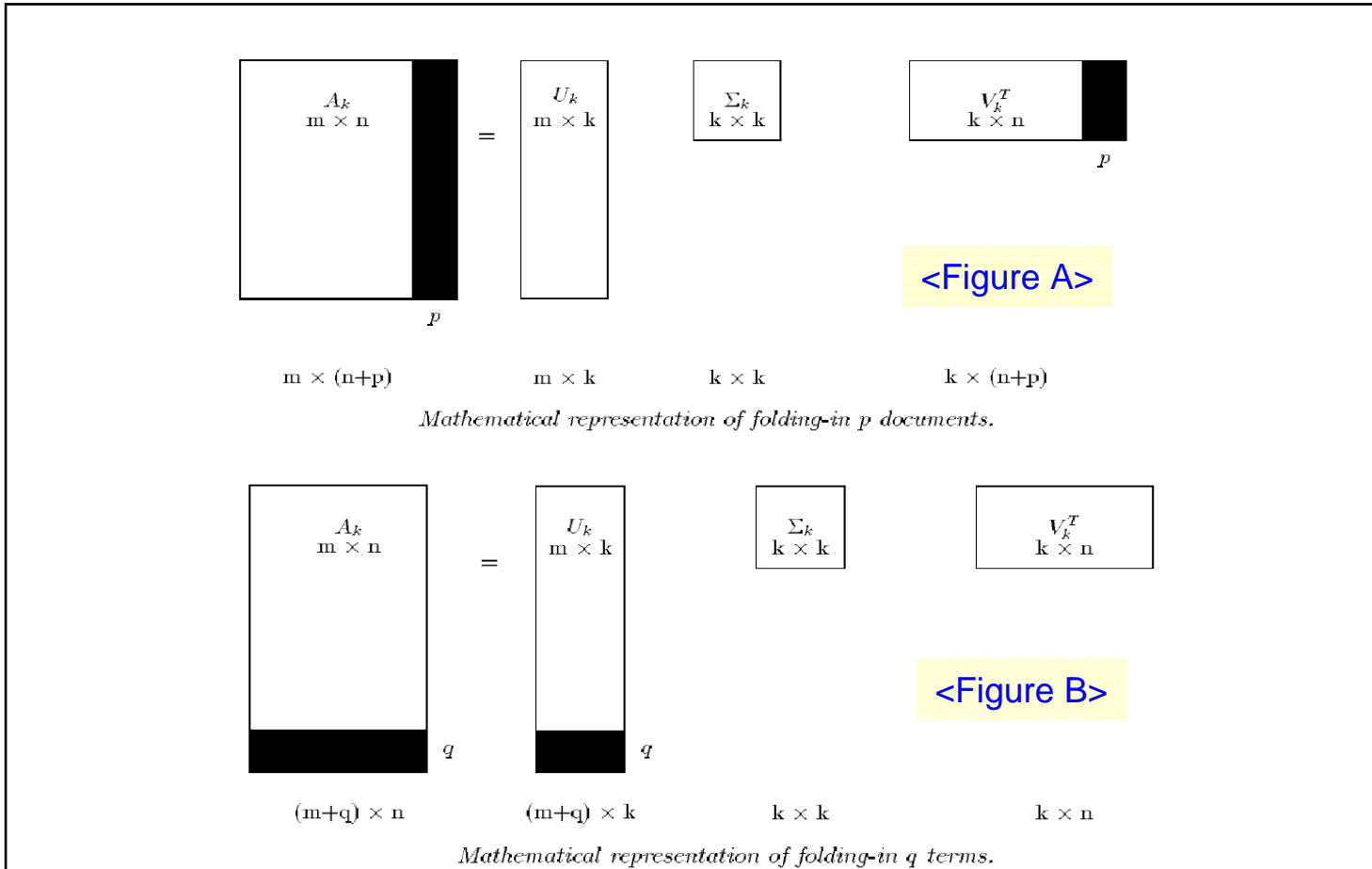
row vectors

LSA: Theoretical Foundation

- Fold-in a new 1 x n term vector

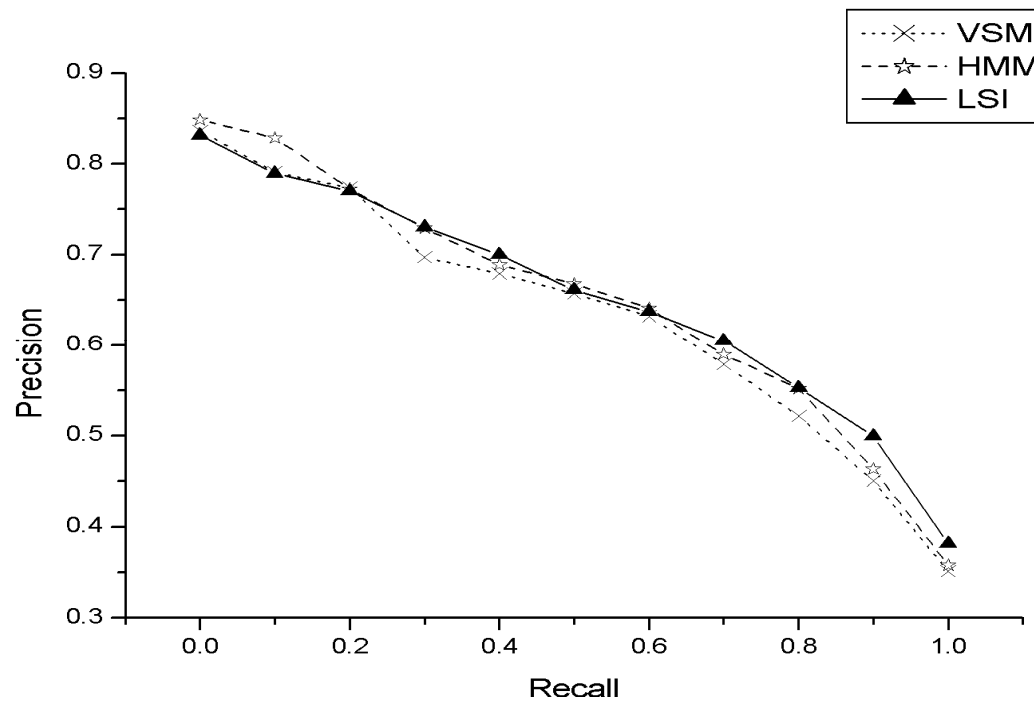
$$\hat{t}_{1 \times k} = t_{1 \times n} V_{n \times k} \Sigma_{k \times k}^{-1} \mathbf{1}_{k \times k}$$

See Figure B below



LSA: A Simple IR Evaluation

- Experimental results
 - HMM is consistently better than VSM at all recall levels
 - LSA 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)

LSA: Pro and Con (1/2)

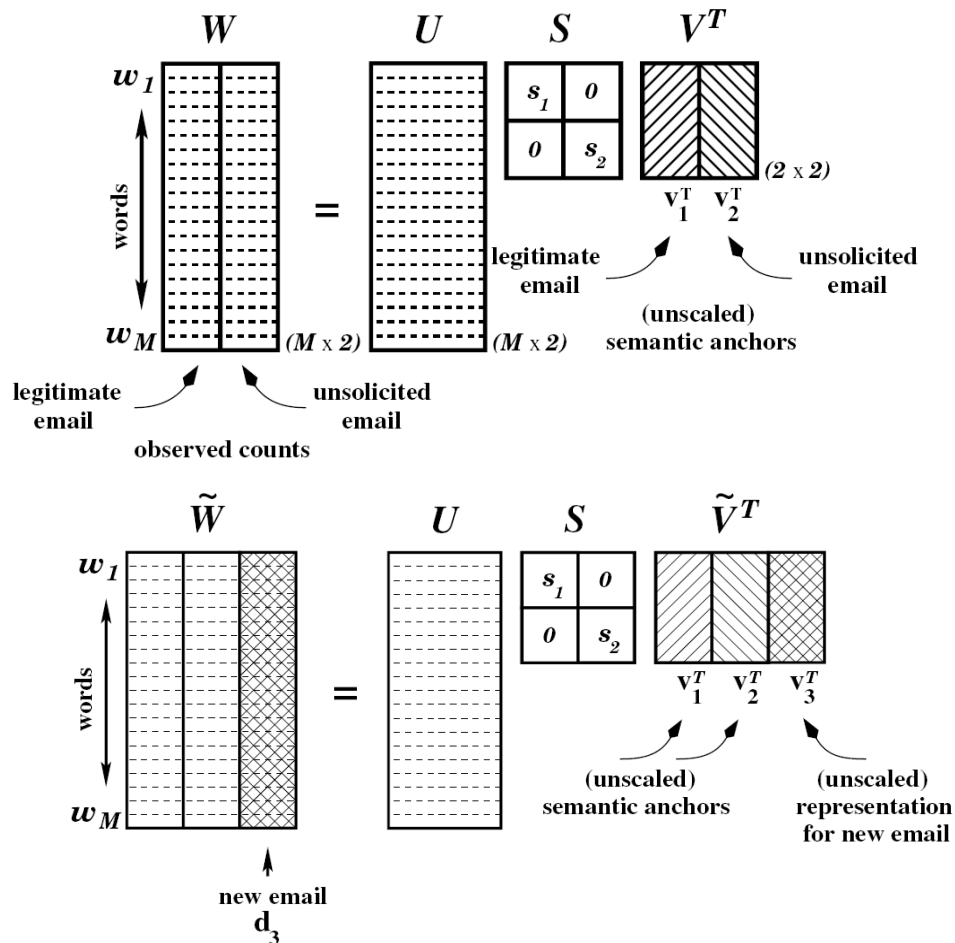
- Pro (Advantages)
 - A clean formal framework and a clearly defined optimization criterion (least-squares)
 - Conceptual simplicity and clarity
 - Handle synonymy problems (“heterogeneous vocabulary”)
 - Replace individual terms as the descriptors of documents by independent “artificial concepts” that can be specified by any one of several terms (or documents) or combinations
 - Good results for high-recall search
 - Take term co-occurrence into account

LSA: Pro and Con (2/2)

- Disadvantages
 - High computational complexity (e.g., SVD decomposition)
 - Exhaustive comparison of a query against all stored documents is needed (cannot make use of inverted files ?)
 - LSA offers only a partial solution to polysemy (e.g. bank, bass,...)
 - Every term is represented as just one point in the latent space (represented as weighted average of different meanings of a term)

LSA: Junk E-mail Filtering

- One vector represents the centroid of all e-mails that are of interest to the user, while the other the centroid of all e-mails that are not of interest



LSA: Dynamic Language Model Adaptation (1/4)

- Let w_q denote the word about to be predicted, and H_{q-1} the admissible LSA history (context) for this particular word
 - The vector representation of H_{q-1} is expressed by \tilde{d}_{q-1}
 - Which can be then projected into the latent semantic space

LSA representation $\tilde{\tilde{v}}_{q-1} = \tilde{v}_{q-1} S = \tilde{d}_{q-1}^T U$ [change of notation : $S = \Sigma$]

- Iteratively update \tilde{d}_{q-1} and $\tilde{\tilde{v}}_{q-1}$ as the decoding evolves

VSM representation $\tilde{d}_q = \frac{n_q - 1}{n_q} \tilde{d}_{q-1} + \frac{1 - \varepsilon_i}{n_q} [0 \dots 1 \dots 0]^T$

LSA representation $\tilde{\tilde{v}}_q = \tilde{v}_q S = d_{q-1}^T U = \frac{1}{n_q} [(n_q - 1) \tilde{\tilde{v}}_{q-1} + \underline{(1 - \varepsilon_i) u_i}]$

or $= \frac{1}{n_q} [\lambda \cdot (n_q - 1) \tilde{\tilde{v}}_{q-1} + (1 - \varepsilon_i) u_i]$

with exponential decay

LSA: Dynamic Language Model Adaptation (2/4)

- Integration of LSA with N-grams

$$\Pr(w_q | H_{q-1}^{(n+l)}) = \Pr(w_q | H_{q-1}^{(n)}, H_{q-1}^{(l)})$$

where H_{q-1} denotes some suitable history for word w_q ,

and the superscripts $^{(n)}$ and $^{(l)}$ refer to the n -gram component ($w_{q-1}w_{q-2}\dots w_{q-n+1}$, with $n > 1$), the LSA

component (\tilde{d}_{q-1}):

This expression can be rewritten as :

$$\Pr(w_q | H_{q-1}^{(n+l)}) = \frac{\Pr(w_q, H_{q-1}^{(l)} | H_{q-1}^{(n)})}{\sum_{w_i \in V} \Pr(w_i, H_{q-1}^{(l)} | H_{q-1}^{(n)})}$$

LSA: Dynamic Language Model Adaptation (3/4)

- Integration of LSA with N-grams (cont.)

$$\Pr(w_q, H_{q-1}^{(l)} | H_{q-1}^{(n)}) =$$

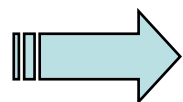
Assume the probability of the document history given the current word is not affected by the immediate context preceding it

$$\Pr(w_q | H_{q-1}^{(n)}) \cdot \Pr(H_{q-1}^{(l)} | w_q, H_{q-1}^{(n)})$$

$$= \Pr(w_q | w_{q-1} w_{q-2} \cdots w_{q-n+1}) \cdot \Pr(\tilde{d}_{q-1} | w_q \underline{w_{q-1} w_{q-2} \cdots w_{q-n+1}})$$

$$= \Pr(w_q | w_{q-1} w_{q-2} \cdots w_{q-n+1}) \cdot \Pr(\tilde{d}_{q-1} | w_q)$$

$$= \Pr(w_q | w_{q-1} w_{q-2} \cdots w_{q-n+1}) \cdot \frac{\Pr(w_q | \tilde{d}_{q-1}) \Pr(\tilde{d}_{q-1})}{\Pr(w_q)}$$



$$\Pr(w_q | H_{q-1}^{(n+l)}) =$$

$$\frac{\Pr(w_q | w_{q-1} w_{q-2} \cdots w_{q-n+1}) \cdot \frac{\Pr(w_q | \tilde{d}_{q-1})}{\Pr(w_q)}}{\sum_{w_i \in V} \Pr(w_i | w_{q-1} w_{q-2} \cdots w_{q-n+1}) \cdot \frac{\Pr(w_i | \tilde{d}_{q-1})}{\Pr(w_i)}}$$

LSA: Dynamic Language Model Adaptation (4/4)

Intuitively, $\Pr(w_q | \tilde{d}_{q-1})$ reflects the "relevance" of word w_q to the admissible history, as observed through \tilde{d}_{q-1} :

$$\begin{aligned} & \Pr(w_q | \tilde{d}_{q-1}) \\ & \approx K(w_q, \tilde{d}_{q-1}) \\ & = \cos(u_q S^{1/2}, \tilde{v}_{q-1} S^{1/2}) = \frac{u_q S \tilde{v}_{q-1}^T}{\|u_q S^{1/2}\| \|\tilde{v}_{q-1} S^{1/2}\|} \end{aligned}$$

As such, it will be highest for words whose meaning aligns most closely with the semantic fabric of \tilde{d}_{q-1} (i.e., relevant "content" words), and lowest for words which do not convey any particular information about this fabric (e.g., "function" words like "*the*").

LSA: Cross-lingual Language Model Adaptation (1/2)

- Assume that a document-aligned (instead of sentence-aligned) Chinese-English bilingual corpus is provided

$$\begin{array}{c}
 \begin{array}{cccc}
 & \mathbf{W} & & \\
 \hline
 d_1^E & d_2^E & \dots & d_N^E \\
 \hline
 d_1^C & d_2^C & \dots & d_N^C \\
 \hline
 & \mathbf{M} \times \mathbf{N} & &
 \end{array}
 =
 \begin{array}{c}
 \begin{array}{cc}
 \mathbf{U} & \mathbf{S} \\
 \hline
 & \\
 \hline
 & \\
 \hline
 \mathbf{M} \times \mathbf{R} & \mathbf{R} \times \mathbf{R}
 \end{array}
 \mathbf{X}
 \begin{array}{c}
 \mathbf{V}^T \\
 \hline
 \\
 \hline
 \\
 \hline
 \mathbf{R} \times \mathbf{N}
 \end{array}
 \end{array}$$

SVD of a word-document matrix for CL-LSA.

$$\begin{array}{c}
 \begin{array}{cccc}
 & \bar{\mathbf{W}} & & \\
 \hline
 \bar{d}_1^E & \bar{d}_2^E & \dots & \bar{d}_P^E \\
 \hline
 0 & 0 & \dots & 0 \\
 \hline
 & \mathbf{M} \times \mathbf{P} & &
 \end{array}
 =
 \begin{array}{c}
 \begin{array}{cc}
 \mathbf{U} & \mathbf{S} \\
 \hline
 \phantom{\bar{d}_1^E} & \phantom{\bar{d}_2^E} \\
 \hline
 \phantom{\bar{d}_1^E} & \phantom{\bar{d}_2^E} \\
 \hline
 \mathbf{M} \times \mathbf{R} & \mathbf{R} \times \mathbf{R}
 \end{array}
 \mathbf{X}
 \begin{array}{c}
 \bar{\mathbf{V}}^T \\
 \hline
 \phantom{\bar{d}_1^E} \\
 \hline
 \phantom{\bar{d}_1^E} \\
 \hline
 \mathbf{R} \times \mathbf{P}
 \end{array}
 \end{array}$$

Folding-in a monolingual corpus into LSA.

LSA: Cross-lingual Language Model Adaptation (2/2)

- CL-LSA adapted Language Model

$$P_{\text{Adapt}} \left(c_k | c_{k-1}, c_{k-2}, d_i^E \right) \quad d_i^E \text{ is a relevant English doc of the Mandarin } d_i^C \text{ doc being transcribed, obtained by CL-IR}$$
$$\approx \lambda \cdot P_{\text{CL-LSA-Unigram}} \left(c_k | d_i^E \right) + (1 - \lambda) \cdot P_{\text{BG}} \left(c_k | c_{k-1}, c_{k-2} \right)$$

$$P_{\text{CL-LSA-Unigram}} \left(c_k | d_i^E \right) = \sum_e P_T(c|e) P(e | d_i^E)$$

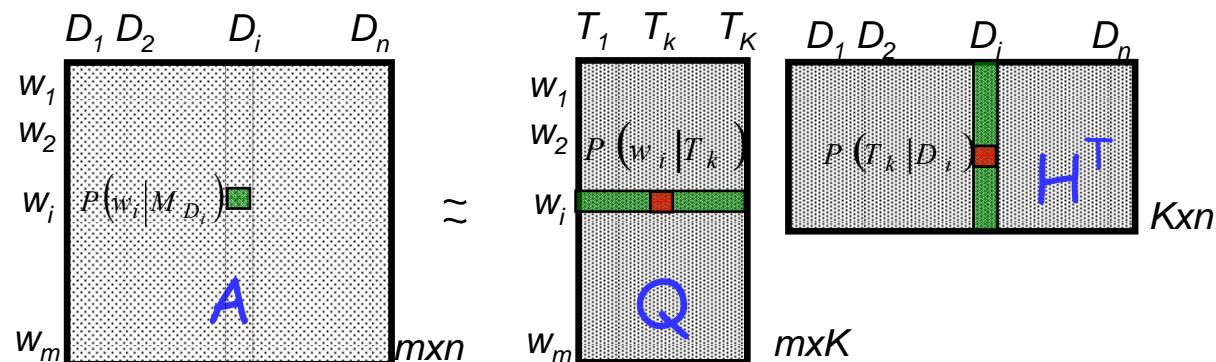
$$P_T(c|e) \approx \frac{\text{sim}(\vec{c}, \vec{e})^\gamma}{\sum_{c'} \text{sim}(\vec{c}', \vec{e})^\gamma} \quad (\gamma \gg 1)$$

Probabilistic Latent Semantic Analysis (PLSA)

- PLSA models the co-occurrence of word and documents and evaluates the relevance in a low dimensional semantic/topic space
 - Each document D is treated as a document model M_D

$$P_{\text{PLSA}}(w_i | M_D) = \sum_{k=1}^K P(w_i | T_k) P(T_k | M_D)$$

- PLSA can be viewed as a nonnegative factorization of a “word-document” matrix consisting probability entries
 - A procedure similar to the SVD performed by its algebraic counterpart- LSA



PLSA: Information Retrieval (1/3)

- The relevance measure between a query and a document can be expressed by

$$P_{\text{PLSA}}(Q|M_D) = \prod_{w_i \in Q} \left[\sum_{k=1}^K P(w_i|T_k)P(T_k|M_D) \right]^{c(w_i,Q)}$$

- Relevance measure is not obtained based on the frequency of a respective query term occurring in a document, but instead based on the frequency of the term and document in the latent topics
- A query and a document thus may have a high relevance score even if they do not share any terms in common

PLSA: Information Retrieval (2/3)

- **Unsupervised training:** The model parameters are trained beforehand using a set of text documents
 - Maximize the log-likelihood of entire collection \mathbf{D}

$$\log L_{\mathbf{D}} = \sum_{D \in \mathbf{D}} \log P_{PLSA}(D | M_D) = \sum_{D \in \mathbf{D}} \sum_{w_i \in D} c(w_i, D) \log P_{PLSA}(w_i | M_D)$$

- **Supervised training:** The model parameters are trained using a training set of query exemplars and the associated query-document relevance information
 - Maximize the log-likelihood of the training set of query exemplars generated by their relevant documents

$$\begin{aligned} \log L_{\mathbf{Q}_{TrainSet}} &= \sum_{Q \in \mathbf{Q}_{TrainSet}} \sum_{D \in \mathbf{D}_{R \text{ to } Q}} \log P_{PLSA}(Q | M_D) \\ &= \sum_{Q \in \mathbf{Q}_{TrainSet}} \sum_{D \in \mathbf{D}_{R \text{ to } Q}} \sum_{w_i \in Q} c(w_i, Q) \log P(w_i | M_D) \end{aligned}$$

PLSA: Information Retrieval (3/3)

- Example: most probable words form 4 latent topics

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.

PLSA: Dynamic Language Model Adaptation

- The search history can be treated as a pseudo-document which is varying during the speech recognition process

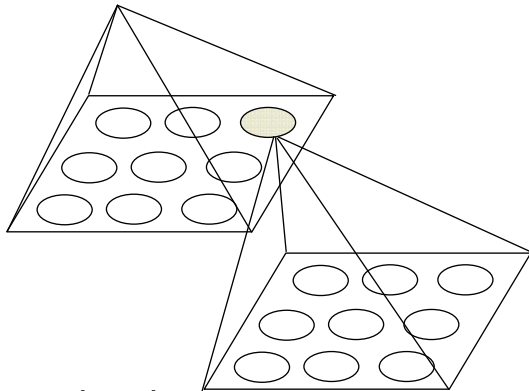
$$P_{\text{PLSA}}(w_i | H_{w_i}) = \sum_{k=1}^K P(w_i | T_k) P(T_k | H_{w_i})$$

- The topic unigrams $P(w_i | T_k)$ are kept unchanged
- The history's probability distribution over the latent topics is gradually updated
- The topic mixture weights $P(T_k | H_{w_i})$ are estimated on the fly
 - It would be time-consuming

PLSA: Document Organization (1/3)

- Each document is viewed as a document model to generate itself
 - Additional transitions between topical mixtures have to do with the topological relationships between topical classes on a 2-D map map

$$P_{\text{PLSA}}(w_i | M_D) = \sum_{k=1}^K P(T_k | M_D) \left[\sum_{l=1}^K P(T_l | T_k) P(w_i | T_l) \right]$$



Two-dimensional
Tree Structure
for Organized Topics

$$E(T_l, T_k) = \frac{1}{\sqrt{2\pi\sigma}} \exp \left[-\frac{\text{dist}(T_k, T_l)^2}{2\sigma^2} \right]$$

$$P(T_l | T_k) = \frac{E(T_l, T_k)}{\sum_{s=1}^K E(T_s, T_k)}$$

PLSA: Document Organization (2/3)

- Document models can be trained in an unsupervised way by maximizing the total log-likelihood of the document collection

$$L_T = \sum_{j=1}^n \sum_{i=1}^V c(w_i, D_j) \log P(w_i | D_j)$$

- Each topical class can be labeled by words selected using the following criterion

$$\text{Sig}(w_i, T_k) = \frac{\sum_{j=1}^n c(w_i, D_j) P(T_k | D_j)}{\sum_{i=1}^n c(w_i, D_j) [1 - P(T_k | D_j)]}$$

PLSA: Document Organization (3/3)

- Spoken Document Retrieval and Browsing System developed by NTU (Prof. Lin-shan Lee)

廣播新聞搜尋瀏覽系統
Broadcast News Retrieval/Browsing System

[國外政治 \[International Political News\]](#) Topic Map
[國內政治 \[Local Political News\]](#) Topic Map
[國外財經 \[International Business\]](#) Topic Map
[國內財經 \[Local Business\]](#) Topic Map
[國外影劇 \[International Entertainment\]](#) Topic Map
[國內影劇 \[Local Entertainment\]](#) Topic Map
[國外體育 \[International Sports\]](#) Topic Map
[國內體育 \[Local Sports\]](#) Topic Map

(a) (b)

伊拉克 巴格達 美軍 陸戰隊	以色列 阿拉法特 巴勒斯坦 迦薩市
國土安全部 民航機 蓋達組織 中情局	聯合國 安理會 武檢人員 武器

(c)

go to Level-1

(d)

阿拉法特 阿巴斯 雷馬拉 任命	以色列 夏隆 約旦河 美國
中東 鮑爾 和平 路線	巴格達 炸彈 自殺 巴士

(e)

go to Level-2

1] 以色列結束對阿拉法特總部的包圍 [sum.] 02.09.20
 2] 阿拉法特反對以色列保所提結束包圍條件 [sum.] 02.09.20
 3] 以色列部隊進攻阿拉法特總部後撤軍 [sum.] 02.10.22
 4] 以色列結束對阿拉法特總部的包圍 [sum.] 02.10.01
 5] 以色列坦克撤出阿拉法特辦公區 [sum.] 02.09.21
 6] 以色列與巴勒斯坦展開安全問題會議 [sum.] 02.11.23
 7] 以色列在加薩擊斃一名回教聖戰組織領袖 [sum.] 02.06.05
 8] 以色列巴勒斯坦就伯利恆撤軍達成協議 [sum.] 02.02.12
 9] 以色列坦克闖入加薩難民營 兩人喪生 [sum.] 02.04.20

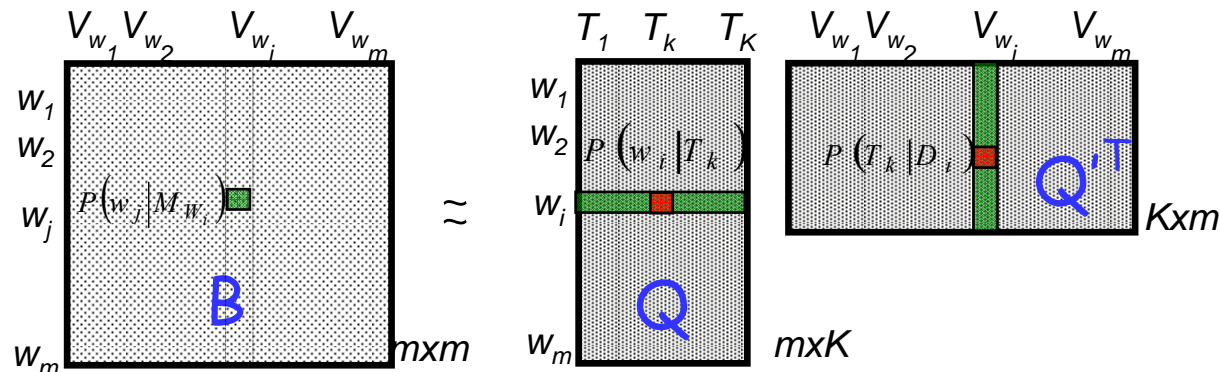
阿拉法特原則接受歐盟所提中東和平計畫 [summary] (May 03/02/12:00)
 英美就解決阿拉法特所受包圍與巴方展開談判 [summary] (May 06/02/12:00)
 阿拉法特反對以色列保所提結束包圍條件 [summary] (Sep 20/02/12:00)
 阿拉法特宣布新內閣引發巴勒斯坦國會激辯 [summary] (Oct 30/02/12:00)
 阿拉伯人支持阿拉法特及巴勒斯坦人正當抵抗 [summary] (Nov 02/02/12:00)

Word Topical Mixture Models (WTMM)

- Each word of language are treated as a word topical mixture model for predicting the occurrences of other words

$$P_{\text{WTMM}}(w_i | M_{w_j}) = \sum_{k=1}^K P(w_i | T_k) P(T_k | M_{w_j})$$

- WTMM also can be viewed as a nonnegative factorization of a “word-word” matrix consisting probability entries
 - Each column encodes the vicinity information of all occurrences of a certain type of word



WTMM: Information Retrieval (1/2)

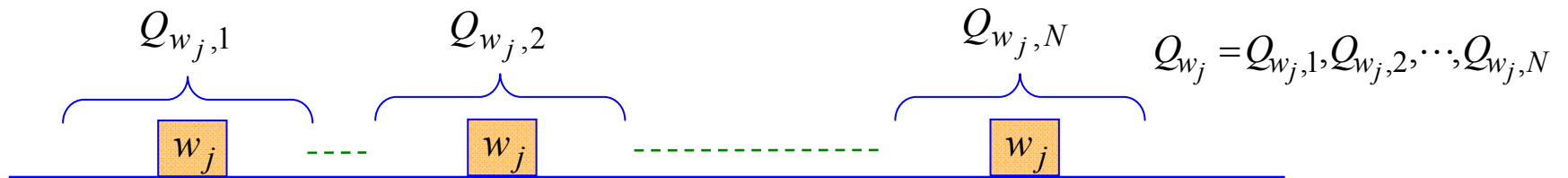
- The relevance measure between a query and a document can be expressed by

$$P_{\text{WTMM}}(Q|D) = \prod_{w_i \in Q} \left[\sum_{w_j \in D} \alpha_{j,D} \sum_{k=1}^K P(w_i|T_k) P(T_k|M_{w_j}) \right]^{c(w_i,Q)}$$

- Unsupervised training**

- The WTMM of each word can be trained by concatenating those words occurring within a context window of size around each occurrence of the word, which are postulated to be relevant to the word

$$\log L_{\mathbf{w}} = \sum_{w_j \in \mathbf{w}} \log P_{\text{WTMM}}(O_{w_j} | M_{w_j}) = \sum_{w_j \in \mathbf{w}} \sum_{w_i \in Q_{w_j}} c(w_i, O_{w_j}) \log P_{\text{WTMM}}(w_i | M_{w_j})$$



WTMM: Information Retrieval (2/2)

- **Supervised training:** The model parameters are trained using a training set of query exemplars and the associated query-document relevance information
 - Maximize the log-likelihood of the training set of query exemplars generated by their relevant documents

$$\log L_{\mathbf{Q}_{TrainSet}} = \sum_{Q \in \mathbf{Q}_{TrainSet}} \sum_{D \in \mathbf{D}_{R \text{ to } Q}} \log P_{WTMM}(Q|D)$$

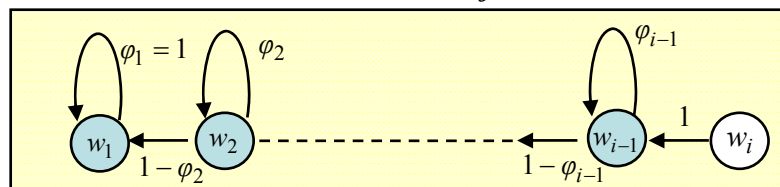
- The detailed training formulas of WTMM are a little bit complicated !

WTMM: Dynamic Language Model Adaptation (1/3)

- For a decoded word w_i , we can again interpret it as a (single-word) query; while for each of its search histories, expressed by $H_{w_i} = w_1, w_2, \dots, w_{i-1}$, we can linearly combine the associated TMM models of the words occurring in H_{w_i} to form a composite WTMM model

$$P_{\text{WTMM}}(w_i | M_{H_{w_i}}) = \sum_{j=1}^{i-1} \beta_j P_{\text{WTMM}}(w_i | M_{w_j}) = \sum_{j=1}^{i-1} \beta_j \sum_{k=1}^K P(w_i | T_k) P(T_k | M_{w_j})$$

$$\beta_j = \varphi_j \prod_{s=1}^{i-j-1} (1 - \varphi_{j+s})$$



- $\beta_j = \varphi_j$ are nonnegative weighting coefficients which empirically set to be exponentially decayed as the word is being apart from w_i
- φ_j is set to a fixed value (between 0 and 1) for $j = 2, \dots, i-1$, and set to 1 for $j = 1$

WTMM: Dynamic Language Model Adaptation (2/3)

- For our speech recognition test data, it was experimentally observed that the language model access time of WTMM was approximately 1/30 of that of PLSA for language model adaptation, as the iteration number of the online EM estimation of $P(T_k | H_{w_i})$ for PLSA was set to 5

WTMM: Dynamic Language Model Adaptation (3/3)

- An Alternative Formulation of WTMM

$$\begin{aligned} P_{\text{WTMM}}(w_i | H_{w_i}) &= \frac{P_{\text{WTMM}}(H_{w_i} | w_i) P_{\text{Unigram}}(w_i)}{P(H_{w_i})} \\ &= \frac{P_{\text{WTMM}}(H_{w_i} | M_{w_i}) P_{\text{Unigram}}(w_i)}{\sum_{w_j} P_{\text{WTMM}}(H_{w_i} | M_{w_j}) P_{\text{Unigram}}(w_j)} \end{aligned}$$

$$\text{where } P_{\text{WTMM}}(H_{w_i} | M_{w_i}) = \prod_{w \in H_{w_i}} \left[\sum_k P(w | T_k) P(T_k | M_{w_i}) \right]^{c(w, H_{w_i})}$$

- It will be a bit more time-consuming !

Comparison: PLSA vs. WTMM in Language Modeling

	WTMM	PLSA
Modeling Relationship	Words	Word and History
Model Estimation	Offline	On-the-fly
Topic Modeling	Explicit	Explicit
Parameters	$V \times K \times 2$	$V \times K + K \times D$
Prediction Ability	Yes	No

V : Vocabulary size; K : Topic number; D : Number of documents used for training

- Topic Modeling: Model topics with explicit or implicit probability distribution
- Prediction Ability: The prediction of the decoded word given the search history

Experimental Results: Information Retrieval (1/3)

- Supervised Training of PLSA and WTMM

Table II. Retrieval results on the TDT-2 development set, achieved, respectively, with the WTMM and the PLSA trained in a supervised manner.

No. Latent Topic		2	4	8	16	32	64	128	256
WTMM-S	TD	0.6505	0.6630	0.6887	0.7177	0.7351	0.7532	0.7672	0.7852
	SD	0.5731	0.5962	0.6186	0.6730	0.6864	0.7387	0.7558	0.7858
PLSA-S	TD	0.6362	0.6721	0.6750	0.6769	0.6823	0.6930	0.7243	0.7794
	SD	0.5759	0.5894	0.5918	0.5988	0.6255	0.6528	0.6652	0.6591

Experimental Results: Information Retrieval (2/3)

- Unsupervised Training of PLSA and WTMM

Table III: Retrieval results on the TDT-2 development set, achieved, respectively, with WTMM and PLSA trained in an unsupervised manner.

No. Latent Topic		2	4	8	16	32	64	128	256
WTMM-U	TD	0.6336	0.6350	0.6359	0.6368	0.6382	0.6386	0.6395	0.6287
	SD	0.5693	0.5723	0.5734	0.5733	0.5739	0.5767	0.5737	0.5652
PLSA-U	TD	0.6277	0.6332	0.6266	0.5973	0.5949	0.6267	0.6041	0.5878
	SD	0.5545	0.5659	0.5681	0.5503	0.5534	0.5664	0.5484	0.5831

Experimental Results: Information Retrieval (3/3)

- Other Retrieval Models
 - HMMs are trained with supervision

Table IV. Retrieval results on the TDT-2 development set, achieved with HMM, VSM and LSA, respectively.

Retrieval Model	HMM/Unigram	HMM/Bigram	VSM	LSA
TD	0.6327	0.5427	0.5548	0.5510
SD	0.5658	0.4803	0.5122	0.5310

Experimental Results: Language Model Adaptation (1/2)

- Experiments were conducted on the MATBN Broadcast News Corpus

Table VII. CER (%) and perplexity (PP) results, achieved by using WTMM and PLSA, respectively, for language model adaptation.

		CER (%)		PP	
Baseline (Background Trigram Model)		15.22%		752.49	
		WTMM		PLSA	
Adaptation Corpus	No. Latent Topic	CER (%)	PP	CER (%)	PP
Texts	16	14.77	566.10	14.83	588.51
	32	14.69	553.88	14.73	571.46
	64	14.60	540.62	14.58	552.80
	128	14.44	524.15	14.53	527.41
	256	14.38	508.29	14.47	510.20
Automatic Transcripts	16	14.87	574.60	14.99	591.21
	32	14.90	568.76	14.92	580.80
	64	14.85	564.56	14.82	569.93
	128	14.81	563.25	14.87	562.45
	256	14.96	567.53	14.92	565.85

Experimental Results: Language Model Adaptation (2/2)

- Hybrid of PLSA and WTMM

Table IX. CER (%) and perplexity (PP) results, achieved by combining WTMM with PLSA.

		WTMM+ PLSA	
Adaptation Corpus	No. Latent Topic	CER (%)	PP
Texts	16	14.78	551.62
	32	14.61	530.00
	64	14.47	506.19
	128	14.34	474.87
	256	14.21	449.09
Automatic Transcripts	16	14.95	546.54
	32	14.97	531.40
	64	14.82	516.82
	128	14.81	504.77
	256	14.94	502.06
Texts + Automatic Transcripts	256	14.10	441.07

- No apparent CER improvement is observed !


LSA: SVDLIBC

- Doug Rohde's SVD C Library version 1.3 is based on the [SVDPACKC](#) library
- Download it at <http://tedlab.mit.edu/~dr/>

LSA: Exercise (1/4)

- Given a sparse term-document matrix
 - E.g., 4 terms and 3 docs

	Doc		
Term	2.3	0.0	4.2
	0.0	1.3	2.2
	3.8	0.0	0.5
	0.0	0.0	0.0



Row #Tem	Col. # Doc	Nonzero entries
4	3	6
2		2 nonzero entries at Col 0
0	2.3	Col 0, Row 0
2	3.8	Col 0, Row 2
1		1 nonzero entry at Col 1
1	1.3	Col 1, Row 1
3		3 nonzero entry at Col 2
0	4.2	Col 2, Row 0
1	2.2	Col 2, Row 1
2	0.5	Col 2, Row 2

- Each entry can be weighted by *TFxIDF* score

- Perform SVD to obtain term and document vectors represented in the latent semantic space
- Evaluate the information retrieval capability of the LSA approach by using varying sizes (e.g., 100, 200, ..., 600 etc.) of LSA dimensionality

LSA: Exercise (2/4)

- Example: term-document matrix

```
Indexing      Nonzero
Term no. Doc no. entries
51253 2265 218852
77
508 7.725771
596 16.213399
612 13.080868
709 7.725771
713 7.725771
744 7.725771
1190 7.725771
1200 16.213399
1259 7.725771
.....
```

- SVD command (IR_svd.bat)

```
svd -r st -o LSA100 -d 100 Term-Doc-Matrix
```

Annotations for the command:

- Red arrow from "st" to "sparse matrix input"
- Red arrow from "LSA100" to "prefix of output files"
- Red arrow from "100" to "No. of reserved eigenvectors"
- Red arrow from "Term-Doc-Matrix" to "name of sparse matrix input"

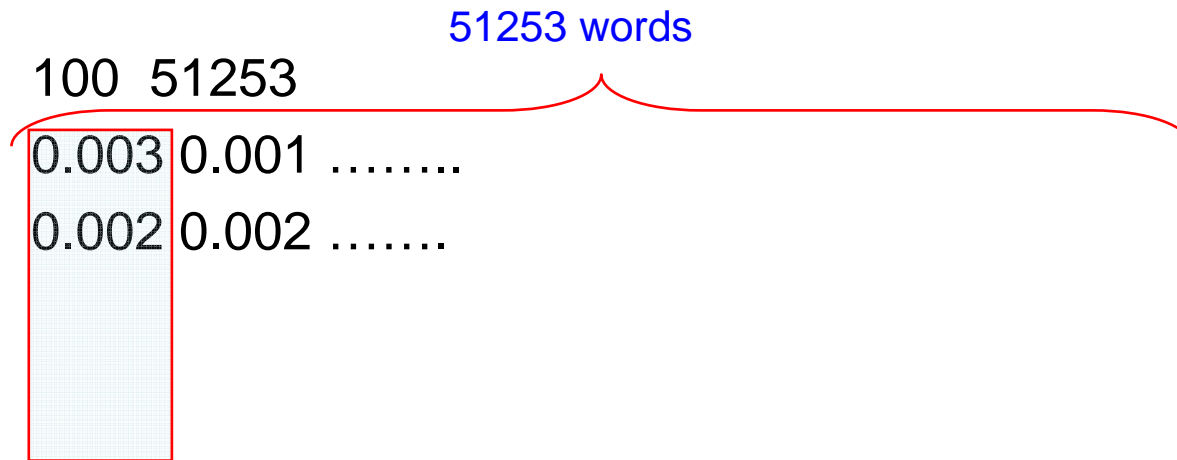
A large blue arrow points from the command to the output files.

output

- LSA100-Ut
- LSA100-S
- LSA100-Vt

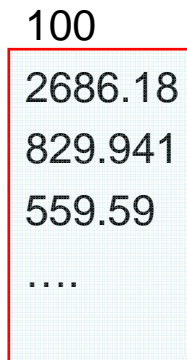
LSA: Exercise (3/4)

- **LSA100-Ut**



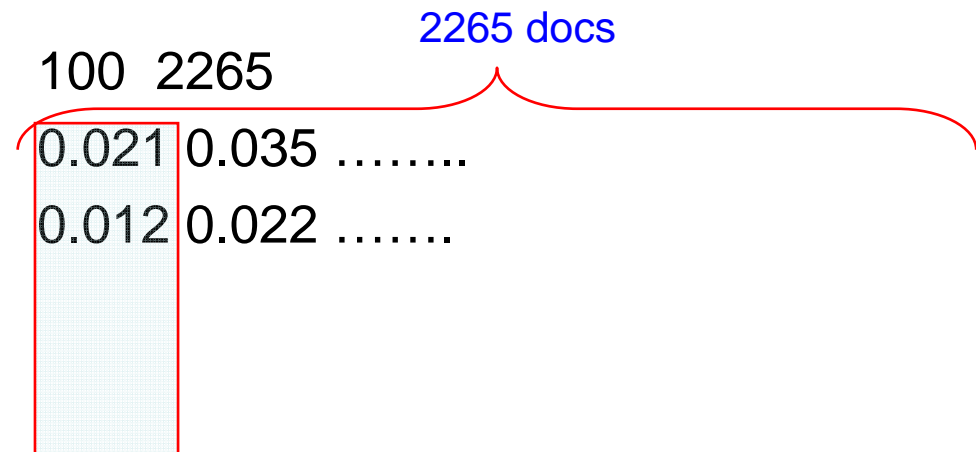
word vector (u^T): 1x100

- **LSA100-S**



100 eigenvalues

- **LSA100-Vt**



doc vector (v^T): 1x100

LSA: Exercise (4/4)

- Fold-in a new $m \times 1$ query vector

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

Just like a row of V

Query represented by the weighted sum of its constituent term vectors

The separate dimensions are differentially weighted

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

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

References

- G.W.Furnas, S. Deerwester, S.T. Dumais, T.K. Landauer, R. Harshman, L.A. Streeter, K.E. Lochbaum, “*Information Retrieval using a Singular Value Decomposition Model of Latent Semantic Structure,*” ACM SIGIR Conference on R&D in Information Retrieval , 1988
- J.R. Bellegarda, “*Latent semantic mapping,*” IEEE Signal Processing Magazine, September 2005
- J.R. Bellegarda. *Latent Semantic Mapping: Principles and Applications.* Morgan and Claypool, 2007
- T. Hofmann, “*Unsupervised learning by probabilistic latent semantic analysis,*” Machine Learning 42, 2001
- H.-S. Chiu, B. Chen, “*Word topical mixture models for dynamic language model adaptation,*” ICASSP2007
- D.M. Blei, A.Y.Ng, M. I. Jordan, “*Latent Dirichlet Allocation,*” Journal of Machine Learning Research, 2003
- W. Kim, S. Khudanpur, “*Lexical triggers and latent semantic analysis for cross-lingual language model adaptation,*” ACM TALIP, 3(2), 2004
- D. Gildea, T. Hofmann, “Topic-based language models using EM,” Eurospeech1999
- L. K. Saul and F. C. N. Pereira, “*Aggregate and mixed-order Markov models for statistical language processing,*” EMNLP1997

