

Mixture Language Models

ChengXiang Zhai

*Department of Computer Science
University of Illinois, Urbana-Champaign*

Central Questions to Ask about a LM: “ADMI”

- **Application:** Why do you need a LM? For what purpose?



Evaluation metric for a LM

Topic mining and analysis

- **Data:** What kind of data do you want to model?



Data set for estimation & evaluation

Text documents & context

- **Model:** How do you define the model?



Assumptions to be made

Mixture of unigram LMs

- **Inference:** How do you infer/estimate the parameters?



Inference/Estimation algorithm

EM algorithm, Bayesian

Outline

- Motivation
- Mining one topic
- Two-component mixture model
- EM algorithm
- Probabilistic Latent Semantic Analysis (PLSA)
- Extensions of PLSA

Topic Mining and Analysis: Motivation

- Topic \approx main idea discussed in text data
 - Theme/subject of a discussion or conversation
 - Different granularities (e.g., topic of a sentence, an article, etc.)
- Many applications require discovery of topics in text
 - What are Twitter users talking about today?
 - What are the current research topics in data mining? How are they different from those 5 years ago?
 - What do people like about the iPhone 6? What do they dislike?
 - What were the major topics debated in 2012 presidential election?

Topics As Knowledge About the World

Knowledge about the world

Non-Text Data

Text Data

+ Context
Time
Location
...

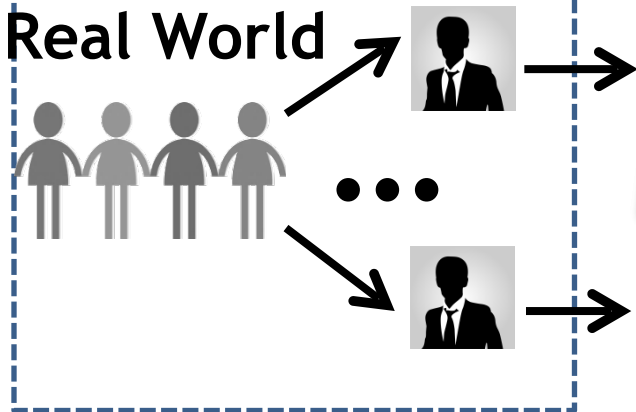
Topic 1

Topic 2

...

Topic k

Real World



Tasks of Topic Mining and Analysis

Task 2: Figure out which documents cover which topics →

Text Data



Topic 1

Topic 2

...

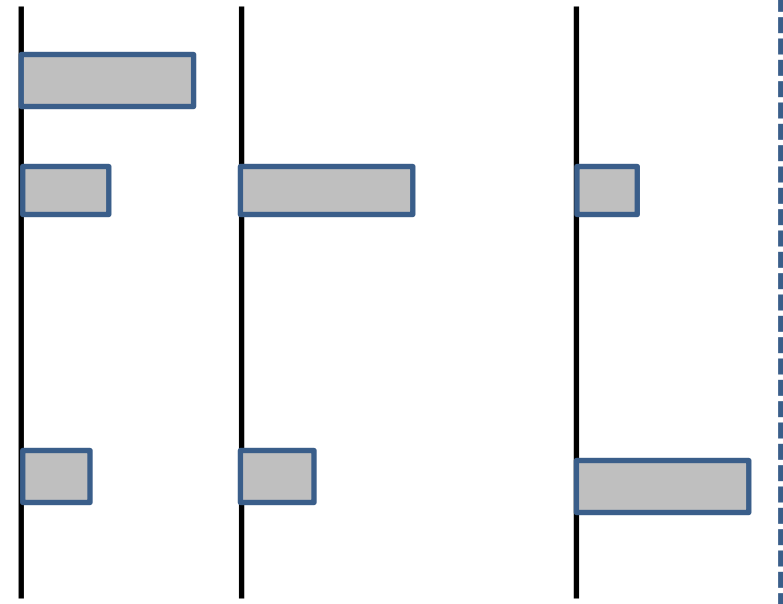
Topic k

Doc 1

Doc 2

...

Doc N



Task 1: Discover k topics

Formal Definition of Topic Mining and Analysis

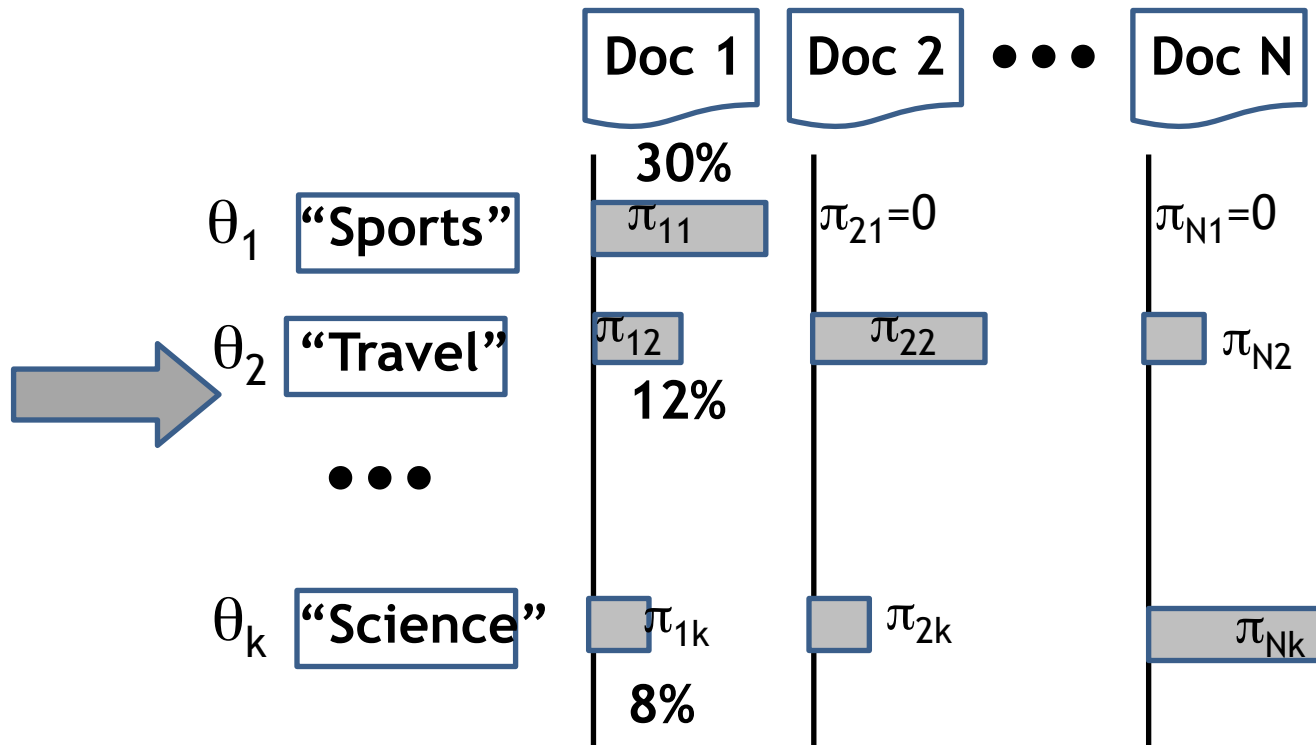
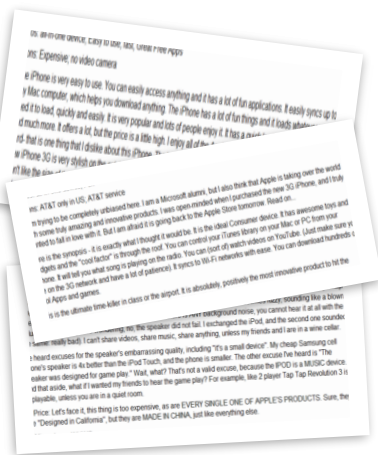
- Input
 - A collection of **N** text documents **$C = \{d_1, \dots, d_N\}$**
 - Number of topics: **k**
- Output
 - k topics: **$\{\theta_1, \dots, \theta_k\}$**
 - Coverage of topics in each d_i : **$\{\pi_{i1}, \dots, \pi_{ik}\}$**
 - π_{ij} = prob. of d_i covering topic θ_j

$$\sum_{j=1}^k \pi_{ij} = 1$$

How to define θ_i ?

Initial Idea: Topic = Term

Text Data



Mining k Topical Terms from Collection C

- Parse text in C to obtain candidate terms (e.g., term = word).
- Design a scoring function to measure how good each term is as a topic.
 - Favor a representative term (high frequency is favored)
 - Avoid words that are too frequent (e.g., “the”, “a”).
 - TF-IDF weighting from retrieval can be very useful.
 - Domain-specific heuristics are possible (e.g., favor title words, hashtags in tweets).
- Pick k terms with the highest scores but try to minimize redundancy.
 - If multiple terms are very similar or closely related, pick only one of them and ignore others.

Computing Topic Coverage: π_{ij}

Doc d_i

θ_1 “Sports” π_{i1} count(“sports”, d_i)=4

θ_2 “Travel” π_{i2} count(“travel”, d_i) =2

...

θ_k “Science” π_{ik} count(“science”, d_i)=1

$$\pi_{ij} = \frac{\text{count}(\theta_j, d_i)}{\sum_{L=1}^k \text{count}(\theta_L, d_i)}$$

How Well Does This Approach Work?

Doc d_i

Cavaliers vs. Golden State Warriors: NBA playoff
finals ... basketball game ... **travel** to Cleveland ...
star ...

θ_1 "Sports"

$$\pi_{i1} \propto c(\text{"sports"}, d_i) = 0$$

1. Need to count
related words also!

θ_2 "Travel"

$$\pi_{i2} \propto c(\text{"travel"}, d_i) = 1 > 0$$

...

2. "Star" can be ambiguous (e.g., star in the sky).

θ_k "Science"

$$\pi_{ik} \propto c(\text{"science"}, d_i) = 0$$

3. Mine complicated
topics?

Problems with “Term as Topic”

- Lack of expressive power → **Topic = {Multiple Words}**
 - Can only represent simple/general topics
 - Can't represent complicated topics
- Incompleteness in vocabulary coverage **+ weights on words**
 - Can't capture variations of vocabulary (e.g., related words)
- Word sense ambiguity → **Split an ambiguous word**
 - A topical term or related term can be ambiguous (e.g., basketball star vs. star in the sky)

A probabilistic topic model can do all these!

Improved Idea: Topic = Word Distribution

θ_1 **“Sports”**

$P(w | \theta_1)$

sports	0.02
game	0.01
basketball	0.005
football	0.004
play	0.003
star	0.003
...	
nba	0.001
...	
travel	0.0005
...	

θ_2 **“Travel”**

$P(w | \theta_2)$

travel	0.05
attraction	0.03
trip	0.01
flight	0.004
hotel	0.003
island	0.003
...	
culture	0.001
...	
play	0.0002
...	

...

θ_k **“Science”**

$P(w | \theta_k)$

science	0.04
scientist	0.03
spaceship	0.006
telescope	0.004
genomics	0.004
star	0.002
...	
genetics	0.001
...	
travel	0.00001
...	

$$\sum_{w \in V} p(w | \theta_i) = 1$$

Vocabulary Set: $V = \{w_1, w_2, \dots\}$

Probabilistic Topic Mining and Analysis

- Input

- A collection of N text documents $C=\{d_1, \dots, d_N\}$
- Vocabulary set: $V=\{w_1, \dots, w_M\}$
- Number of topics: k

- Output

- k topics, each a word distribution: $\{ \theta_1, \dots, \theta_k \}$
- Coverage of topics in each d_i : $\{ \pi_{i1}, \dots, \pi_{ik} \}$
- π_{ij} =prob. of d_i covering topic θ_j

$$\sum_{w \in V} p(w | \theta_i) = 1$$

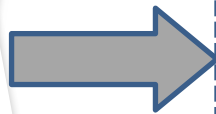
$$\sum_{j=1}^k \pi_{ij} = 1$$

The Computation Task

INPUT: C, k, V

OUTPUT: $\{ \theta_1, \dots, \theta_k \}, \{ \pi_{i1}, \dots, \pi_{ik} \}$

Text Data



θ_1

sports 0.02
game 0.01
basketball 0.005
football 0.004
...

Doc 1

30%

π_{11}

Doc 2

$\pi_{21}=0\%$

Doc N

$\pi_{N1}=0\%$

θ_2

travel 0.05
attraction 0.03
trip 0.01
...

12%

π_{12}

π_{22}

π_{N2}

...

θ_k

science 0.04
scientist 0.03
spaceship 0.006
...

8%

π_{1k}

π_{2k}

π_{Nk}

Generative Model for Text Mining

INP

Modeling of Data Generation: $P(\text{Data} \mid \text{Model}, \Lambda)$

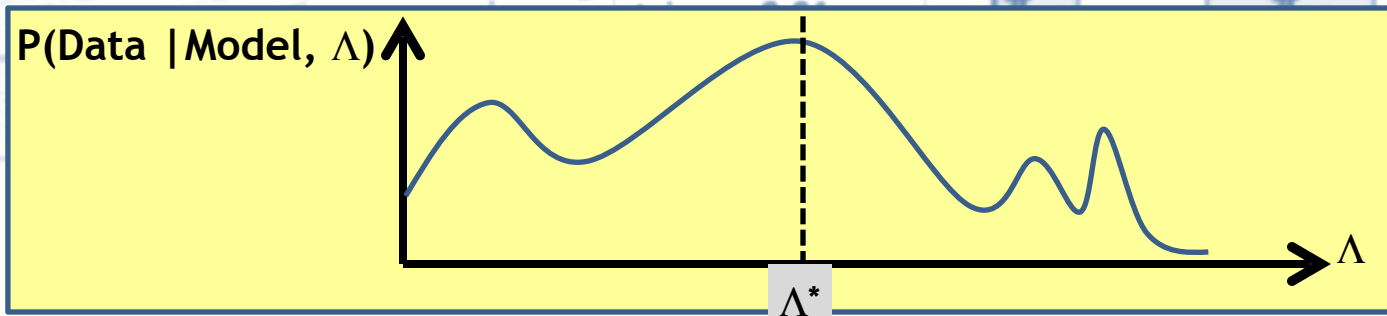
$\Lambda = (\{ \theta_1, \dots, \theta_k \}, \{ \pi_{11}, \dots, \pi_{1k} \}, \dots, \{ \pi_{N1}, \dots, \pi_{Nk} \})$

Text Data

How many parameters in total?

Parameter Estimation/ Inferences

$$\Lambda^* = \operatorname{argmax}_{\Lambda} p(\text{Data} \mid \text{Model}, \Lambda)$$



General Ideas of Generative Models for Text Mining

- **Model data generation:** $P(\text{Data} \mid \text{Model}, \Lambda)$
- **Infer the most likely parameter values Λ^* given a particular data set:** $\Lambda^* = \operatorname{argmax}_{\Lambda} p(\text{Data} \mid \text{Model}, \Lambda)$
- **Take Λ^* as the “knowledge”** to be mined for the text mining problem
- **Adjust** the design of the model to discover different knowledge

Simplest Case of Topic Model: Mining One Topic

INPUT: $C=\{d\},$
 V

Text Data



OUTPUT: $\{ \theta \}$

$P(w | \theta)$

θ

text ?
mining ?
association ?
database ?

...
query ?
...

Doc d

100%



Language Model Setup

- **Data:** Document $d = x_1 x_2 \dots x_{|d|}$, $x_i \in V = \{w_1, \dots, w_M\}$ is a word
- **Model:** Unigram LM θ (=topic) : $\{\theta_i = p(w_i | \theta)\}$, $i=1, \dots, M$; $\theta_1 + \dots + \theta_M = 1$
- **Likelihood function:** $p(d | \theta) = p(x_1 | \theta) \times \dots \times p(x_{|d|} | \theta)$
$$= p(w_1 | \theta)^{c(w_1, d)} \times \dots \times p(w_M | \theta)^{c(w_M, d)}$$
$$= \prod_{i=1}^M p(w_i | \theta)^{c(w_i, d)} = \prod_{i=1}^M \theta_i^{c(w_i, d)}$$
- **ML estimate:** $(\hat{\theta}_1, \dots, \hat{\theta}_M) = \arg \max_{\theta_1, \dots, \theta_M} p(d | \theta) = \arg \max_{\theta_1, \dots, \theta_M} \prod_{i=1}^M \theta_i^{c(w_i, d)}$

Computation of Maximum Likelihood Estimate

Maximize $p(d | \theta)$ $(\hat{\theta}_1, \dots, \hat{\theta}_M) = \arg \max_{\theta_1, \dots, \theta_M} p(d | \theta) = \arg \max_{\theta_1, \dots, \theta_M} \prod_{i=1}^M \theta_i^{c(w_i, d)}$

Max. Log-Likelihood $(\hat{\theta}_1, \dots, \hat{\theta}_M) = \arg \max_{\theta_1, \dots, \theta_M} \log[p(d | \theta)] = \arg \max_{\theta_1, \dots, \theta_M} \sum_{i=1}^M c(w_i, d) \log \theta_i$

Subject to constraint: $\sum_{i=1}^M \theta_i = 1$

Use Lagrange multiplier approach

Lagrange function: $f(\theta | d) = \sum_{i=1}^M \alpha(w_i, d) \log \theta_i + \lambda (\sum_{i=1}^M \theta_i - 1)$

$$\frac{\partial f(\theta | d)}{\partial \theta_i} = \frac{\alpha(w_i, d)}{\theta_i} + \lambda = 0 \rightarrow \theta_i = -\frac{\alpha(w_i, d)}{\lambda}$$

$$\sum_{i=1}^M -\frac{\alpha(w_i, d)}{\lambda} = 1 \rightarrow \lambda = -\sum_{i=1}^M \alpha(w_i, d) \rightarrow \hat{\theta}_i = p(w_i | \hat{\theta}) = \frac{\alpha(w_i, d)}{\sum_{i=1}^M \alpha(w_i, d)} = \frac{\alpha(w_i, d)}{|d|}$$

**Normalized
Counts**



What Does the Topic Look Like?

d

Text mining
paper

$p(w | \theta)$

the 0.031
a 0.018
...
text 0.04
mining 0.035
association 0.03
clustering 0.005
computer 0.0009
...
food 0.000001
...

Can we get rid of
these common words?

Factoring out Background Words

d

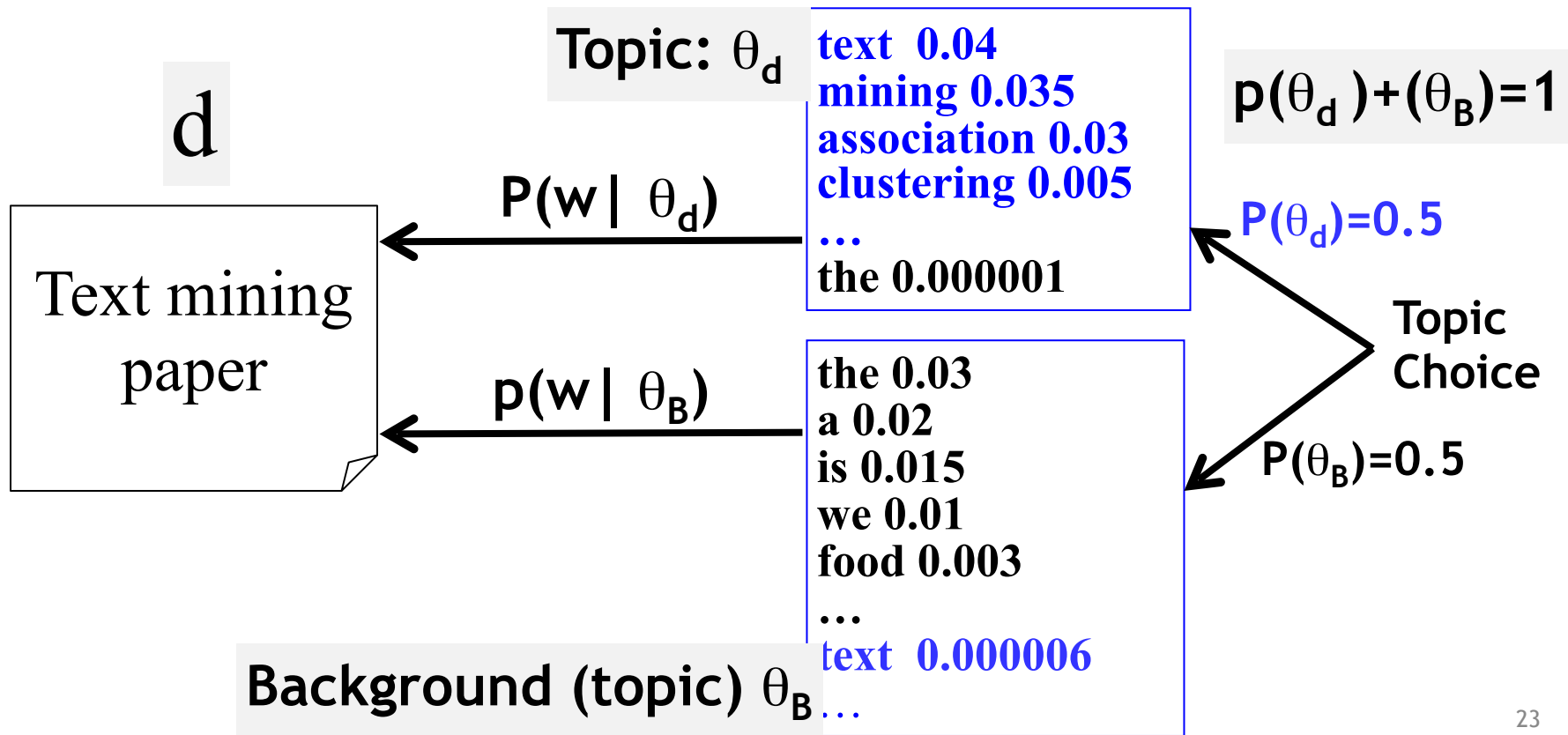
Text mining
paper

$p(w | \theta)$

the 0.031
a 0.018
...
text 0.04
mining 0.035
association 0.03
clustering 0.005
computer 0.0009
...
food 0.000001
...

How can we get rid of
these common words?

Generate d Using Two Word Distributions



What's the probability of observing a word w ?

d

Topic: θ_d text 0.04
mining 0.035

$p(\theta_d) + p(\theta_B) = 1$

$$P(\text{"the"}) = p(\theta_d)p(\text{"the"} | \theta_d) + p(\theta_B)p(\text{"the"} | \theta_B)$$

$$= 0.5 * 0.000001 + 0.5 * 0.03$$

$$= 0.5$$

Topic choice

$$P(\text{"text"}) = p(\theta_d)p(\text{"text"} | \theta_d) + p(\theta_B)p(\text{"text"} | \theta_B)$$

$$= 0.5 * 0.04 + 0.5 * 0.000006$$

$$0.5$$

... 0.000

...

text 0.000006

Background (topic) θ_B ...

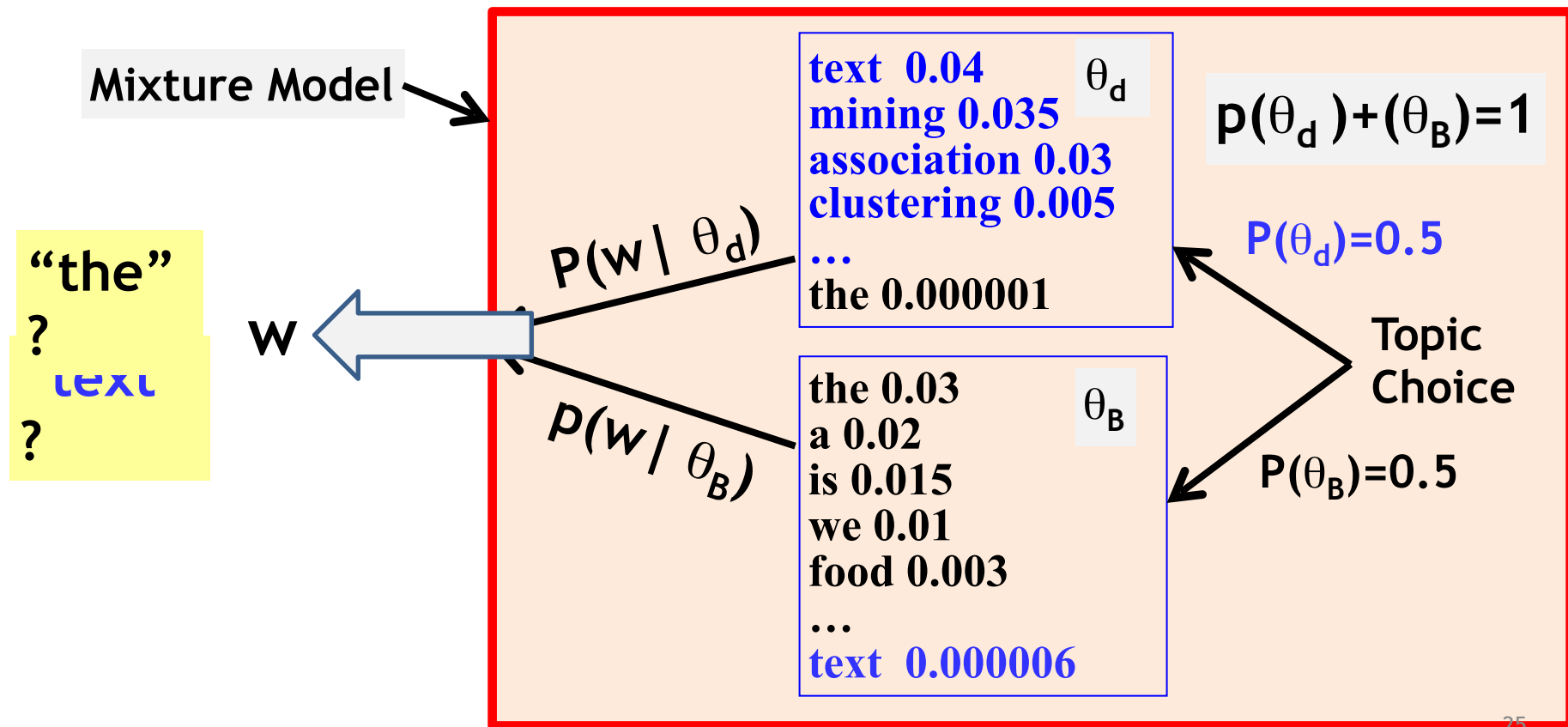
"the"

?

text

?

The Idea of a Mixture Model



As a Generative Model...

text 0.04 θ_d
mining 0.035
association 0.03
clustering 0.005

$$p(\theta_d) + p(\theta_B) = 1$$

Formally defines the following generative model:

$$p(w) = p(\theta_d)p(w | \theta_d) + p(\theta_R)p(w | \theta_R)$$

Estimate of the model “discovers”

two topics + topic coverage

What if $p(\theta_d)=1$ or $p(\theta_B)=1$?

W

[illegible]

Mixture of Two Unigram Language Models

- **Data:** Document d
- **Mixture Model:** parameters $\Lambda = (\{p(w | \theta_d)\}, \{p(w | \theta_B)\}, p(\theta_B), p(\theta_d))$
 - Two unigram LMs: θ_d (the topic of d); θ_B (background topic)
 - Mixing weight (topic choice): $p(\theta_d) + p(\theta_B) = 1$

- **Likelihood function:**

$$\begin{aligned} p(d | \Lambda) &= \prod_{i=1}^{|d|} p(x_i | \Lambda) = \prod_{i=1}^{|d|} [p(\theta_d)p(x_i | \theta_d) + p(\theta_B)p(x_i | \theta_B)] \\ &= \prod_{i=1}^M [p(\theta_d)p(w_i | \theta_d) + p(\theta_B)p(w_i | \theta_B)]^{c(w_i, d)} \end{aligned}$$

- **ML Estimate:** $\Lambda^* = \arg \max_{\Lambda} p(d | \Lambda)$

$$\text{Subject to } \sum_{i=1}^M p(w_i | \theta_d) = \sum_{i=1}^M p(w_i | \theta_B) = 1 \quad p(\theta_d) + p(\theta_B) = 1$$

Back to Factoring out Background Words

Text Mining Paper
d

... text
mining... is...
clustering...
we.... Text.. the

$$P(w | \theta_d)$$

text 0.04
mining 0.035
association 0.03
clustering 0.005
...
the 0.000001

θ_d

$$p(\theta_d) + p(\theta_B) = 1$$

$$P(\theta_d) = 0.5$$

Topic
Choice

$$P(\theta_B) = 0.5$$

$$p(w | \theta_B)$$

the 0.03
a 0.02
is 0.015
we 0.01
food 0.003
...

θ_B

text 0.000006

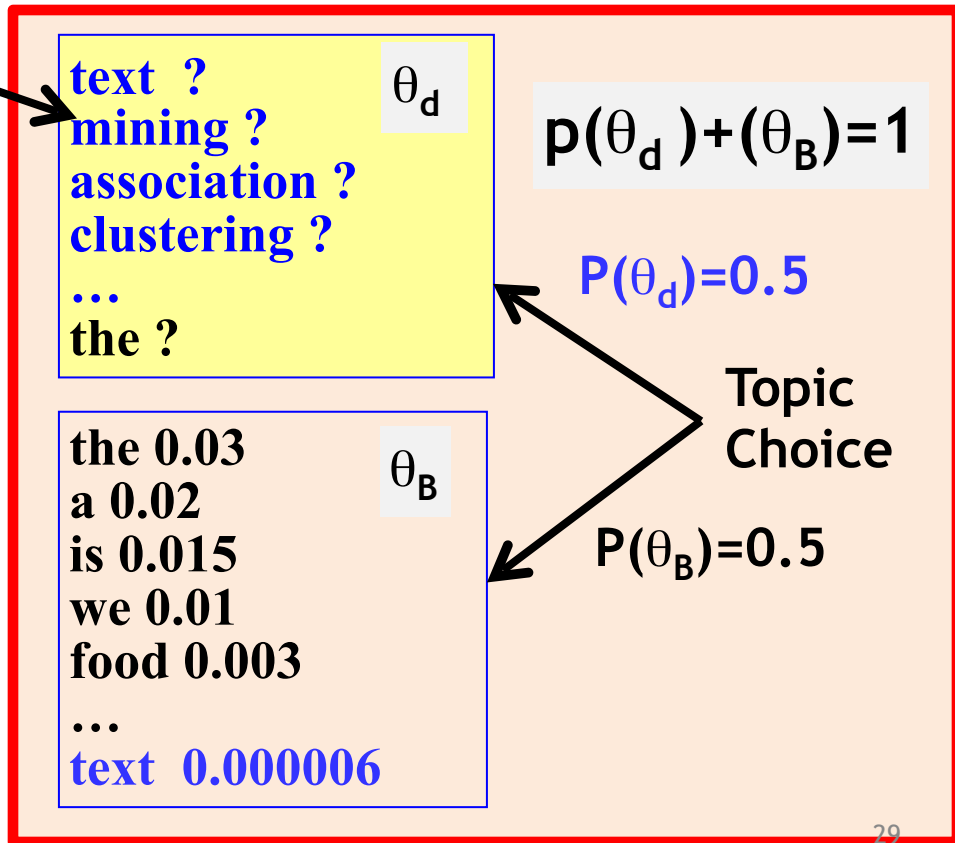
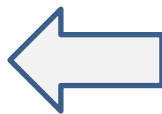
Estimation of One Topic: $P(w | \theta_d)$

Adjust θ_d to maximize $p(d | \Lambda)$
(all other parameters are known)

Would the ML estimate demote
background words in θ_d ?

d

... text
mining... is...
clustering...
we.... Text.. the



Behavior of a Mixture Model

$d =$ text

Likelihood:

the

$$P(\text{"text"}) = p(\theta_d)p(\text{"text"} | \theta_d) +$$

$$p(\theta_B)p(\text{"text"} | \theta_B)$$

$$P(\text{"the"}) = 0.5 * p(\text{"text"} | \theta_d) + 0.5 * 0.1$$

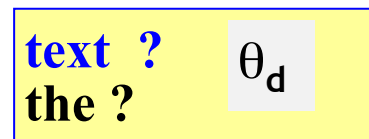
$$p(d = \text{"text"}) = 0.5 * 0.9$$

$$p(d = \text{"the"}) = 0.5 * 0.1$$

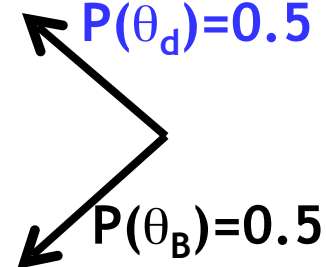
$$p(d | \Lambda) = p(\text{"text"} | \Lambda) p(\text{"the"} | \Lambda)$$

$$= [0.5 * p(\text{"text"} | \theta_d) + 0.5 * 0.1] \times$$

$$[0.5 * p(\text{"the"} | \theta_d) + 0.5 * 0.9]$$



$P(\theta_d) = 0.5$



$P(\theta_B) = 0.5$

How can we set $p(\text{"text"} | \theta_d)$ & $p(\text{"the"} | \theta_d)$ to maximize it?

Note that $p(\text{"text"} | \theta_d) + p(\text{"the"} | \theta_d) = 1$

“Collaboration” and “Competition” of θ_d and θ_B

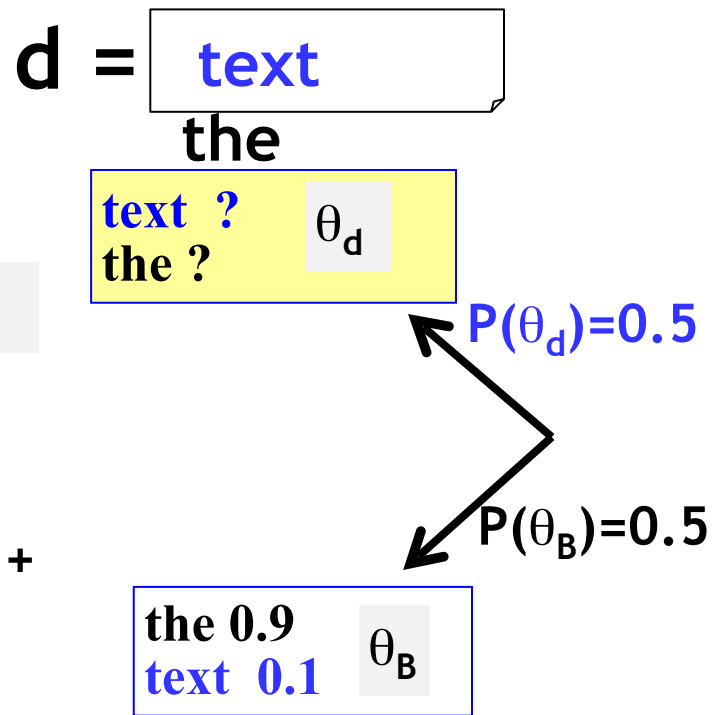
$$\begin{aligned} p(d|\Lambda) &= p(\text{“text”}|\Lambda) p(\text{“the”}|\Lambda) \\ &= [0.5 \cdot p(\text{“text”}|\theta_d) + 0.5 \cdot 0.1] \times \\ &\quad [0.5 \cdot p(\text{“the”}|\theta_d) + 0.5 \cdot 0.9] \end{aligned}$$

Note that $p(\text{“text”}|\theta_d) + p(\text{“the”}|\theta_d) = 1$

If , then reaches maximum when

$$\begin{aligned} 0.5 \cdot p(\text{“text”}|\theta_d) + 0.5 \cdot 0.1 &= 0.5 \cdot p(\text{“the”}|\theta_d) + \\ 0.5 \cdot 0.9 & \\ \rightarrow p(\text{“text”}|\theta_d) &= 0.9 \gg p(\text{“the”}|\theta_d) \\ &= 0.1 ! \end{aligned}$$

Behavior 1: if $p(w1|\theta_B) > p(w2|\theta_B)$, then $p(w1|\theta_d) < p(w2|\theta_d)$



Response to Data Frequency

$d =$ text
the

$$p(d|\Lambda) = [0.5 \cdot p(\text{"text"}|\theta_d) + 0.5 \cdot 0.1]$$

$$\rightarrow p(\text{"text"}|\theta_d) = 0.9 \gg p(\text{"the"}|\theta_d) = 0.1$$

$$p(d'|\Lambda) = [0.5 \cdot p(\text{"text"}|\theta_d) + 0.5 \cdot 0.1]$$

$$\times [0.5 \cdot p(\text{"the"}|\theta_d) + 0.5 \cdot 0.9] \dots$$

$d' =$ text the
the the
the ...the

What if we increase $p(\theta_B)$?

What's the optimal solution now? $p(\text{"the"}|\theta_d) > 0.1$? or $p(\text{"the"}|\theta_d)$

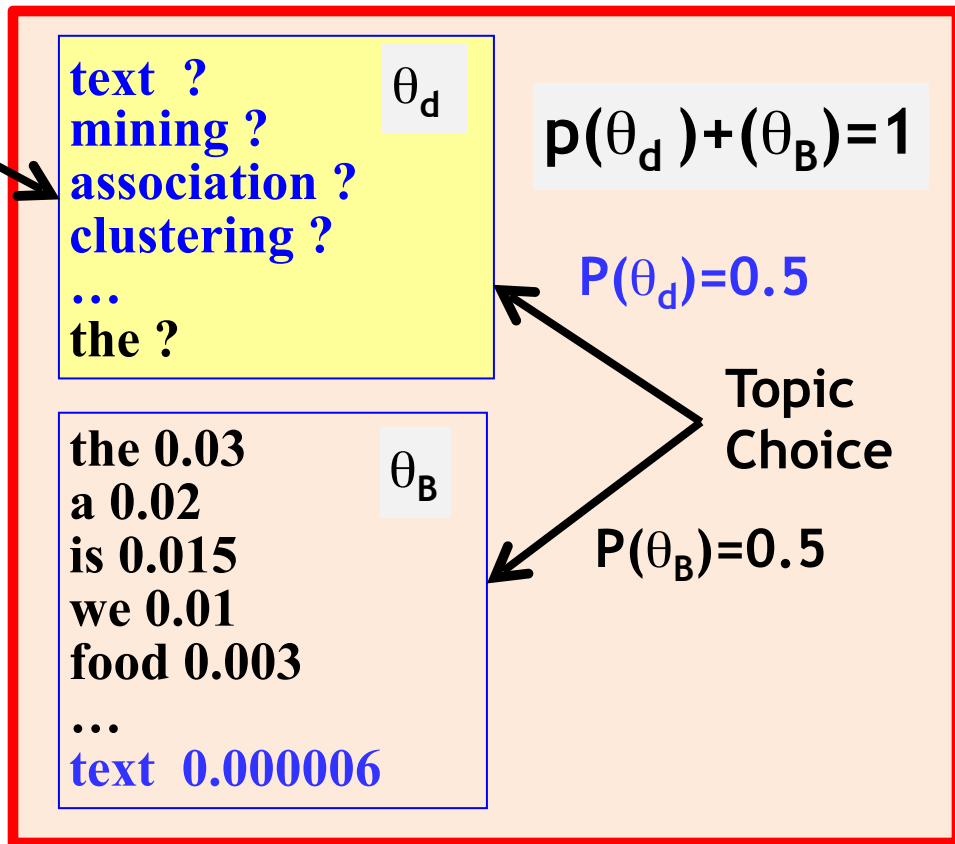
Behavior 2: high frequency words get higher $p(w|\theta_d)$

Estimation of One Topic: $P(w \mid \theta_d)$

How to set θ_d to maximize $p(d | \Lambda)$?
(all other parameters are known)

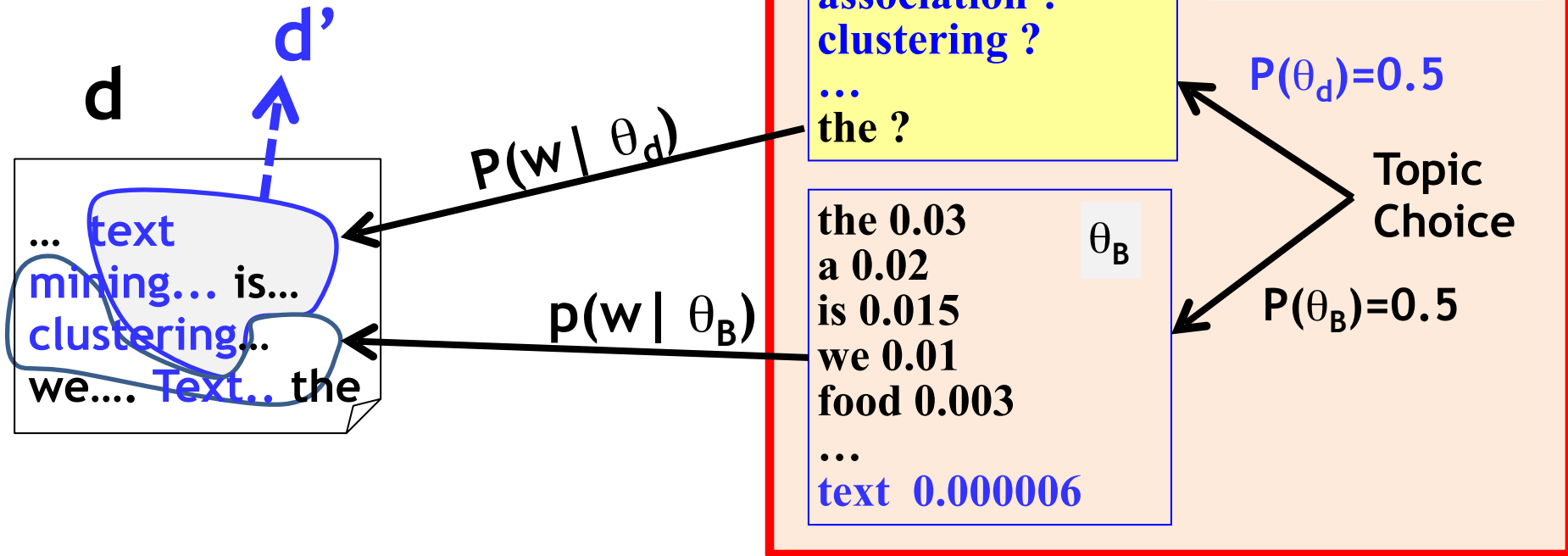
d

... text
mining... is...
clustering...
we... Text.. the



If we know which word is from which distribution...

$$p(w_i | \theta_d) = \frac{c(w_i, d')}{\sum_{w' \in V} c(w', d')}$$



Given all the parameters, infer the distribution a word is from...

Is “**text**” more likely from θ_d or θ_B ?

From θ_d
($Z=0$)?

From θ_B ($Z=1$)?

text | θ_d

$P(w | \theta_d)$

text 0.04 θ_d
mining 0.035
association 0.03
clustering 0.005
 ...
 the 0.000001

$p(\theta_d) + p(\theta_B) = 1$

$P(\theta_d) = 0.5$

Topic
Choice

$P(\theta_B) = 0.5$

$p(\theta_B)p(\text{“text”} | \theta_B)$
 $p(w | \theta_B)$

the 0.03 θ_B
 a 0.02
 is 0.015
 we 0.01
 food 0.003
 ...
text 0.000006

$p(z = 0 | w = \text{“text”}) =$

$$\frac{p(\theta_d)p(\text{“text”} | \theta_d)}{p(\theta_d)p(\text{“text”} | \theta_d) + p(\theta_B)p(\text{“text”} | \theta_B)}$$

The Expectation-Maximization (EM) Algorithm

Hidden Variable:

$$z \in \{0, 1\}$$

	z
the	1
paper	1
presents	1
a	1
text	0
mining	0
algorithm	0
for	1
clustering	0
...	...

Initialize $p(w|\theta_d)$ with random values.

Then iteratively improve it using E-step & M-step.
Stop when likelihood doesn't change.

$$p^{(n)}(z = 0 | w) = \frac{p(\theta_d)p^{(n)}(w | \theta_d)}{p(\theta_d)p^{(n)}(w | \theta_d) + p(\theta_B)p(w | \theta_B)}$$

E-step

How likely w is from θ_d

$$p^{(n+1)}(w | \theta_d) = \frac{c(w, d)p^{(n)}(z = 0 | w)}{\sum_{w' \in V} c(w', d)p^{(n)}(z = 0 | w')}$$

M-step

EM Computation in Action

E-step
$$p^{(n)}(z = 0 | w) = \frac{p(\theta_d)p^{(n)}(w | \theta_d)}{p(\theta_d)p^{(n)}(w | \theta_d) + p(\theta_B)p(w | \theta_B)}$$

M-step
$$p^{(n+1)}(w | \theta_d) = \frac{c(w, d)p^{(n)}(z = 0 | w)}{\sum_{w' \in V} c(w', d)p^{(n)}(z = 0 | w')}$$

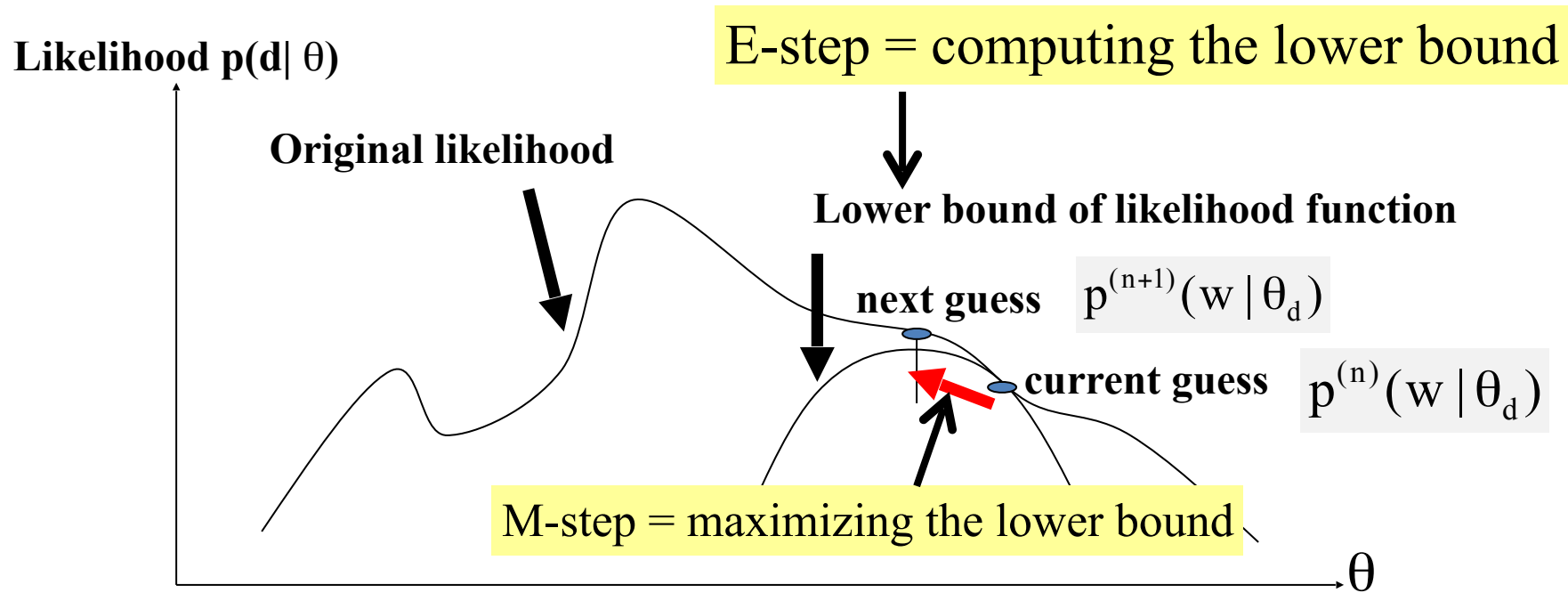
Assume $p(\theta_d) = p(\theta_B) = 0.5$
and $p(w | \theta_B)$ is known

Word	#	$p(w \theta_B)$	Iteration 1		Iteration 2		Iteration 3	
			$P(w \theta)$	$p(z=0 w)$	$P(w \theta)$	$P(z=0 w)$	$P(w \theta)$	$P(z=0 w)$
The	4	0.5	0.25	0.33	0.20	0.29	0.18	0.26
Paper	2	0.3	0.25	0.45	0.14	0.32	0.10	0.25
Text	4	0.1	0.25	0.71	0.44	0.81	0.50	0.93
Mining	2	0.1	0.25	0.71	0.22	0.69	0.22	0.69
Log-Likelihood			-16.96		-16.13		-16.02	

Likelihood increasing

“By products”: Are they also useful?

EM As Hill-Climbing → Converge to Local Maximum



A General Introduction to EM

Data: X (observed) + H (hidden) Parameter: θ

“Incomplete” likelihood: $L(\theta) = \log p(X | \theta)$

“Complete” likelihood: $L_c(\theta) = \log p(X, H | \theta)$

EM tries to iteratively maximize the incomplete likelihood:

Starting with an initial guess $\theta^{(0)}$,

1. E-step: compute the expectation of the complete likelihood

$$Q(\theta; \theta^{(n-1)}) = E_{\theta^{(n-1)}} [L_c(\theta) | X] = \sum_{h_i} p(H = h_i | X, \theta^{(n-1)}) \log P(X, h_i)$$

2. M-step: compute $\theta^{(n)}$ by maximizing the Q-function

$$\theta^{(n)} = \arg \max_{\theta} Q(\theta; \theta^{(n-1)}) = \arg \max_{\theta} \sum_{h_i} p(H = h_i | X, \theta^{(n-1)}) \log P(X, h_i)$$

Convergence Guarantee

Goal: maximizing “Incomplete” likelihood: $L(\theta) = \log p(X|\theta)$

I.e., choosing $\theta^{(n)}$, so that $L(\theta^{(n)}) - L(\theta^{(n-1)}) \geq 0$

Note that, since $p(X, H|\theta) = p(H|X, \theta) P(X|\theta)$, $L(\theta) = L_c(\theta) - \log p(H|X, \theta)$

$$L(\theta^{(n)}) - L(\theta^{(n-1)}) = L_c(\theta^{(n)}) - L_c(\theta^{(n-1)}) + \log [p(H|X, \theta^{(n-1)}) / p(H|X, \theta^{(n)})]$$

Taking expectation w.r.t. $p(H|X, \theta^{(n-1)})$,

$$\underbrace{L(\theta^{(n)}) - L(\theta^{(n-1)})}_{\text{Doesn't contain H}} = \underbrace{Q(\theta^{(n)}; \theta^{(n-1)}) - Q(\theta^{(n-1)}; \theta^{(n-1)})}_{\text{EM chooses } \theta^{(n)} \text{ to maximize } Q} + \underbrace{D(p(H|X, \theta^{(n-1)}) || p(H|X, \theta^{(n)}))}_{\text{KL-divergence, always non-negative}}$$

Doesn't contain H

EM chooses $\theta^{(n)}$ to maximize Q

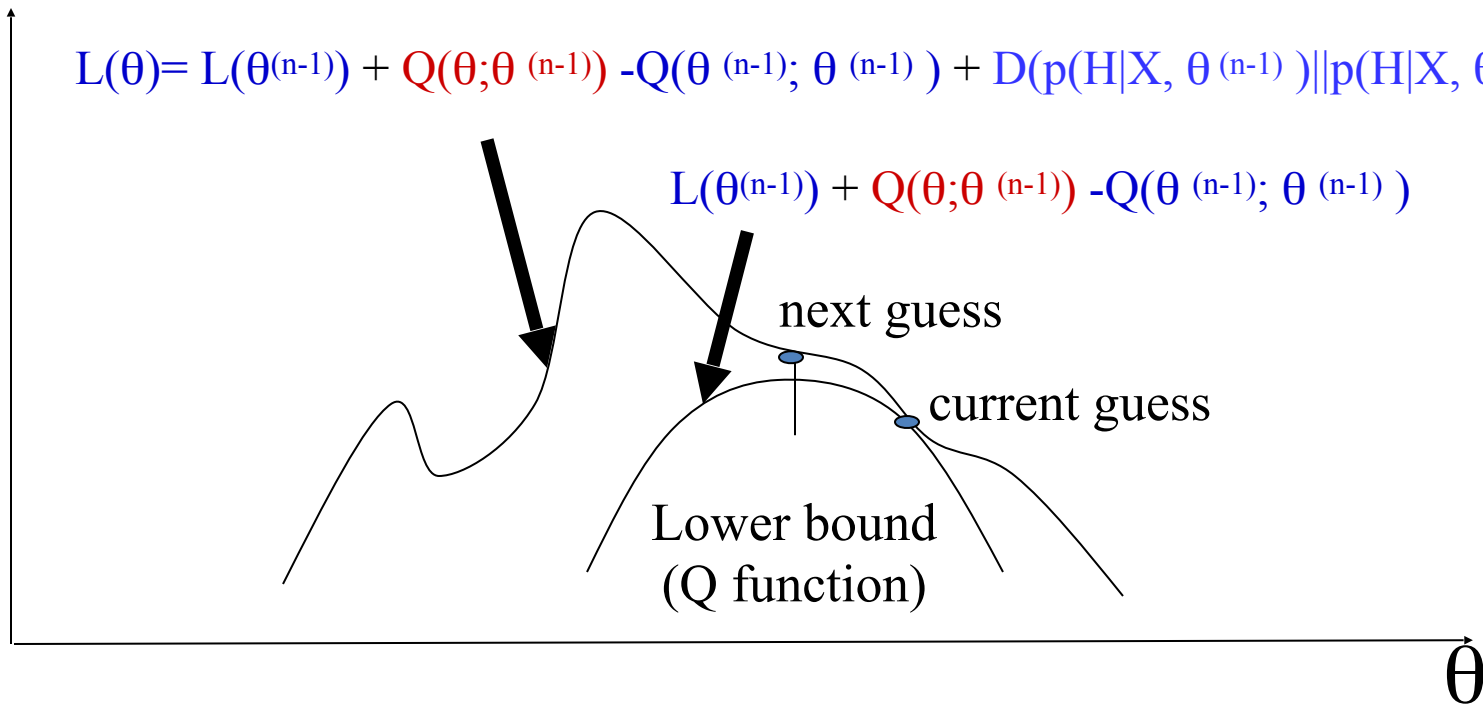
KL-divergence, always non-negative

Therefore, $L(\theta^{(n)}) \geq L(\theta^{(n-1)})!$

EM as Hill-Climbing: converging to a local maximum

Likelihood $p(X|\theta)$

$$L(\theta) = L(\theta^{(n-1)}) + \underbrace{Q(\theta; \theta^{(n-1)}) - Q(\theta^{(n-1)}; \theta^{(n-1)})}_{\text{Lower bound}} + D(p(H|X, \theta^{(n-1)}) || p(H|X, \theta))$$



E-step = computing the lower bound
M-step = maximizing the lower bound

Document as a Sample of Mixed Topics

Topic θ_1

government 0.3
response 0.2

...

Topic θ_2

city 0.2
new 0.1
orleans 0.05

...

...

Topic θ_k

donate 0.1
relief 0.05
help 0.02

...

Background θ_B

the 0.04
a 0.03

...

Blog article about “Hurricane Katrina”

[Criticism of government response to the hurricane primarily consisted of criticism of its response to the approach of the storm and its aftermath, specifically in the delayed response] to the [flooding of New Orleans. ... 80% of the 1.3 million residents of the greater New Orleans metropolitan area evacuated] ... [Over seventy countries pledged monetary donations or other assistance]. ...

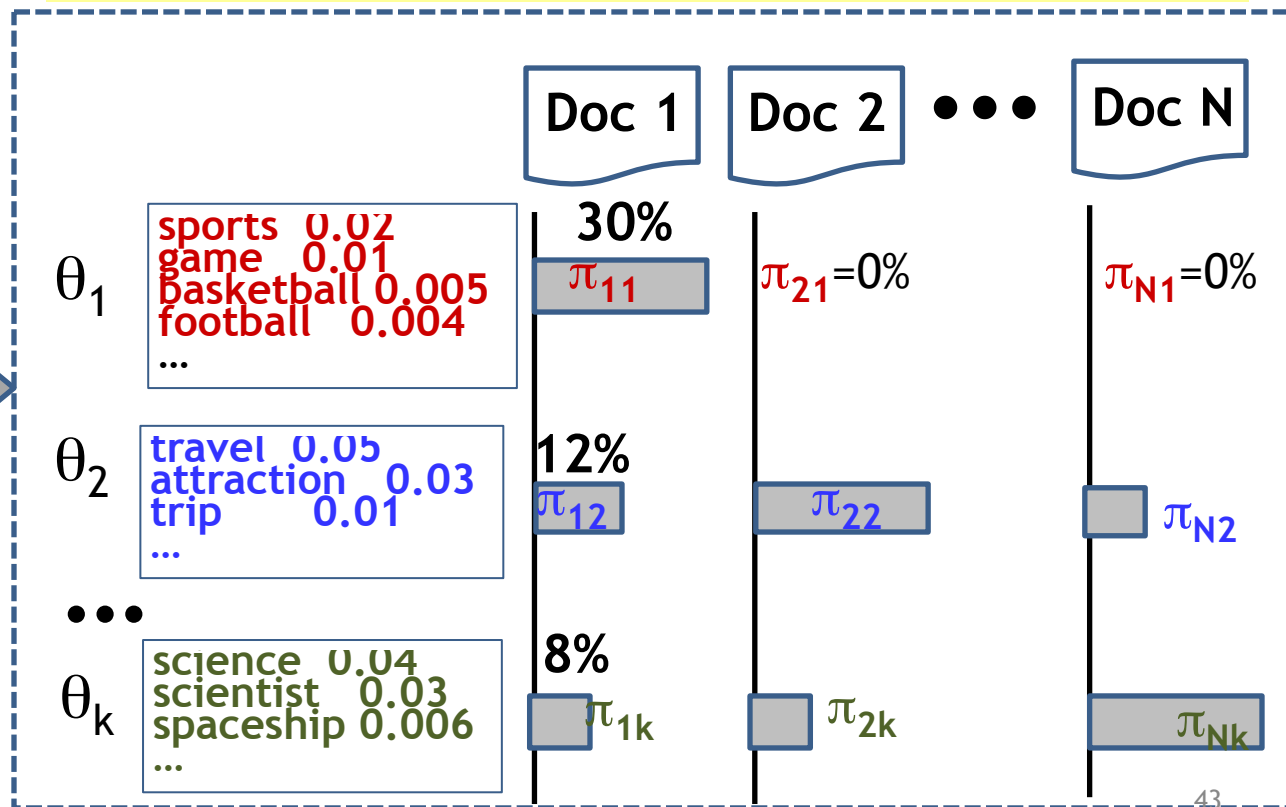
Many applications are possible if we can “decode” the topics in text...

Mining Multiple Topics from Text

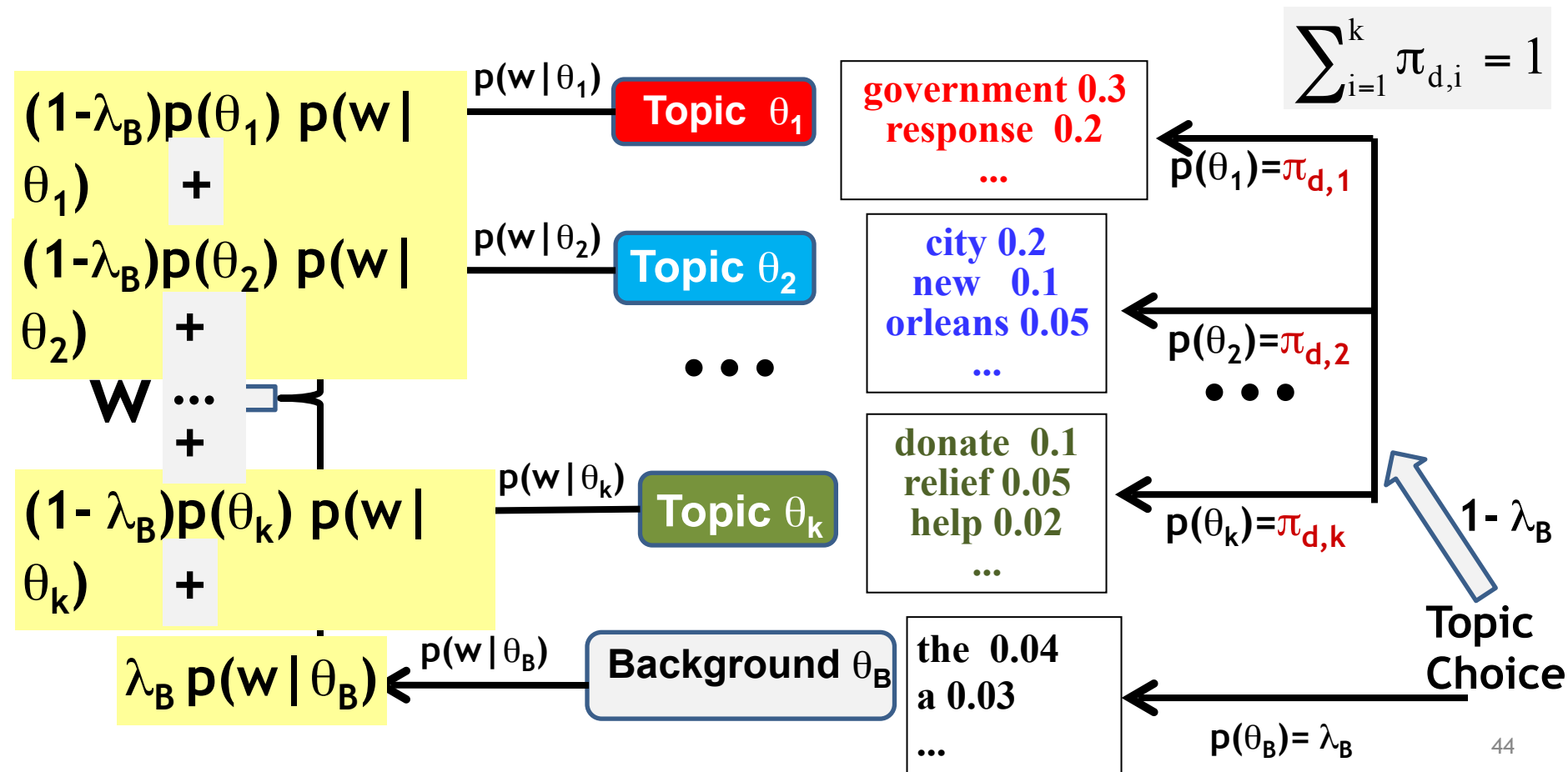
OUTPUT: $\{ \theta_1, \dots, \theta_k \}, \{ \pi_{i1}, \dots, \pi_{ik} \}$

INPUT: C, k, V

Text Data



Generating Text with Multiple Topics: $p(w)=?$



Probabilistic Latent Semantic Analysis (PLSA)

Percentage of
background words
(known)

Background
LM (known)

Coverage of topic θ_j in doc d

Prob. of word w in topic θ_j

$$p_d(w) = \lambda_B p(w | \theta_B) + (1 - \lambda_B) \sum_{j=1}^k \pi_{d,j} p(w | \theta_j)$$

$$\log p(d) = \sum_{w \in V} c(w, d) \log [\lambda_B p(w | \theta_B) + (1 - \lambda_B) \sum_{j=1}^k \pi_{d,j} p(w | \theta_j)]$$

$$\log p(C | \Lambda) = \sum_{d \in C} \sum_{w \in V} c(w, d) \log [\lambda_B p(w | \theta_B) + (1 - \lambda_B) \sum_{j=1}^k \pi_{d,j} p(w | \theta_j)]$$

Unknown Parameters: $\Lambda = (\{\pi_{d,j}\}, \{\theta_j\})$, $j=1, \dots, k$

How many unknown parameters are there in

ML Parameter Estimation

$$p_d(w) = \lambda_B p(w | \theta_B) + (1 - \lambda_B) \sum_{j=1}^k \pi_{d,j} p(w | \theta_j)$$

$$\log p(d) = \sum_{w \in \mathcal{V}} c(w, d) \log[\lambda_B p(w | \theta_B) + (1 - \lambda_B) \sum_{j=1}^k \pi_{d,j} p(w | \theta_j)]$$

$$\log p(C | \Lambda) = \sum_{d \in \mathcal{C}} \sum_{w \in \mathcal{V}} c(w, d) \log[\lambda_B p(w | \theta_B) + (1 - \lambda_B) \sum_{j=1}^k \pi_{d,j} p(w | \theta_j)]$$

Constrained Optimization: $\Lambda^* = \arg \max_{\Lambda} p(C | \Lambda)$

$$\forall j \in [1, k], \sum_{i=1}^M p(w_i | \theta_j) = 1$$

$$\forall d \in \mathcal{C}, \sum_{j=1}^k \pi_{d,j} = 1$$

EM Algorithm for PLSA: E-Step

Hidden Variable (=topic indicator): $z_{d,w} \in \{B, 1, 2, \dots, k\}$

Probability that w in doc d is generated from topic θ_j

Use of Bayes Rule

$$p(z_{d,w} = j) = \frac{\pi_{d,j}^{(n)} p^{(n)}(w | \theta_j)}{\sum_{j'=1}^k \pi_{d,j'}^{(n)} p^{(n)}(w | \theta_{j'})}$$

$$p(z_{d,w} = B) = \frac{\lambda_B p(w | \theta_B)}{\lambda_B p(w | \theta_B) + (1 - \lambda_B) \sum_{j=1}^k \pi_{d,j}^{(n)} p^{(n)}(w | \theta_j)}$$

Probability that w in doc d is generated from background θ_B

EM Algorithm for PLSA: M-Step

Hidden Variable (=topic indicator): $z_{d,w} \in \{B, 1, 2, \dots, k\}$

Re-estimated probability of doc d covering topic θ_j based on
 “allocated” word counts to topic θ_j

$$\pi_{d,j}^{(n+1)} = \frac{\sum_{w \in V} c(w, d)(1 - p(z_{d,w} = B))p(z_{d,w} = j)}{\sum_{j'} \sum_{w \in V} c(w, d)(1 - p(z_{d,w} = B))p(z_{d,w} = j')}$$

$$p^{(n+1)}(w | \theta_j) = \frac{\sum_{d \in D} c(w, d)(1 - p(z_{d,w} = B))p(z_{d,w} = j)}{\sum_{w' \in V} \sum_{d \in D} c(w', d)(1 - p(z_{d,w'} = B))p(z_{d,w'} = j)}$$

Re-estimated probability of word w for topic θ_j

Computation of the EM Algorithm

- Initialize all unknown parameters randomly
- Repeat until likelihood converges

– E-step $p(z_{d,w} = j) \propto \pi_{d,j}^{(n)} p^{(n)}(w | \theta_j)$ $\sum_{j=1}^k p(z_{d,w} = j) = 1$

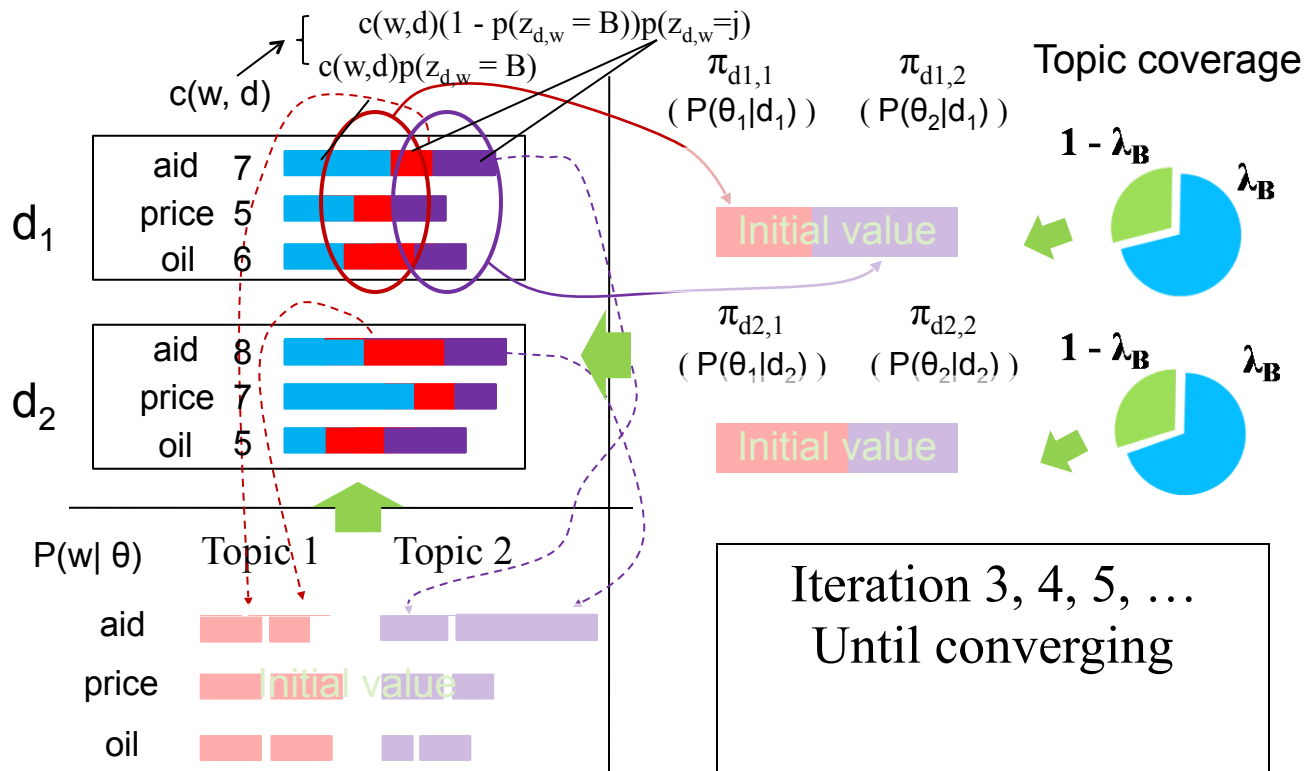
– M-step $p(z_{d,w} = B) \propto \lambda_B p(w | \theta_B) \leftarrow$ What's the normalizer for this one?

$$\pi_{d,j}^{(n+1)} \propto \sum_{w \in V} c(w, d) (1 - p(z_{d,w} = B)) p(z_{d,w} = j) \quad \forall d \in C, \sum_{j=1}^k \pi_{d,j} = 1$$

$$p^{(n+1)}(w | \theta_j) \propto \sum_{d \in C} c(w, d) (1 - p(z_{d,w} = B)) p(z_{d,w} = j) \quad \forall j \in [1, k], \sum_{w \in V} p(w | \theta_j) = 1$$

In general, accumulate counts, and then normalize

Illustration of EM Algorithm for PLSA



Applications of Mixture Models for Text Mining

Likelihood:

$$p(d | \theta_1 \oplus \theta_2) = \prod_{w \in V} [\lambda p(w | \theta_1) + (1 - \lambda) p(w | \theta_2)]^{c(w, d)}$$
$$\log p(d | \theta_1 \oplus \theta_2) = \sum_{w \in V} c(w, d) \log [\lambda p(w | \theta_1) + (1 - \lambda) p(w | \theta_2)]$$

Application Scenarios:

- $p(w | \theta_1)$ & $p(w | \theta_2)$ are known; estimate λ

↙ The doc is about text mining and food nutrition, how much percent is about text mining?

↙ 30% of the doc is about text mining, what's the rest about?

- $p(w | \theta_1)$ & λ are known; estimate $p(w | \theta_2)$

↙ The doc is about text mining, is it also about some other topic, and if so to what extent?

- $p(w | \theta_1)$ is known; estimate λ & $p(w | \theta_2)$

↙ 30% of the doc is about one topic and 70% is about another, what are these two topics?

- λ is known; estimate $p(w | \theta_1)$ & $p(w | \theta_2)$

↙ The doc is about two subtopics, find out what these two subtopics are and to what extent the doc covers each.

Use PLSA for Text Mining

- PLSA would be able to generate
 - Topic coverage in each document: $p(\pi_d = j)$
 - Word distribution for each topic: $p(w | \theta_j)$
 - Topic assignment at the word level for each document
 - The number of topics must be given in advance
- These probabilities can be used in many different ways
 - θ_j naturally serves as a word cluster
 - $\pi_{d,j}$ can be used for document clustering $j^* = \arg \max_j \pi_{d,j}$
 - Contextual text mining: Make these parameters conditioned on context, e.g.,
 - $p(\theta_j | \text{time})$, from which we can compute/plot $p(\text{time} | \theta_j)$
 - $p(\theta_j | \text{location})$, from which we can compute/plot $p(\text{loc} | \theta_j)$

How to Help Users Interpret a Topic Model? [Mei et al. 07b]

- Use top words
 - automatic, but hard to make sense
- Human generated labels
 - Make sense, but don't scale up

Term, relevance,
weight, feedback

term	0.16
relevance	0.08
weight	0.07
feedback	0.04
independence	0.03
model	0.03
frequent	0.02
probabilistic	0.02
document	0.02
...	

insulin
foraging
foragers
collected
grains
loads
collection
nectar
...

Retrieval Models

Question: Can we automatically generate understandable labels for topics?

What is a Good Label?

Retrieval models

term	0.1599
relevance	0.0752
weight	0.0660
feedback	0.0372
independence	0.0311
model	0.0310
frequent	0.0233
probabilistic	0.0188
document	0.0173
...	

A topic from
[Mei & Zhai 06b]

- Semantically close (**relevance**)
- **Understandable** - phrases?
- High **coverage** inside topic
- **Discriminative** across topics
- ... ~~iPod Nano~~

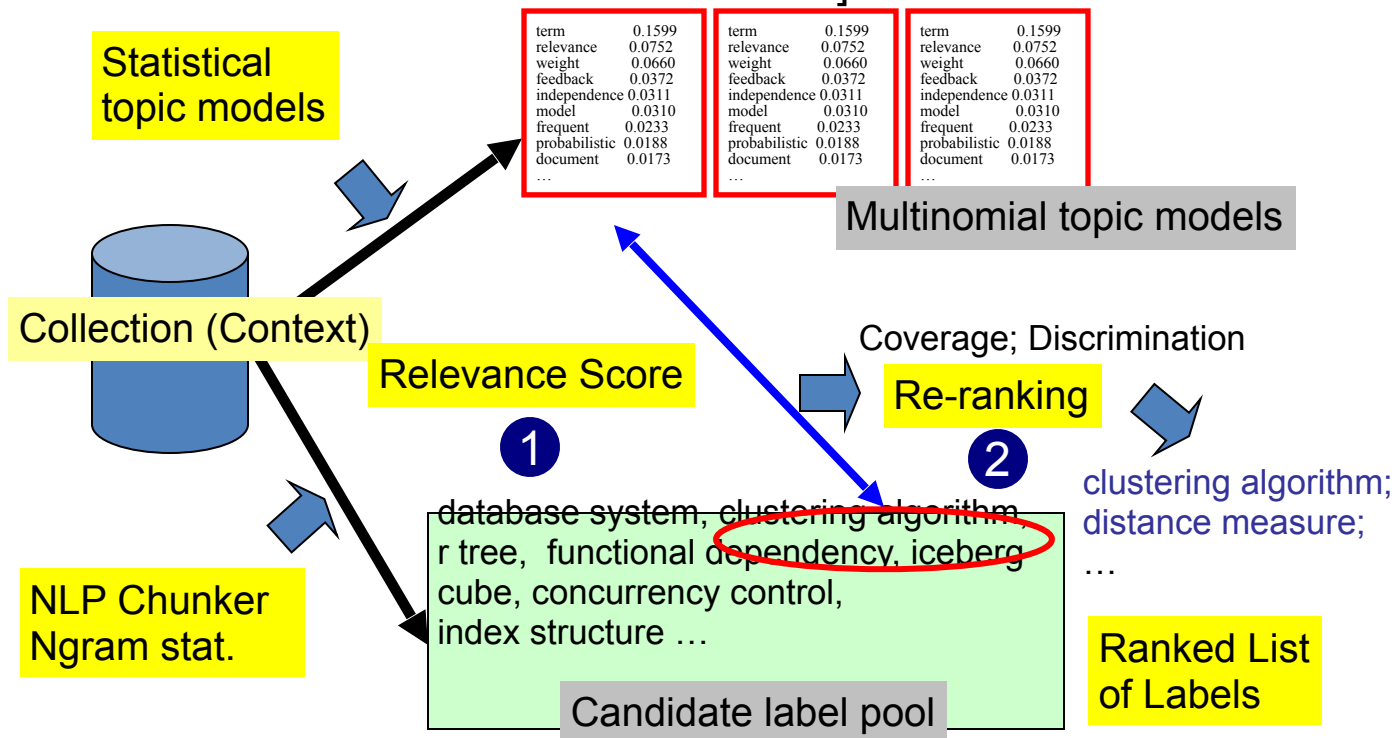
~~じょうほうけんさく~~

~~Pseudo-feedback~~

~~Information Retrieval~~

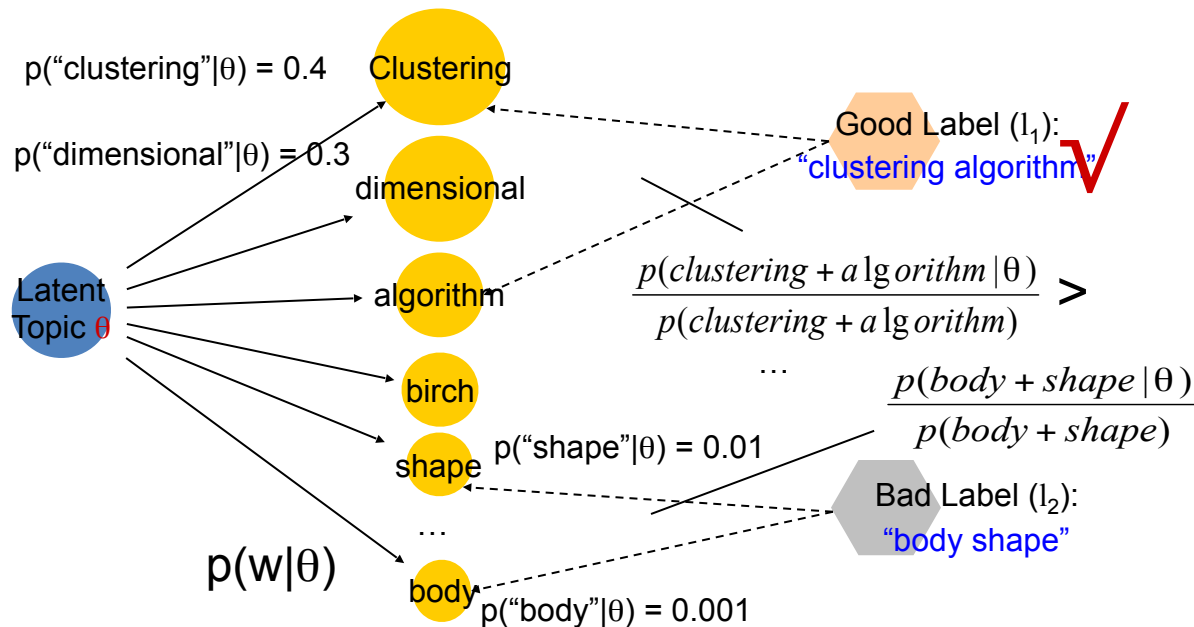
Automatic Labeling of Topics [Mei et

al. 07b]



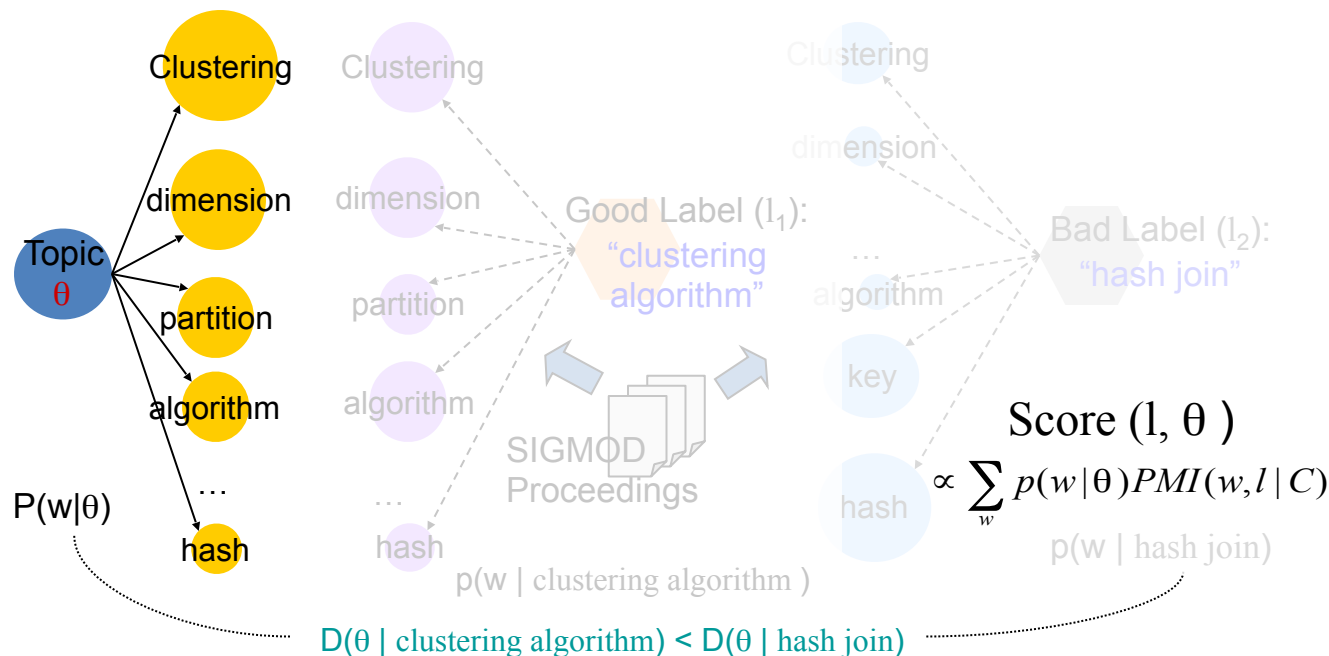
Relevance: the Zero-Order Score

- Intuition: prefer phrases well covering top words

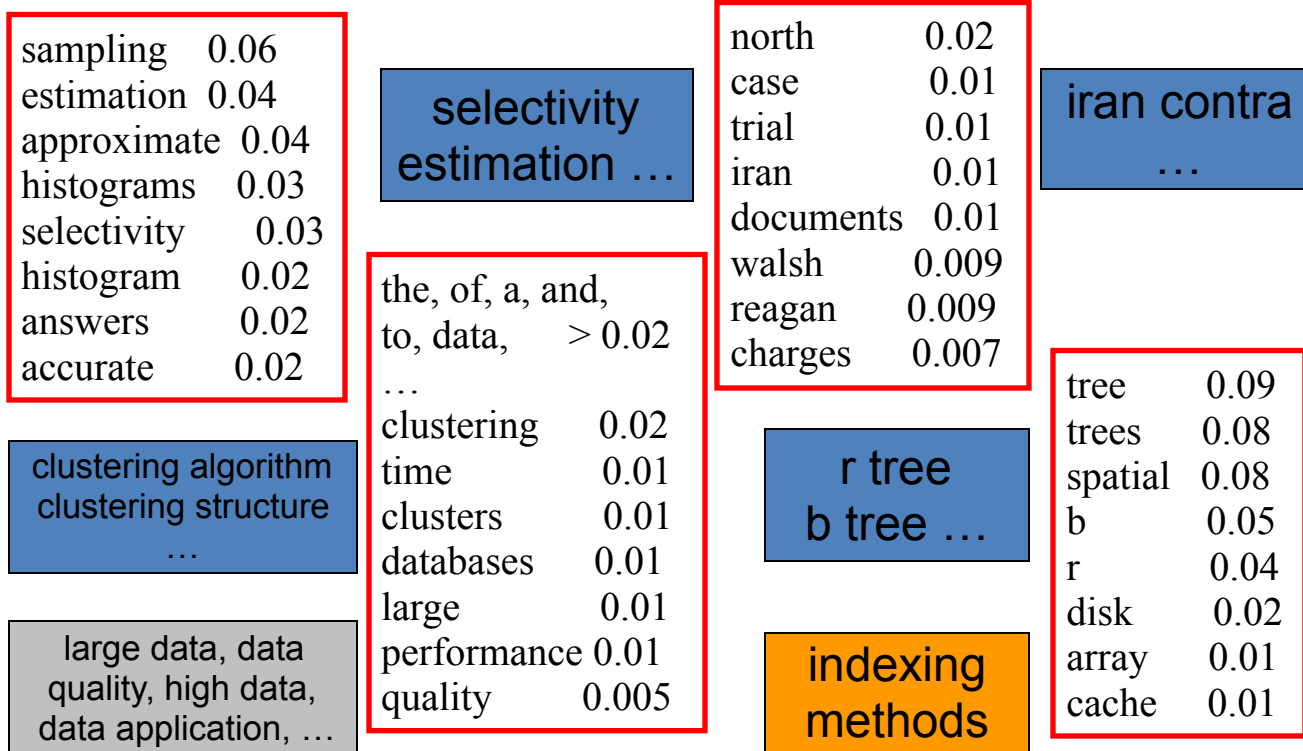


Relevance: the First-Order Score

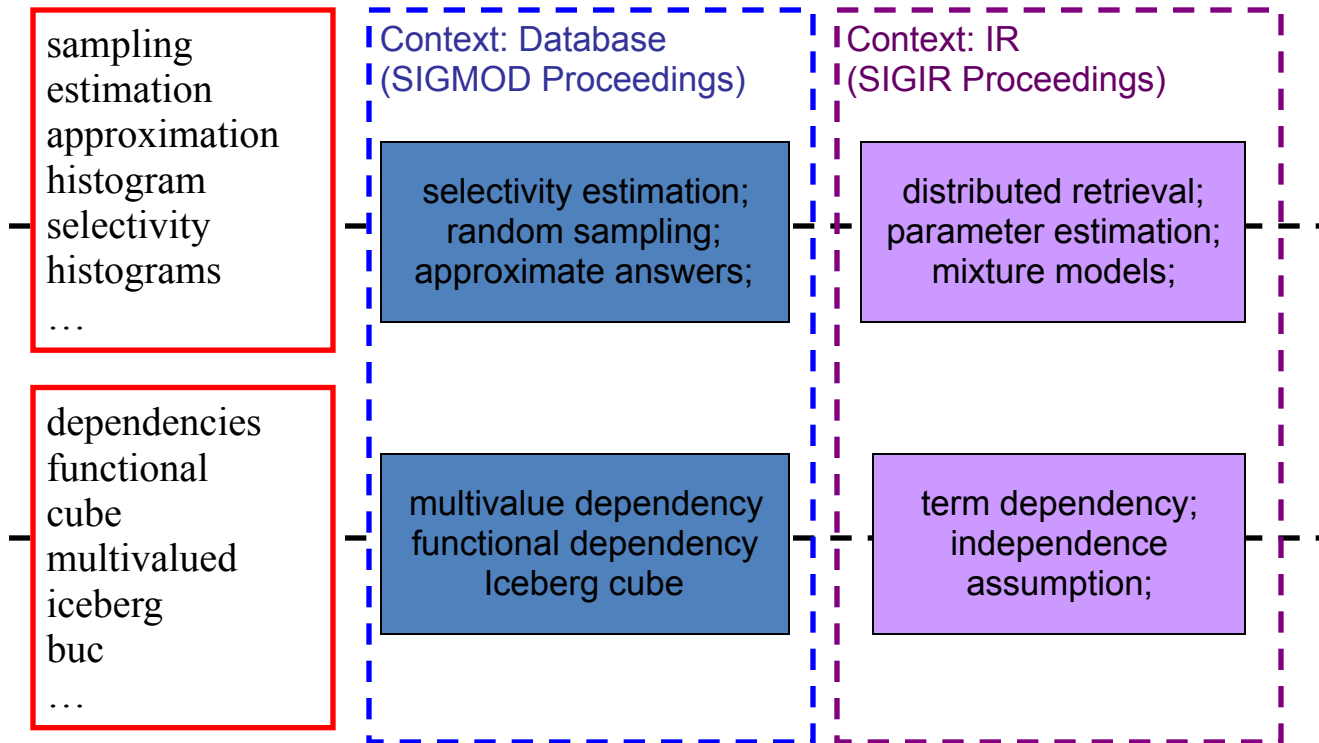
- Intuition: prefer phrases with similar context (distribution)



Results: Sample Topic Labels



Results: Contextual-Sensitive Labeling



Extensions of PLSA

- PLSA with prior knowledge → User-controlled PLSA
- PLSA for text data with context → Contextualized PLSA
- PLSA as a generative model → Latent Dirichlet Allocation (Bayesian inference for mixture models, covered later)

PLSA with Prior Knowledge

- Users may have expectations about which topics to analyze:
 - We expect to see “retrieval models” as a topic in IR
 - We want to see aspects such as “battery” and “memory” for opinions about a laptop
- Users may have knowledge about what topics are (or are NOT) covered in a document
 - Tags = topics → A doc can only be generated using topics corresponding to the tags assigned to the document
- We can incorporate such knowledge as priors of PLSA model

Maximum a Posteriori (MAP) Estimate

$$\Lambda^* = \arg \max_{\Lambda} p(\Lambda) p(Data | \Lambda)$$

- We may use $p(\Lambda)$ to encode all kinds of preferences and constraints, e.g.,
 - $p(\Lambda) > 0$ if and only if one topic is precisely “background”: $p(w | \theta_B)$
 - $p(\Lambda) > 0$ if and only if for a particular doc d , $\pi_{d,3} = 0$ and $\pi_{d,1} = 1/2$
 - $p(\Lambda)$ favors a Λ with topics that assign high probabilities to some particular words
- The MAP estimate (with conjugate prior) can be computed using a similar EM algorithm to the ML estimate with smoothing to reflect prior preferences

EM Algorithm with Conjugate Prior on $p(w | \theta_j)$

$$p(z_{d,w} = j) = \frac{\pi_{d,j}^{(n)} p^{(n)}(w | \theta_j)}{\sum_{j'=1}^k \pi_{d,j'}^{(n)} p^{(n)}(w | \theta_{j'})}$$

Prior: $p(w | \theta'_j)$

battery 0.5
life 0.5

$$p(z_{d,w} = B) = \frac{\lambda_B p(w | \theta_B)}{\lambda_B p(w | \theta_B) + (1 - \lambda_B) \sum_{j=1}^k \pi_{d,j}^{(n)} p^{(n)}(w | \theta_j)}$$

$$\pi_{d,j}^{(n+1)} = \frac{\sum_{w \in V} c(w, d) (1 - p(z_{d,w} = B)) p(z_{d,w} = j)}{\sum_{j'} \sum_{w \in V} c(w, d) (1 - p(z_{d,w} = B)) p(z_{d,w} = j')}$$

Pseudo counts of w
from prior θ'

$$p^{(n+1)}(w | \theta_j) = \frac{\sum_{d \in C} c(w, d) (1 - p(z_{d,w} = B)) p(z_{d,w} = j) + \mu p(w | \theta'_j)}{\sum_{w' \in V} \sum_{d \in C} c(w', d) (1 - p(z_{d,w'} = B)) p(z_{d,w'} = j) + \mu}$$

What if $\mu=0$? What if $\mu=+\infty$?

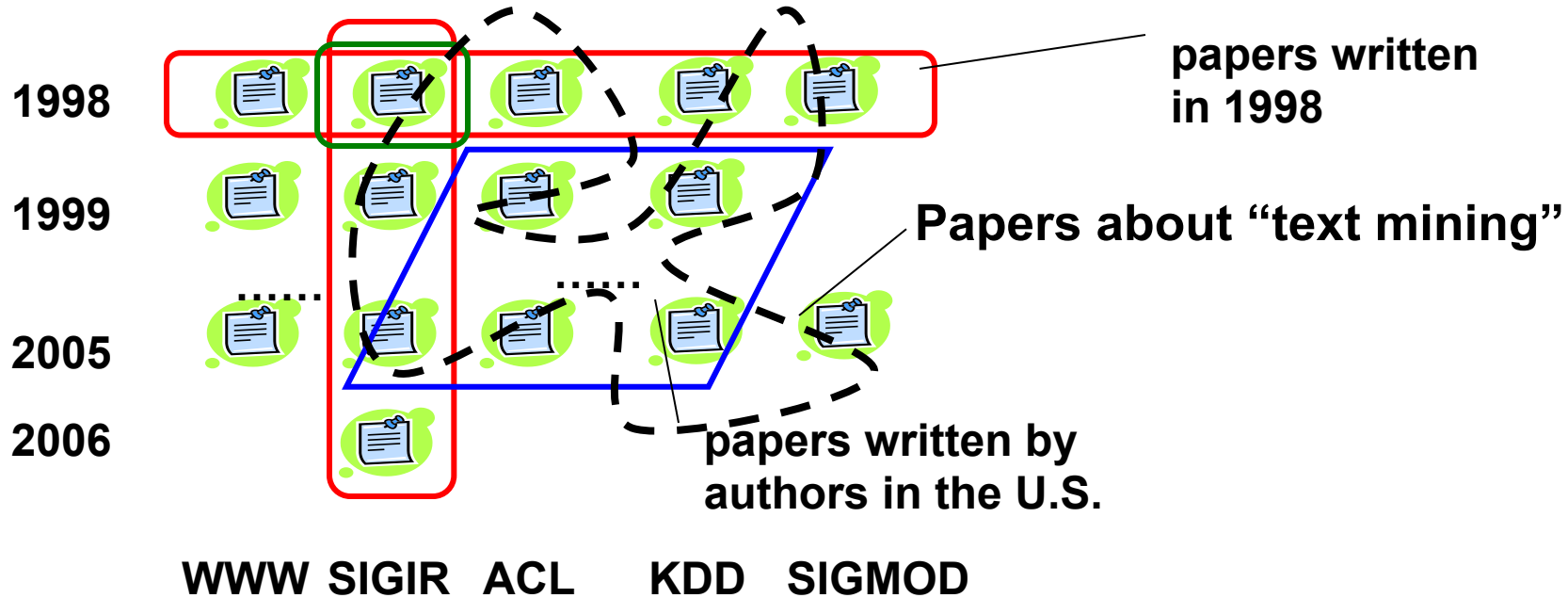
Sum of all pseudo counts

We may also set any parameter to a constant (including 0) as needed

Contextual Text Mining: Motivation

- Text often has rich context information
 - Direct context (Meta-Data): time, location, authors, source, ...
 - Indirect context (additional data related to meta-data): social network of the author, author's age, other text from the same source, etc.
 - Any related data can be regarded as context
- Context can be used to
 - Partition text data for comparative analysis
 - Provide meaning to the discovered topics

Context = Partitioning of Text



Enables discovery of knowledge associated with different context as needed

Many Interesting Questions Require Contextual Text Mining

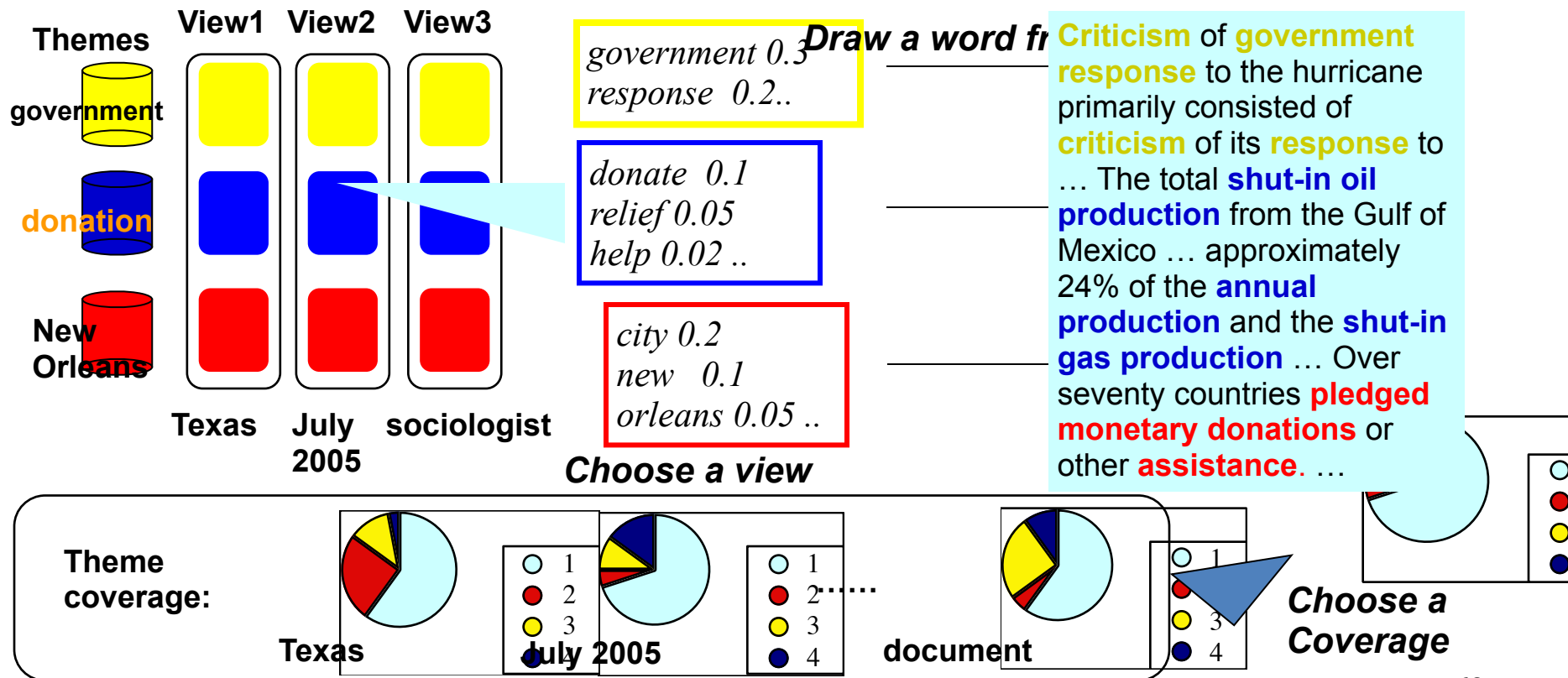
- What topics have been gaining increasing attention recently in data mining research? (time as context)
- Is there any difference in the responses of people in different regions to the event? (location as context)
- What are the common research interests of two researchers? (authors as context)
- Is there any difference in the research topics published by authors in the USA and those outside? (author's affiliation and location as context)
- Is there any difference in the opinions about a topic expressed on one social network and another? (social network of authors and topic as context)
- Are there topics in news data that are correlated with sudden changes in stock prices? (time series as context)
- What issues “mattered” in the 2012 presidential election? (time series as context)

Contextual Probabilistic Latent Semantic Analysis (CPLSA)

[Mei & Zhai 06]

- General idea:
 - Explicitly add interesting context variables into a generative model (→ enable discovery contextualized topics)
 - Context influences both coverage and content variation of topics
- As an extension of PLSA
 - Model the conditional likelihood of text given context
 - Assume context-dependent views of a topic
 - Assume context-dependent topic coverage
 - EM algorithm can still be used for parameter estimation
 - Estimated parameters naturally contain context variables, enabling contextual text mining

Generation Process of CPLSA



Comparing News Articles [Zhai et al. 04]

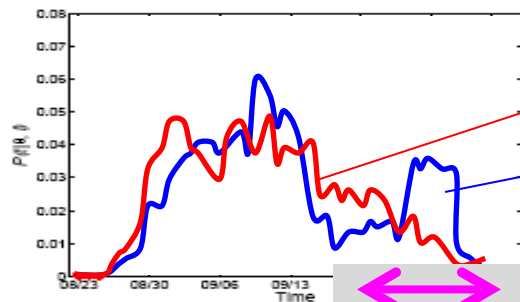
Iraq War (30 articles) vs. Afghan War (26 articles)

The common theme indicates that “United Nations” is involved in both wars

	Cluster 1	Cluster 2	Cluster 3
Common Theme	united 0.042 nations 0.04 ...	killed 0.035 month 0.032 deaths 0.023
Iraq Theme	n 0.03 Weapons 0.024 Inspections 0.023 ...	troops 0.016 hoon 0.015 sanches 0.012
Afghan Theme	Northern 0.04 alliance 0.04 kabul 0.03 taleban 0.025 aid 0.02	taleban 0.026 rumsfeld 0.02 hotel 0.012 front 0.011

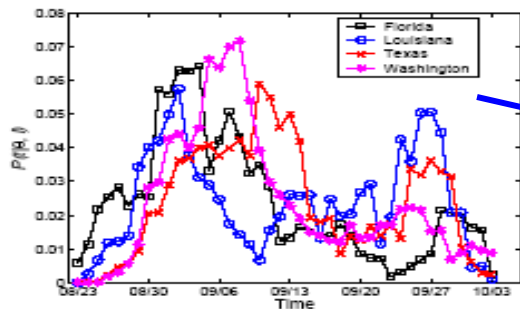
Collection-specific themes indicate different roles of “United Nations” in the two wars

Theme Life Cycles in Blog Articles About “Hurricane Katrina” [Mei et al. 06]



(a) Theme life cycles in Texas (Hurricane Katrina)

Hurricane Rita



(b) Theme “New Orleans” over states (Hurricane Katrina)

Oil Price

New Orleans

price 0.0772
oil 0.0643
gas 0.0454
increase 0.0210
product 0.0203
fuel 0.0188
company 0.0182
 ...

city 0.0634
orleans 0.0541
new 0.0342
louisiana 0.0235
flood 0.0227
evacuate 0.0211
storm 0.0177
 ...

Spatial Distribution of the Topic “Government Response” in Blog Articles About Hurricane Katrina [Mei et al. 06]



(a) Week1: 08/23-08/29



(b) Week Two: 08/30-09/05



(c) Week Three: 09/06-09/12

Theme 1
Government Response
bush 0.0716374
president 0.0610942
federal 0.0514114
govern 0.0476977
fema 0.0474692
administrate 0.0233903
response 0.0208351
brown 0.0199573
blame 0.0170033
governor 0.0142153



(d) Week Four: 09/13-09/19



(e) Week Five: 09/20-09/26

Event Impact Analysis: IR Research [Mei & Zhai 06]

Topic: retrieval models

<i>term</i>	0.1599
<i>relevance</i>	0.0752
<i>weight</i>	0.0660
<i>feedback</i>	0.0372
<i>independence</i>	0.0311
<i>model</i>	0.0310
<i>frequent</i>	0.0233
<i>probabilistic</i>	0.0188
<i>document</i>	0.0173
...	

<i>vector</i>	0.0514
<i>concept</i>	0.0298
<i>extend</i>	0.0297
<i>model</i>	0.0291
<i>space</i>	0.0236
<i>boolean</i>	0.0151
<i>function</i>	0.0123
<i>feedback</i>	0.0077
...	

<i>xml</i>	0.0678
<i>email</i>	0.0197
<i>model</i>	0.0191
<i>collect</i>	0.0187
<i>judgment</i>	0.0102
<i>rank</i>	0.0097
<i>subtopic</i>	0.0079
...	

SIGIR papers

A seminal paper [Croft & Ponte 98]

1992

Star

<i>probabilist</i>	0.0778
<i>model</i>	0.0432
<i>logic</i>	0.0404
<i>ir</i>	0.0338
<i>boolean</i>	0.0281
<i>algebra</i>	0.0200
<i>estimate</i>	0.0119
<i>weight</i>	0.0111
...	

1998

<i>model</i>	0.1687
<i>language</i>	0.0753
<i>estimate</i>	0.0520
<i>parameter</i>	0.0281
<i>distribution</i>	0.0268
<i>probable</i>	0.0205
<i>smooth</i>	0.0198
<i>markov</i>	0.0137
<i>likelihood</i>	0.0059
...	

year

Topic Analysis with Network Context

- The **context** of a text article can form a **network**, e.g.,
 - Authors of research articles may form **collaboration networks**
 - Authors of social media content form **social networks**
 - Locations associated with text can be connected to form a **geographic network**
- **Benefit of joint analysis** of text and its network context
 - Network imposes **constraints** on topics in text (**authors connected in a network tend to write about similar topics**)
 - Text helps **characterize** the content associated with each subnetwork (e.g., difference in opinions expressed in two subnetworks?)

Network Supervised Topic Modeling: General Idea [Mei et al. 08]

- Probabilistic topic modeling as optimization: maximize likelihood

$$\Lambda^* = \arg \max_{\Lambda} p(\text{TextData} \mid \Lambda)$$

- Main idea: network imposes constraints on model parameters Λ
 - The text at two adjacent nodes of the network tends to cover similar topics
 - Topic distributions are smoothed over adjacent nodes
 - Add network-induced regularizers to the likelihood objective function

Any generative model

Any network

$$\Lambda^* = \arg \max_{\Lambda} f(p(\text{TextData} \mid \Lambda), r(\Lambda, \text{Network}))$$

Any way to combine

Any regularizer

Instantiation: NetPLSA [Mei et al. 08]

Network-induced prior: Neighbors have similar topic distribution

Modified objective function

PLSA log-likelihood

Text collection

$$O(C, G) = (1 - \lambda) \cdot \left(\sum_d \sum_w c(w, d) \log \sum_{j=1}^k p(\theta_j | d) p(w | \theta_j) \right)$$

Network graph

$$+ \lambda \cdot \left(-\frac{1}{2} \sum_{\langle u, v \rangle \in E} \frac{w(u, v)}{\sum_{j=1}^k (p(\theta_j | u) - p(\theta_j | v))^2} \right)$$

Influence of network constraint

Weight of edge (u,v)

Quantify the difference in the topic coverage at node u and v

Mining 4 Topical Communities: Results of PLSA

Can't uncover the 4 communities (IR, DM, ML, Web)

Topic 1		Topic 2		Topic 3		Topic 4	
term	0.02	peer	0.02	visual	0.02	interface	0.02
question	0.02	patterns	0.01	analog	0.02	towards	0.02
protein	0.01	mining	0.01	neurons	0.02	browsing	0.02
training	0.01	clusters	0.01	vlsi	0.01	xml	0.01
weighting	0.01	stream	0.01	motion	0.01	generation	0.01
multiple	0.01	frequent	0.01	chip	0.01	design	0.01
recognition	0.01	e	0.01	natural	0.01	engine	0.01
relations	0.01	page	0.01	cortex	0.01	service	0.01
library	0.01	gene	0.01	spike	0.01	social	0.01

Mining 4 Topical Communities: Results of NetPLSA

Uncovers the 4 communities well

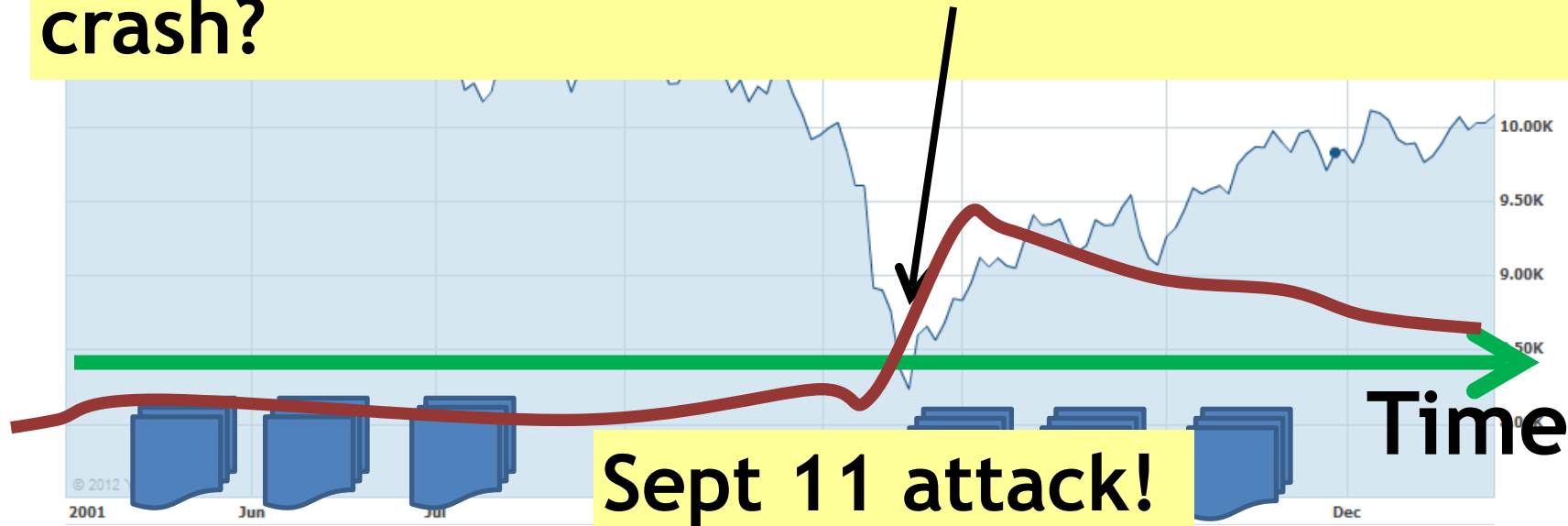
Information Retrieval		Data Mining		Machine Learning		Web	
retrieval	0.13	mining	0.11	neural	0.06	web	0.05
information	0.05	data	0.06	learning	0.02	services	0.03
document	0.03	discovery	0.03	networks	0.02	semantic	0.03
query	0.03	databases	0.02	recognition	0.02	services	0.03
text	0.03	rules	0.02	analog	0.01	peer	0.02
search	0.03	association	0.02	vlsi	0.01	ontologies	0.02
evaluation	0.02	patterns	0.02	neurons	0.01	rdf	0.02
user	0.02	frequent	0.01	gaussian	0.01	management	0.01
relevance	0.02	streams	0.01	network	0.01	ontology	0.01

Text Information Network

- In general, we can view text data that naturally “lives” in a rich information network with all other related data
- Text data can be associated with
 - Nodes of the network
 - Edges of the network
 - Paths of the network
 - Subnetworks
 - ...
- Analysis of text should be using the entire network!

Text Mining for Understanding Time Series

What might have caused the stock market crash?



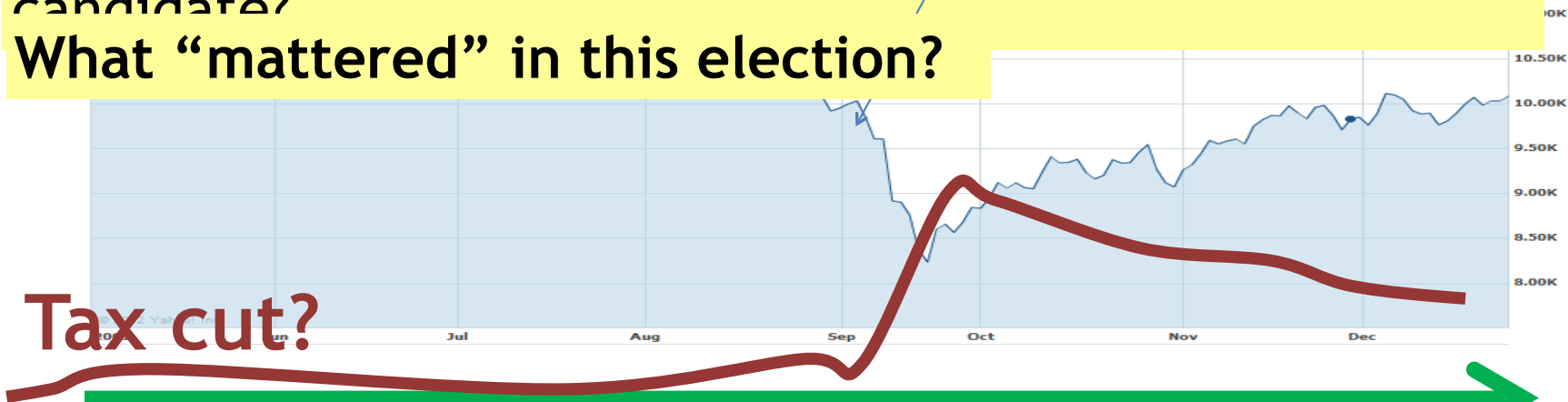
Any clues in the companion news stream?

Dow Jones Industrial Average [Source: Yahoo Finance]

Analysis of Presidential Prediction Markets

What might have caused the sudden drop of price for this candidate?

What “mattered” in this election?



Tax cut?



...

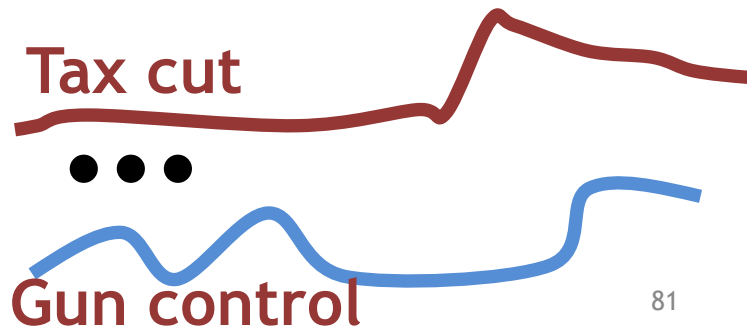


Time

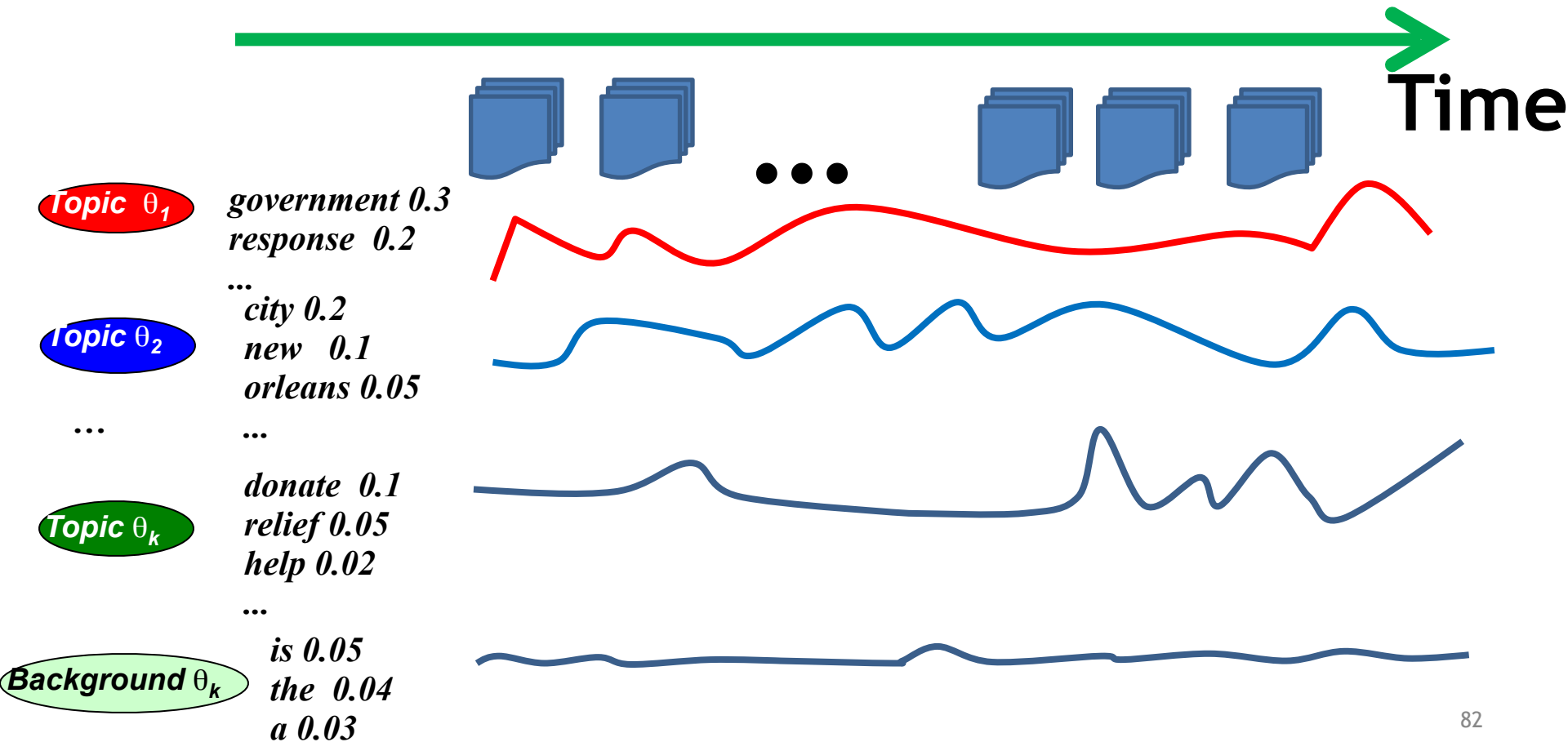
Any clues in the companion news stream?

Joint Analysis of Text and Time Series to Discover “Causal Topics”

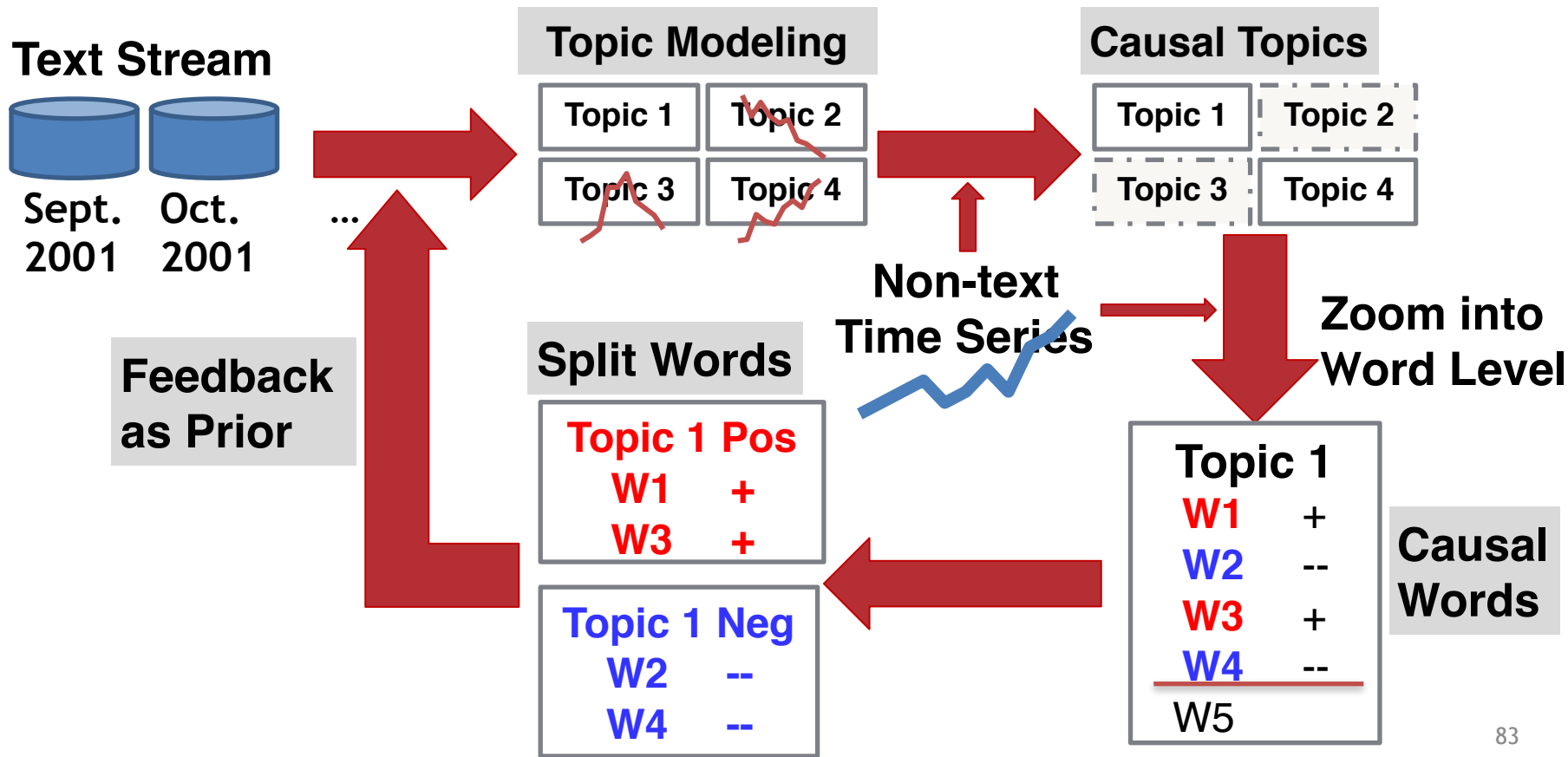
- Input:
 - Time series
 - Text data produced in a similar time period (text stream)
- Output
 - Topics whose coverage in the text stream has strong correlations with the time series (“causal” topics)



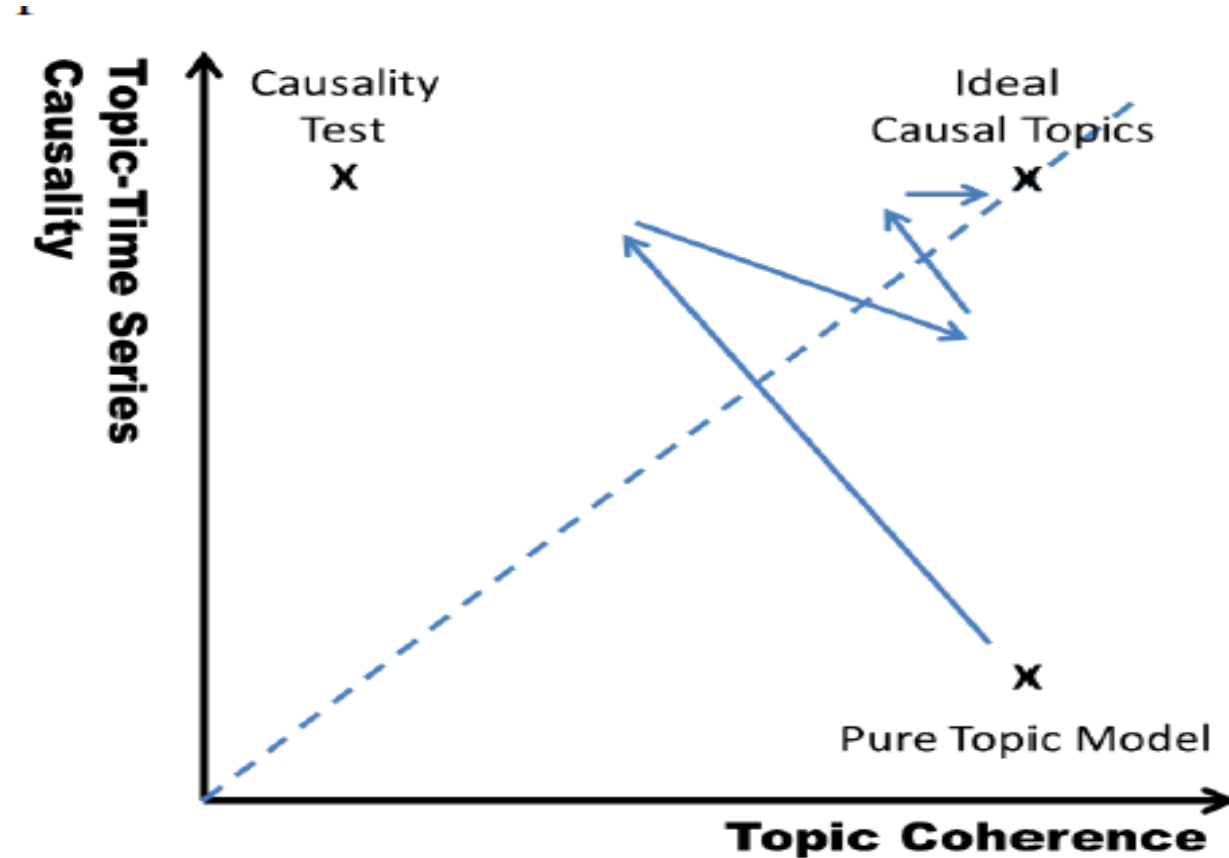
When a Topic Model Applied to Text Stream



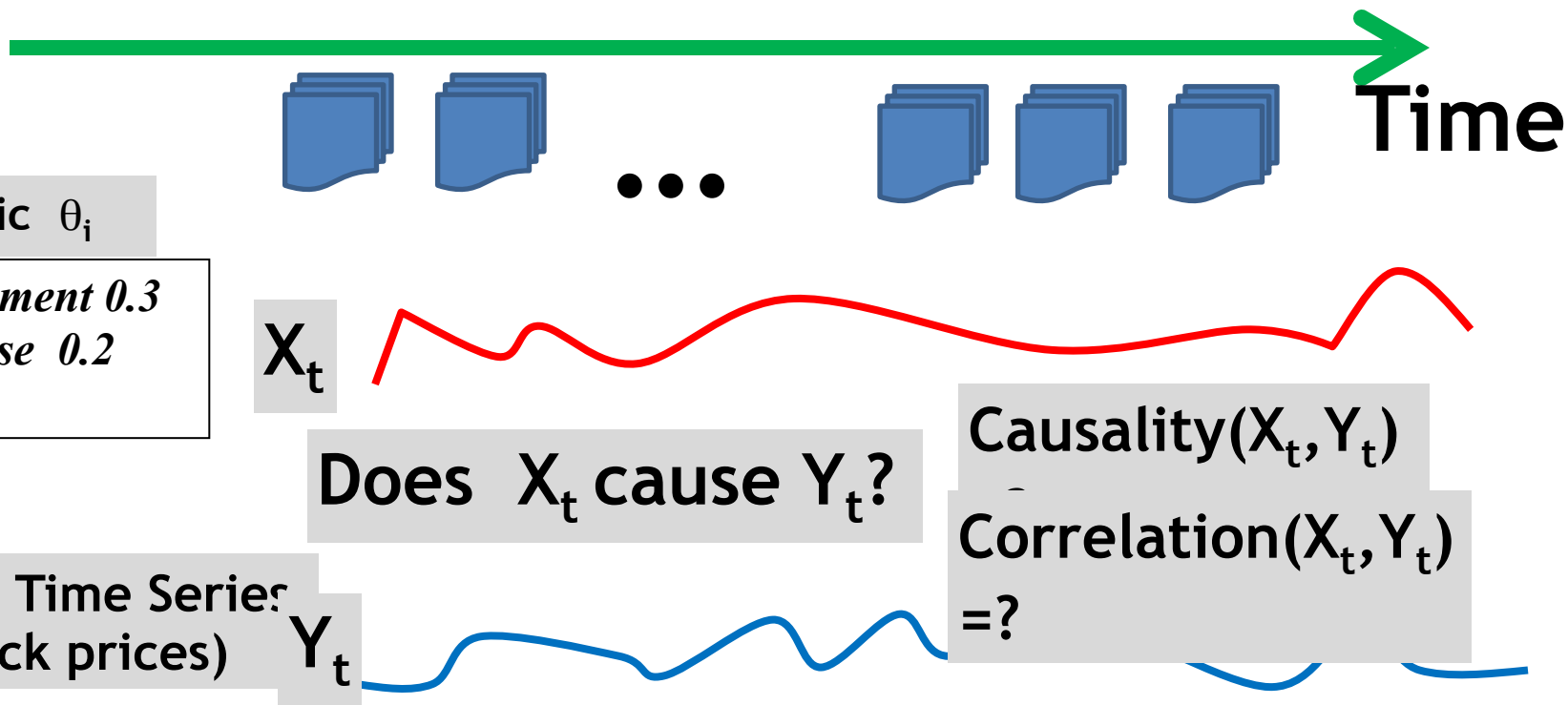
Iterative Causal Topic Modeling [Kim et al. 13]



Heuristic Optimization of Causality + Coherence



Measuring Causality (Correlation)



Granger Causality Test is often useful [Seth 07]

Topics in NY Times Correlated with Stocks

[Kim et al. 13]: June 2000 ~ Dec. 2011

AAMRQ (American Airlines)	AAPL (Apple)
<p>russia russian putin europe european germany bush gore presidential police court judge <u>airlines airport air</u> <u>united trade terrorism</u> food foods cheese nets scott basketball tennis williams open awards gay boy moss minnesota chechnya</p>	<p>paid notice st russia russian europe olympic games olympics she her ms oil ford prices black fashion blacks <u>computer technology software</u> <u>internet com web</u> football giants jets japan japanese plane</p>

Topics are biased toward each time series

Major Topics in 2000 Presidential Election

[Kim et al. 13]

Top Three Words in Significant Topics from NY Times

tax cut 1

screen pataki guiliani

enthusiasm door symbolic

oil energy prices

news w top

pres al vice

love tucker presented

partial abortion privatization

court supreme abortion

gun control nra

Text: NY Times (May 2000 - Oct.
2000)
Series: Iowa Electronic
Market

<http://tippie.uiowa.edu/iem/>

Issues known to be
important in the
2000 presidential election

What You Should Know

- What is a mixture language model, particularly mixture of unigram language model?
- What is the general form of the likelihood function of a mixture model?
- Why can the two-component mixture model with a background component language model factor out common words?
- What is PLSA and how does it work?
- How does EM algorithm work for PLSA?
- Why do we want to add a prior to PLSA and how can we modify the EM algorithm to incorporate a conjugate prior?
- Why is contextual topic modeling interesting? Can you give multiple examples of applications of contextual topic modeling?