

Language Models for Information Retrieval

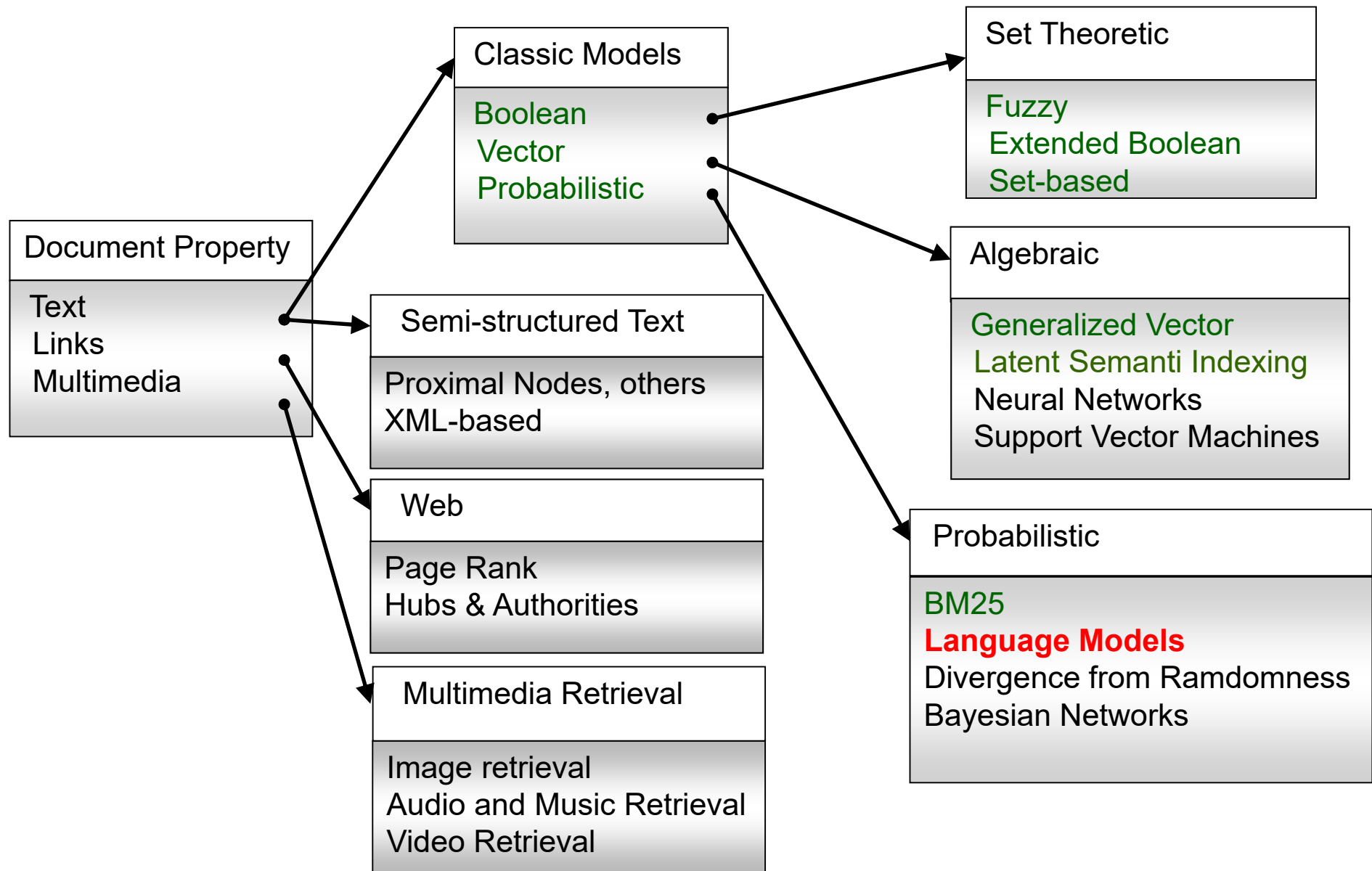
Berlin Chen

Department of Computer Science & Information Engineering
National Taiwan Normal University

References:

1. W. B. Croft and J. Lafferty (Editors). Language Modeling for Information Retrieval. July 2003
2. X. Liu and W.B. Croft, Statistical Language Modeling For Information Retrieval, the Annual Review of Information Science and Technology, vol. 39, 2005
3. Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, Introduction to Information Retrieval, Cambridge University Press, 2008. (Chapter 12)
4. D. A. Grossman, O. Frieder, Information Retrieval: Algorithms and Heuristics, Springer, 2004 (Chapter 2)
5. C.X. Zhai. Statistical Language Models for Information Retrieval (Synthesis Lectures Series on Human Language Technologies). Morgan & Claypool Publishers, 2008

Taxonomy of Classic IR Models



Statistical Language Models (1/2)

- A probabilistic mechanism for “generating” a piece of text
 - Define a distribution over all possible word sequences

$$W = w_1 w_2 \dots w_L$$

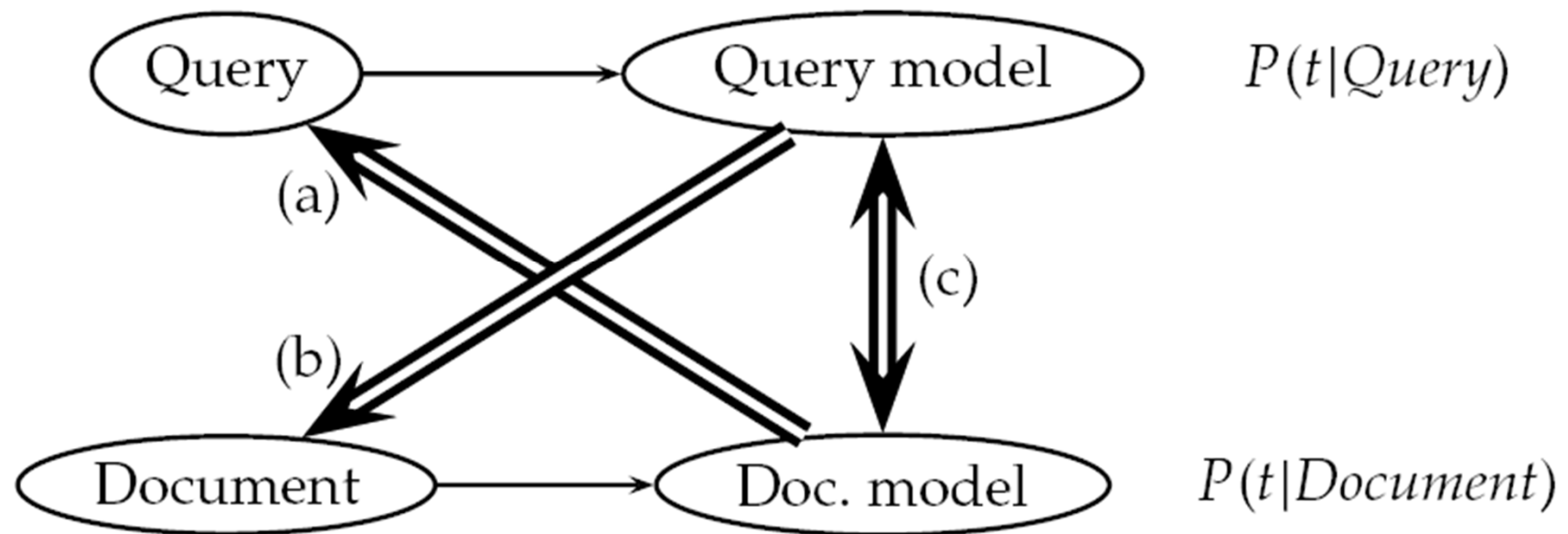
$$P(W) = ?$$

- Used LM to quantify the acceptability of a given word sequence
- What is LM Used for ?
 - Speech recognition
 - Spelling correction
 - Handwriting recognition
 - Optical character recognition
 - Machine translation
 - Document classification and routing
 - Information retrieval ...

Statistical Language Models (2/2)

- (Statistical) language models (LM) have been widely used for speech recognition and language (machine) translation for more than thirty years
- However, their use for information retrieval started only in 1998 [Ponte and Croft, SIGIR 1998]
 - Basically, a query is considered generated from an “ideal” document that satisfies the information need
 - The system’s job is then to estimate the likelihood of each document in the collection being the ideal document and rank then accordingly (in decreasing order)

Three Ways of Developing LM Approaches for IR



(a) Query likelihood

(b) Document likelihood

(c) Model comparison

literal term matching
or concept matching

Query-Likelihood Language Models

- Criterion: Documents are ranked based on Bayes (decision) rule

$$P(D|Q) = \frac{P(Q|D)P(D)}{P(Q)}$$

- $P(Q)$ is the same for all documents, and can be ignored
- $P(D)$ might have to do with authority, length, genre, etc.
 - There is no general way to estimate it
 - Can be treated as uniform across all documents
- Documents can therefore be ranked based on $P(Q|D)$ (or denoted as $P(Q|M_D)$) document model
 - The user has a prototype (ideal) document in mind, and generates a query based on words that appear in this document
 - A document D is treated as a model M_D to predict (generate) the query

Another Criterion: Maximum Mutual Information

- Documents can be ranked based their mutual information with the query (in decreasing order)

$$\begin{aligned} MI(Q, D) &= \log \frac{P(Q, D)}{P(Q)P(D)} \\ &= \log P(Q|D) - \underbrace{\log P(Q)} \end{aligned}$$

being the same for all documents,
and hence can be ignored

- Document ranking by mutual information (MI) is equivalent that by likelihood

$$\arg \max_D MI(Q, D) \stackrel{\text{rank}}{=} \arg \max_D P(Q|D)$$

Yet Another Criterion: Minimum KL Divergence

- Documents are ranked by Kullback-Leibler (KL) divergence (in increasing order)

$$\begin{aligned}
 KL(Q||D) &= \sum_w P(w|Q) \log \frac{P(w|Q)}{P(w|D)} \\
 &= \sum_w P(w|Q) \log P(w|Q) - \sum_w P(w|Q) \log P(w|D)
 \end{aligned}$$

Query model Document model

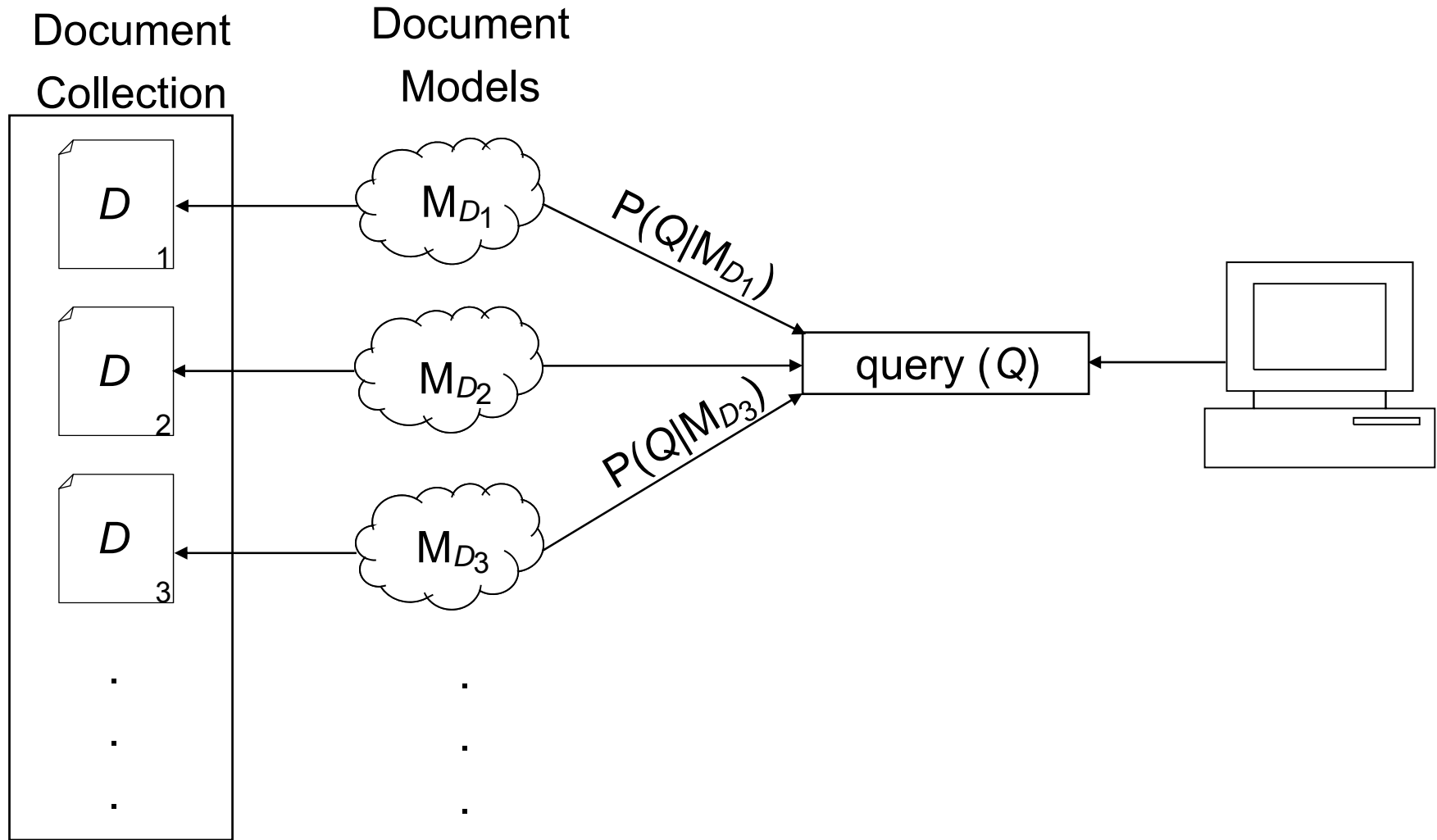
The same for all document
=> can be disregarded

Cross entropy between the language models of a query and a document

Equivalent to ranking **in decreasing order** of

$$\begin{aligned}
 &\sum_w P(w|Q) \log P(w|D) \quad \text{Relevant documents are deemed to have lower cross entropies} \\
 &= \sum_w^{\text{rank}} c(w, Q) \log P(w|D) = P(Q|D)
 \end{aligned}$$

Schematic Depiction for Query-Likelihood Approach



Building Document Models: n -grams

- Multiplication (Chain) rule

$$P(w_1 w_2 \dots w_L) = P(w_1) P(w_2 | w_1) P(w_3 | w_1 w_2) \dots P(w_L | w_1 w_2 \dots w_{L-1})$$

- Decompose the probability of a sequence of events into the probability of each successive events conditioned on earlier events

- n -gram assumption

- **Unigram**

$$P(w_1 w_2 \dots w_L) = P(w_1) P(w_2) P(w_3) \dots P(w_L)$$

- Each word occurs independently of the other words
- The so-called “bag-of-words” model (e.g., how to distinguish “street market” from “market street”)

- Bigram

$$P(w_1 w_2 \dots w_L) = P(w_1) P(w_2 | w_1) P(w_3 | w_2) \dots P(w_L | w_{L-1})$$

- Most language-modeling work in IR has used unigram models
 - IR does not directly depend on the structure of sentences

Unigram Model (1/4)

- The likelihood of a query $Q = w_1 w_2 \dots w_L$ given a document D

$$\begin{aligned} P(Q|M_D) &= P(w_1|M_D)P(w_2|M_D)\cdots P(w_L|M_D) \\ &= \prod_{i=1}^L P(w_i|M_D) \end{aligned}$$

- **Words are conditionally independent of each other given the document**

- How to estimate the probability of a (query) word given the document $P(w|M_D)$?

- Assume that words follow a **multinomial distribution** given the document

permutation is considered here

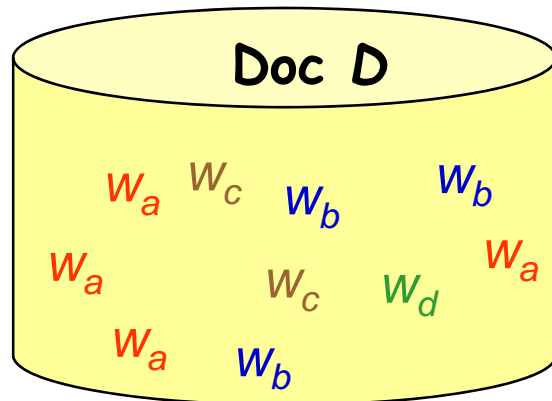
$$P(c(w_1), \dots, c(w_V) | M_D) = \frac{(\sum_{j=1}^V c(w_j))!}{\prod_{i=1}^V (c(w_i))!} \prod_{i=1}^V \lambda_{w_i}^{c(w_i)}$$

where $c(w_i)$: the number of times a word occurs

$$\lambda_{w_i} = P(w_i | M_D), \quad \sum_{i=1}^V \lambda_{w_i} = 1$$

Unigram Model (2/4)

- Use each document itself a sample for estimating its corresponding unigram (multinomial) model
 - If Maximum Likelihood Estimation (MLE) is adopted



$$\hat{P}(w_i | M_D) = \frac{c(w_i, D)}{|D|}$$

where

$c(w_i, D)$: number of times w_i occurs in D

$|D|$: length of D , $\sum_i c(w_i, D) = |D|$

$$P(w_b | M_D) = 0.3$$

$$P(w_c | M_D) = 0.2$$

$$P(w_d | M_D) = 0.1$$

$$P(w_e | M_D) = 0.0$$

$$P(w_f | M_D) = 0.0$$

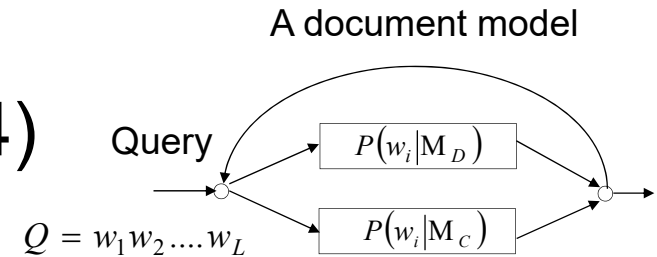
The zero-probability problem

If w_e and w_f do not occur in D

then $P(w_e | M_D) = P(w_f | M_D) = 0$

This will cause a problem in predicting the query likelihood (See the equation for the query likelihood in the preceding slide)

Unigram Model (3/4)



- Smooth the document-specific unigram model with a collection model (**two states, or a mixture of two multinomials**)

$$P(Q|M_D) = \prod_{i=1}^L [\lambda \cdot P(w_i|M_D) + (1 - \lambda) \cdot P(w_i|M_C)]$$

- The role of the collection unigram model $P(w_i|M_C)$
 - Help to solve zero-probability problem
 - Help to differentiate the contributions of different missing terms in a document (**global information like IDF ?**)

$$P(w_i|M_C) = \frac{c(w_i, Collection)}{\sum_{w_l} c(w_l, Collection)} \quad \text{or} \quad \frac{\frac{N}{n_i}}{\sum_{w_l} \frac{N}{n_l}}$$

N : number of doc in the collection
 n_i : number of doc in the collection containing w_i
Normalized doc freq

- The collection unigram model can be estimated in a similar way as what we do for the document-specific unigram model

Unigram Model (4/4)

- An evaluation on the Topic Detection and Tracking (TDT) corpora
 - Language Model

mAP		Unigram	Unigram+Bigram
TDT2	TQ/TD	0.6327	0.5427
	TQ/SD	0.5658	0.4803
TDT3	TQ/TD	0.6569	0.6141
	TQ/SD	0.6308	0.5808

- Vector Space Model

mAP		Unigram	Unigram+Bigram
TDT2	TQ/TD	0.5548	0.5623
	TQ/SD	0.5122	0.5225
TDT3	TQ/TD	0.6505	0.6531
	TQ/SD	0.6216	0.6233

$$P_{Unigram}(Q|M_D) = \prod_{i=1}^L [\lambda \cdot P(w_i|M_D) + (1 - \lambda) \cdot P(w_i|M_C)]$$

$$P_{Unigram+Bigram}(Q|M_D) = \prod_{i=1}^L [\lambda_1 \cdot P(w_i|M_D) + \lambda_2 \cdot P(w_i|M_C) + \lambda_3 \cdot P(w_i|w_{i-1},M_D) + (1 - \lambda_1 - \lambda_2 - \lambda_3) \cdot P(w_i|w_{i-1},M_C)]$$

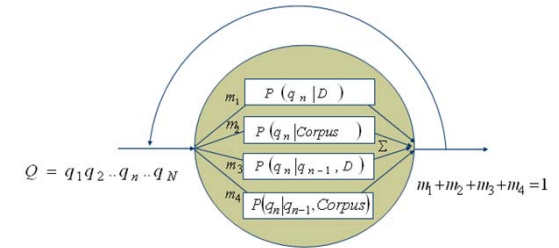
- Consideration of contextual information (Higher-order language models, e.g., bigrams) will not always lead to improved performance

Training Mixture Weights of LMs

- Expectation-Maximization (EM) Training

- The weights are tied among the documents

- E.g. m_1 of Type I HMM can be trained using the following equation:



$$\hat{m}_1 = \frac{\sum_{Q \in [TrainSet]_Q} \sum_{D \in [Doc]_{R \text{ to } Q}} \sum_{q_n \in Q} \left[\frac{m_1 P(q_n | D)}{m_1 P(q_n | D) + m_2 P(q_n | Corpus)} \right]}{\sum_{Q \in [TrainSet]_Q} |Q| \cdot |[Doc]_{R \text{ to } Q}|}$$

Annotations in the diagram:

- the new weight**: Points to \hat{m}_1 .
- the old weight**: Points to $m_1 P(q_n | D)$ in the numerator.

- 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 $|[Doc]_{R \text{ to } Q}|$ is the total number of docs relevant to the query Q

Discriminative Training of LMs (1/5)

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

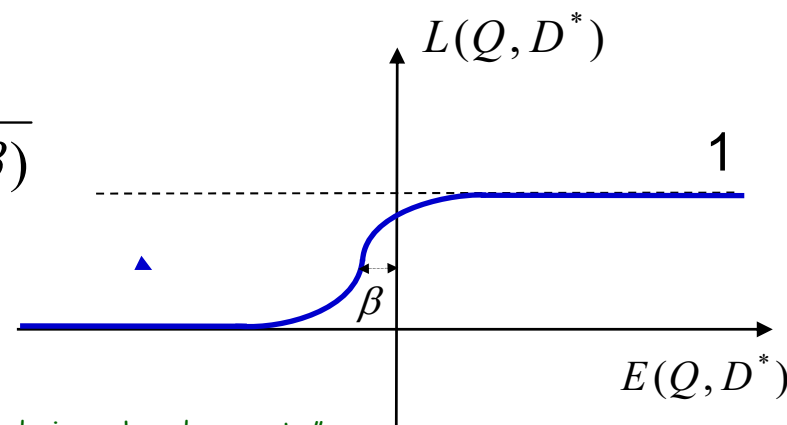
$$E(Q, D^*) = \frac{1}{|Q|} \left[-\log P(Q|D^* \text{ is } R) + \max_{D'} \log P(Q|D' \text{ is not } R) \right]$$

“>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



Discriminative Training of LMs (2/5)

- 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 Transformation, (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)

Gradient descent

$$\tilde{m}_1(i+1) = \tilde{m}_1(i) - \varepsilon(i) \cdot \frac{\partial L(Q, D^*)}{\partial \tilde{m}_1} \Big|_{D^* = D^*(i)}$$

- Where,

$$\nabla_{D^*, \tilde{m}_1} = \varepsilon(i) \cdot \frac{\partial L(Q, D^*)}{\partial \tilde{m}_1}$$

$$= \varepsilon(i) \cdot \frac{\partial L(Q, D^*)}{\partial E(Q, D^*)} \cdot \frac{\partial E(Q, D^*)}{\partial \tilde{m}_1},$$

$$\frac{\partial L(Q, D^*)}{\partial E(Q, D^*)} = \alpha \cdot L(Q, D^*) \cdot [1 - L(Q, D^*)]$$

Discriminative Training of LMs (3/5)

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

$$\begin{aligned}
 \frac{\partial E(Q, D^*)}{\partial \tilde{m}_1} &= \frac{-1}{|Q|} \frac{\partial \left\{ \sum_{q_n \in Q} \log \left[\frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | D^*) + \frac{e^{\tilde{m}_2}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | Corpus) \right] \right\}}{\partial \tilde{m}_1} \\
 &= \frac{-1}{|Q|} \sum_{q_n \in Q} \left\{ \frac{\frac{-e^{\tilde{m}_1}}{(e^{\tilde{m}_1} + e^{\tilde{m}_2})^2} [e^{\tilde{m}_1} P(q_n | D^*) + e^{\tilde{m}_2} P(q_n | Corpus)] + \frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | D^*)}{\frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | D^*) + \frac{e^{\tilde{m}_2}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | Corpus)}} \right\} \\
 &= \frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} - \frac{1}{|Q|} \sum_{q_n \in Q} \left\{ \frac{\frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | D^*)}{\frac{e^{\tilde{m}_1}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | D^*) + \frac{e^{\tilde{m}_2}}{e^{\tilde{m}_1} + e^{\tilde{m}_2}} P(q_n | Corpus)}} \right\} \\
 &= - \left[-m_1 + \frac{1}{|Q|} \sum_{q_n \in Q} \frac{m_1 P(q_n | D^*)}{m_1 P(q_n | D^*) + m_2 P(q_n | Corpus)} \right],
 \end{aligned}$$

Note :

$$[\log f(x)]' = \frac{1}{f(x)} f'(x)$$

$$[f(x)g(x)]' = f'(x)g(x) + f(x)g'(x)$$

$$\left[\frac{f(x)}{g(x)} \right]' = \frac{f'(x)g(x) - f(x)g'(x)}{g^2(x)}$$

Discriminative Training of LMs (4/5)

- 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

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

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

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

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

the old weight

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

Discriminative Training of LMs (5/5)

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

$$\nabla_{D^*, \tilde{m}_1}(i) = -\varepsilon(i) \cdot \alpha \cdot L(Q, D^*) \cdot [1 - L(Q, D^*)] \cdot \left[-m_1(i) + \frac{1}{|Q|} \sum_{q_n \in Q} \frac{m_1(i)P(q_n|D^*)}{m_1(i)P(q_n|D^*) + m_2(i)P(q_n|Corpus)} \right]$$

$$m_1(i+1) = \frac{m_1(i) \cdot e^{-\nabla_{D^*, \tilde{m}_1}(i)}}{m_1(i) \cdot e^{-\nabla_{D^*, \tilde{m}_1}(i)} + m_2(i) \cdot e^{-\nabla_{D^*, \tilde{m}_2}(i)}}$$

- m_2 can be updated in the similar way

Statistical Translation Model (1/2)

Berger & Lafferty (1999)

- A query Q is viewed as a translation or distillation from a document D
 - That is, the similarity measure is computed by estimating the probability that the query would have been generated as a translation of that document

$$\text{sim}(Q, D) = P(Q|D) = \prod_{q \in Q} P_{\text{Trans}}(q|D)^{c(q,Q)} = \prod_{q \in Q} \sum_{w \in D} \underbrace{[P(q|w)P(w|D)]}_{\text{word-to-word translation}}^{c(q,Q)}$$

- Assumption of context-independence (the ability to handle the ambiguity of word senses is limited)
- However, it has the capability of handling the issues of synonymy (multiple terms having similar meaning) and polysemy (the same term having multiple meanings)

$$\hat{P}(Q|D) = \prod_{q \in Q} [\lambda P_{\text{Trans}}(q|D) + (1 - \lambda)(q|M_c)]^{c(q,Q)}$$

Statistical Translation Model (2/2)

- **Weakness** of the statistical translation model
 - The need of a large collection of training data for estimating translation probabilities, and inefficiency for ranking documents
- Jin et al. (2002) proposed a “**Title Language Model**” approach to capture the intrinsic document to query translation patterns
 - Queries are more like titles than documents (queries and titles both tend to be very short and concise descriptions of information, and created through a similar generation process)
 - Train the statistical translation model based on the document-title pairs in the whole collection

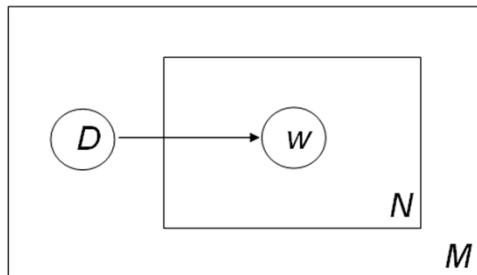
$$M^* = \arg \max_M \prod_{j=1}^N P_M(T_j | D_j) = \arg \max_M \prod_{j=1}^N \prod_{t \in T_j} P_M(t | D_j)$$
$$\approx \arg \max_M \prod_{j=1}^N \prod_{t \in T_j} \left(\sum_{w \in D} [P_M(t | w) P(w | D)]^{c(t, T_j)} \right)$$

Probabilistic Latent Semantic Analysis (PLSA)

Hofmann (1999)

- Also called The Aspect Model, Probabilistic Latent Semantic Indexing (PLSI)
 - Graphical Model Representation (a kind of Bayesian Networks)

Language (unigram) model



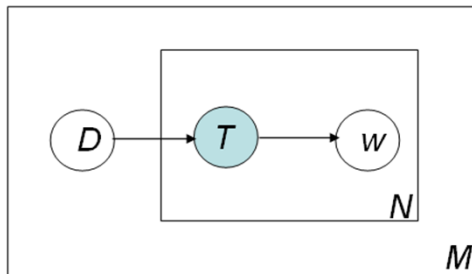
$$\text{sim}(Q, D) = P(D|Q) = \frac{P(Q|D)P(D)}{P(Q)}$$

$$\propto P(Q|D)P(D)$$

$$\approx P(Q|D)$$

$$= \prod_{w \in Q} [\lambda \cdot P(w|M_D) + (1 - \lambda) \cdot P(w|M_C)]^{c(w,Q)}$$

PLSA



$$\text{sim}(Q, D) = P(Q|D) = \prod_{w \in Q} P(w|D)^{c(w,Q)}$$

$$= \prod_{w \in Q} \left[\sum_{k=1}^K P(w, T_k | D) \right]^{c(w,Q)}$$

$$= \prod_{w \in Q} \left[\sum_{k=1}^K P(w|T_k)P(T_k|D) \right]^{c(w,Q)}$$

N : number of distinct in the vocabulary

M : number of documents in the collection

○ : observed variable

● : latent variable

The latent variables
 \Rightarrow The unobservable class variables T_k
 (topics or domains)

PLSA: Formulation

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

PLSA: Assumptions

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

$$\text{sim}(Q, D) = P(Q|D) = \prod_w P(w|D)^{c(w, Q)}$$

- **Conditional independent:** the doc D and word w are independent conditioned on the state of the associated latent variable T_k

$$P(w, D|T_k) \approx P(w|T_k)P(D|T_k)$$

$$\begin{aligned} P(w|D) &= \sum_{k=1}^K P(w, T_k|D) = \sum_{k=1}^K \frac{P(w, D, T_k)}{P(D)} = \sum_{k=1}^K \frac{P(w, D|T_k)P(T_k)}{P(D)} \\ &= \sum_{k=1}^K \frac{P(w|T_k)P(D|T_k)P(T_k)}{P(D)} = \sum_{k=1}^K \frac{P(w|T_k)P(T_k, D)}{P(D)} \\ &= \sum_{k=1}^K P(w|T_k)P(T_k|D) \end{aligned}$$

PLSA: Training (1/2)

- Probabilities are estimated by maximizing the collection likelihood using the Expectation-Maximization (EM) algorithm

$$\begin{aligned} L_C &= \sum_D \sum_w c(w, D) \log P(w|D) \\ &= \sum_D \sum_w c(w, D) \log \left[\sum_{T_k} P(w|T_k) P(T_k|D) \right] \end{aligned}$$

EM tutorial:

- Jeff A. Bilmes ["A Gentle Tutorial of the EM Algorithm and its Application to Parameter Estimation for Gaussian Mixture and Hidden Markov Models,"](#) U.C. Berkeley TR-97-021

PLSA: Training (2/2)

- E (expectation) step

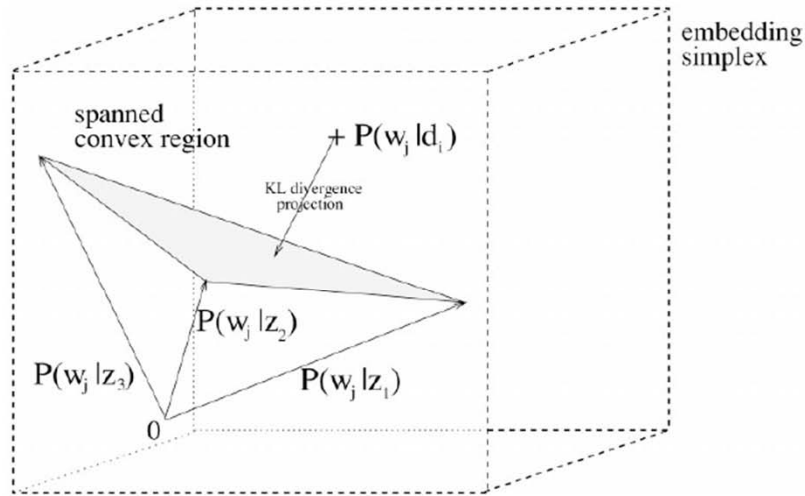
$$P(T_k | w, D) = \frac{P(w | T_k)P(T_k | D)}{\sum_{T_k} P(w | T_k)P(T_k | D)}$$

- M (Maximization) step

$$\hat{P}(w | T_k) = \frac{\sum_D c(w, D)P(T_k | w, D)}{\sum_w \sum_D c(w, D)P(T_k | w, D)}$$

$$\hat{P}(T_k | D) = \frac{\sum_w c(w, D)P(T_k | w, D)}{\sum_{w'} c(w', D)}$$

PLSA: Latent Probability Space (1/2)



Dimensionality $K=128$ (latent classes)

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

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

medical imaging

image sequence

context of contour

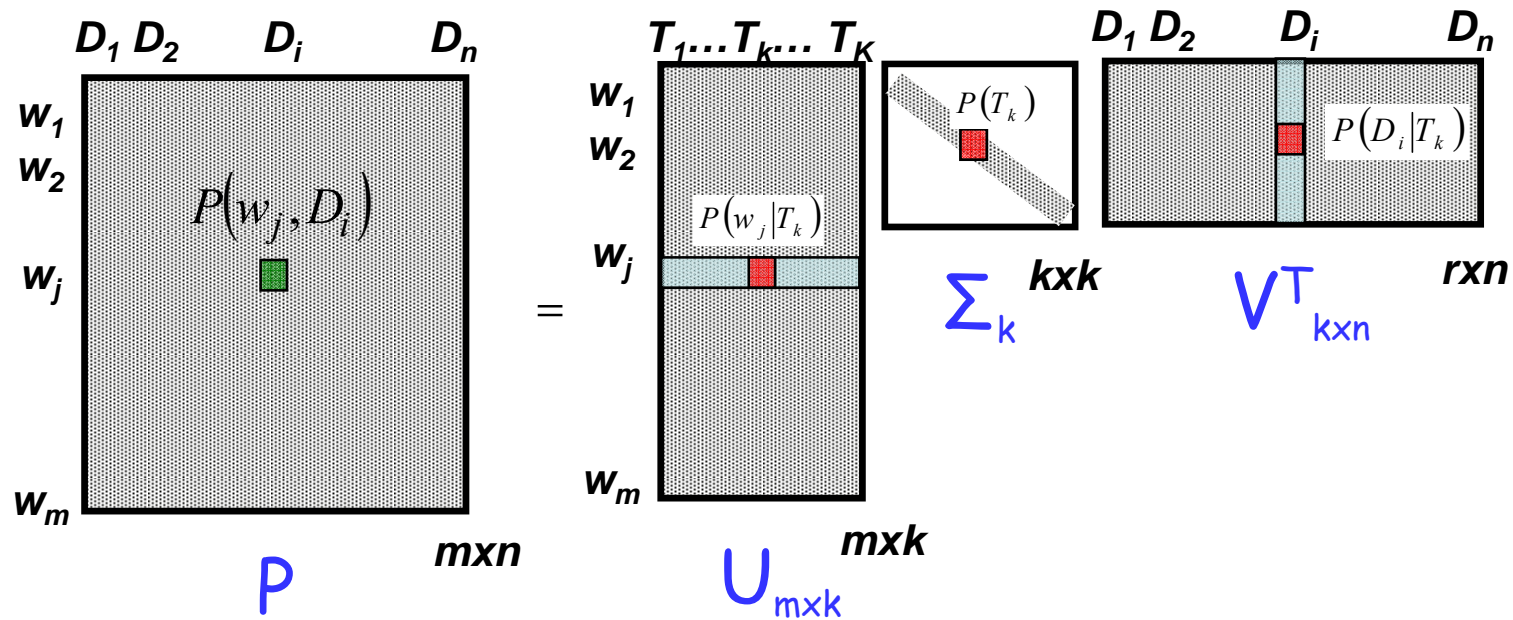
phonetic

analysis

boundary detection segmentation

$$\begin{aligned}
 P(w_j, D_i) &= \sum_{T_k} P(w_j, T_k, D_i) = \sum_{T_k} P(w_j | T_k, D_i) P(T_k, D_i) \\
 &= \sum_{T_k} P(w_j | T_k) P(T_k) P(D_i | T_k) \\
 \mathbf{P}(\mathbf{W}, \mathbf{D}) &= \hat{\mathbf{U}} : (P(w_j | T_k))_{j,k} \cdot \hat{\mathbf{\Sigma}} : \text{diag}(P(T_k))_k \cdot \hat{\mathbf{V}} : (P(D_i | T_k))_{i,k}
 \end{aligned}$$

PLSA: Latent Probability Space (2/2)



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

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.

PLSA: Experiment Results (1/4)

- Experimental Results

- Two ways to smoothen empirical distribution with PLSA

- Combine the cosine score with that of the vector space model (so does LSA)

PLSA-U* (See next slide)

- Combine the multinomials individually

PLSA-Q*

$$P_{PLSA}(w|D) = \sum_{k=1}^K P(w|T_k)P(T_k|D)$$

$$P_{PLSA-Q^*}(w|D) = \lambda \cdot P_{Empirical}(w|D) + (1-\lambda) \cdot P_{PLSA}(w|D)$$

$$P_{Empirical}(w|D) = \frac{c(w,D)}{c(D)}$$

$$P_{PLSA-Q^*}(Q|D) = \prod_{w \in Q} \left(\lambda \cdot P_{Empirical}(w|D) + (1-\lambda) \cdot P_{PLSA}(w|D) \right)^{c(w,D)}$$

Both provide almost identical performance

- The performance of PLSA ($P_{PLSA}(w|D)$) is not promising when being used alone

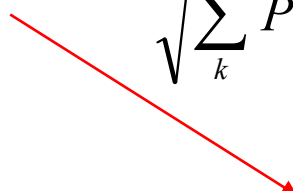
PLSA: Experiment Results (2/4)

PLSA-U*

- Use the low-dimensional representation $P(T_k | Q)$ and $P(T_k | D)$ (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 re-weight the different model components (dimensions) by

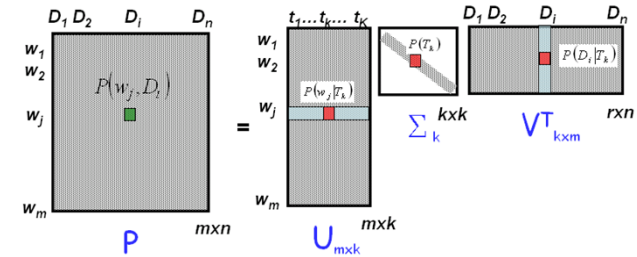
$$R_{PLSA-U^*}(Q, D) = \frac{\sum_k P(T_k | Q)P(T_k | D)}{\sqrt{\sum_k P(T_k | Q)^2} \sqrt{\sum_k P(T_k | D)^2}} \quad , \text{ where } P(T_k | Q) = \frac{\sum_{w \in Q} c(w, Q) P(T_k | w, Q)}{\sum_{w' \in Q} c(w', Q)}$$

online folded-in


$$\tilde{R}_{PLSA-U^*}(Q, D) = \lambda \cdot R_{PLSA-U^*}(Q, D) + (1 - \lambda) \cdot R_{VSM}(\vec{Q}, \vec{D})$$

PLSA: Experiment Results (3/4)

• **Why** $R_{PLSI-Q^*}(Q, D_i) = \frac{\sum_k P(T_k|Q)P(T_k|D_i)}{\sqrt{\sum_k P(T_k|Q)^2} \sqrt{\sum_k P(T_k|D_i)^2}}$?



– Reminder that in LSA, the relations between any two docs can be formulated as



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

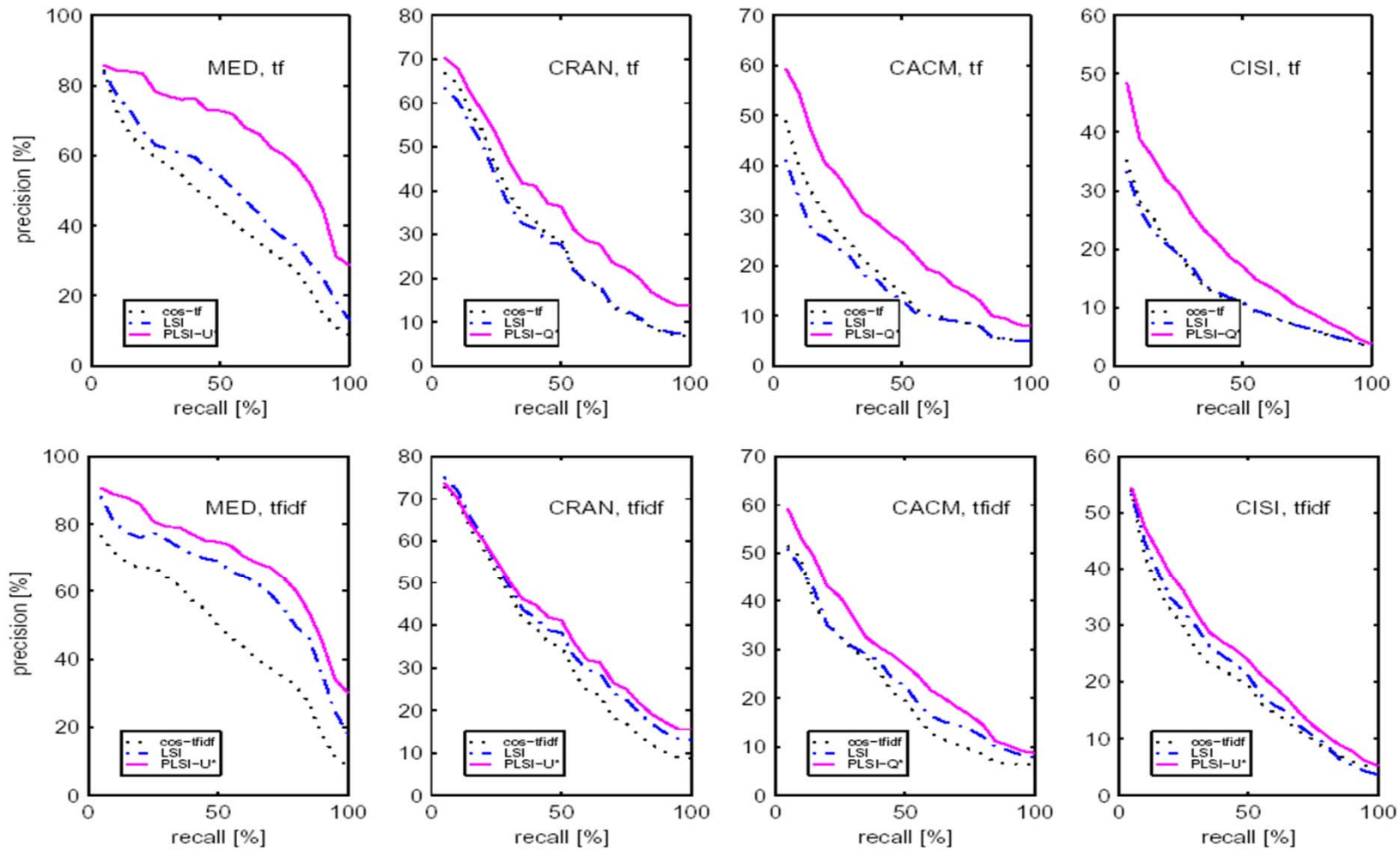
$$sim(D_i, D_s) = coine(\hat{D}_i \Sigma, \hat{D}_s \Sigma) = \frac{\hat{D}_i \Sigma^2 \hat{D}_s^T}{|\hat{D}_i \Sigma| |\hat{D}_s \Sigma|}$$

– **PLSA mimics LSA in similarity measure** \hat{D}_i and \hat{D}_s are row vectors

$$\begin{aligned} R_{PLSI-Q^*}(D_i, D_s) &= \frac{\sum_k P(D_i|T_k)P(T_k)P(T_k)P(D_s|T_k)}{\sqrt{\sum_k [P(D_i|T_k)P(T_k)]^2} \sqrt{\sum_k [P(D_s|T_k)P(T_k)]^2}} \\ &= \frac{\sum_k P(T_k|D_i)P(D_i)P(T_k|D_s)P(D_s)}{\sqrt{\sum_k [P(T_k|D_i)P(D_i)]^2} \sqrt{\sum_k [P(T_k|D_s)P(D_s)]^2}} \\ &= \frac{\sum_k P(T_k|D_i)P(T_k|D_s)}{\sqrt{\sum_k P(T_k|D_i)^2} \sqrt{\sum_k P(T_k|D_s)^2}} \end{aligned}$$

$\curvearrowright P(D_i|T_k)P(T_k) = P(T_k|D_i)P(D_i)$

PLSA: Experiment Results (4/4)



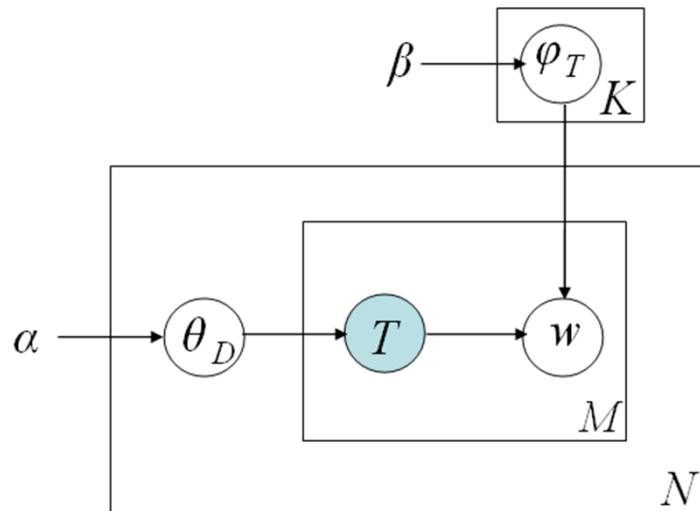
PLSA vs. LSA

- Decomposition/Approximation
 - **LSA**: **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
 - LSA: SVD decomposition
 - PLSA: EM training, is time-consuming for iterations ?
 - The model complexity of Both LSA and PLSA grows linearly with the number of training documents
 - There is no general way to estimate or predict the vector representation (of LSA) or the model parameters (of PLSA) for a newly observed document

Latent Dirichlet Allocation (LDA) (1/2)

Blei et al. (2003)

- The basic generative process of LDA closely resembles PLSA; however,
 - In PLSA, the topic mixture $P(T_k|D)$ is conditioned on each document ($P(T_k|D)$ is fixed, unknown)
 - While in LDA, the topic mixture $P(T_k|D)$ is drawn from a Dirichlet distribution, so-called the conjugate prior, ($P(T_k|D)$ is unknown and follows a probability distribution)



Process of generating a corpus with LDA

- 1) Pick a multinomial distribution φ_T for each topic T from a Dirichlet distribution with parameter β
- 2) Pick a multinomial distribution θ_D for each document D from a Dirichlet distribution with parameter α
- 3) Pick a topic $T \in \{1, 2, \dots, K\}$ from a multinomial distribution with parameter θ_D
- 4) Pick a word from a multinomial distribution with parameter φ_T

Latent Dirichlet Allocation (2/2)

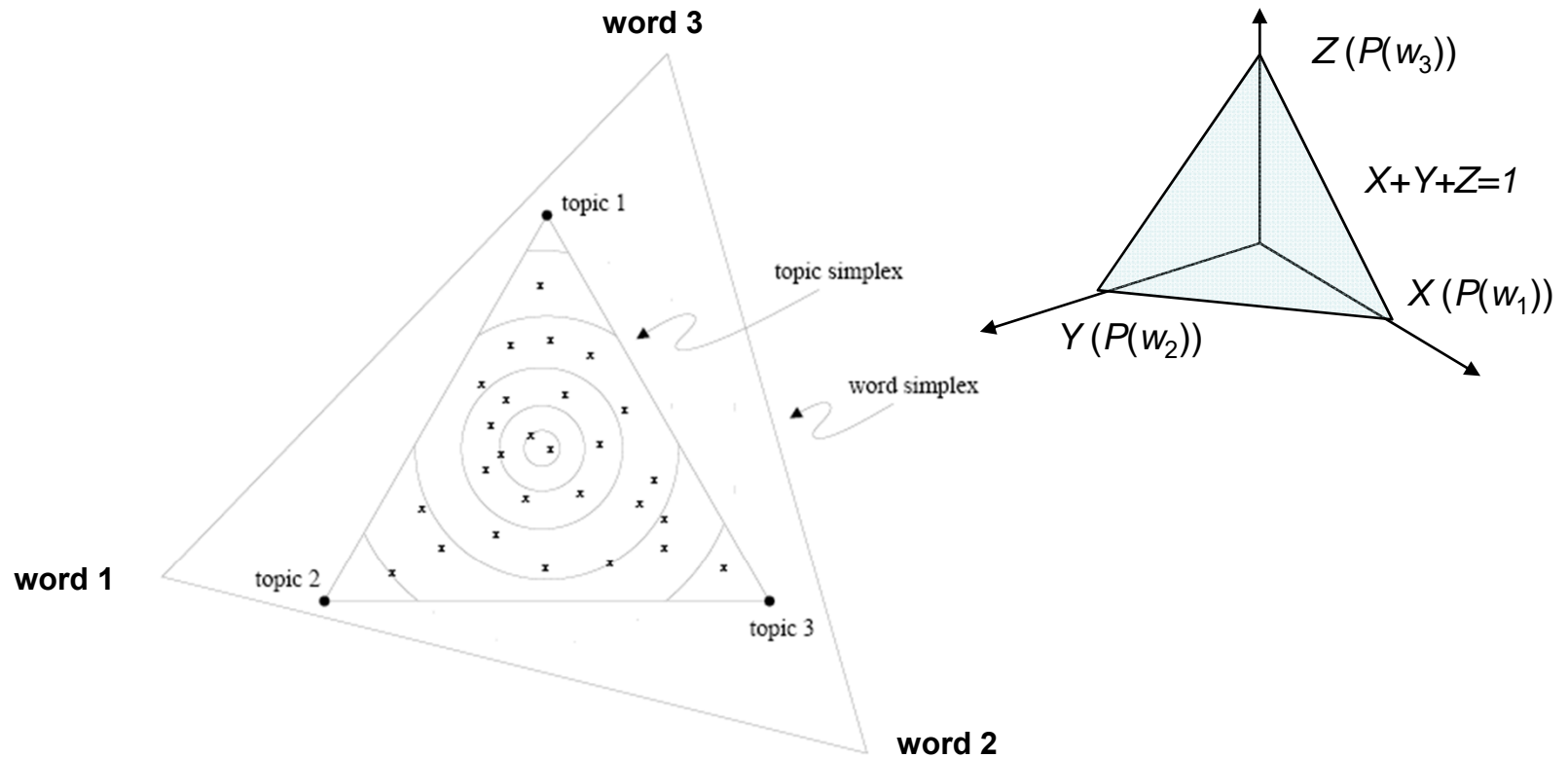


Figure 4: The topic simplex for three topics embedded in the word simplex for three words. The corners of the word simplex correspond to the three distributions where each word (respectively) has probability one. The three points of the topic simplex correspond to three different distributions over words. The mixture of unigrams places each document at one of the corners of the topic simplex. The pLSI model induces an empirical distribution on the topic simplex denoted by x . LDA places a smooth distribution on the topic simplex denoted by the contour lines.

$$L_{\text{LDA}} = \iint \prod_{k=1}^K P(\phi_T | \beta) \prod_{D \in \mathbf{D}} p(\theta_D | \alpha) \left(\prod_{i=1}^{|D|} \sum_{k=1}^K P(w_i | T_k, \phi_z) P(T_k | \theta_D) \right) d\theta d\phi$$

Word Topic Models (WTM)

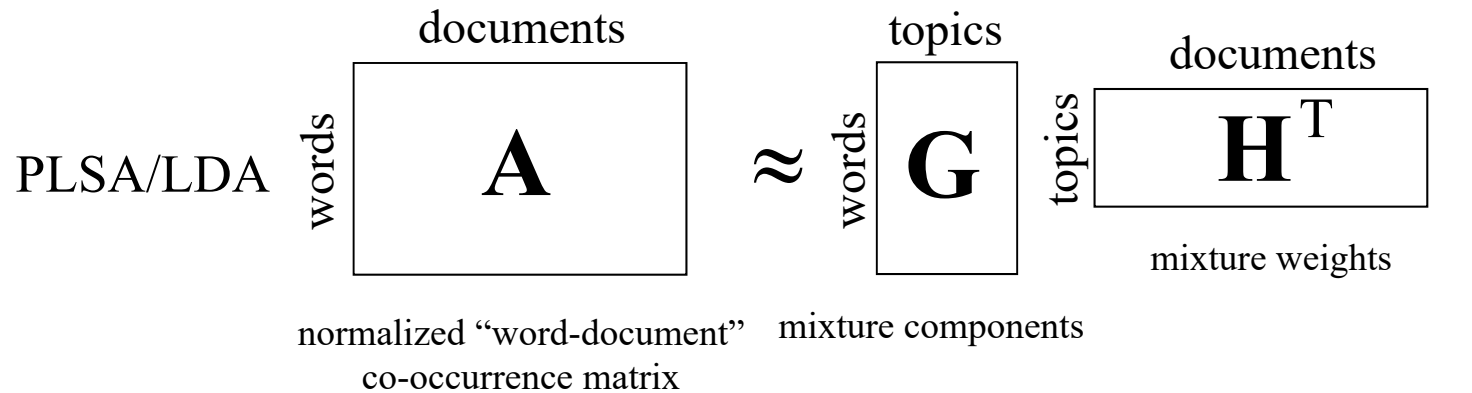
- Each word of language are treated as a word topical mixture model for predicting the occurrences of other words

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

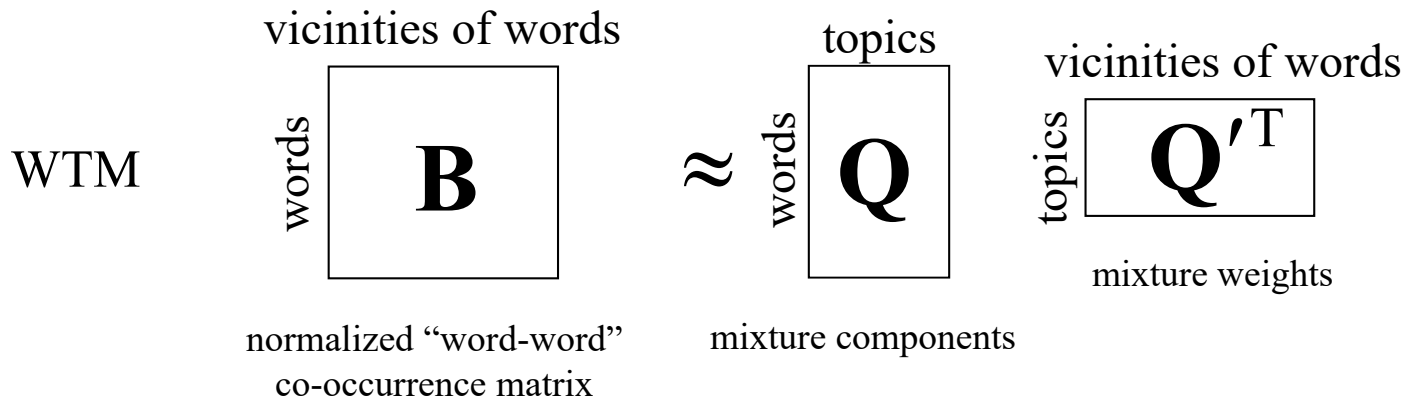
- WTM also can be viewed as a nonnegative factorization of a “word-word” matrix consisting probability entries
 - Each column encodes the vicinity information of all occurrences of a distinct word in the document collection

Comparison of WTM and PLSA/LDA

- A schematic comparison for the matrix factorizations of PLSA/LDA and WTM



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



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

WTM: Information Retrieval (1/4)

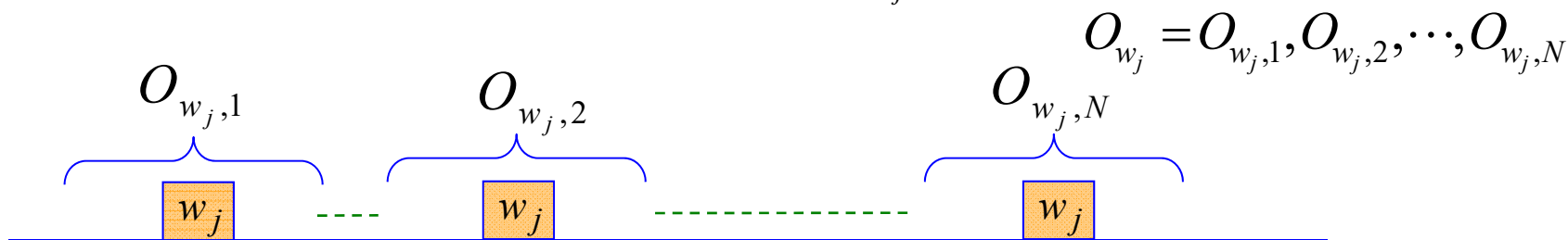
- The relevance measure between a query and a document can be expressed by

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

- Unsupervised training**

- The WTM 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

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



WTM: Information Retrieval (2/4)

- **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

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

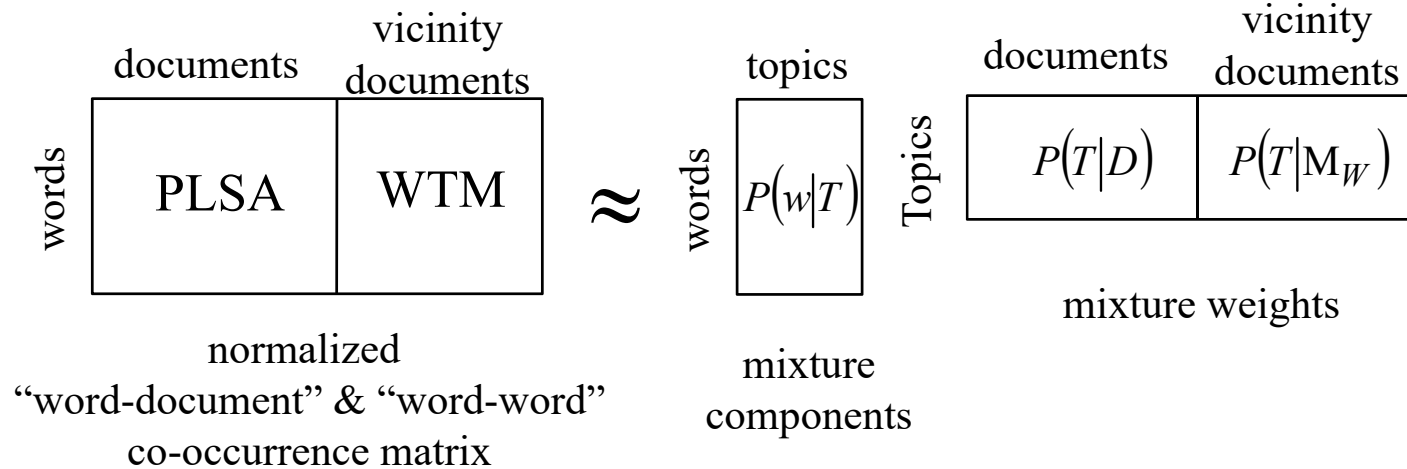
WTM: Information Retrieval (3/4)

- Example topic distributions of WTM

Topic 13		Topic 14		Topic 23	
word	weight	word	weight	word	weight
Vena (靜脈)	1.202	Land tax (土地稅)	0.704	Cholera (霍亂)	0.752
Resection (切除)	0.674	Tobacco and alcohol tax law (菸酒稅法)	0.489	Colorectal cancer (大腸直腸癌)	0.681
Myoma (肌瘤)	0.668	Tax (財稅)	0.457	Salmonella enterica (沙門氏菌)	0.471
Cephalitis (腦炎)	0.618	Amend drafts (修正草案)	0.446	Apthae epizooticae (口蹄疫)	0.337
Uterus (子宮)	0.501	Acquisition (購併)	0.396	Thyroid (甲狀腺)	0.303
Bronchus (支氣管)	0.500	Insurance law (保險法)	0.373	Gastric cancer (胃癌)	0.298

WTM: Information Retrieval (4/4)

- Pairing of PLSA and WTM
 - Sharing the same set of latent topics



Applying Relevance Feedback to LM Framework (1/2)

- There is still no formal mechanism to incorporate relevance feedback (judgments) into the language modeling framework (**especially for the query-likelihood approach**)
 - The query is a fixed sample while focusing on estimating accurate estimation of document language models $P(w|D)$
- Ponte (1998) proposed a limited way to incorporate blind reference feedback into the LM framework
 - Think of example relevant documents $D \in \tilde{R}$ as examples of what the query might have been, and re-sample (or expand) the query by adding k highly descriptive words from the these documents (blind reference feedback)

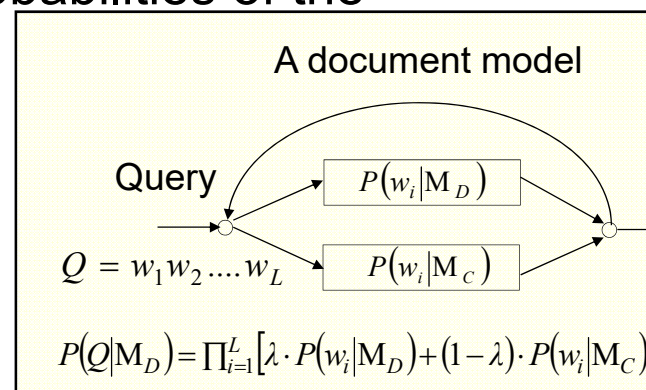
$$w^* = \arg \max_w \sum_{D \in \tilde{R}} \log \frac{P(w|M_D)}{P(w|M_C)}$$

Applying Relevance Feedback to LM Framework (2/2)

- Miller et al. (1999) propose two relevance feedback approach
 - **Query expansion**: add those words to the initial query that appear in two or more of the top m retrieved documents
 - **Document model re-estimation**: use a set of outside training query exemplars to train the transition probabilities of the document models

$$\hat{\lambda} = \frac{\sum_{Q \in [\text{TrainSet}]_Q} \sum_{D \in [\text{Doc}]_{R \text{ to } Q}} \sum_{q_n \in Q} \left[\frac{\lambda P(q_n | M_D)}{\lambda P(q_n | M_D) + (1 - \lambda) P(q_n | M_C)} \right]}{\sum_{Q \in [\text{TrainSet}]_Q} |Q| \cdot |[\text{Doc}]_{R \text{ to } Q}}$$

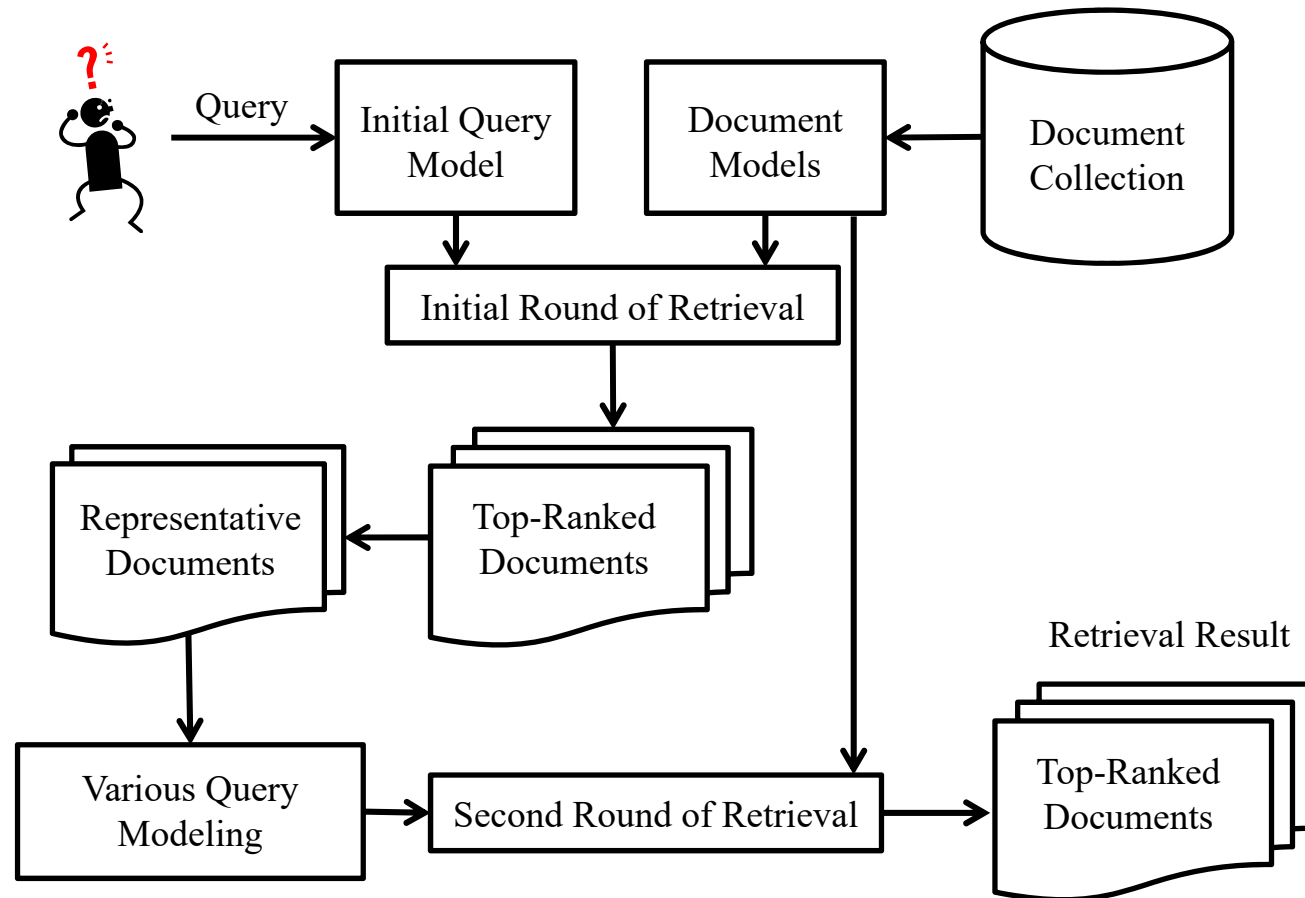
Annotations: "the new weight" points to $\hat{\lambda}$; "819 queries" points to the first sum; " ≤ 2265 docs" points to the second sum; "the old weight" points to $\lambda P(q_n | M_D)$.



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

Relevance Modeling (1/3)

- A schematic illustration of the information retrieval system with relevance feedback for improved query modeling



Relevance Modeling (2/3)

- Use the top-ranked pseudo-relevant documents to approximate the relevance class of each query
 - The joint probability of a query Q and any word w being generated by the relevance class of is computed as follows, on the basis of the top-ranked list of pseudo-relevant documents obtained from the initial round of retrieval:

$$P_{\text{RM}}(Q, w) = \sum_{m=1}^M P(D_m) P(q_1, q_2, \dots, q_L, w | D_m)$$

- If we further assume that words are conditionally independent given D_m and their order is of no importance

$$P_{\text{RM}}(Q, w) = \sum_{m=1}^M P(D_m) P(w | D_m) \prod_{l=1}^L P(q_l | D_m)$$

Relevance Modeling (3/3)

- The enhanced query model , therefore, can be expressed by

$$P_{\text{RM}}(w|Q) = \frac{P_{\text{RM}}(Q, w)}{P_{\text{RM}}(Q)} = \frac{\sum_{m=1}^M P(D_m) P(w | D_m) \prod_{l=1}^L P(q_l | D_m)}{\sum_{m=1}^M P(D_m) \prod_{l=1}^L P(q_l | D_m)}$$

- Incorporation of Latent Topic Information

$$\tilde{P}(w | D_m) = \sum_{k=1}^K P(w | T_k) P(T_k | D_m)$$

$$P_{\text{TRM}}(Q, w) = \sum_{m=1}^M \sum_{k=1}^K P(D_m) P(T_k | D_m) P(w | T_k) \prod_{l=1}^L P(q_l | T_k)$$

Leveraging Non-relevant Information

- Hypothesize that the low-ranked (or pseudo non-relevant) documents can provide useful cues as well to boost the retrieval effectiveness of a given query
 - For this idea to work, we may estimate a non-relevance model $P(w|NR_Q)$ for each test query based on those selected pseudo non-relevant documents
- The similarity measure between a query and document thus can be computed as follows

$$SIM(Q, D) = -KL(Q||D) + \alpha \cdot KL(NR_Q||D)$$

Incorporating Prior Knowledge into LM Framework

- Several efforts have been paid to using prior knowledge for the LM framework, especially modeling the document prior $P(D)$
 - Document length
 - Document source
 - Average word-length
 - Aging (time information/period)
 - URL
 - Page links

$$P(D|Q) = \frac{P(Q|D)P(D)}{P(Q)}$$

Implementation Notes: Probability Manipulation

- For language modeling approaches to IR, many conditional probabilities are usually multiplied. This can result in a “floating point underflow”

$$P(Q|M_D) = \prod_{i=1}^L [\lambda \cdot P(w_i|M_D) + (1 - \lambda) \cdot P(w_i|M_C)]$$

- It is better to perform the computation by “adding” logarithms of probabilities instead
 - The logarithm function is monotonic (order-preserving)

$$\log P(Q|M_D) = \sum_{i=1}^L \log [\lambda \cdot P(w_i|M_D) + (1 - \lambda) \cdot P(w_i|M_C)]$$

- We also should avoid the problem of “zero probabilities (or estimates)” owing to sparse data, by using appropriate probability smoothing techniques

Implementation Notes: Converting to *tf-idf*-like Weighting

- The query likelihood retrieval model

$$P(Q|M_D) = \prod_{i=1}^L [\lambda \cdot P(w_i|M_D) + (1-\lambda) \cdot P(w_i|M_C)]$$



$$\log P(Q|M_D) = \sum_{i=1}^L \log(\lambda \cdot P(w_i|M_D) + (1-\lambda) \cdot P(w_i|M_C))$$



Logarithm is a monotonic
(rank-preserving) transformation

$$= \sum_{i, c(w_i, D) > 0} \log \left(\lambda \cdot \frac{c(w_i, D)}{|D|} + (1-\lambda) \cdot \frac{c(w_i, C)}{|C|} \right) + \sum_{i, c(w_i, D) = 0} \log \left((1-\lambda) \cdot \frac{c(w_i, C)}{|C|} \right)$$

$$= \sum_{i, c(w_i, D) > 0} \left[\log \left(\lambda \cdot \frac{c(w_i, D)}{|D|} + (1-\lambda) \cdot \frac{c(w_i, C)}{|C|} \right) - \log \left((1-\lambda) \cdot \frac{c(w_i, C)}{|C|} \right) \right] + \sum_{i=1}^L \log \left((1-\lambda) \cdot \frac{c(w_i, C)}{|C|} \right)$$

$$= \sum_{i, c(w_i, D) > 0} \log \left(\frac{\lambda \cdot \frac{c(w_i, D)}{|D|} + (1-\lambda) \cdot \frac{c(w_i, C)}{|C|}}{(1-\lambda) \cdot \frac{c(w_i, C)}{|C|}} \right) + \sum_{i=1}^L \log \left((1-\lambda) \cdot \frac{c(w_i, C)}{|C|} \right)$$

Being the same for all
documents
=> can be discarded!

$$\text{rank} = \sum_{i, c(w_i, D) > 0} \log \left(\frac{\lambda \cdot \frac{c(w_i, D)}{|D|}}{(1-\lambda) \cdot \frac{c(w_i, C)}{|C|}} + 1 \right)$$

Therefore, the similarity score is directly proportional
to the document frequency and inversely proportional
to the collection frequency.

=> Can be efficiently implemented with inverted files

(To be discussed later on!)