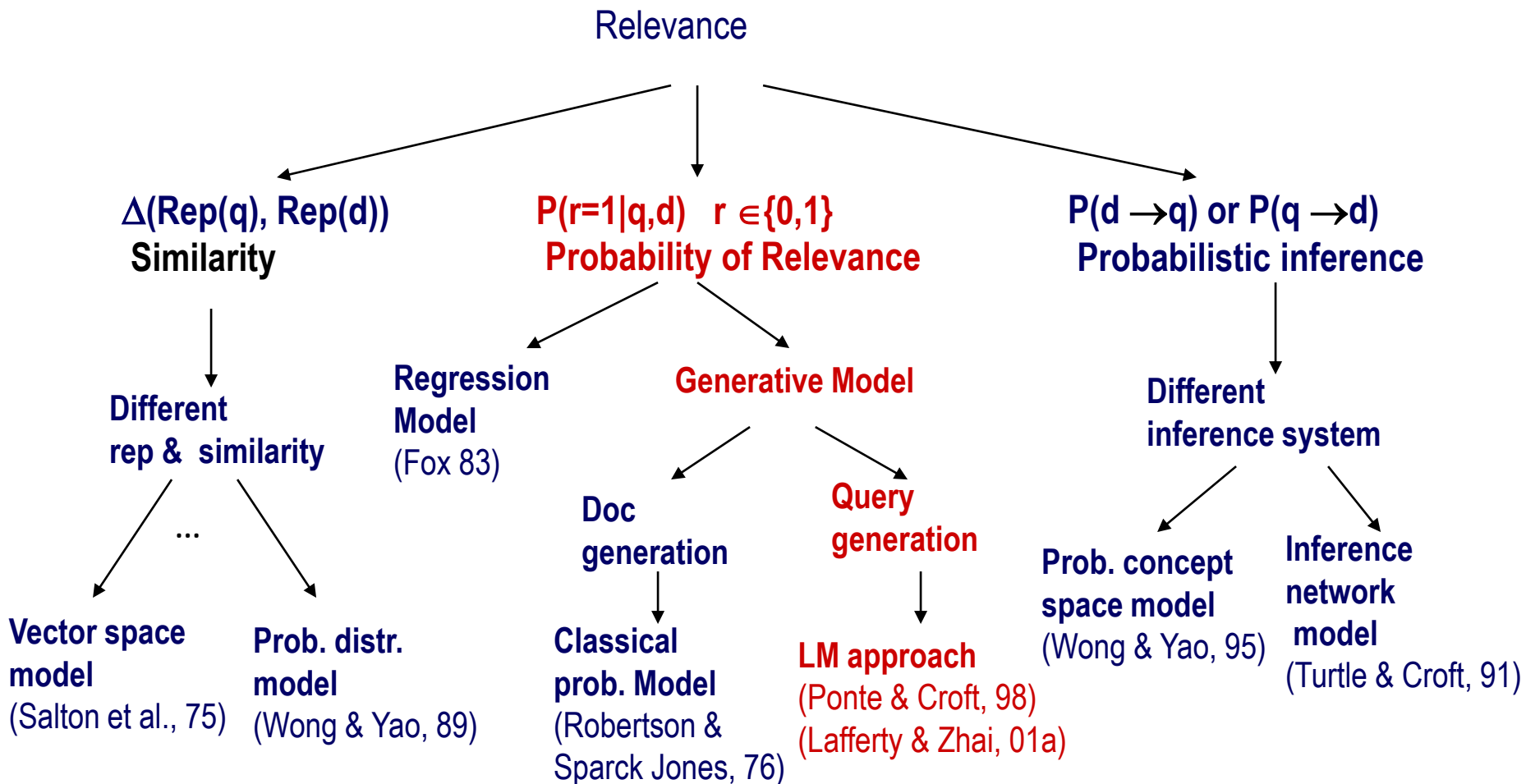


Language Models

Hongning Wang

CS@UVa

Notion of Relevance



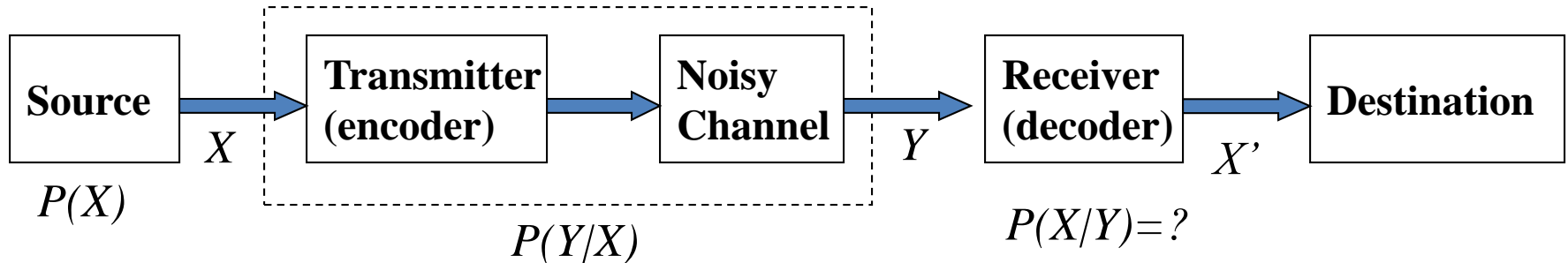
What is a statistical LM?

- A model specifying probability distribution over word sequences
 - $p(\textit{“Today is Wednesday”}) \approx 0.001$
 - $p(\textit{“Today Wednesday is”}) \approx 0.0000000000000001$
 - $p(\textit{“The eigenvalue is positive”}) \approx 0.00001$
- It can be regarded as a probabilistic mechanism for “generating” text, thus also called a “generative” model

Why is a LM useful?

- Provides a principled way to quantify the uncertainties associated with natural language
- Allows us to answer questions like:
 - Given that we see “*John*” and “*feels*”, how likely will we see “*happy*” as opposed to “*habit*” as the next word?
(speech recognition)
 - Given that we observe “baseball” three times and “game” once in a news article, how likely is it about “sports”?
(text categorization, information retrieval)
 - Given that a user is interested in sports news, how likely would the user use “baseball” in a query?
(information retrieval)

Source-Channel framework [Shannon 48]



$$\hat{X} = \arg \max_X p(X | Y) = \arg \max_X p(Y | X) p(X) \quad (\text{Bayes Rule})$$

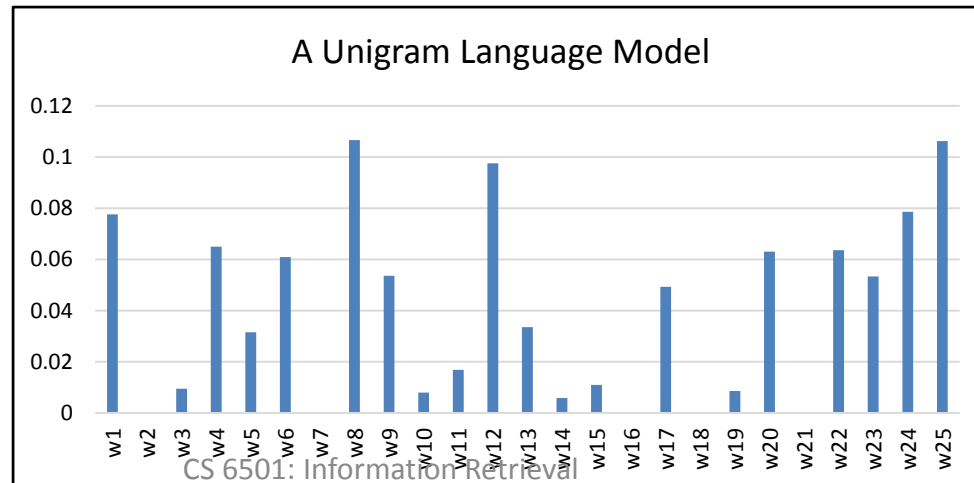
When X is text, $p(X)$ is a language model

Many Examples:

Speech recognition:	X =Word sequence	Y =Speech signal
Machine translation:	X =English sentence	Y =Chinese sentence
OCR Error Correction:	X =Correct word	Y = Erroneous word
Information Retrieval:	X =Document	Y =Query
Summarization:	X =Summary	Y =Document

Unigram language model

- Generate a piece of text by generating each word independently
 - $p(w_1 w_2 \dots w_n) = p(w_1)p(w_2) \dots p(w_n)$
 - $s. t. \{p(w_i)\}_{i=1}^N, \sum_i p(w_i) = 1, p(w_i) \geq 0$
- Essentially a multinomial distribution over the vocabulary



The simplest and most popular choice!

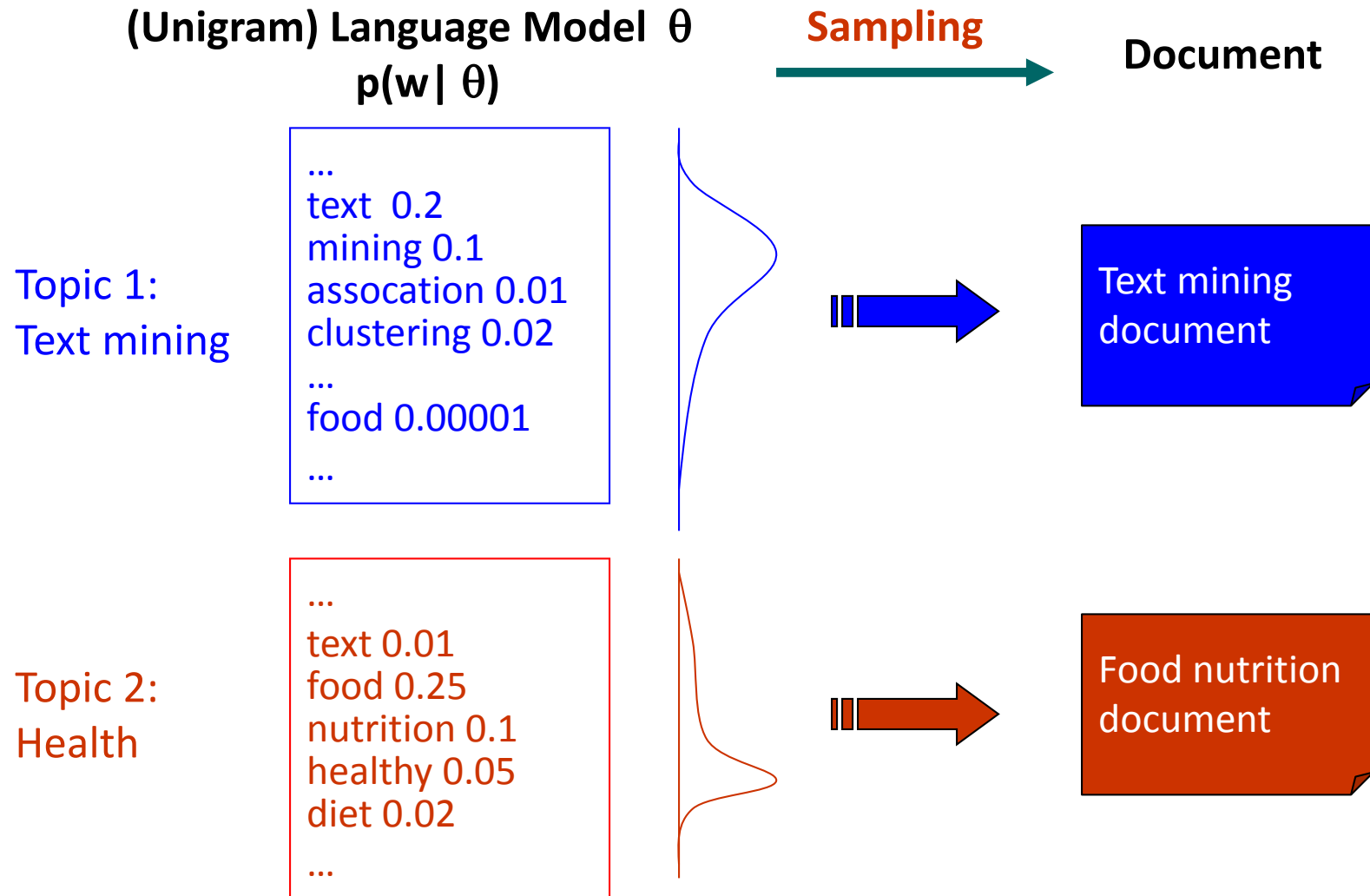
More sophisticated LMs

- N-gram language models
 - In general, $p(w_1 w_2 \dots w_n) = p(w_1)p(w_2|w_1) \dots p(w_n|w_1 \dots w_{n-1})$
 - n-gram: conditioned only on the past n-1 words
 - E.g., bigram: $p(w_1 \dots w_n) = p(w_1)p(w_2|w_1) p(w_3|w_2) \dots p(w_n|w_{n-1})$
- Remote-dependence language models (e.g., Maximum Entropy model)
- Structured language models (e.g., probabilistic context-free grammar)

Why just unigram models?

- Difficulty in moving toward more complex models
 - They involve more parameters, so need more data to estimate
 - They increase the computational complexity significantly, both in time and space
- Capturing word order or structure may not add so much value for “topical inference”
- But, using more sophisticated models can still be expected to improve performance ...

Generative view of text documents



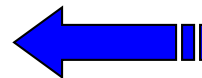
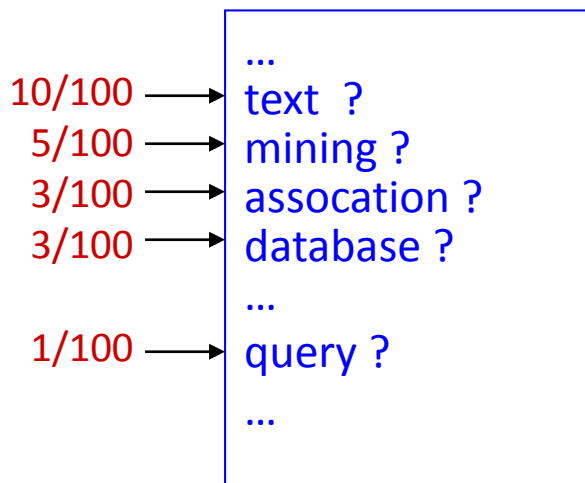
Estimation of language models

- Maximum likelihood estimation

~~Unigram~~ Language Model θ
 $p(w | \theta) = ?$

Estimation

Document

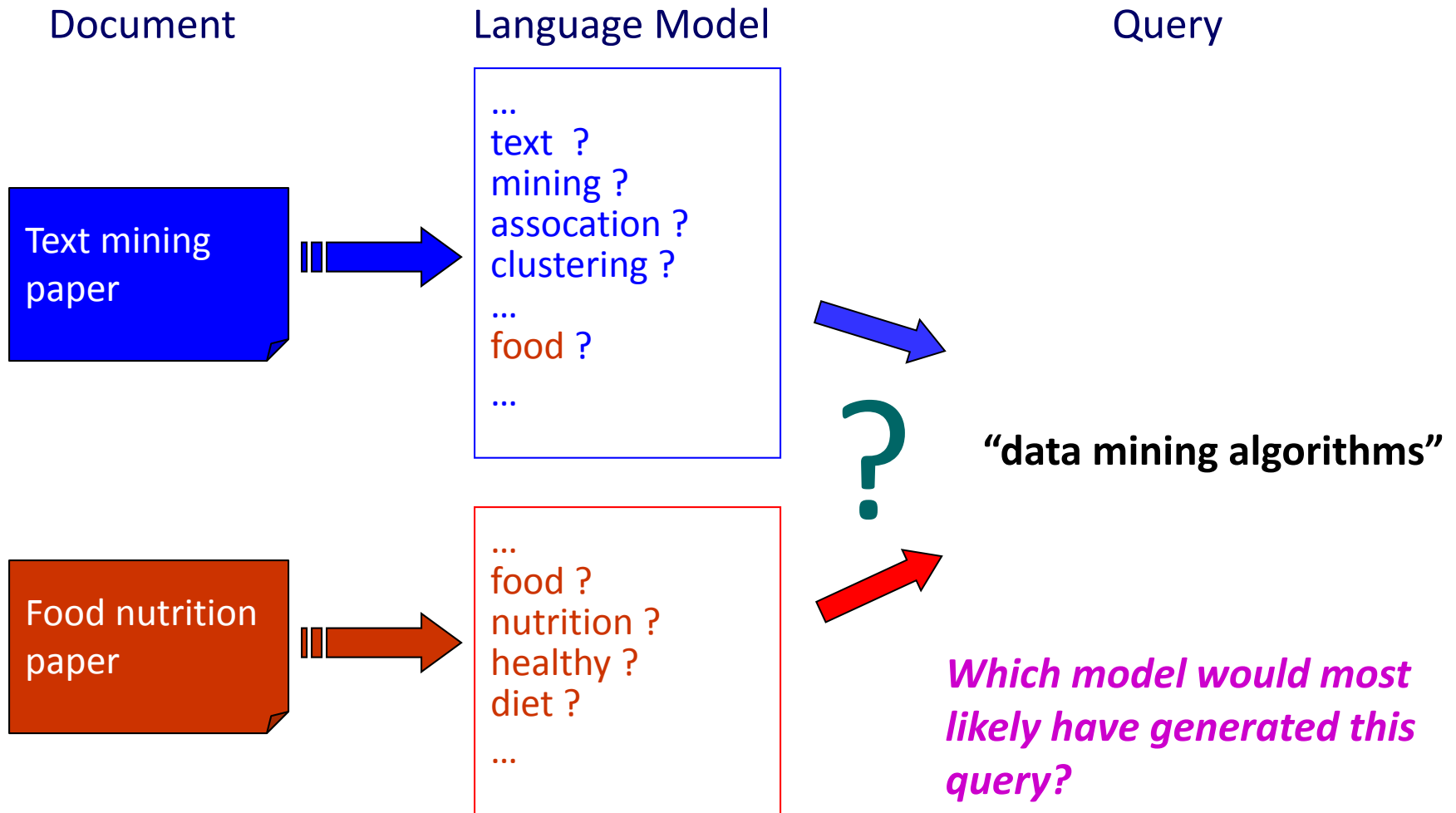


text 10
mining 5
association 3
database 3
algorithm 2
...
query 1
efficient 1

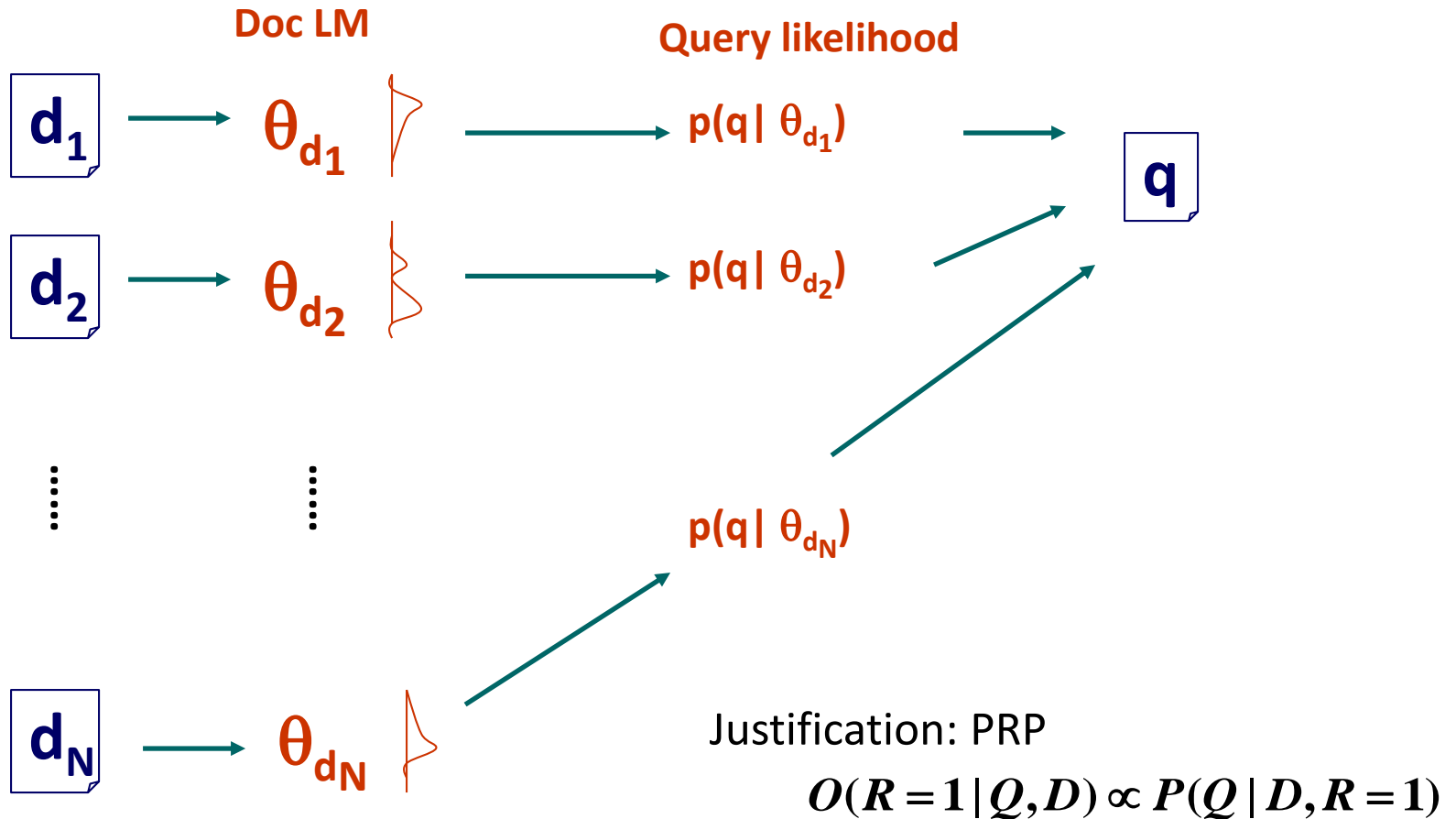
A “text mining” paper
(total #words=100)

Language models for IR

[Ponte & Croft SIGIR'98]



Ranking docs by query likelihood



Justification from PRP

$$\begin{aligned} O(R=1|Q,D) &\propto \frac{P(Q,D|R=1)}{P(Q,D|R=0)} \\ &= \frac{P(Q|D,R=1)P(D|R=1)}{P(Q|D,R=0)P(D|R=0)} \\ &\propto \underbrace{P(Q|D,R=1)}_{\text{Query likelihood } p(q|\theta_d)} \underbrace{\frac{P(D|R=1)}{P(D|R=0)}}_{\text{Document prior}} \quad (\text{Assume } P(Q|D,R=0) \approx P(Q|R=0)) \end{aligned}$$

Query generation

Assuming uniform document prior, we have

$$O(R=1|Q,D) \propto P(Q|D,R=1)$$

Retrieval as language model estimation

- Document ranking based on *query likelihood*

$$\log p(q | d) = \sum_i \log p(w_i | d)$$

where, $q = w_1 w_2 \dots w_n$

Document language model

- Retrieval problem \approx Estimation of $p(w_i | d)$
- Common approach
 - Maximum likelihood estimation (MLE)

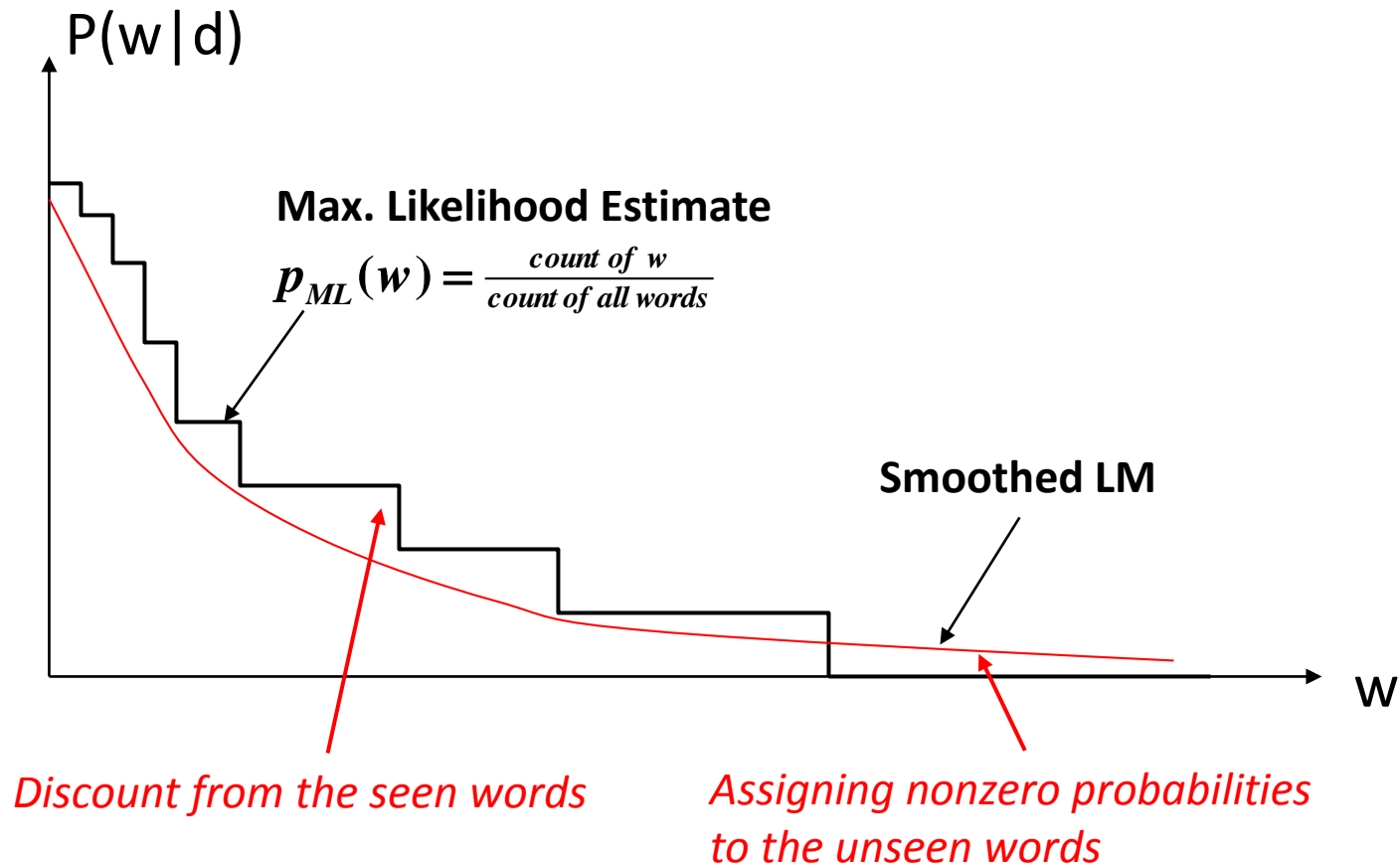
Problem with MLE

- What probability should we give a word that has not been observed in the document?
 - $\log 0$?
- If we want to assign non-zero probabilities to such words, we'll have to discount the probabilities of observed words
- This is so-called “smoothing”

General idea of smoothing

- All smoothing methods try to
 1. Discount the probability of words seen in a document
 2. Re-allocate the extra counts such that unseen words will have a non-zero count

Illustration of language model smoothing



Smoothing methods

- Method 1: Additive smoothing
 - Add a constant δ to the counts of each word

The diagram shows the formula for Laplace smoothing: $p(w|d) = \frac{c(w,d) + 1}{|d| + |V|}$. The entire formula is highlighted in yellow. Annotations with arrows point to specific parts: 'Counts of w in d' points to the numerator $c(w,d)$; '“Add one”, Laplace smoothing' points to the '+1' in the numerator; 'Vocabulary size' points to $|V|$ in the denominator; and 'Length of d (total counts)' points to $|d|$ in the denominator.

Counts of w in d

$p(w|d) = \frac{c(w,d) + 1}{|d| + |V|}$

“Add one”, Laplace smoothing

Vocabulary size

Length of d (total counts)

- Problems?
 - Hint: all words are equally important?

Refine the idea of smoothing

- Should all unseen words get equal probabilities?
- We can use a reference model to discriminate unseen words

$$p(w | d) = \begin{cases} p_{seen}(w | d) & \text{if } w \text{ is seen in } d \\ \alpha_d p(w | REF) & \text{otherwise} \end{cases}$$

Discounted ML estimate

Reference language model


$$\alpha_d = \frac{1 - \sum_{w \text{ is seen}} p_{seen}(w | d)}{\sum_{w \text{ is unseen}} p(w | REF)}$$

Smoothing methods (cont)

- Method 2: Absolute discounting
 - Subtract a constant δ from the counts of each word

$$p(w|d) = \frac{\max(c(w;d) - \delta, 0) + \delta |d|_u p(w|REF)}{|d|}$$

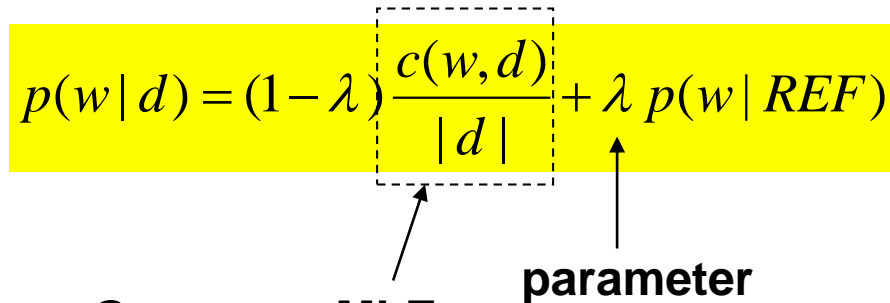
uniq words



- Problems?
 - Hint: varied document length?

Smoothing methods (cont)

- Method 3: Linear interpolation, Jelinek-Mercer
 - “Shrink” uniformly toward $p(w | REF)$

$$p(w | d) = (1 - \lambda) \frac{c(w, d)}{|d|} + \lambda p(w | REF)$$


The equation is displayed on a yellow background. A dashed box encloses the term $\frac{c(w, d)}{|d|}$. An arrow points from the label "MLE" below to this dashed box. Another arrow points from the label "parameter" below to the λ term in the equation.

- Problems?
 - Hint: what is missing?

Smoothing methods (cont)

- Method 4: Dirichlet Prior/Bayesian
 - Assume pseudo counts $\mu p(w | REF)$

$$p(w | d) = \frac{c(w; d) + \mu p(w | REF)}{|d| + \mu} = \frac{|d|}{|d| + \mu} \frac{c(w, d)}{|d|} + \frac{\mu}{|d| + \mu} p(w | REF)$$

parameter

- Problems?

Dirichlet prior smoothing

- Bayesian estimator
 - Posterior of LM: $p(\theta|d) \propto p(d|\theta)p(\theta)$
- Conjugate prior
 - Posterior will be in the same form as prior
 - Prior can be interpreted as “extra”/“pseudo” data
- Dirichlet distribution is a conjugate prior for multinomial distribution

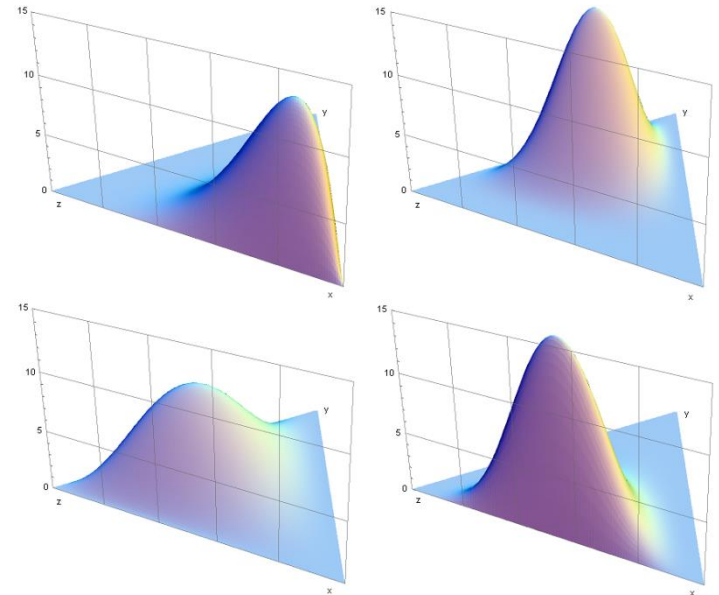
$$Dir(\theta | \underline{\alpha_1, \dots, \alpha_N}) = \frac{\Gamma(\alpha_1 + \dots + \alpha_N)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_N)} \prod_{i=1}^N \theta_i^{\alpha_i - 1}$$

“extra”/“pseudo” word counts, we set $\alpha_i = \mu p(w_i | \text{REF})$

Some background knowledge

- Conjugate prior
 - Posterior dist in the same family as prior
- Dirichlet distribution
 - Continuous
 - Samples from it will be the parameters in a multinomial distribution

Gaussian \rightarrow Gaussian
Beta \rightarrow Binomial
Dirichlet \rightarrow Multinomial



Dirichlet prior smoothing (cont)

Posterior distribution of parameters:

$$p(\theta | d) = \text{Dir}(\theta | c(w_1) + \alpha_1, \dots, c(w_N) + \alpha_N)$$

Property : If $\theta \sim \text{Dir}(\theta | \alpha)$, then $E(\theta) = \{\frac{\alpha_i}{\sum \alpha_i}\}$

The predictive distribution is the same as the mean:

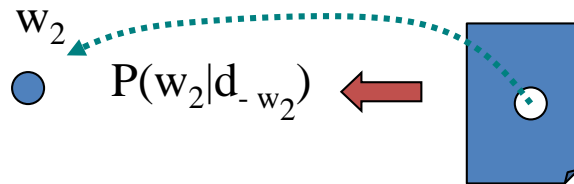
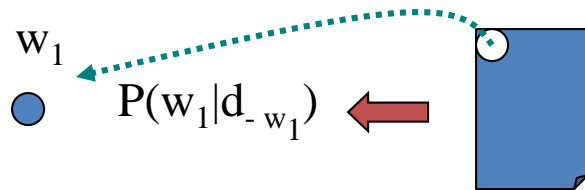
$$\begin{aligned} p(w_i | \hat{\theta}) &= \int p(w_i | \theta) \text{Dir}(\theta | \alpha) d\theta \\ &= \frac{c(w_i) + \alpha_i}{|d| + \sum_{i=1}^N \alpha_i} = \boxed{\frac{c(w_i) + \mu p(w_i | \text{REF})}{|d| + \mu}} \end{aligned}$$



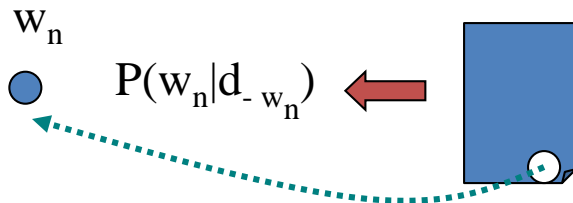
Dirichlet prior smoothing

Estimating μ using leave-one-out

[Zhai & Lafferty 02]



...



log-likelihood

$$L_{-1}(\mu | C) = \sum_{i=1}^N \sum_{w \in V} c(w, d_i) \log \left(\frac{c(w, d_i) - 1 + \mu p(w | C)}{|d_i| - 1 + \mu} \right)$$

Leave-one-out

Maximum Likelihood Estimator

$$\hat{\mu} = \operatorname{argmax}_{\mu} L_{-1}(\mu | C)$$

Why would “leave-one-out” work?

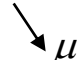
20 word by author1

abc abc ab c d d
abc cd d d
abd ab ab ab ab
cd d e cd e

Now, suppose we leave “e” out...

μ doesn't have to be big


$$p_{ml}("e" | author1) = \frac{1}{19}$$

$$p_{smooth}("e" | author1) = \frac{20}{20 + \mu} \frac{1}{19} + \frac{\mu}{20 + \mu} p("e" | REF)$$


20 word by author2

abc abc ab c d d
abe cb e f
acf fb ef aff abef
cdc db ge f s

$$p_{ml}("e" | author2) = \frac{0}{19}$$

$$p_{smooth}("e" | author2) = \frac{20}{20 + \mu} \frac{0}{19} + \frac{\mu}{20 + \mu} p("e" | REF)$$


μ must be big! more smoothing

The amount of smoothing is closely related to the underlying vocabulary size

Understanding smoothing

		Content words					
	Query = "the	algorithms	for	data	mining"		
$p_{ML}(w d1):$	0.04	0.001	0.02	0.002	0.003	4.8×10^{-12}	
$p_{ML}(w d2):$	0.02	0.001	0.01	0.003	0.004	2.4×10^{-12}	

$p(\text{"algorithms"}|d1) = p(\text{"algorithm"}|d2)$

$p(\text{"data"}|d1) < p(\text{"data"}|d2)$

$p(\text{"mining"}|d1) < p(\text{"mining"}|d2)$

Intuitively, d2 should have a higher score, but $p(q|d1) > p(q|d2)$...

So we should make $p(\text{"the"})$ and $p(\text{"for"})$ **less different** for all docs, 2.35×10^{-13}
and smoothing helps to achieve this goal... 4.53×10^{-13}

After smoothing with $p(w|d) = 0.1p_{DML}(w|d) + 0.9p(w|REF)$, $p(q|d1) < p(q|d2)$!

Query	= "the	algorithms	for	data	mining"
$P(w REF)$	0.2	0.00001	0.2	0.00001	0.00001
Smoothed $p(w d1):$	0.184	0.000109	0.182	0.000209	0.000309
Smoothed $p(w d2):$	0.182	0.000109	0.181	0.000309	0.000409

Understanding smoothing

$$\alpha_d = \frac{1 - \sum_{w \text{ is seen}} p_{\text{seen}}(w | d)}{\sum_{w \text{ is unseen}} p(w | REF)}$$

- Plug in the general smoothing scheme to the query likelihood retrieval formula, we obtain

$$\log p(q | d) = \sum_{w_i \in d \cap q} \left[\log \frac{p_{\text{seen}}(w_i | d)}{\alpha_d p(w_i | C)} \right] + n \log \alpha_d + \boxed{\sum_{w_i \in q} \log p(w_i | C)}$$

TF weighting (points to $p_{\text{seen}}(w_i | d)$)
IDF weighting (points to $p(w_i | C)$)
Doc length normalization
 (longer doc is expected to have a smaller α_d) (points to α_d)
Ignore for ranking (points to the boxed term)

- Smoothing with $p(w/C) \approx \text{TF-IDF} + \text{doc-length normalization}$**

Smoothing & TF-IDF weighting

Retrieval formula using the
general smoothing scheme

$$p(w|d) = \begin{cases} p_{\text{Seen}}(w|d) & \text{if } w \text{ is seen in } d \\ \alpha_d p(w|C) & \text{otherwise} \end{cases}$$

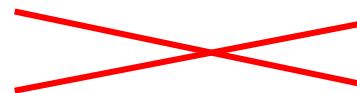
Smoothed ML estimate

Reference language model



$$\alpha_d = \frac{1 - \sum_{w \text{ is seen}} p_{\text{Seen}}(w|d)}{\sum_{w \text{ is unseen}} p(w|C)}$$

$$\log p(q|d) = \sum_{w \in V, c(w,q) > 0} c(w,q) \log p(w|d)$$



Key rewriting step (where did we see it before?)

Similar rewritings are very common when
using probabilistic models for IR...

What you should know

- How to estimate a language model
- General idea and different ways of smoothing
- Effect of smoothing