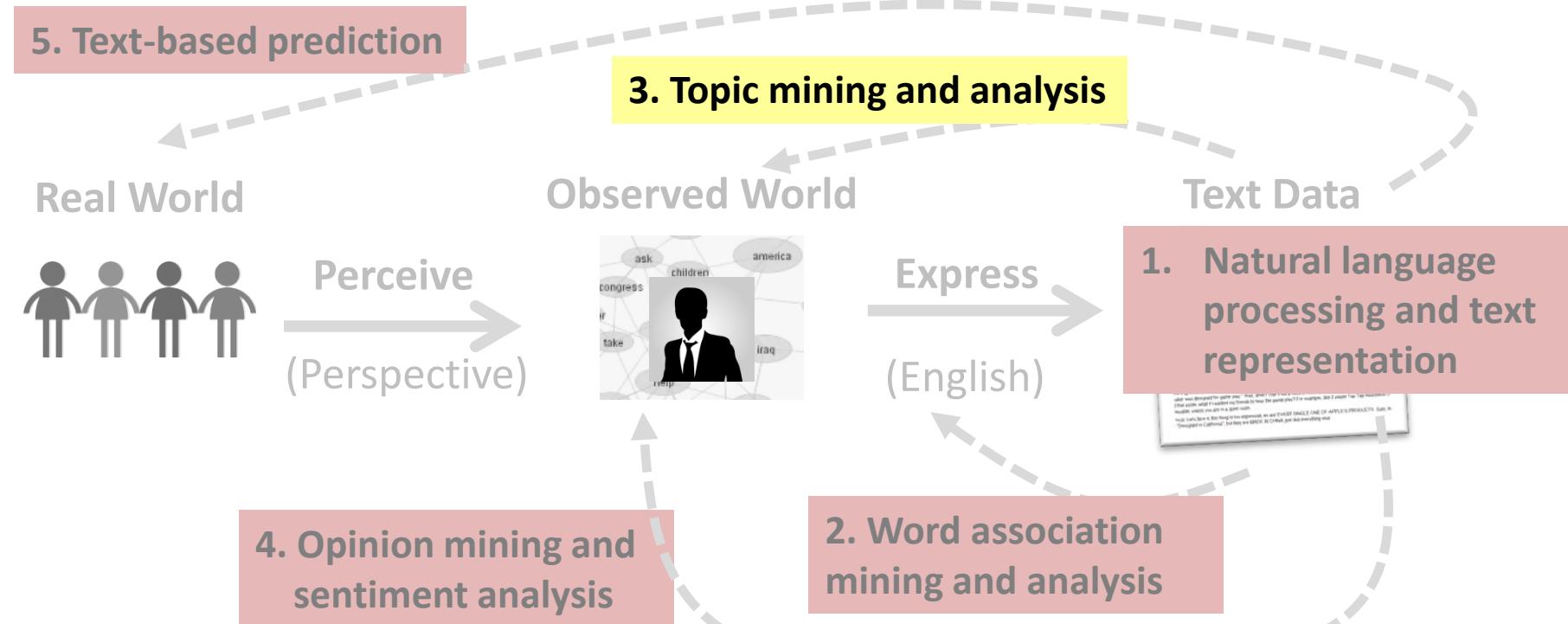


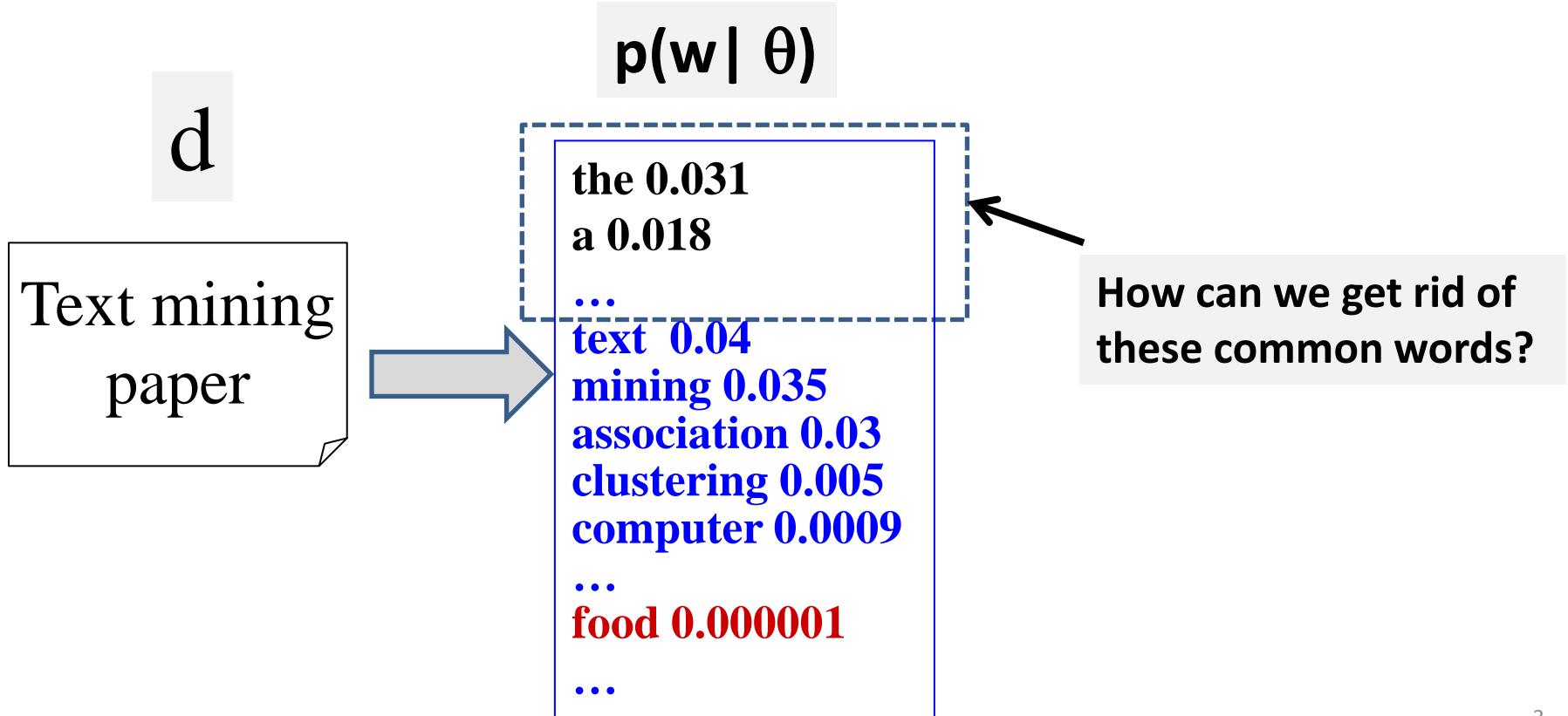
# Probabilistic Topic Models: Mixture of Unigram Language Models

ChengXiang “Cheng” Zhai  
Department of Computer Science  
University of Illinois at Urbana-Champaign

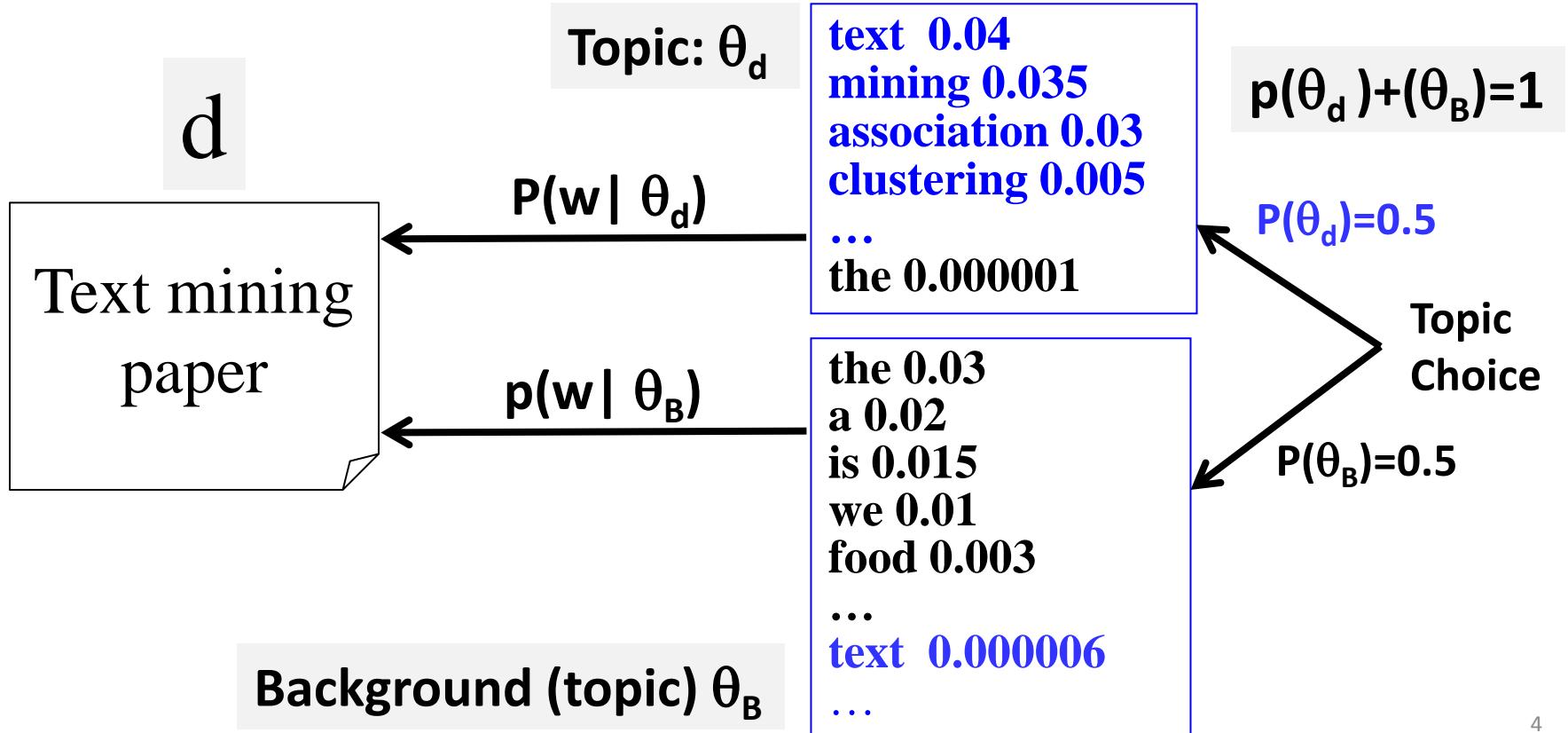
# Probabilistic Topic Models: Mixture of Unigram LMs



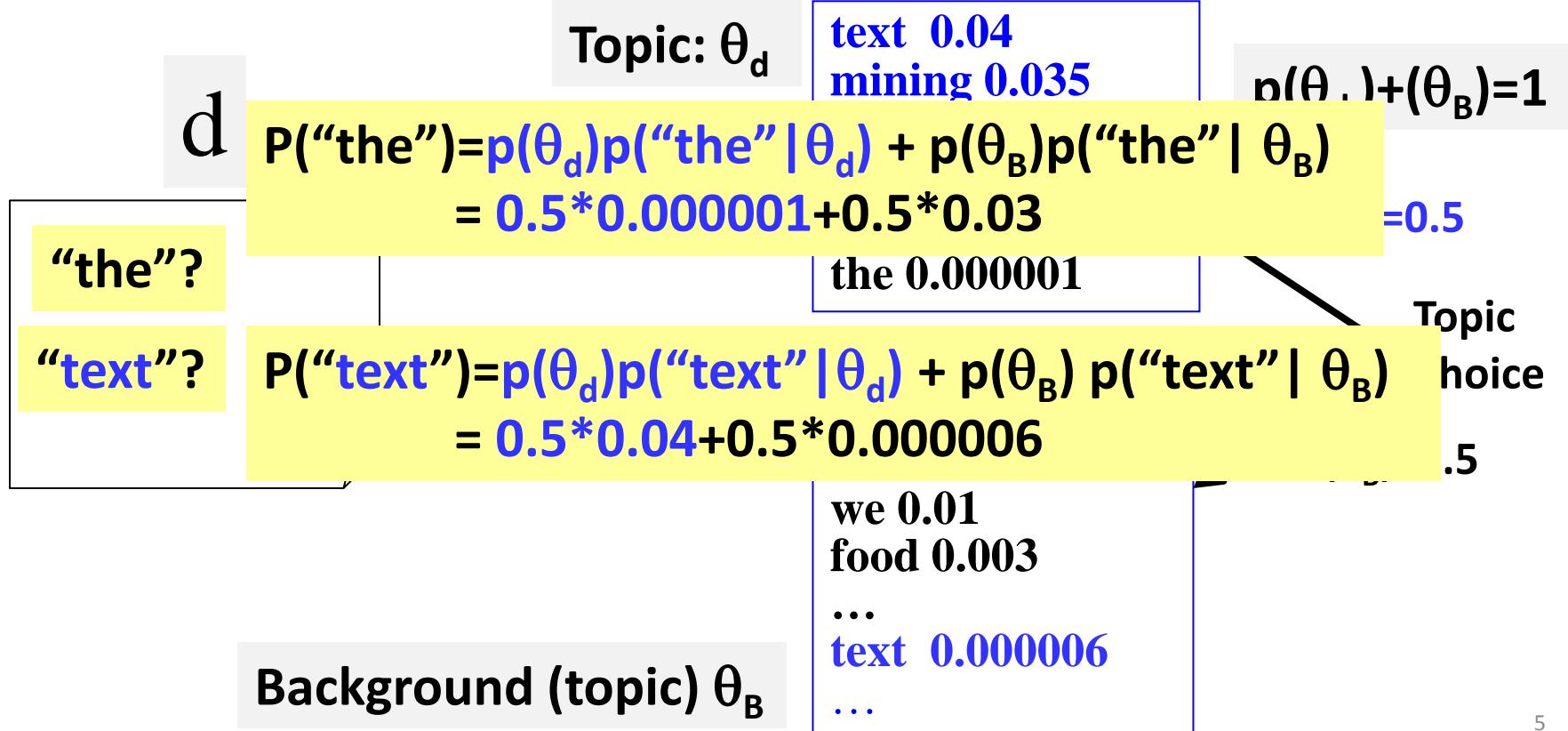
# Factoring out Background Words



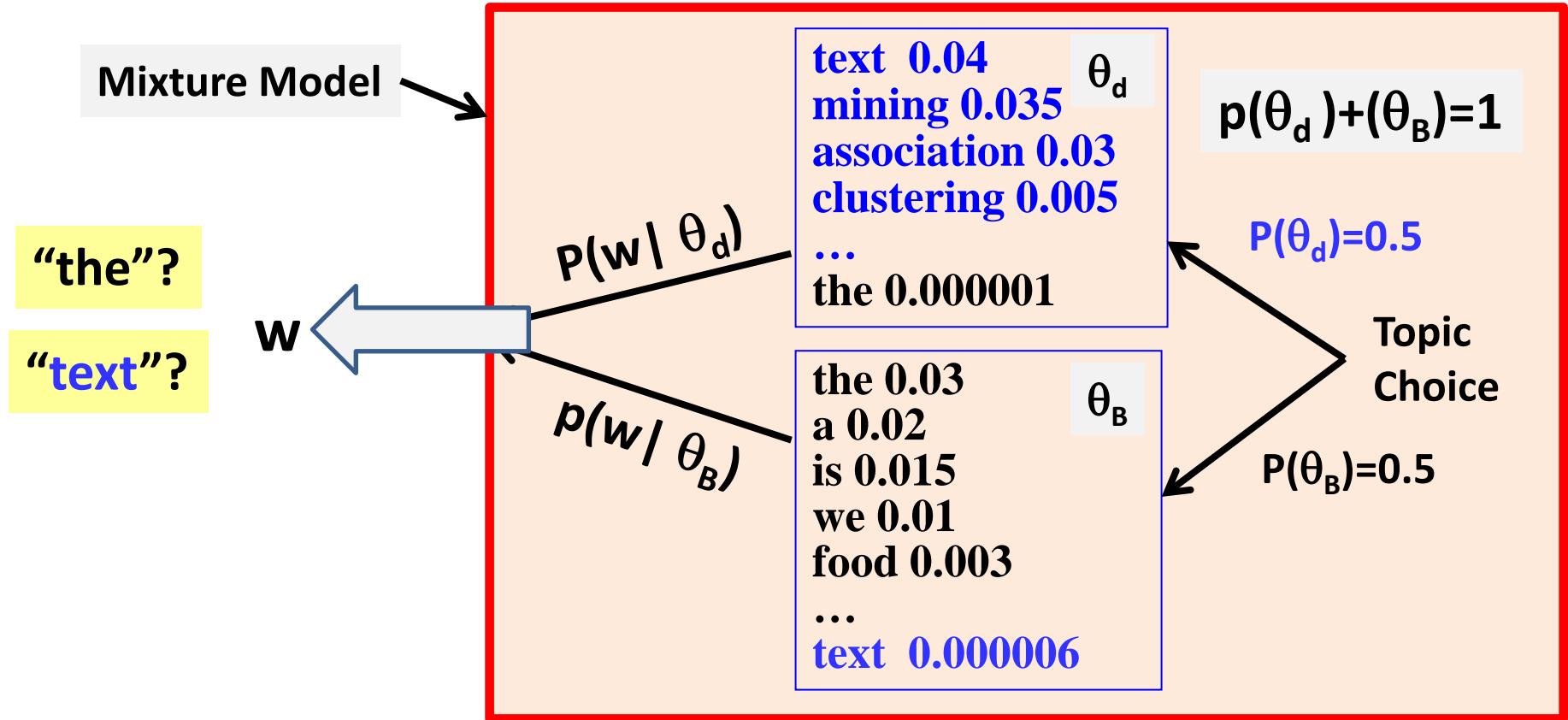
# Generate d Using Two Word Distributions



# What's the probability of observing a word w?



# The Idea of a Mixture Model



# As a Generative Model...

text 0.04  $\theta_d$   
mining 0.035  
association 0.03  
clustering 0.005

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

Formally defines the following generative model:

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

w

Estimate of the model “discovers”  
two topics + topic coverage

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

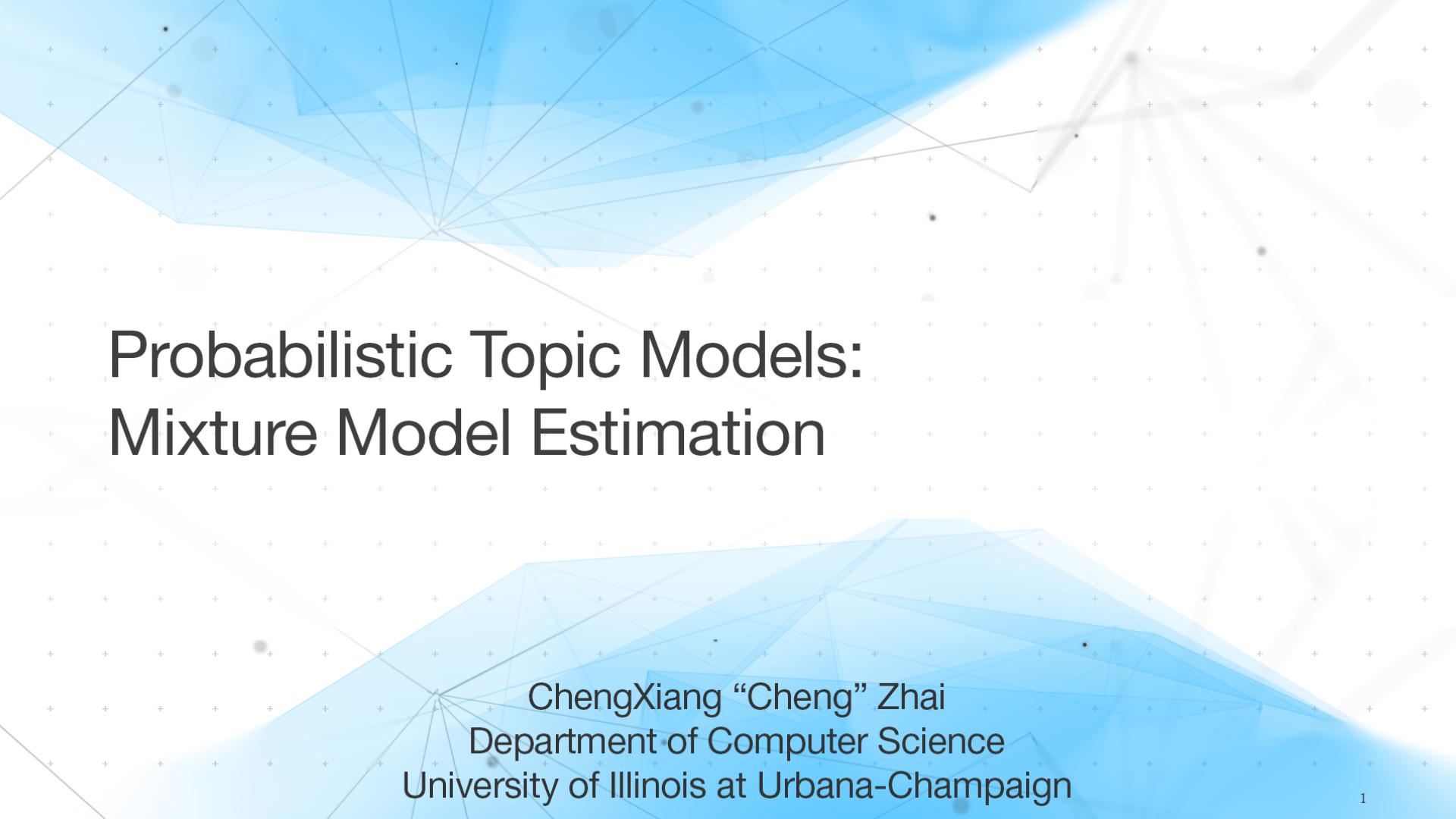
# Mixture of Two Unigram Language Models

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

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

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

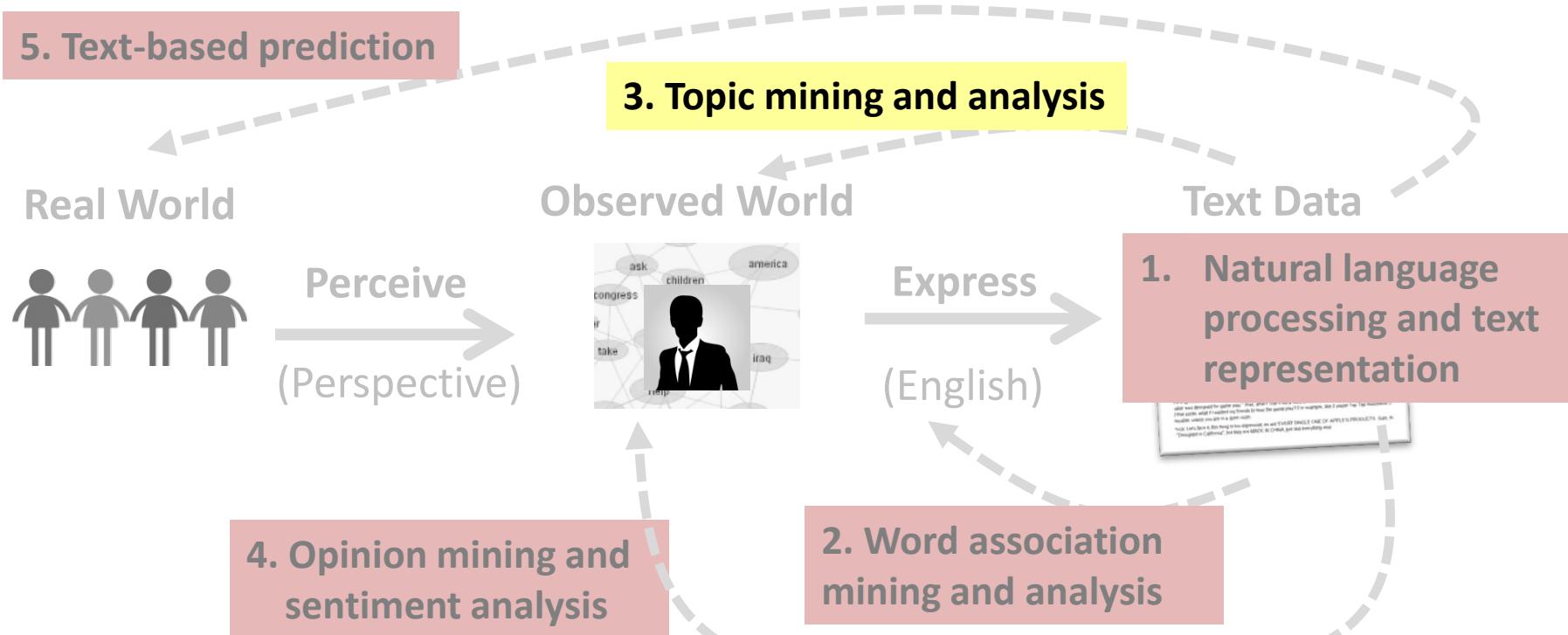
**Subject to**  $\sum_{i=1}^M p(w_i | \theta_d) = \sum_{i=1}^M p(w_i | \theta_B) = 1 \quad p(\theta_d) + p(\theta_B) = 1$



# Probabilistic Topic Models: Mixture Model Estimation

ChengXiang “Cheng” Zhai  
Department of Computer Science  
University of Illinois at Urbana-Champaign

# Probabilistic Topic Models: Mixture Model Estimation



# Back to Factoring out Background Words

Text Mining Paper

d



$$P(w | \theta_d)$$

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

$\theta_d$

$$P(w | \theta_B)$$

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

$\theta_B$

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

$$P(\theta_d) = 0.5$$

Topic  
Choice

$$P(\theta_B) = 0.5$$

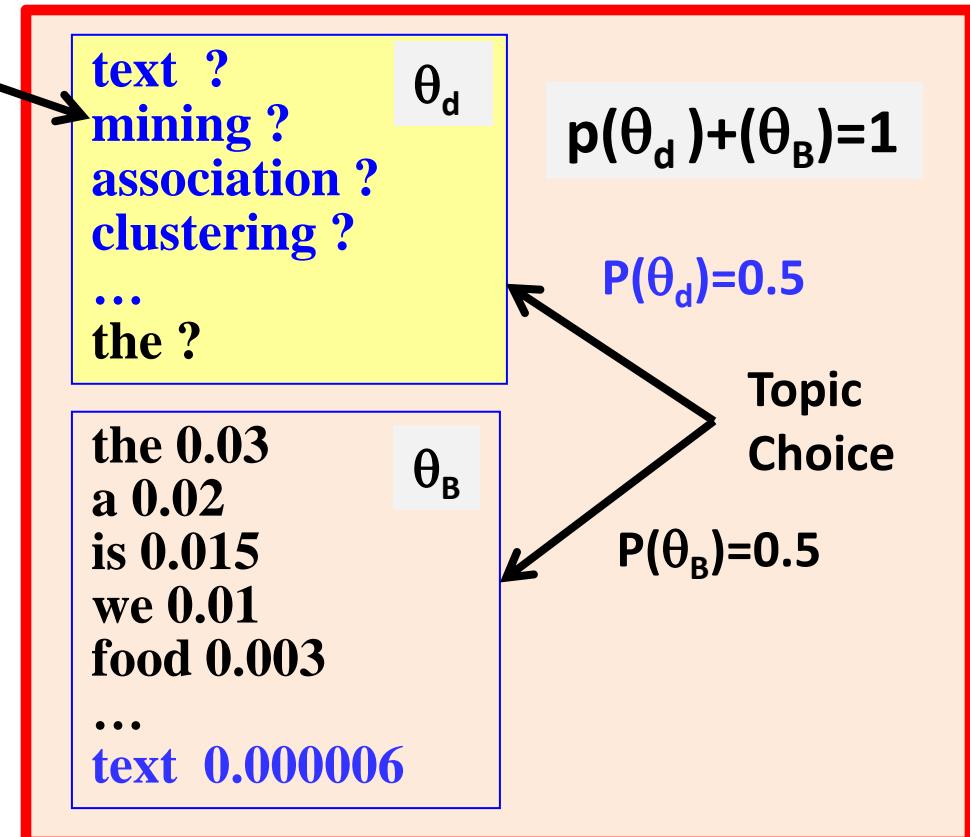
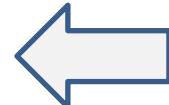
# Estimation of One Topic: $P(w | \theta_d)$

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

Would the ML estimate demote  
background words in  $\theta_d$  ?

$d$

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



# Behavior of a Mixture Model

$d = \boxed{\text{text the}}$

Likelihood:

$$\begin{aligned} P(\text{"text"}) &= p(\theta_d)p(\text{"text"}|\theta_d) + p(\theta_B)p(\text{"text"}|\theta_B) \\ &= 0.5 * p(\text{"text"}|\theta_d) + 0.5 * 0.1 \end{aligned}$$

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

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

$\boxed{\text{text ?} \quad \theta_d}$   
 $\text{the ?}$

$P(\theta_d)=0.5$

$\boxed{\text{the } 0.9 \quad \theta_B}$   
 $\text{text } 0.1$

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 $\theta_d$ and $\theta_B$

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

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

If  $x + y = \text{constant}$ , then  $xy$  reaches maximum when  $x = y$ .

$$0.5 * p(\text{"text"}|\theta_d) + 0.5 * 0.1 = 0.5 * p(\text{"the"}|\theta_d) + 0.5 * 0.9$$

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

$d =$  text the

text ?  $\theta_d$

$P(\theta_d) = 0.5$

$P(\theta_B) = 0.5$

the 0.9 text 0.1  $\theta_B$

**Behavior 1:** if  $p(w_1|\theta_B) > p(w_2|\theta_B)$ , then  $p(w_1|\theta_d) < p(w_2|\theta_d)$

# Response to Data Frequency

$d = \boxed{\text{text the}}$

$$p(d|\Lambda) = [0.5*p(\text{"text"}|\theta_d) + 0.5*0.1] \\ \times [0.5*p(\text{"the"}|\theta_d) + 0.5*0.9]$$

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

$d' = \boxed{\text{text the}} \\ \text{the the} \\ \text{the ...the}}$

$$p(d'|\Lambda) = [0.5*p(\text{"text"}|\theta_d) + 0.5*0.1] \\ \times [0.5*p(\text{"the"}|\theta_d) + 0.5*0.9] \\ \times [0.5*p(\text{"the"}|\theta_d) + 0.5*0.9] \\ \times [0.5*p(\text{"the"}|\theta_d) + 0.5*0.9]$$

...

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

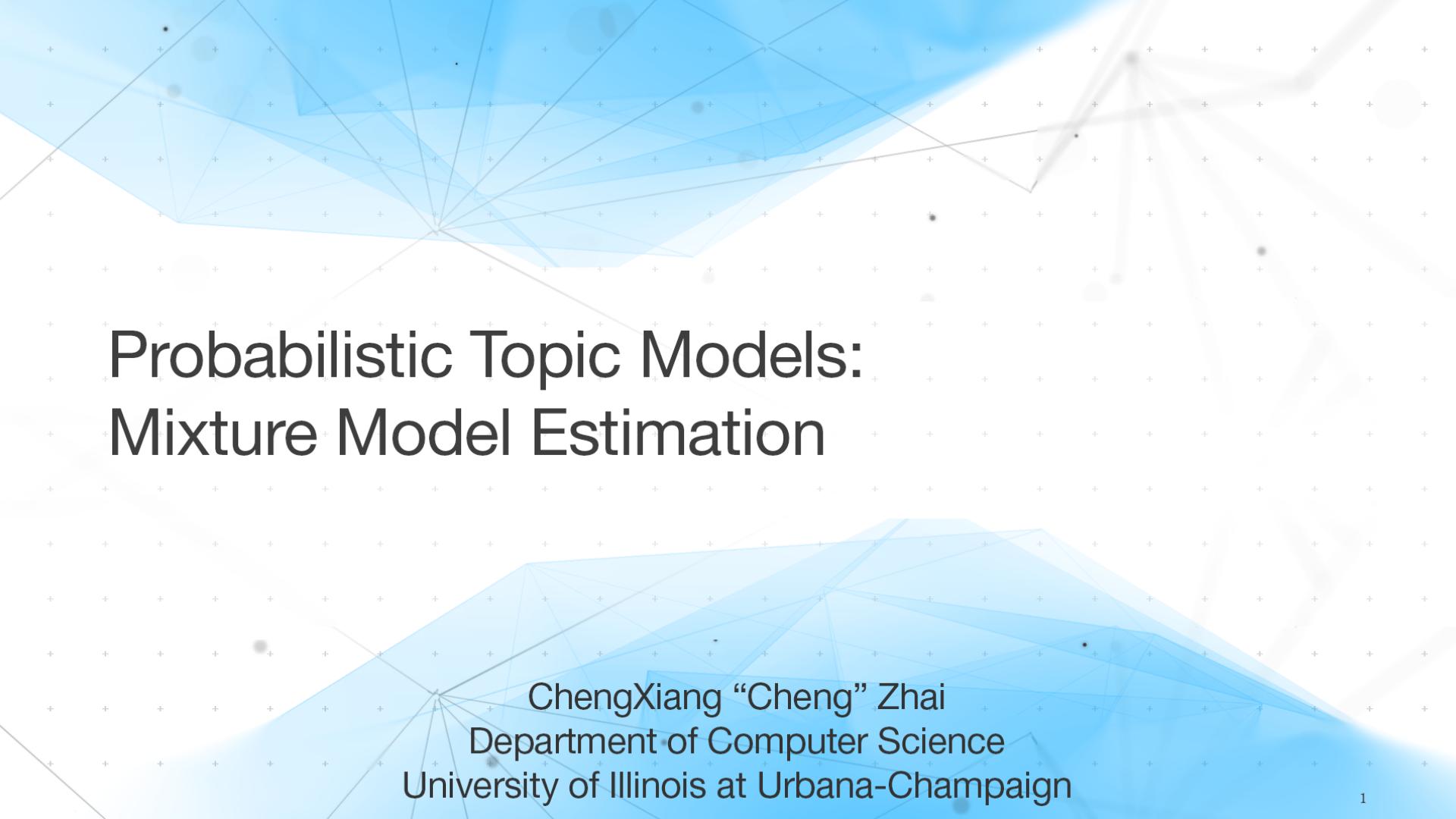
$$\times [0.5*p(\text{"the"}|\theta_d) + 0.5*0.9]$$

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

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

# Summary

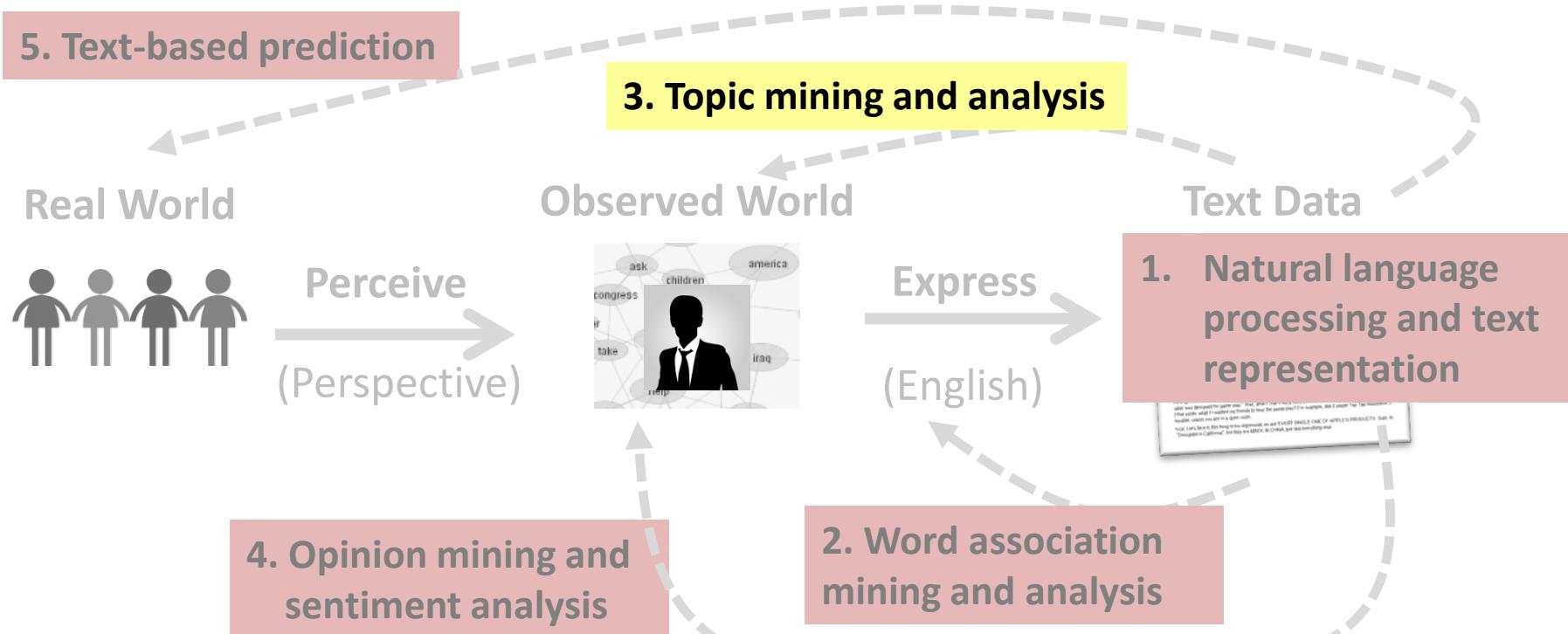
- General behavior of a mixture model:
  - Every component model attempts to assign high probabilities to highly frequent words in the data (to “collaboratively maximize likelihood”)
  - Different component models tend to “bet” high probabilities on different words (to avoid “competition” or “waste of probability”)
  - The probability of choosing each component “regulates” the collaboration/competition between the component models
- Fixing one component to a background word distribution (i.e., background language model):
  - Helps “get rid of background words” in other component
  - Is an example of imposing a prior on the model parameters (prior = one model must be exactly the same as the background LM)



# Probabilistic Topic Models: Mixture Model Estimation

ChengXiang “Cheng” Zhai  
Department of Computer Science  
University of Illinois at Urbana-Champaign

# Probabilistic Topic Models: Mixture Model Estimation



# Back to Factoring out Background Words

Text Mining Paper

d



$$P(w | \theta_d)$$

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

$\theta_d$

$$P(w | \theta_B)$$

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

$\theta_B$

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

$$P(\theta_d) = 0.5$$

Topic  
Choice

$$P(\theta_B) = 0.5$$

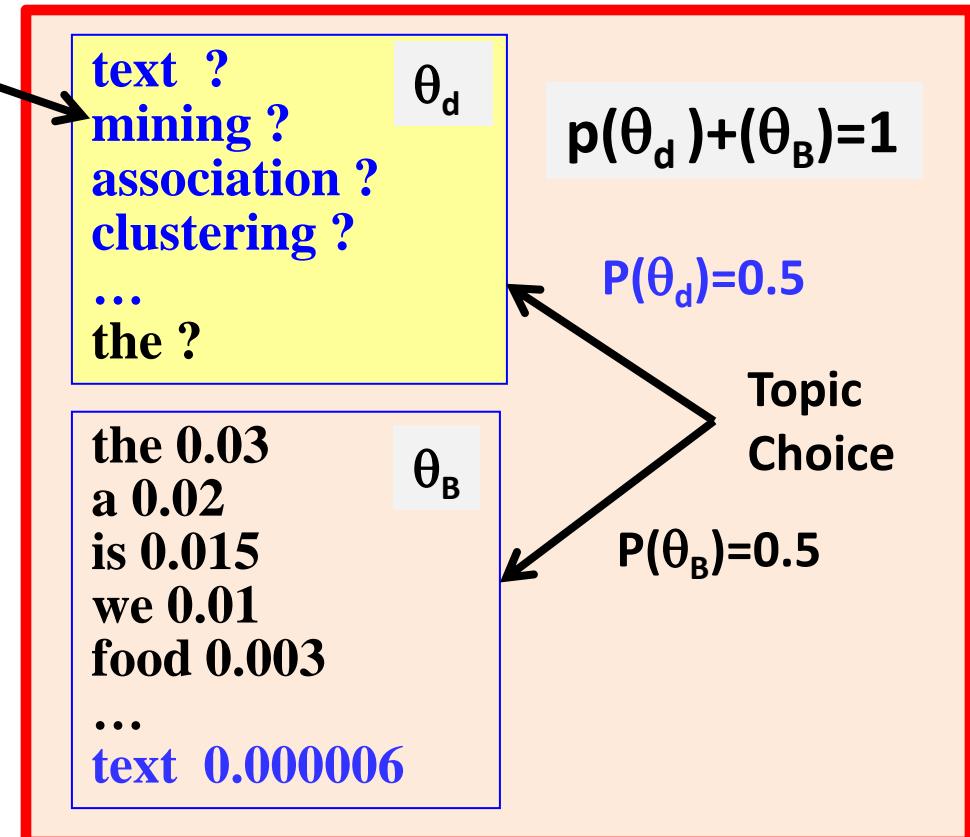
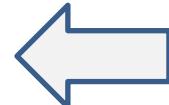
# Estimation of One Topic: $P(w | \theta_d)$

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

Would the ML estimate demote  
background words in  $\theta_d$  ?

$d$

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



# Behavior of a Mixture Model

$d = \boxed{\text{text the}}$

Likelihood:

$$\begin{aligned} P(\text{"text"}) &= p(\theta_d)p(\text{"text"}|\theta_d) + p(\theta_B)p(\text{"text"}|\theta_B) \\ &= 0.5 * p(\text{"text"}|\theta_d) + 0.5 * 0.1 \end{aligned}$$

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

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

$\boxed{\text{text ?} \quad \theta_d}$   
 $\text{the ?}$

$P(\theta_d)=0.5$

$\boxed{\text{the } 0.9 \quad \theta_B}$   
 $\text{text } 0.1$

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 $\theta_d$ and $\theta_B$

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

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

If  $x + y = \text{constant}$ , then  $xy$  reaches maximum when  $x = y$ .

$$0.5 * p(\text{"text"}|\theta_d) + 0.5 * 0.1 = 0.5 * p(\text{"the"}|\theta_d) + 0.5 * 0.9$$

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

$d =$  text the

text ?  $\theta_d$

$P(\theta_d) = 0.5$

$P(\theta_B) = 0.5$

the 0.9 text 0.1  $\theta_B$

**Behavior 1:** if  $p(w_1|\theta_B) > p(w_2|\theta_B)$ , then  $p(w_1|\theta_d) < p(w_2|\theta_d)$

# Response to Data Frequency

$d = \boxed{\text{text the}}$

$$p(d|\Lambda) = [0.5*p(\text{"text"}|\theta_d) + 0.5*0.1] \\ \times [0.5*p(\text{"the"}|\theta_d) + 0.5*0.9]$$

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

$d' = \boxed{\text{text the}} \\ \text{the the} \\ \text{the ...the}}$

$$p(d'|\Lambda) = [0.5*p(\text{"text"}|\theta_d) + 0.5*0.1] \\ \times [0.5*p(\text{"the"}|\theta_d) + 0.5*0.9] \\ \times [0.5*p(\text{"the"}|\theta_d) + 0.5*0.9] \\ \times [0.5*p(\text{"the"}|\theta_d) + 0.5*0.9]$$

...

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

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

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

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

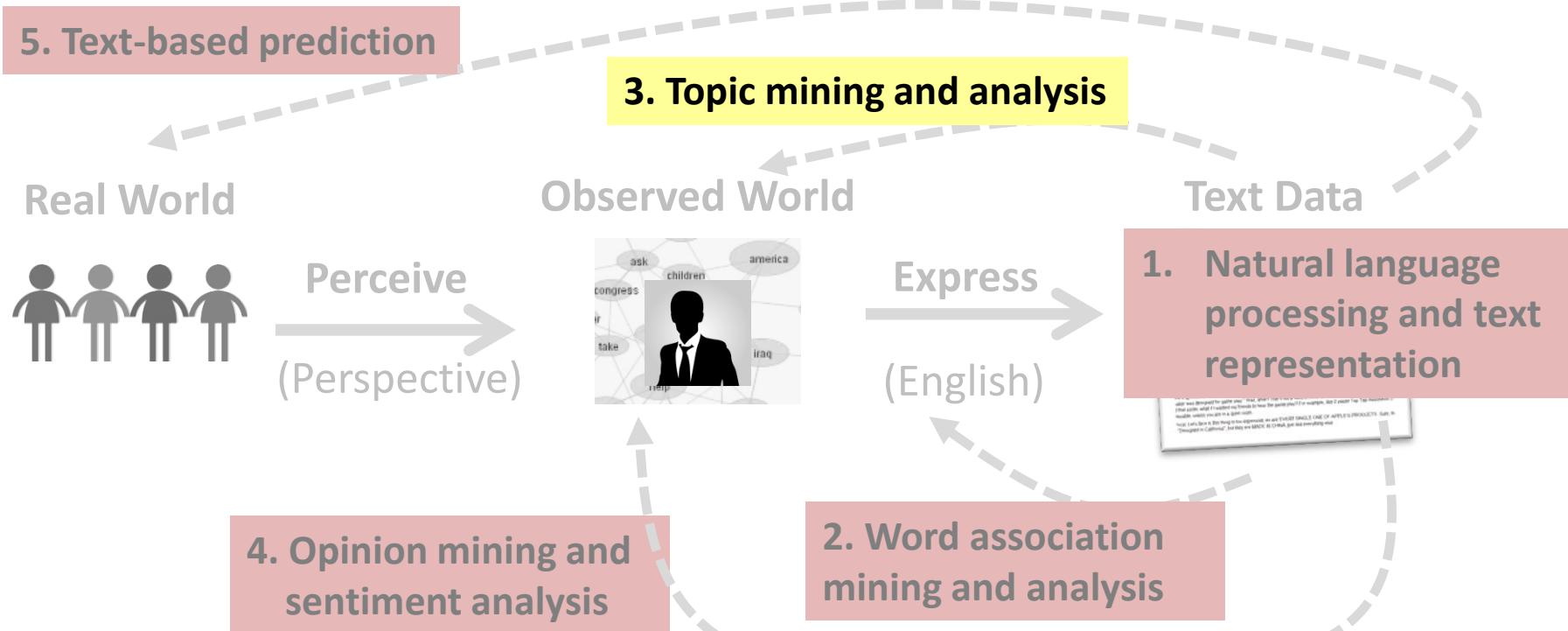
# Summary

- General behavior of a mixture model:
  - Every component model attempts to assign high probabilities to highly frequent words in the data (to “collaboratively maximize likelihood”)
  - Different component models tend to “bet” high probabilities on different words (to avoid “competition” or “waste of probability”)
  - The probability of choosing each component “regulates” the collaboration/competition between the component models
- Fixing one component to a background word distribution (i.e., background language model):
  - Helps “get rid of background words” in other component
  - Is an example of imposing a prior on the model parameters (prior = one model must be exactly the same as the background LM)

# Probabilistic Topic Models: Expectation-Maximization Algorithm

ChengXiang “Cheng” Zhai  
Department of Computer Science  
University of Illinois at Urbana-Champaign

# Probabilistic Topic Models: Expectation-Maximization (EM) Algorithm

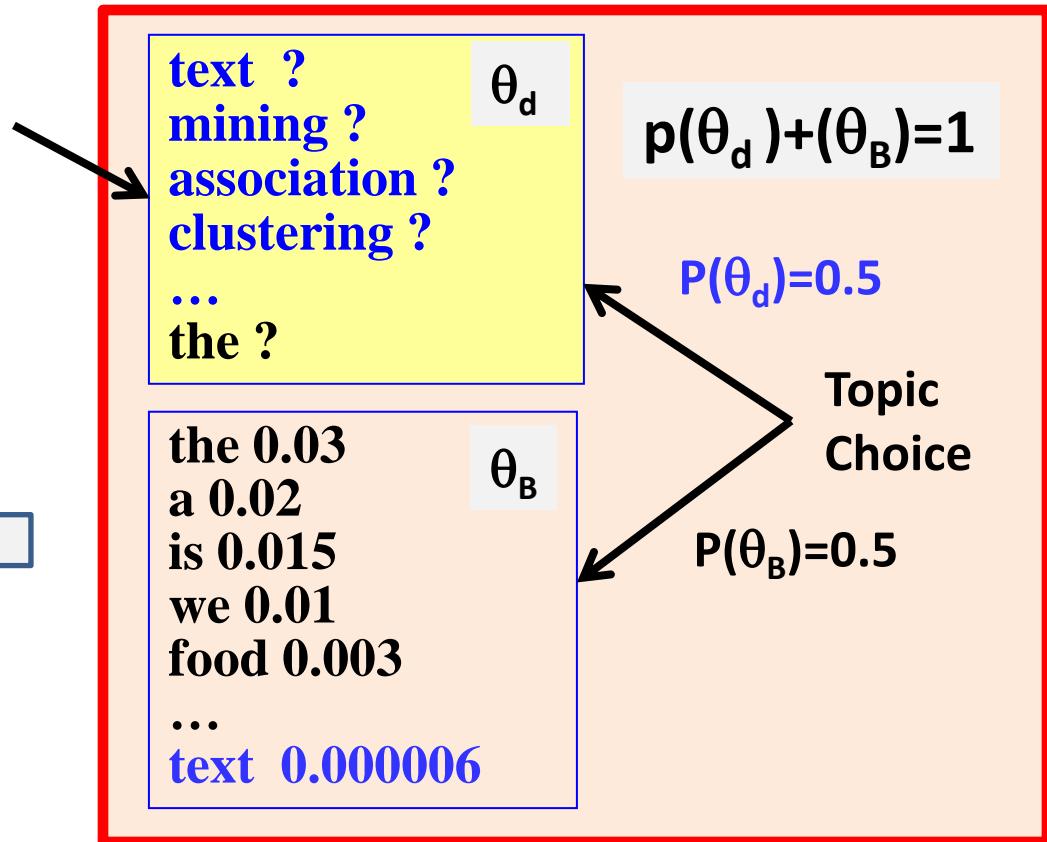
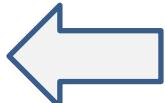


# Estimation of One Topic: $P(w | \theta_d)$

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

d

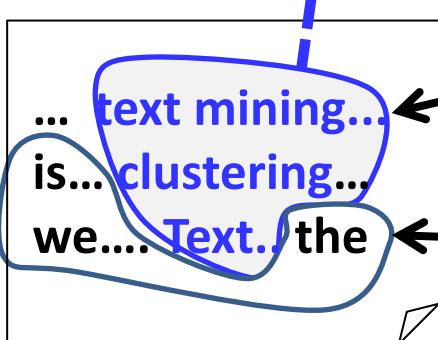
... text mining...  
is... clustering...  
we.... Text.. the



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

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

d  
d'



$$P(w | \theta_d)$$

$$p(w | \theta_B)$$

text ?  
mining ?  
association ?  
clustering ?  
...  
the ?

$\theta_d$

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

$$P(\theta_d) = 0.5$$

Topic  
Choice

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

$\theta_B$

$$P(\theta_B) = 0.5$$

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

Is “text” more likely from  $\theta_d$  or  $\theta_B$  ?

From  $\theta_d$  ( $z=0$ )?

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

From  $\theta_B$  ( $z=1$ )?

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

$$p(z = 0 | w = \text{"text"}) =$$

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

$$P(w | \theta_d)$$

$$p(w | \theta_B)$$

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

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

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

$$P(\theta_d) = 0.5$$

Topic  
Choice

$$P(\theta_B) = 0.5$$

# 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  $\theta_d$

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

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

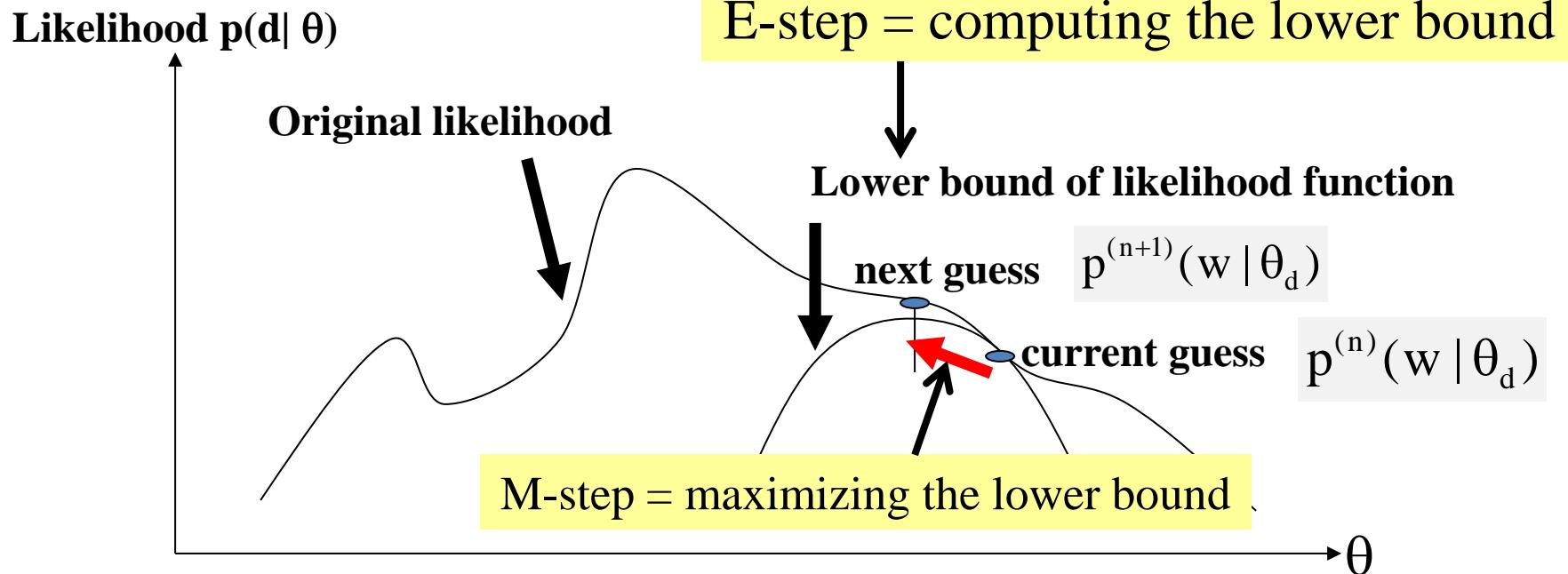
Word	#	$p(w   \theta_B)$	Iteration 1		Iteration 2		Iteration 3	
			$P(w   \theta)$	$p(z=0   w)$	$P(w   \theta)$	$P(z=0   w)$	$P(w   \theta)$	$P(z=0   w)$
The	4	0.5	<b>0.25</b>	0.33	<b>0.20</b>	0.29	<b>0.18</b>	0.26
Paper	2	0.3	<b>0.25</b>	0.45	<b>0.14</b>	0.32	<b>0.10</b>	0.25
Text	4	0.1	<b>0.25</b>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



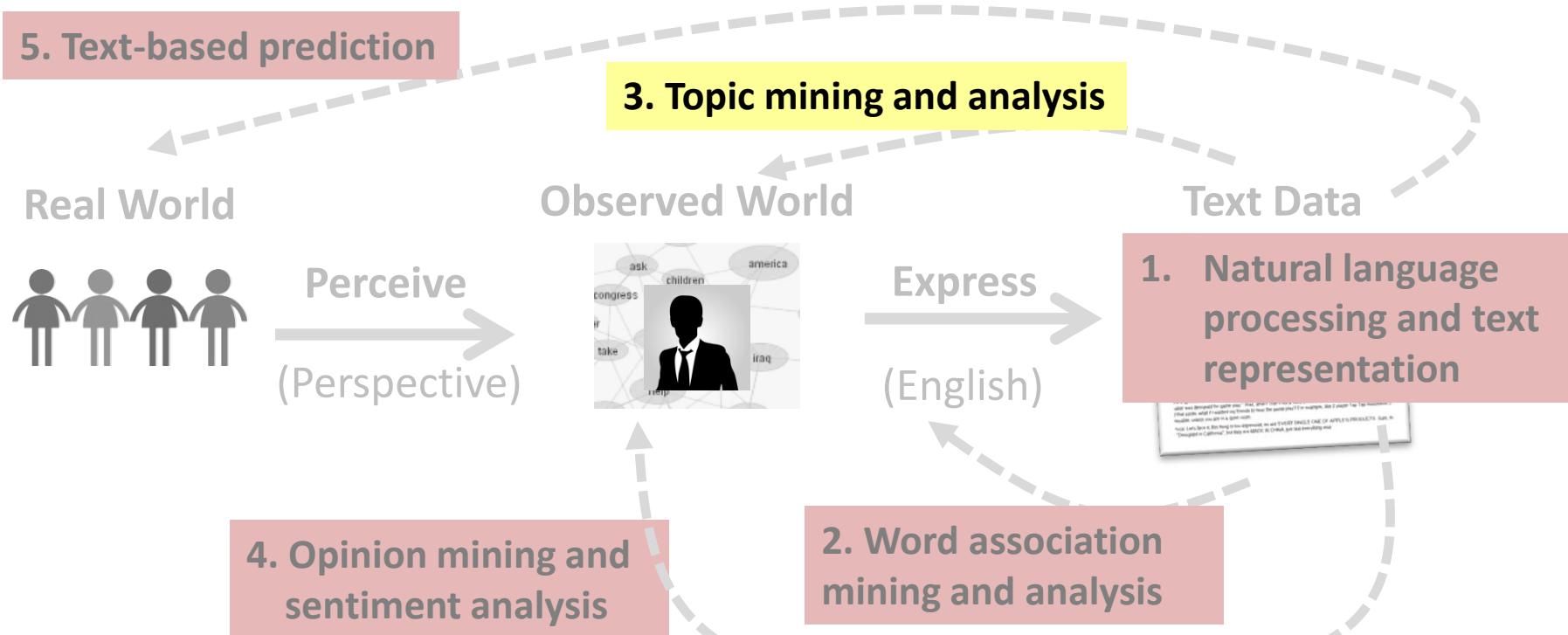
# Summary

- Expectation-Maximization (EM) algorithm
  - General algorithm for computing ML estimate of mixture models
  - Hill-climbing, so can only converge to a local maximum (depending on initial points)
- E-step: “augment” data by predicting values of useful hidden variables
- M-step: exploit the “augmented data” to improve estimate of parameters (“improve” is guaranteed in terms of likelihood)
- “Data augmentation” is probabilistic → Split counts of events probabilistically

# Probabilistic Topic Models: Expectation-Maximization Algorithm

ChengXiang “Cheng” Zhai  
Department of Computer Science  
University of Illinois at Urbana-Champaign

# Probabilistic Topic Models: Expectation-Maximization (EM) Algorithm

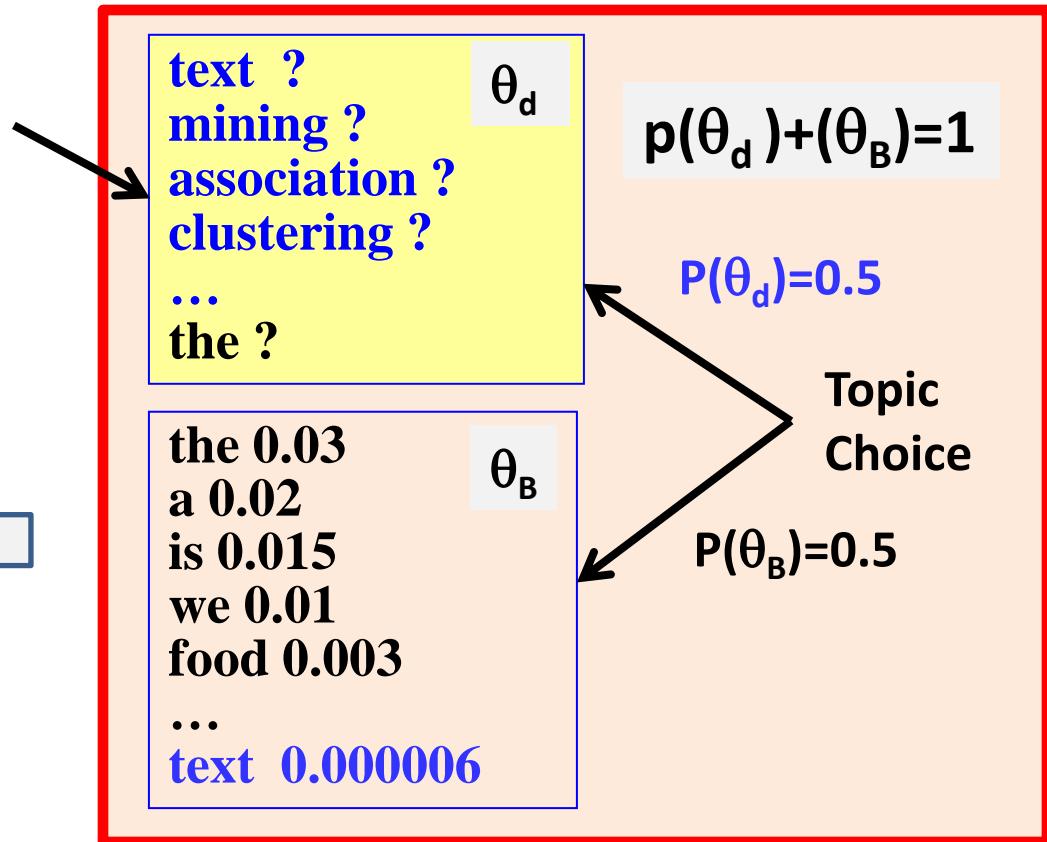
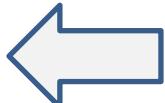


# Estimation of One Topic: $P(w | \theta_d)$

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

d

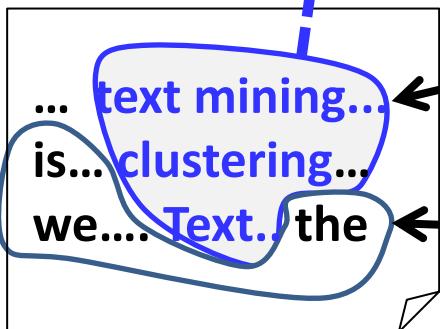
... text mining...  
is... clustering...  
we.... Text.. the



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

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

d  
d'



$$P(w | \theta_d)$$

$$p(w | \theta_B)$$

text ?  
mining ?  
association ?  
clustering ?  
...  
the ?

$\theta_d$

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

$$P(\theta_d) = 0.5$$

Topic  
Choice

$$P(\theta_B) = 0.5$$

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

$\theta_B$

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

Is “text” more likely from  $\theta_d$  or  $\theta_B$  ?

From  $\theta_d$  ( $z=0$ )?

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

From  $\theta_B$  ( $z=1$ )?

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

$$p(z = 0 | w = \text{"text"}) =$$

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

$$P(w | \theta_d)$$

$$p(w | \theta_B)$$

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

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

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

$$P(\theta_d) = 0.5$$

Topic  
Choice

$$P(\theta_B) = 0.5$$

# 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  $\theta_d$

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

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

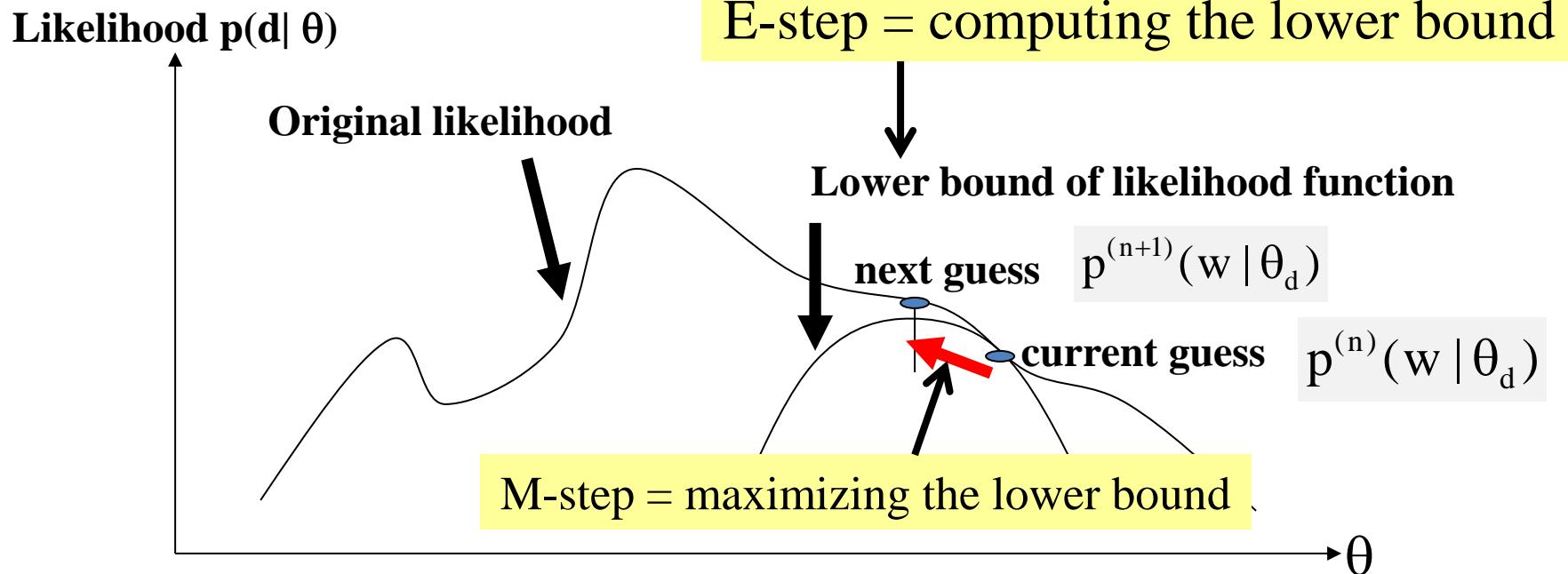
Word	#	$p(w   \theta_B)$	Iteration 1		Iteration 2		Iteration 3	
			$P(w   \theta)$	$p(z=0   w)$	$P(w   \theta)$	$P(z=0   w)$	$P(w   \theta)$	$P(z=0   w)$
The	4	0.5	<b>0.25</b>	0.33	<b>0.20</b>	0.29	<b>0.18</b>	0.26
Paper	2	0.3	<b>0.25</b>	0.45	<b>0.14</b>	0.32	<b>0.10</b>	0.25
Text	4	0.1	<b>0.25</b>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



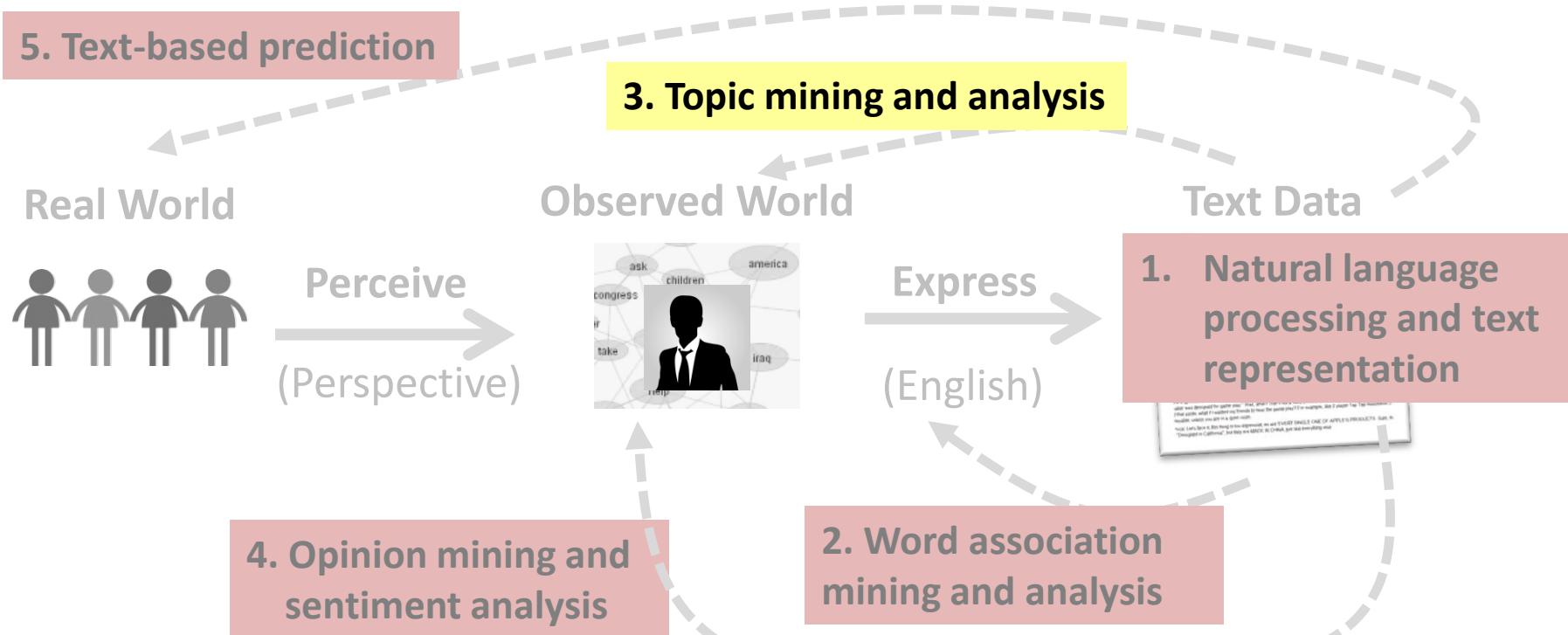
# Summary

- Expectation-Maximization (EM) algorithm
  - General algorithm for computing ML estimate of mixture models
  - Hill-climbing, so can only converge to a local maximum (depending on initial points)
- E-step: “augment” data by predicting values of useful hidden variables
- M-step: exploit the “augmented data” to improve estimate of parameters (“improve” is guaranteed in terms of likelihood)
- “Data augmentation” is probabilistic → Split counts of events probabilistically

# Probabilistic Topic Models: Expectation-Maximization Algorithm

ChengXiang “Cheng” Zhai  
Department of Computer Science  
University of Illinois at Urbana-Champaign

# Probabilistic Topic Models: Expectation-Maximization (EM) Algorithm

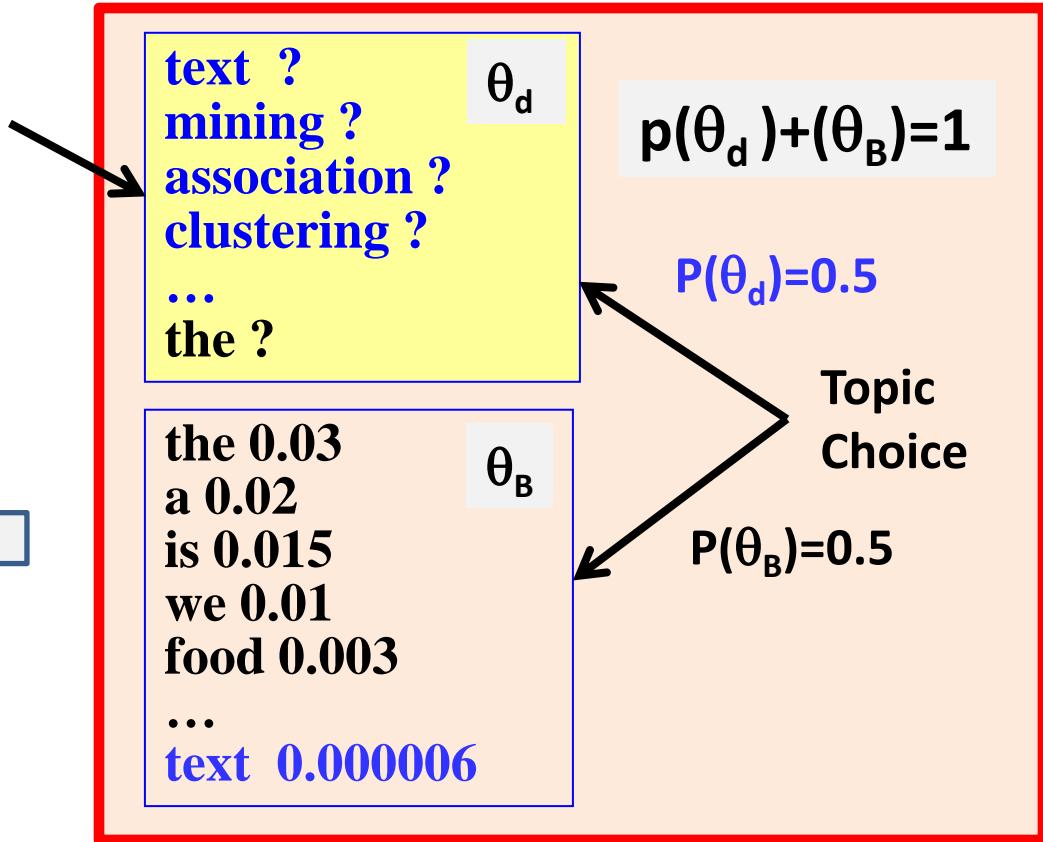
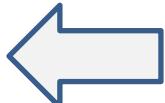


# Estimation of One Topic: $P(w | \theta_d)$

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

d

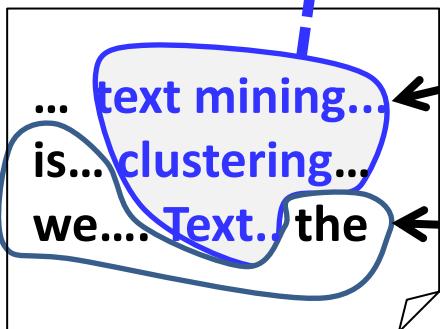
... text mining...  
is... clustering...  
we.... Text.. the



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

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

d  
d'



$$P(w | \theta_d)$$

$$p(w | \theta_B)$$

text ?  
mining ?  
association ?  
clustering ?  
...  
the ?

$\theta_d$

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

$$P(\theta_d) = 0.5$$

Topic  
Choice

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

$\theta_B$

$$P(\theta_B) = 0.5$$

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

Is “text” more likely from  $\theta_d$  or  $\theta_B$  ?

From  $\theta_d$  ( $z=0$ )?

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

From  $\theta_B$  ( $z=1$ )?

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

$$p(z = 0 | w = \text{"text"}) =$$

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

$$P(w | \theta_d)$$

$$p(w | \theta_B)$$

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

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

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

$$P(\theta_d) = 0.5$$

Topic  
Choice

$$P(\theta_B) = 0.5$$

# 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  $\theta_d$

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

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

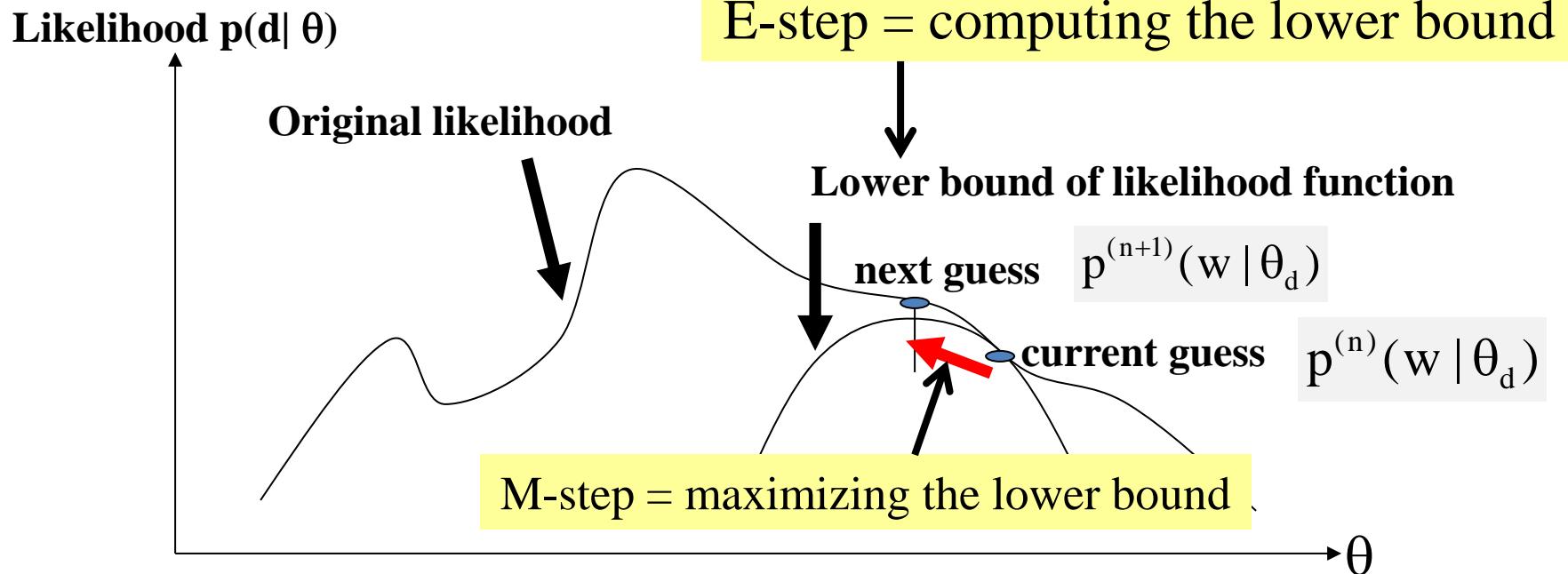
Word	#	$p(w   \theta_B)$	Iteration 1		Iteration 2		Iteration 3	
			$P(w   \theta)$	$p(z=0   w)$	$P(w   \theta)$	$P(z=0   w)$	$P(w   \theta)$	$P(z=0   w)$
The	4	0.5	<b>0.25</b>	0.33	<b>0.20</b>	0.29	<b>0.18</b>	0.26
Paper	2	0.3	<b>0.25</b>	0.45	<b>0.14</b>	0.32	<b>0.10</b>	0.25
Text	4	0.1	<b>0.25</b>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



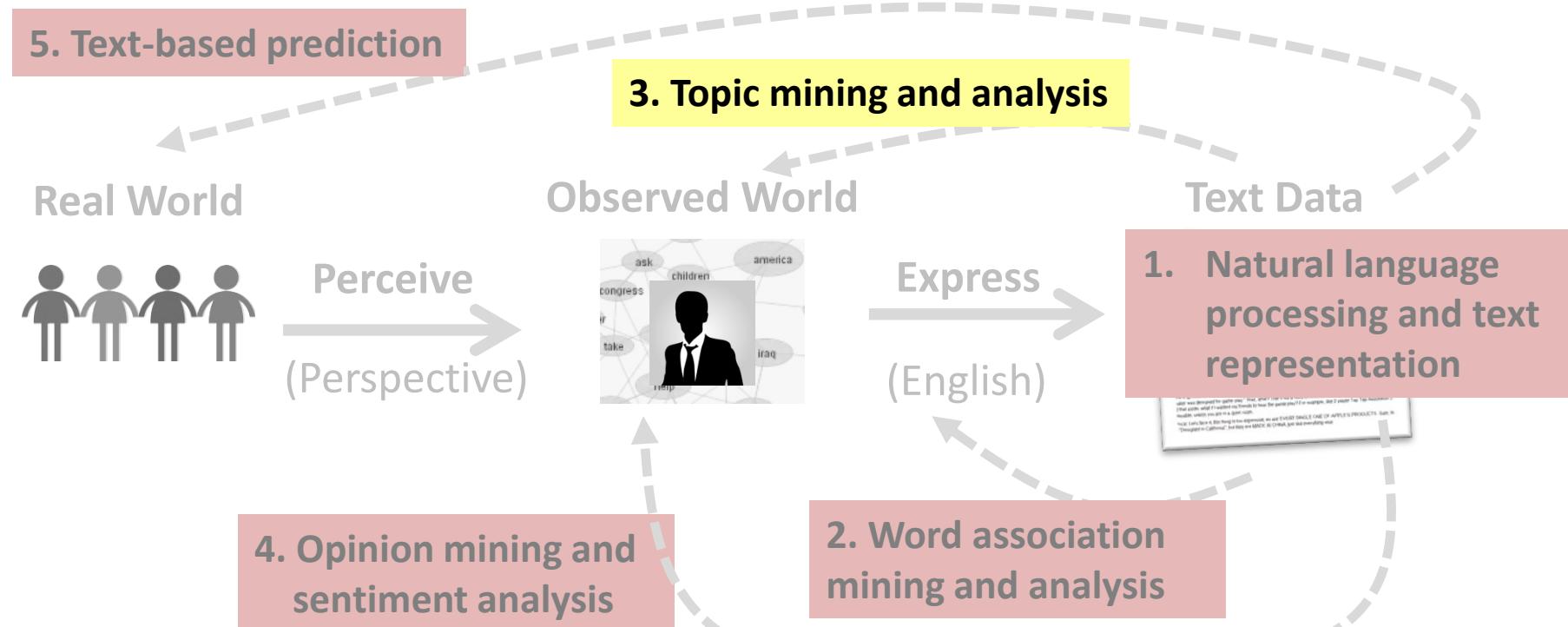
# Summary

- Expectation-Maximization (EM) algorithm
  - General algorithm for computing ML estimate of mixture models
  - Hill-climbing, so can only converge to a local maximum (depending on initial points)
- E-step: “augment” data by predicting values of useful hidden variables
- M-step: exploit the “augmented data” to improve estimate of parameters (“improve” is guaranteed in terms of likelihood)
- “Data augmentation” is probabilistic → Split counts of events probabilistically

# Probabilistic Latent Semantic Analysis (PLSA)

ChengXiang “Cheng” Zhai  
Department of Computer Science  
University of Illinois at Urbana-Champaign

# Probabilistic Latent Semantic Analysis (PLSA)



# Document as a Sample of Mixed Topics

Topic  $\theta_1$

government 0.3  
response 0.2

...

Topic  $\theta_2$

city 0.2  
new 0.1  
orleans 0.05

...

Topic  $\theta_k$

donate 0.1  
relief 0.05  
help 0.02

...

Background  $\theta_B$

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]. ...

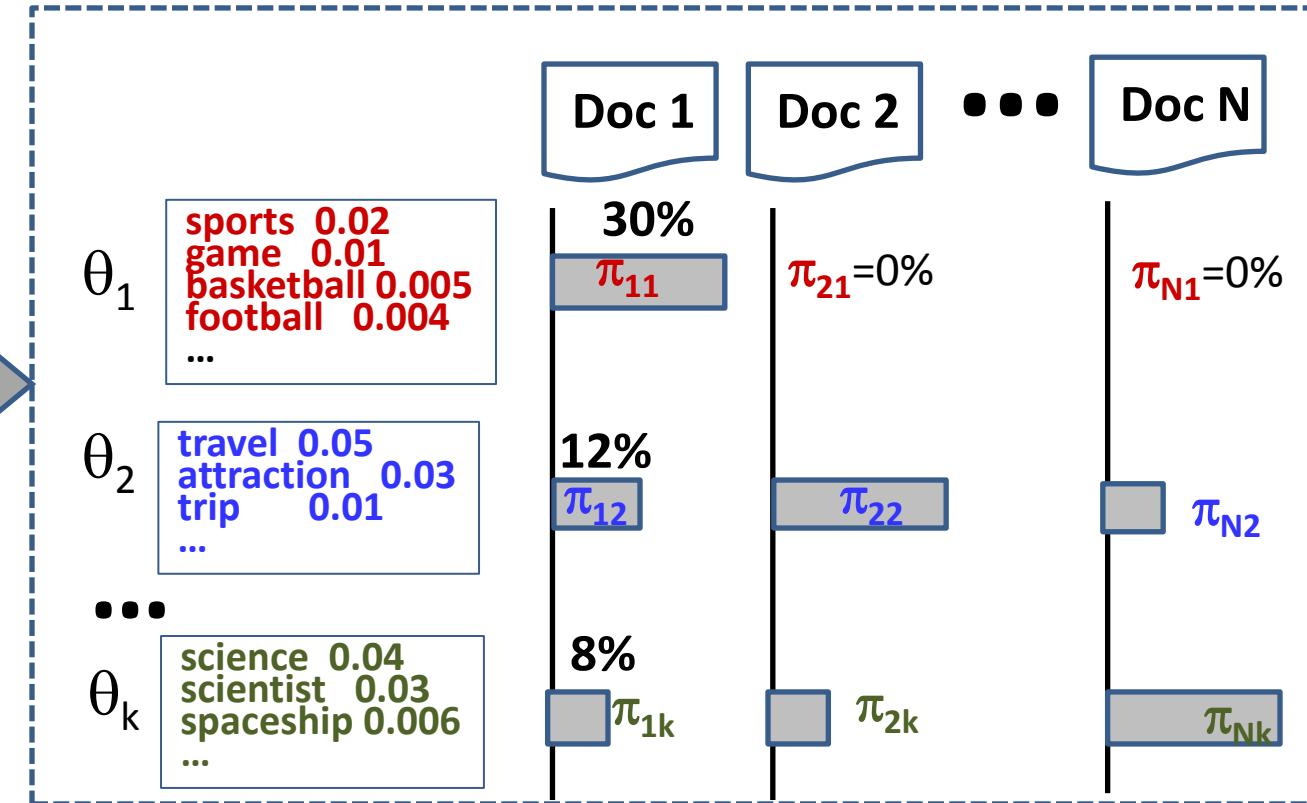
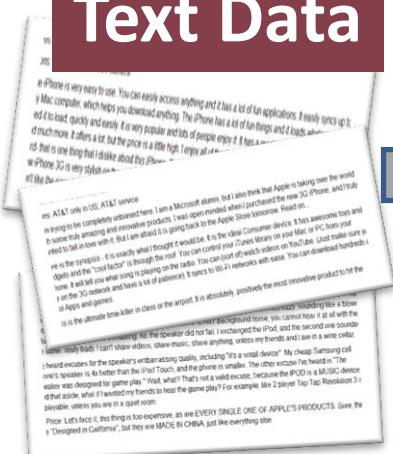
Many applications are possible if we can “decode” the topics in text...

# Mining Multiple Topics from Text

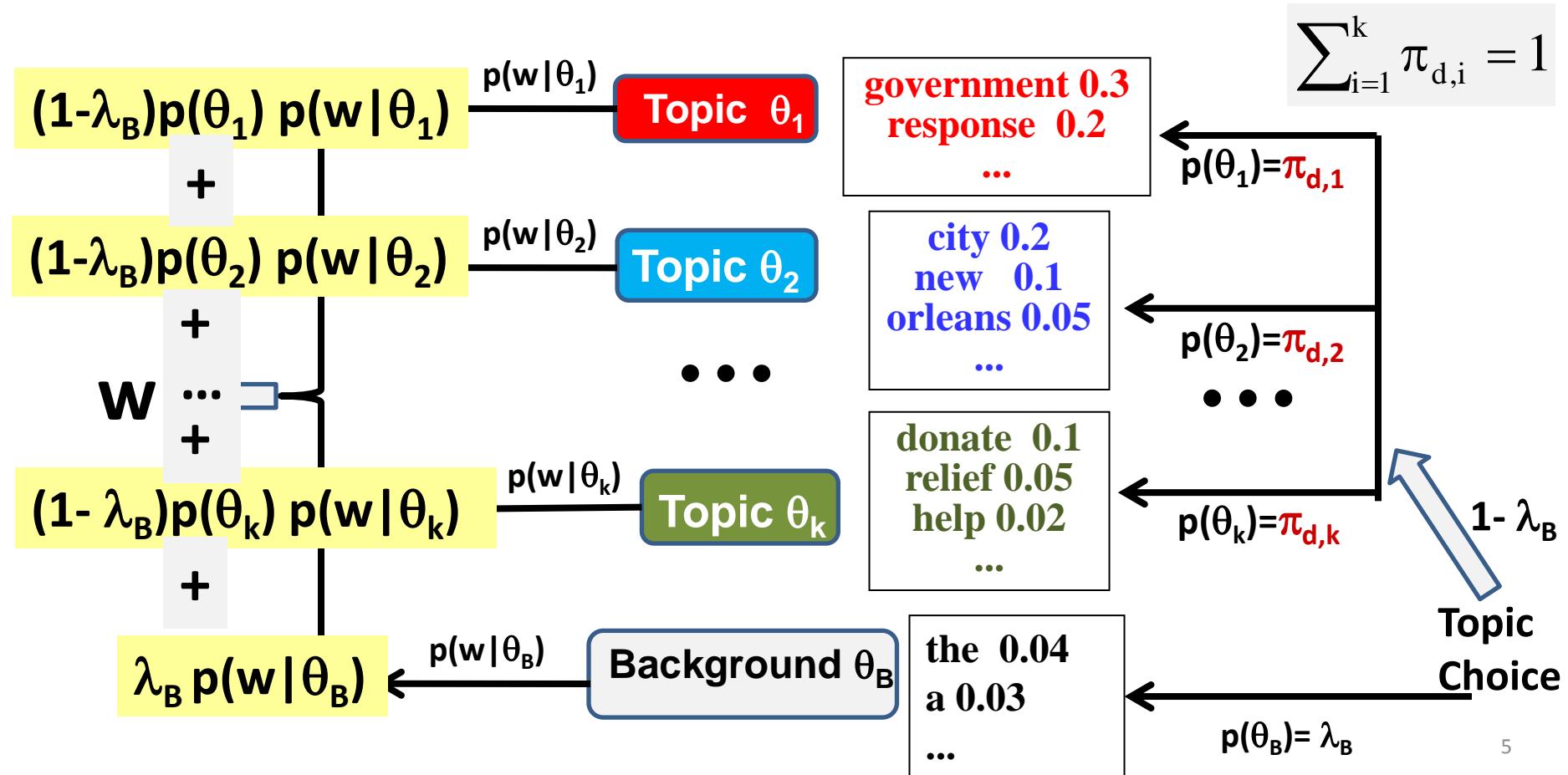
INPUT:  $C, k, V$

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

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  $\theta_j$  in doc d

Prob. of word w in topic  $\theta_j$

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

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

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

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

How many unknown parameters are there in total?

# ML Parameter Estimation

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

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

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

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

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

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

# 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**  $\theta_j$

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

**Use of Bayes Rule**

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

Probability that **w** in doc **d** is generated from **background**  $\theta_B$

# 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  $\theta_j$

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

ML Estimate based on  
“allocated” word  
counts to topic  $\theta_j$

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

↑  
Re-estimated probability of word w for topic  $\theta_j$

# Computation of the EM Algorithm

- Initialize all unknown parameters randomly
- Repeat until likelihood converges

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

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

– M-step

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

In general, accumulate counts, and then normalize

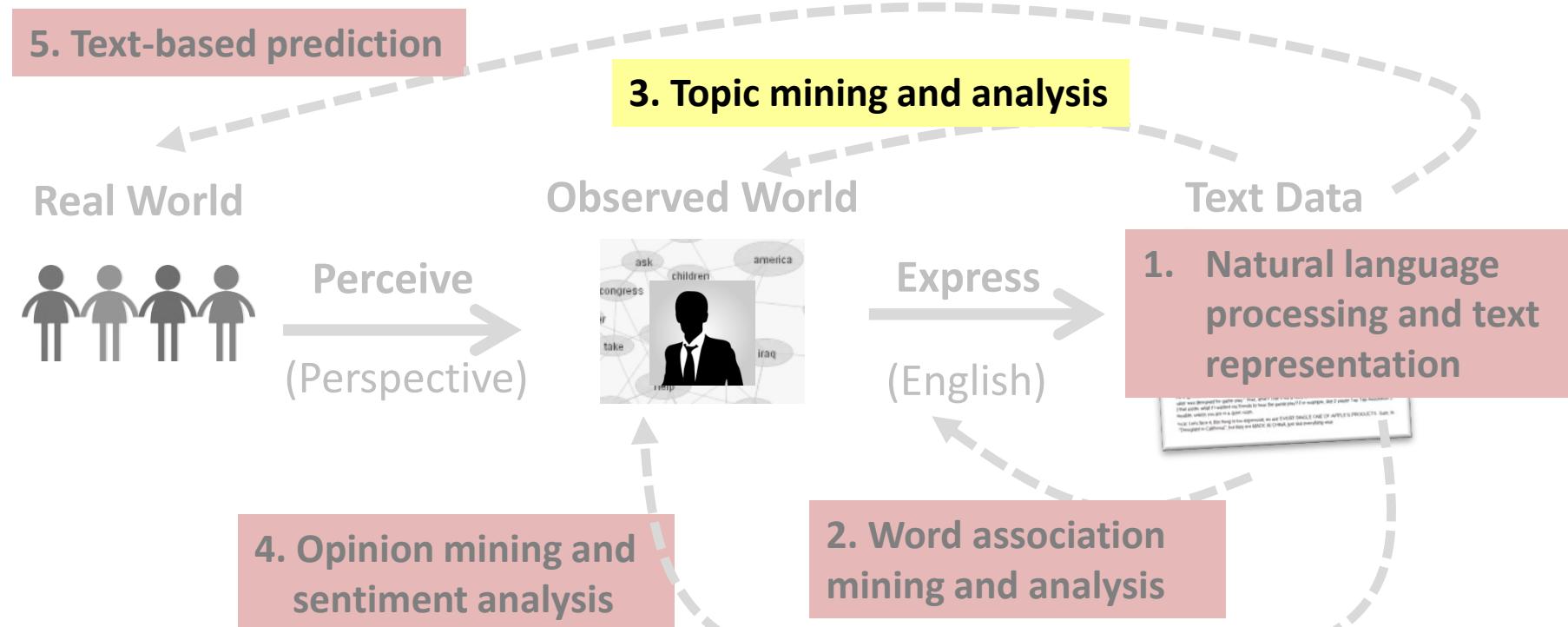
# Summary

- PLSA = mixture model with  $k$  unigram LMs ( $k$  topics)
- Adding a pre-determined background LM helps discover discriminative topics
- ML estimate “discovers” topical knowledge from text data
  - $k$  word distributions ( $k$  topics)
  - proportion of each topic in each document
- The output can enable many applications!
  - Clustering of terms and docs (treat each topic as a cluster)
  - Further associate topics with different contexts (e.g., time periods, locations, authors, sources, etc.)

# Probabilistic Latent Semantic Analysis (PLSA)

ChengXiang “Cheng” Zhai  
Department of Computer Science  
University of Illinois at Urbana-Champaign

# Probabilistic Latent Semantic Analysis (PLSA)



# Document as a Sample of Mixed Topics

Topic  $\theta_1$

government 0.3  
response 0.2

...

Topic  $\theta_2$

city 0.2  
new 0.1  
orleans 0.05

...

Topic  $\theta_k$

donate 0.1  
relief 0.05  
help 0.02

...

Background  $\theta_B$

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]. ...

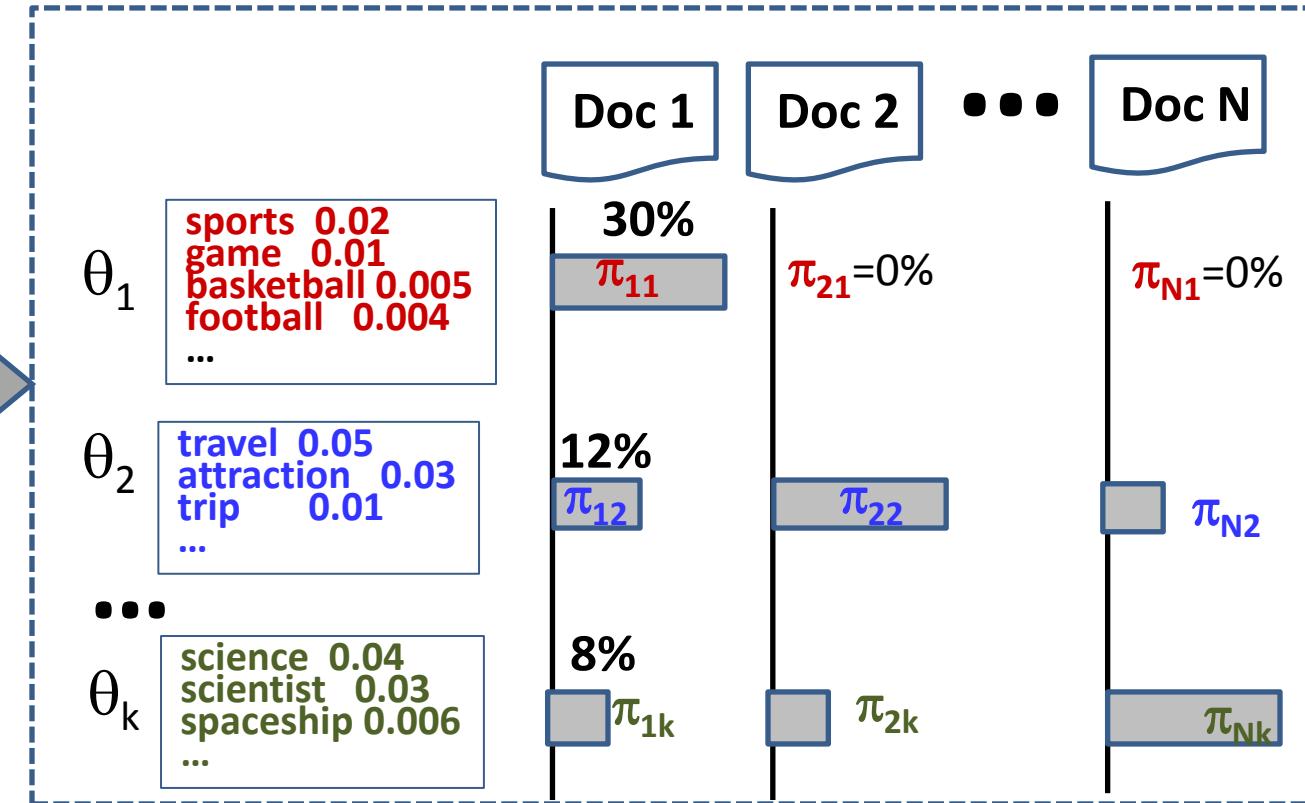
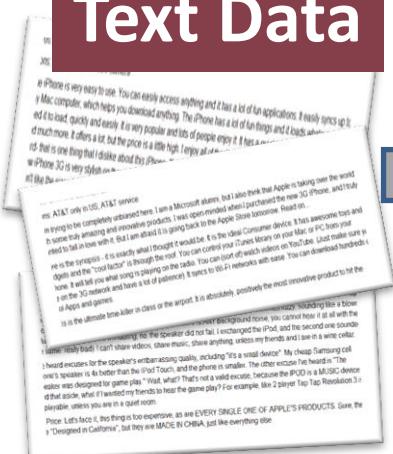
Many applications are possible if we can “decode” the topics in text...

# Mining Multiple Topics from Text

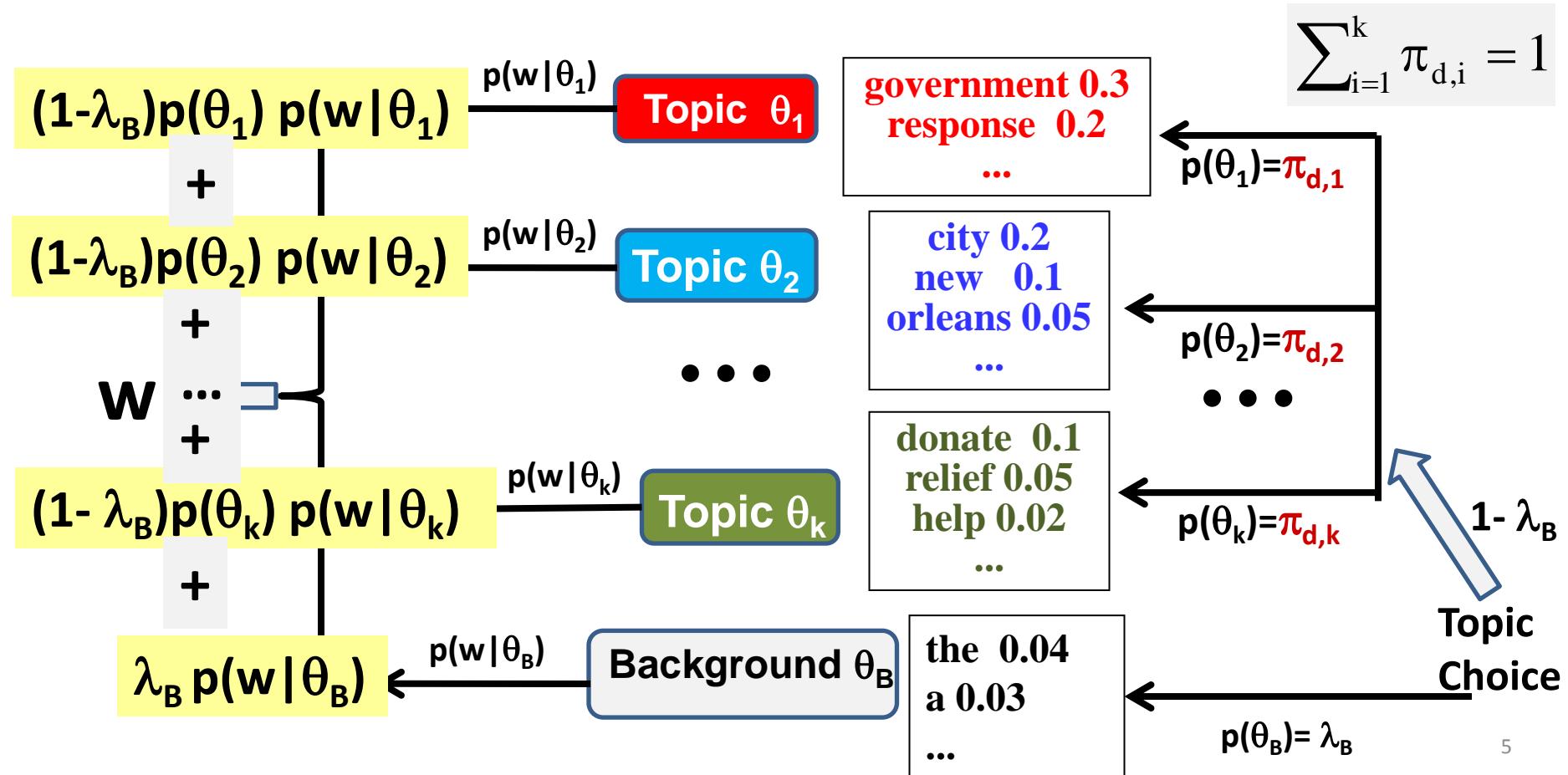
INPUT:  $C, k, V$

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

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  $\theta_j$  in doc d

Prob. of word w in topic  $\theta_j$

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

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

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

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

How many unknown parameters are there in total?

# ML Parameter Estimation

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

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

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

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

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

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

# 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**  $\theta_j$

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

**Use of Bayes Rule**

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

Probability that **w** in doc **d** is generated from **background**  $\theta_B$

# 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  $\theta_j$

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

ML Estimate based on  
“allocated” word  
counts to topic  $\theta_j$

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

↑  
Re-estimated probability of word w for topic  $\theta_j$

# Computation of the EM Algorithm

- Initialize all unknown parameters randomly
- Repeat until likelihood converges

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

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

– M-step

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

In general, accumulate counts, and then normalize

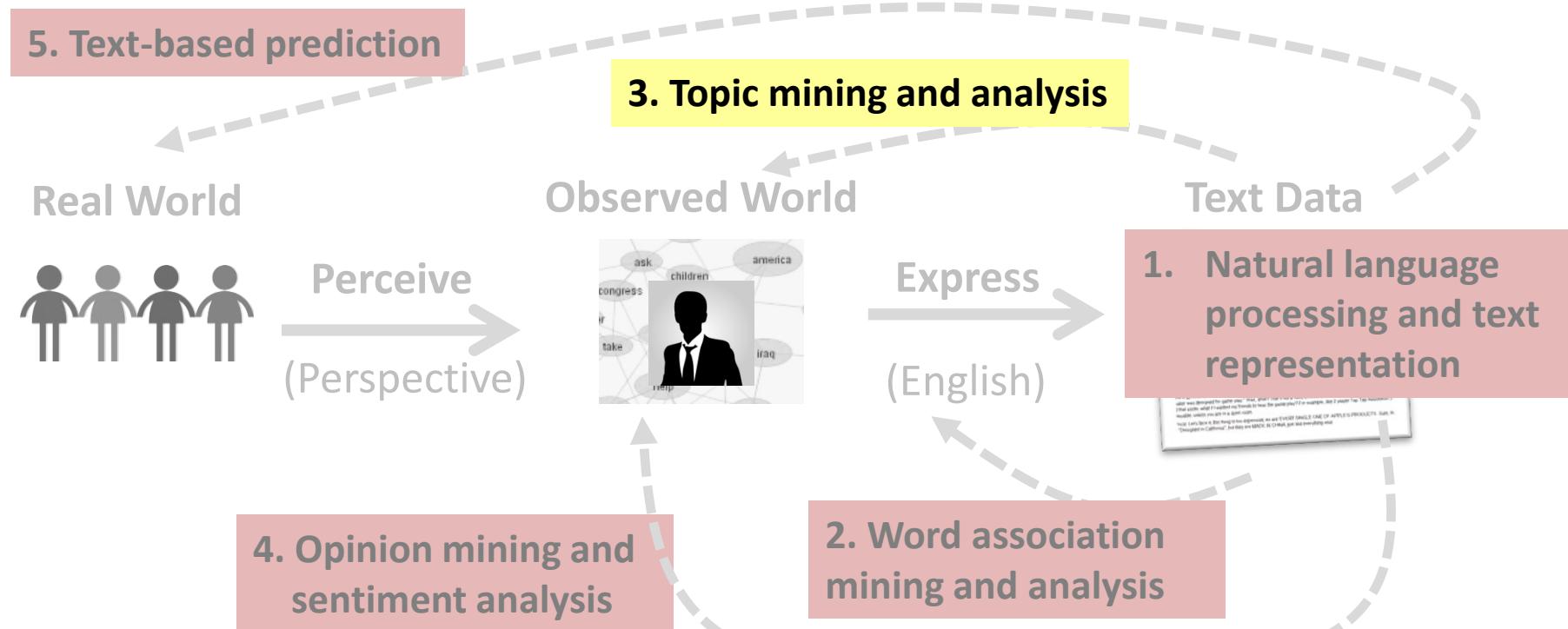
# Summary

- PLSA = mixture model with  $k$  unigram LMs ( $k$  topics)
- Adding a pre-determined background LM helps discover discriminative topics
- ML estimate “discovers” topical knowledge from text data
  - $k$  word distributions ( $k$  topics)
  - proportion of each topic in each document
- The output can enable many applications!
  - Clustering of terms and docs (treat each topic as a cluster)
  - Further associate topics with different contexts (e.g., time periods, locations, authors, sources, etc.)

# Latent Dirichlet Allocation (LDA)

ChengXiang “Cheng” Zhai  
Department of Computer Science  
University of Illinois at Urbana-Champaign

# Latent Dirichlet Allocation (LDA)



# Extensions of PLSA

- PLSA with prior knowledge → User-controlled PLSA
- PLSA as a generative model → Latent Dirichlet Allocation

# 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  $\Lambda$  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 $p(w | \theta_i)$

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

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

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

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

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

battery 0.5  
life 0.5

Pseudo counts of w  
from prior  $\theta'$

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

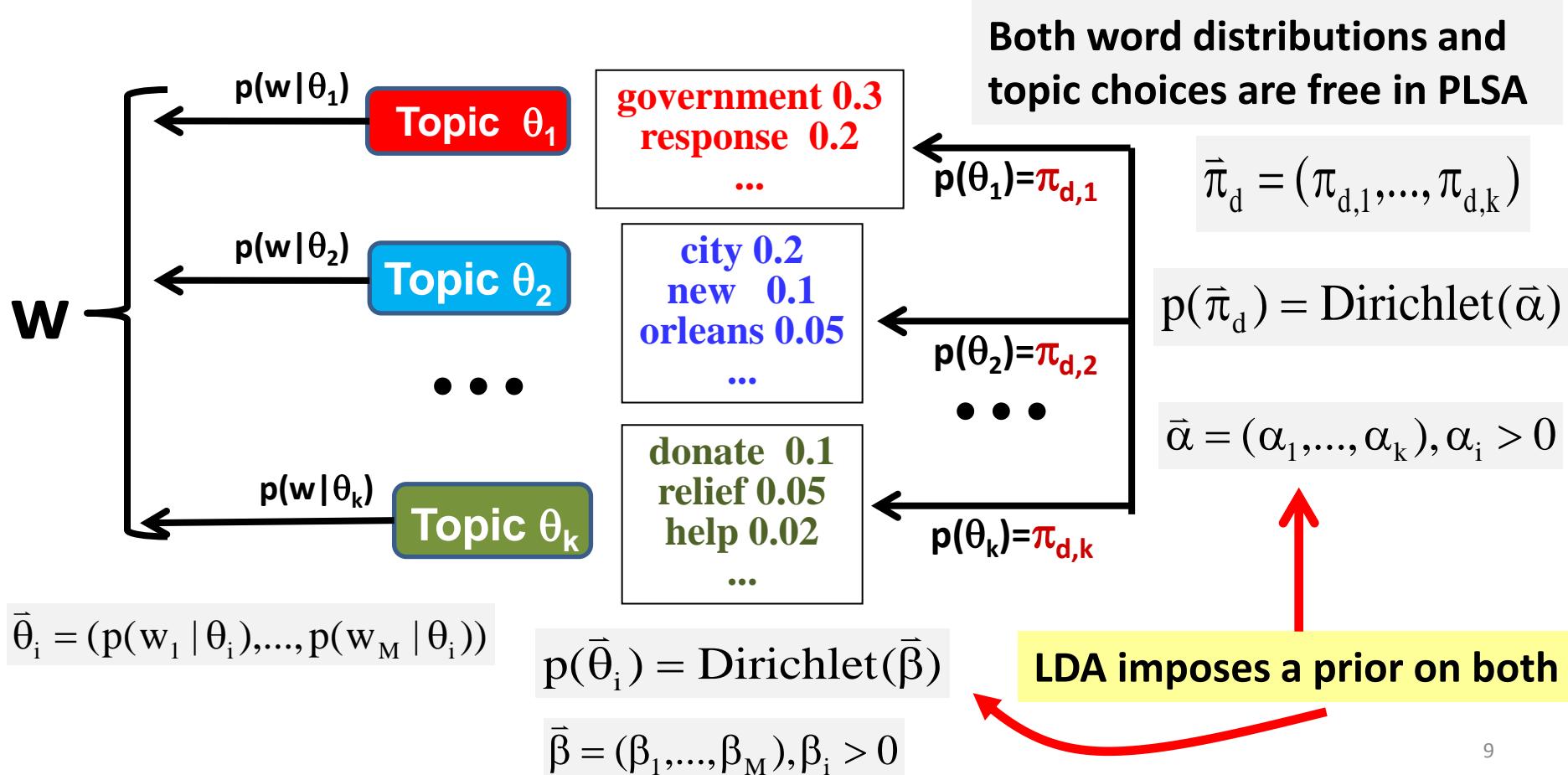
# Deficiency of PLSA

- Not a generative model
  - Can't compute probability of a new document
  - Heuristic workaround is possible, though
- Many parameters → high complexity of models
  - Many local maxima
  - Prone to overfitting
- Not necessarily a problem for text mining (only interested in fitting the “training” documents)

# Latent Dirichlet Allocation (LDA)

- Make PLSA a generative model by imposing a Dirichlet prior on the model parameters →
  - LDA = Bayesian version of PLSA
  - Parameters are regularized
- Can achieve the same goal as PLSA for text mining purposes
  - Topic coverage and topic word distributions can be inferred using Bayesian inference

# PLSA → LDA



# Likelihood Functions for PLSA vs. LDA

PLSA

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

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

$$\log p(C | \{\theta_j\}, \{\pi_{d,j}\}) = \sum_{d \in C} \log p(d | \{\theta_j\}, \{\pi_{d,j}\})$$

Core assumption  
in all topic models

LDA

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

$$\log p(d | \bar{\alpha}, \{\theta_j\}) = \int \left[ \sum_{w \in V} c(w, d) \log [\sum_{j=1}^k \pi_{d,j} p(w | \theta_j)] \right] p(\bar{\pi}_d | \bar{\alpha}) d\bar{\pi}_d$$

$$\log p(C | \bar{\alpha}, \bar{\beta}) = \int \sum_{d \in C} \log p(d | \bar{\alpha}, \{\theta_j\}) \prod_{j=1}^k p(\theta_j | \bar{\beta}) d\theta_1 \dots d\theta_k$$

PLSA component

Added by LDA

# Parameter Estimation and Inferences in LDA

- Parameters can be estimated using ML estimator

$$(\hat{\vec{\alpha}}, \hat{\vec{\beta}}) = \arg \max_{\vec{\alpha}, \vec{\beta}} \log p(C | \vec{\alpha}, \vec{\beta})$$

How many parameters in LDA vs. PLSA?

- However,  $\{\theta_j\}$  and  $\{\pi_{d,j}\}$  must now be computed using posterior inference
  - Computationally intractable
  - Must resort to approximate inference
  - Many different inference methods are available

# Summary of Probabilistic Topic Models

- Probabilistic topic models provide a general principled way of mining and analyzing topics in text with many applications
- Basic task setup:
  - Input: Text data
  - Output:  $k$  topics + proportions of these topics covered in each document
- PLSA is the basic topic model, often adequate for most applications
- LDA improves over PLSA by imposing priors
  - Theoretically more appealing
  - Practically, LDA and PLSA perform similarly for many tasks

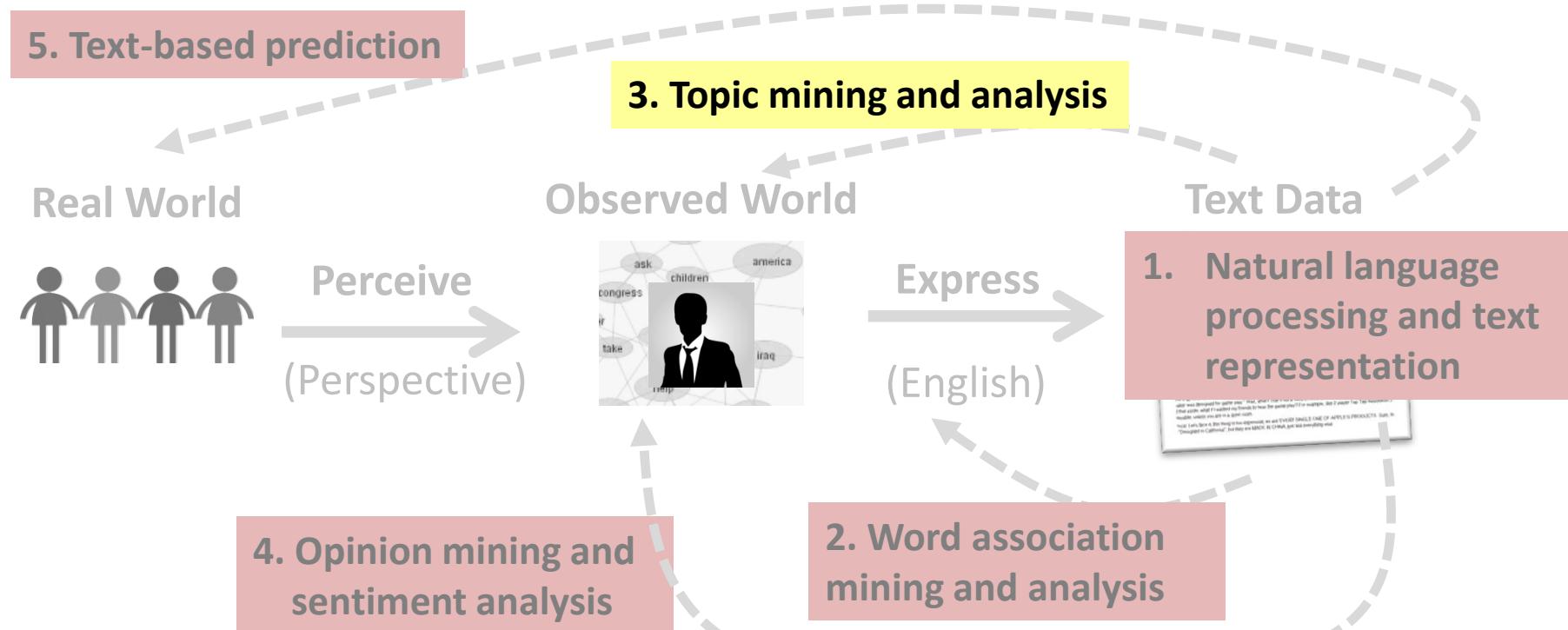
# Suggested Readings

- Blei, D. 2012. “Probabilistic Topic Models.” *Communications of the ACM* 55 (4): 77–84. doi: 10.1145/2133806.2133826.
- Qiaozhu Mei, Xuehua Shen, and ChengXiang Zhai. “Automatic Labeling of Multinomial Topic Models.” *Proceedings of ACM KDD 2007*, pp. 490-499, DOI=10.1145/1281192.1281246.
- Yue Lu, Qiaozhu Mei, and Chengxiang Zhai. 2011. Investigating task performance of probabilistic topic models: an empirical study of PLSA and LDA. *Information Retrieval*, 14, 2 (April 2011), 178-203. DOI=10.1007/s10791-010-9141-9.

# Latent Dirichlet Allocation (LDA)

ChengXiang “Cheng” Zhai  
Department of Computer Science  
University of Illinois at Urbana-Champaign

# Latent Dirichlet Allocation (LDA)



# Extensions of PLSA

- PLSA with prior knowledge → User-controlled PLSA
- PLSA as a generative model → Latent Dirichlet Allocation

# 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  $\Lambda$  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 $p(w | \theta_i)$

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

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

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

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

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

battery 0.5  
life 0.5

Pseudo counts of w  
from prior  $\theta'$

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

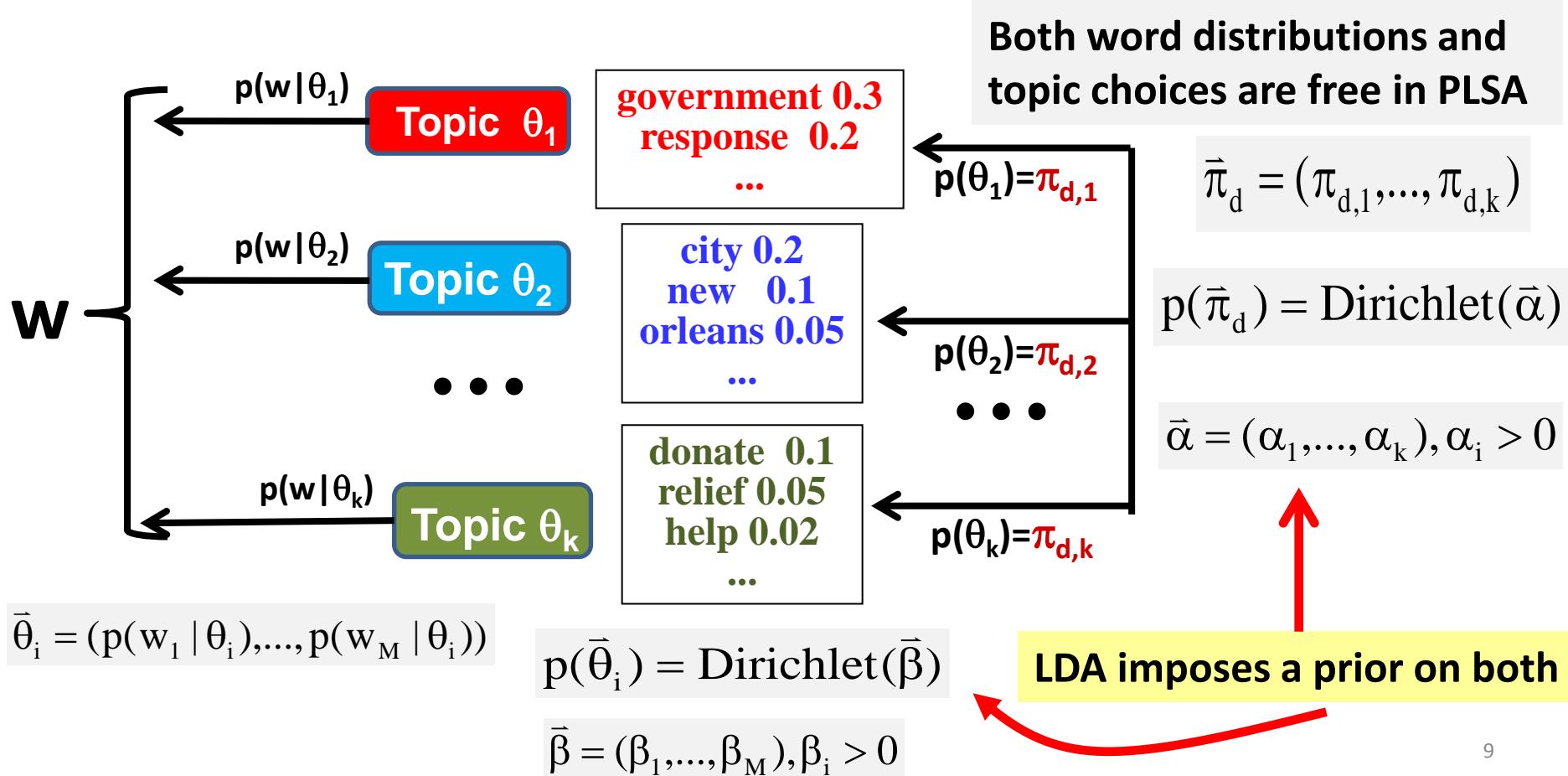
# Deficiency of PLSA

- Not a generative model
  - Can't compute probability of a new document
  - Heuristic workaround is possible, though
- Many parameters → high complexity of models
  - Many local maxima
  - Prone to overfitting
- Not necessarily a problem for text mining (only interested in fitting the “training” documents)

# Latent Dirichlet Allocation (LDA)

- Make PLSA a generative model by imposing a Dirichlet prior on the model parameters →
  - LDA = Bayesian version of PLSA
  - Parameters are regularized
- Can achieve the same goal as PLSA for text mining purposes
  - Topic coverage and topic word distributions can be inferred using Bayesian inference

# PLSA → LDA



# Likelihood Functions for PLSA vs. LDA

PLSA

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

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

$$\log p(C | \{\theta_j\}, \{\pi_{d,j}\}) = \sum_{d \in C} \log p(d | \{\theta_j\}, \{\pi_{d,j}\})$$

Core assumption  
in all topic models

LDA

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

$$\log p(d | \bar{\alpha}, \{\theta_j\}) = \int \left[ \sum_{w \in V} c(w, d) \log [\sum_{j=1}^k \pi_{d,j} p(w | \theta_j)] \right] p(\bar{\pi}_d | \bar{\alpha}) d\bar{\pi}_d$$

$$\log p(C | \bar{\alpha}, \bar{\beta}) = \int \sum_{d \in C} \log p(d | \bar{\alpha}, \{\theta_j\}) \prod_{j=1}^k p(\theta_j | \bar{\beta}) d\theta_1 \dots d\theta_k$$

PLSA component

Added by LDA

# Parameter Estimation and Inferences in LDA

- Parameters can be estimated using ML estimator

$$(\hat{\vec{\alpha}}, \hat{\vec{\beta}}) = \arg \max_{\vec{\alpha}, \vec{\beta}} \log p(C | \vec{\alpha}, \vec{\beta})$$

How many parameters in LDA vs. PLSA?

- However,  $\{\theta_j\}$  and  $\{\pi_{d,j}\}$  must now be computed using posterior inference
  - Computationally intractable
  - Must resort to approximate inference
  - Many different inference methods are available

# Summary of Probabilistic Topic Models

- Probabilistic topic models provide a general principled way of mining and analyzing topics in text with many applications
- Basic task setup:
  - Input: Text data
  - Output:  $k$  topics + proportions of these topics covered in each document
- PLSA is the basic topic model, often adequate for most applications
- LDA improves over PLSA by imposing priors
  - Theoretically more appealing
  - Practically, LDA and PLSA perform similarly for many tasks

# Suggested Readings

- Blei, D. 2012. “Probabilistic Topic Models.” *Communications of the ACM* 55 (4): 77–84. doi: 10.1145/2133806.2133826.
- Qiaozhu Mei, Xuehua Shen, and ChengXiang Zhai. “Automatic Labeling of Multinomial Topic Models.” *Proceedings of ACM KDD 2007*, pp. 490-499, DOI=10.1145/1281192.1281246.
- Yue Lu, Qiaozhu Mei, and Chengxiang Zhai. 2011. Investigating task performance of probabilistic topic models: an empirical study of PLSA and LDA. *Information Retrieval*, 14, 2 (April 2011), 178-203. DOI=10.1007/s10791-010-9141-9.