

Data Mining and Analysis

Text Mining: 2

CSCE 676 :: Fall 2019

Texas A&M University

Department of Computer Science & Engineering

Prof. James Caverlee

Resources

<https://web.stanford.edu/~jurafsky/slp3/> — slides and readings on Sentiment Analysis

http://www.cs.virginia.edu/~hw5x/Course/TextMining-2019Spring/_site/lectures/ — slides from Hongning Wang's course on Text Mining

<http://www.cs.columbia.edu/~blei/topicmodeling.html>

Opinion mining and sentiment analysis by Pang and Lee

Probabilistic Topic Models (CACM) by Blei

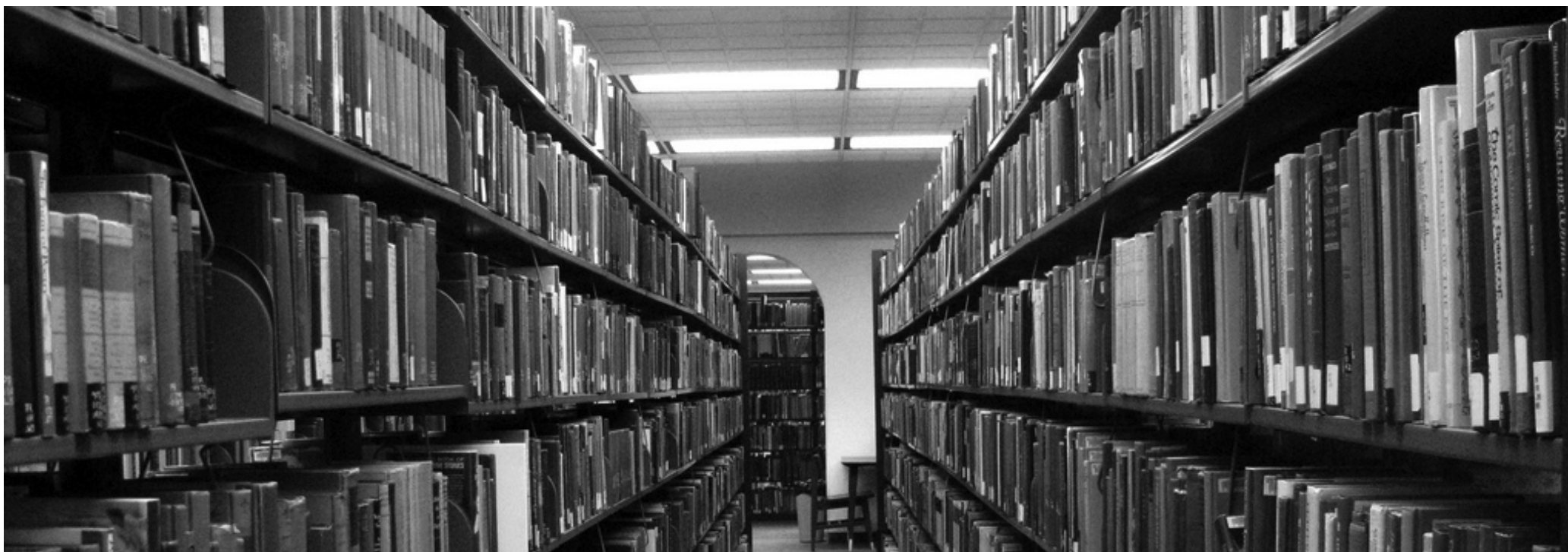
Latent dirichlet allocation by Blei, Ng, and Jordan

Probabilistic latent semantic analysis by Thomas Hofmann

Topic Models

http://www.cs.virginia.edu/~hw5x/Course/TextMining-2019Spring/_site/lectures/

<http://www.cs.columbia.edu/~blei/papers/Blei2012.pdf>



Input: An unorganized collection of documents
Output: An organized collection, and a description of how

What are Topics?

Topic = A broad semantically coherent theme, usually hidden in documents

Examples: politics, sports, technology, etc.

How to Represent Topics?

Typically as a **probability distribution** over **words**

Example for “texas a&m football”:

jimbo	0.020
--------------	--------------

aggies	0.015
---------------	--------------

touchdown	0.011
------------------	--------------

win	0.009
------------	--------------

Remember: **words** could be unigrams, bigrams, phrases, ...

Documents are a mix of topics

Topics

gene	0.04
dna	0.02
genetic	0.01
...	

life	0.02
evolve	0.01
organism	0.01
...	

brain	0.04
neuron	0.02
nerve	0.01
...	

data	0.02
number	0.02
computer	0.01
...	

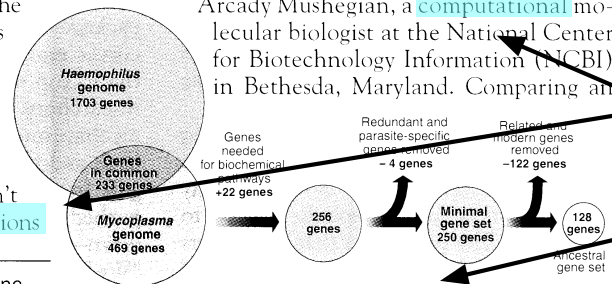
Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic numbers** game, particularly as more and more **genomes** are completely mapped and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

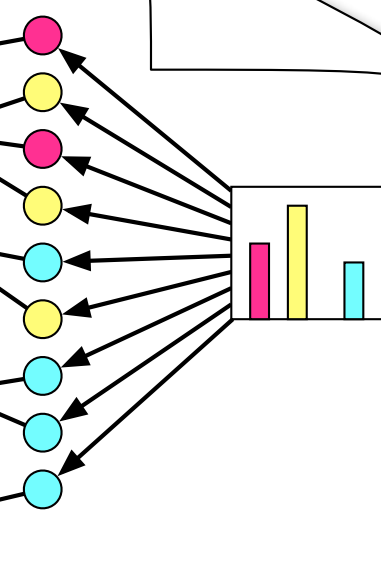


* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. **Computer analysis** yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments



Consider the following “documents”

I like to eat broccoli and bananas.

I ate a banana and spinach smoothie for breakfast.

Chinchillas and kittens are cute.

My sister adopted a kitten yesterday.

Look at this cute hamster munching on a piece of broccoli.



= 30% broccoli, 15% bananas, 10% breakfast, 10% munching, ... (at which point, you could interpret topic A to be about food)



= 20% chinchillas, 20% kittens, 20% cute, 15% hamster, ... (at which point, you could interpret topic B to be about cute animals)

Motivating Question

How can we discover these topics? How about the word distributions?

Many applications would be enabled by discovering such topics

- Summarize themes/aspects

- Facilitate navigation/browsing

- Retrieve documents

- Segment documents

- Many other text mining tasks

Topic Models

Topic: a multinomial distribution over words

Document: a mixture of topics

A document is “generated” by first sampling topics from some prior distribution

Each time, sample a word from a corresponding topic

Many variations of how these topics are mixed

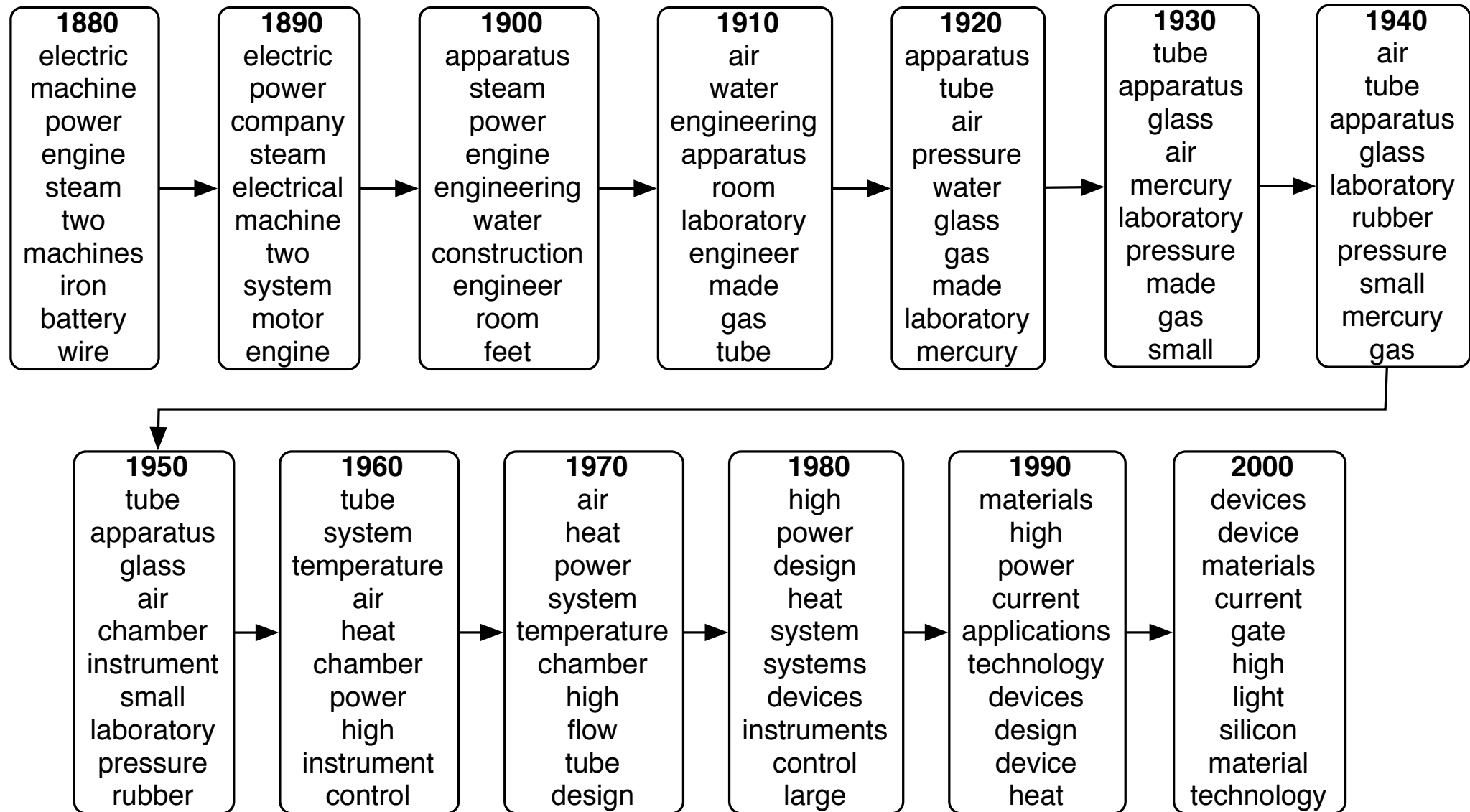
Topic modeling

Fitting the probabilistic model to text

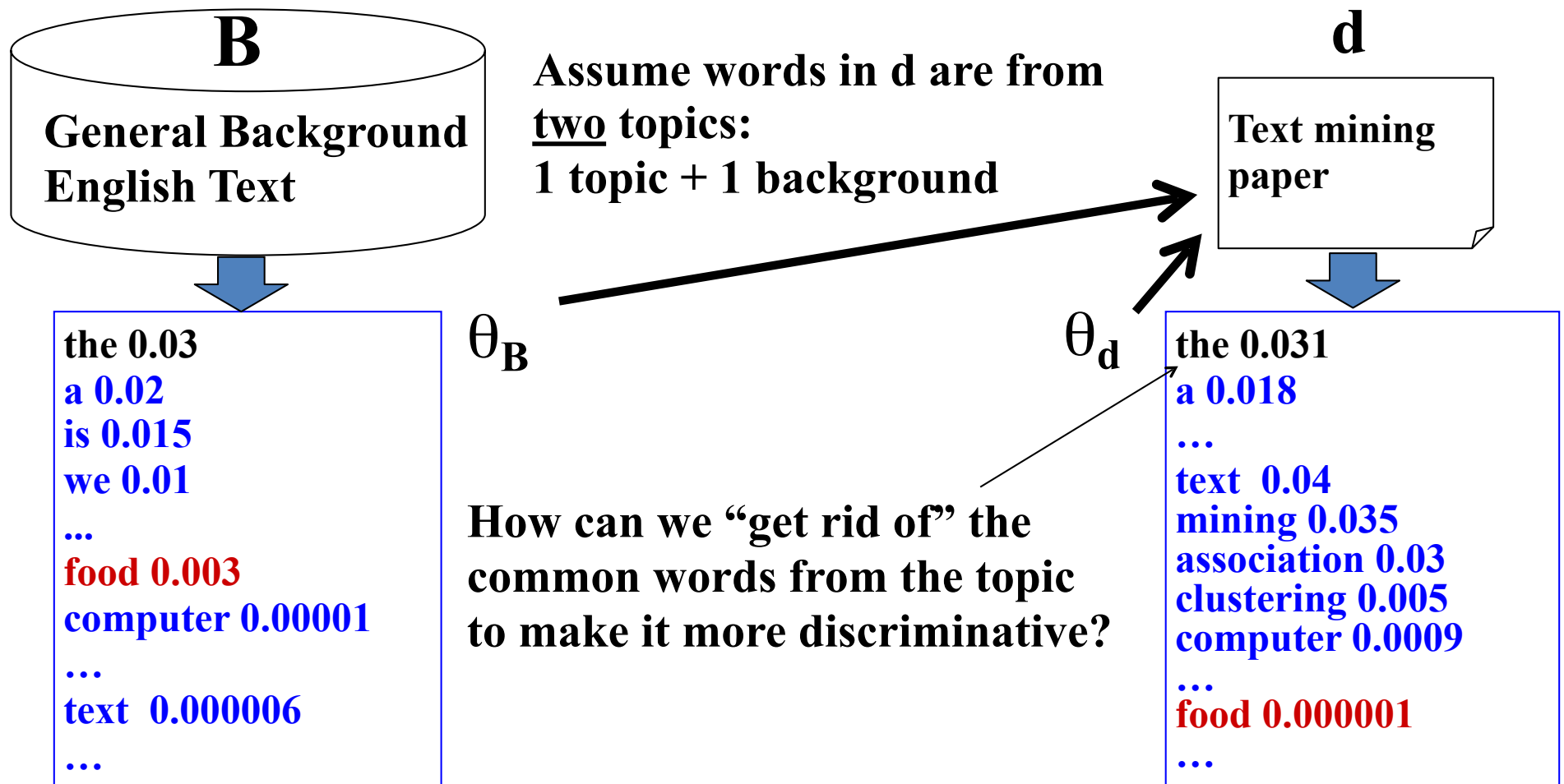
Answer topic-related questions by computing various kinds of posterior distributions

e.g., $p(\text{topic} \mid \text{time})$, $p(\text{sentiment} \mid \text{topic})$

Example: Scientific Topics over Time



Simplest Case: 1 topic + 1 “background”



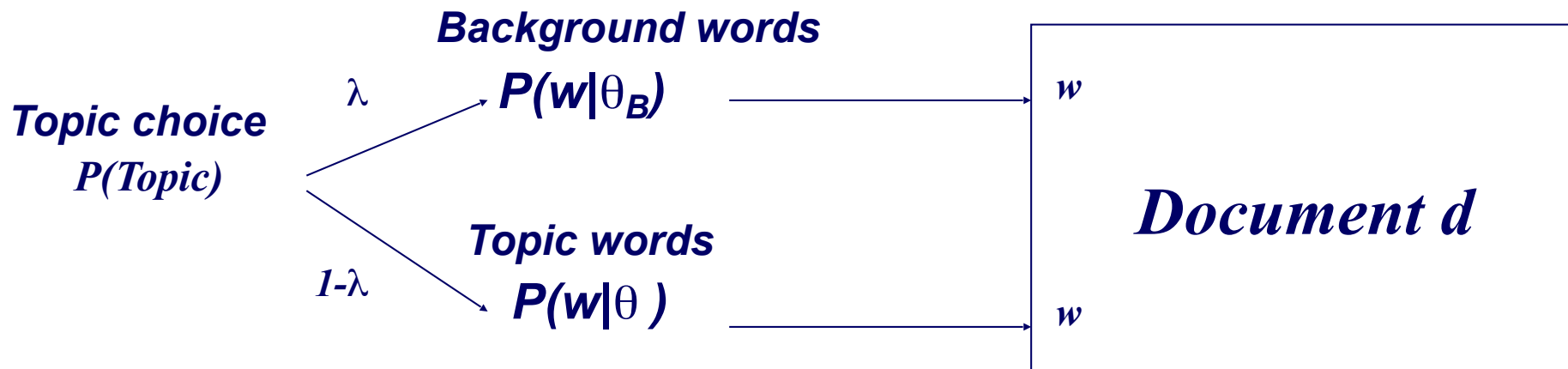
Background Topic: $p(w|\theta_B)$

Document Topic: $p(w|\theta_d)$

The Simplest Case: One Topic + One Background Model

Assume $p(w|\theta_B)$ and λ are *known*

λ = mixing proportion of background topic in d



$$p(w) = \lambda p(w|\theta_B) + (1-\lambda)p(w|\theta)$$

$$\log p(d|\theta) = \sum_{w \in V} c(w, d) \log [\lambda p(w|\theta_B) + (1-\lambda)p(w|\theta)]$$

Expectation Maximization $\hat{\theta} = \arg \max_{\theta} \log p(d|\theta)$

How to Estimate θ ?

**Known
Background
 $p(w|\theta_B)$**

the 0.2
a 0.1
we 0.01
to 0.02
...
text 0.0001
mining 0.00005
...

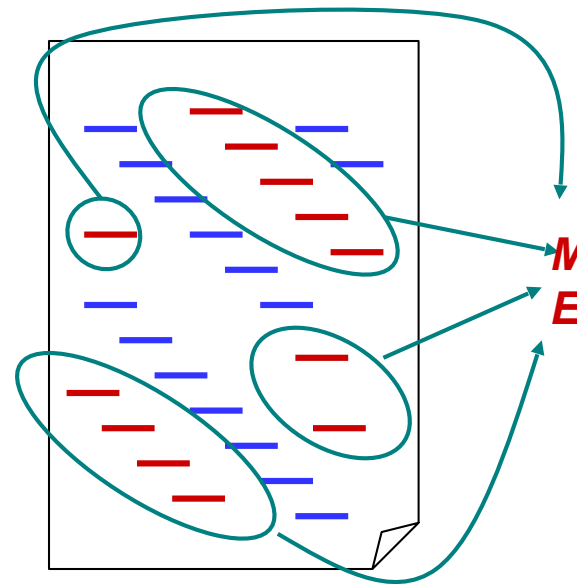
**Unknown
topic $p(w|\theta)$ for
“Text mining”**

...
text =?
mining =?
association =?
word =?
...

$\lambda=0.7$



**Observed
words**



**ML
Estimator**

$\lambda=0.3$



**Suppose we know
the identity/label of each word ...**

**But we
don't!**

We guess the topic assignments

Assignment (“hidden”) variable: $z_i \in \{1 \text{ (background)}, 0 \text{ (topic)}\}$

	z_i
the	1
paper	1
presents	1
a	1
text	0
mining	0
algorithm	0
the	1
paper	0
...	...

Suppose the parameters are all known, what's a reasonable guess of z_i ?

- depends on λ

- depends on $p(w|\theta_B)$ and $p(w|\theta)$

$$p(z_i = 1 | w_i) = \frac{p(z_i = 1)p(w | z_i = 1)}{p(z_i = 1)p(w | z_i = 1) + p(z_i = 0)p(w | z_i = 0)}$$

$$= \frac{\lambda p(w | \theta_B)}{\lambda p(w | \theta_B) + (1 - \lambda) p^{current}(w | \theta)}$$

E-step

$$p^{new}(w_i | \theta) = \frac{c(w_i, d)(1 - p(z_i = 1 | w_i))}{\sum_{w' \in V} c(w', d)(1 - p(z_i = 1 | w'))}$$

M-step

θ_B and θ are competing for explaining words in document d!

Initially, set $p(w | \theta)$ to some random values, then iterate ...



An example of EM computation

$$p^{(n)}(z_i = 1 | w_i) = \frac{\lambda p(w_i | \theta_B)}{\lambda p(w_i | \theta_B) + (1 - \lambda) p^{(n)}(w_i | \theta)}$$
Expectation-Step:
Augmenting data by guessing hidden variables

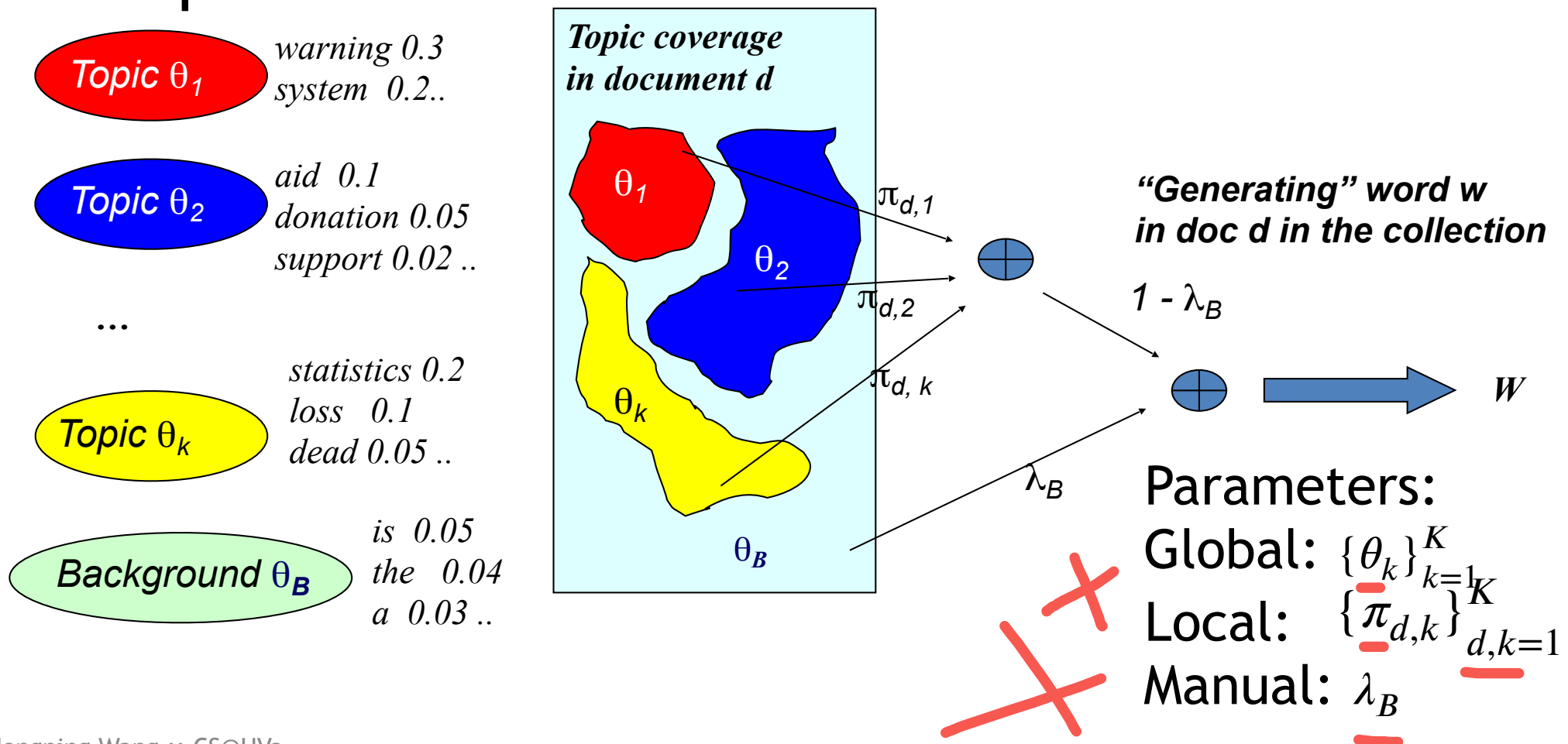
$$p^{(n+1)}(w_i | \theta) = \frac{c(w_i, d)(1 - p^{(n)}(z_i = 1 | w_i))}{\sum_{w_j \in \text{vocabulary}} c(w_j, d)(1 - p^{(n)}(z_j = 1 | w_j))}$$
Maximization-Step
With the “augmented data”, estimate parameters using maximum likelihood

Assume $\lambda=0.5$

Word	#	P(w θ_B)	Iteration 1		Iteration 2		Iteration 3	
			P(w θ)	P(z=1)	P(w θ)	P(z=1)	P(w θ)	P(z=1)
The	4	0.5	0.25	0.67	0.20	0.71	0.18	0.74
Paper	2	0.3	0.25	0.55	0.14	0.68	0.10	0.75
Text	4	0.1	0.25	0.29	0.44	0.19	0.50	0.17
Mining	2	0.1	0.25	0.29	0.22	0.31	0.22	0.31
Log-Likelihood			-16.96		-16.13		-16.02	

Discover multiple topics in a collection

- Generalize the two topic mixture to k topics



Probabilistic Latent Semantic Analysis

[Hofmann 99a, 99b]

- Topic: a multinomial distribution over words
- Document
 - Mixture of k topics
 - Mixing weights reflect the topic coverage
- Topic modeling
 - Word distribution under topic: $p(\underline{w} | \theta)$
 - Topic coverage: $p(\underline{\pi} | d)$



EM for estimating multiple topics



**Known
Background**
 $p(w | \theta_B)$

the 0.2
a 0.1
we 0.01
to 0.02
...

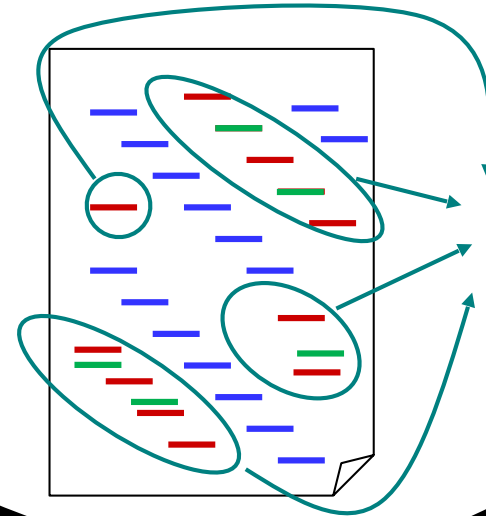
E-Step:
Predict topic labels
using Bayes Rule

Observed Words

**Unknown
topic model**
 $p(w|\theta_1)=?$

“Text mining”

...
text =?
mining =?
association =?
word =?
...

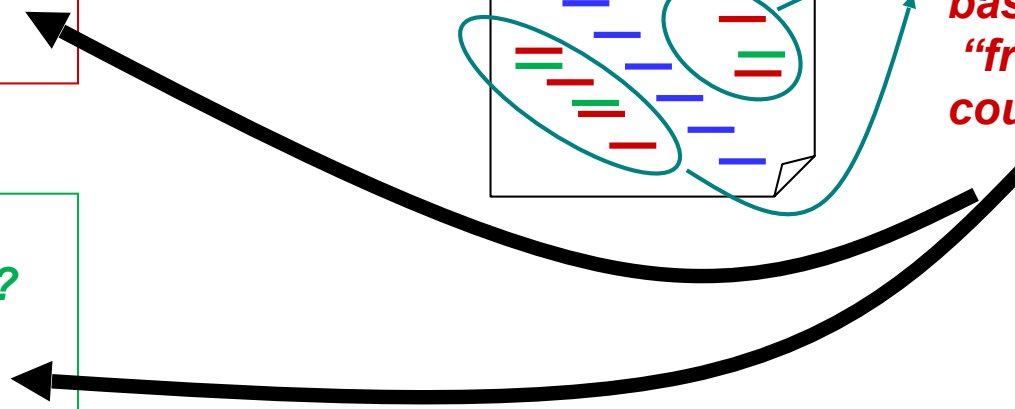


M-Step:
ML Estimator
based on
“fractional
counts”

**Unknown
topic model**
 $p(w|\theta_2)=?$

**“information
retrieval”**

...
information =?
retrieval =?
query =?
document =?
...





Parameter estimation

E-Step:

Word w in doc d is generated

- from topic j
- from background

Posterior: application of Bayes rule

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

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

M-Step:

Re-estimate

- mixing weights
- word-topic distribution

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

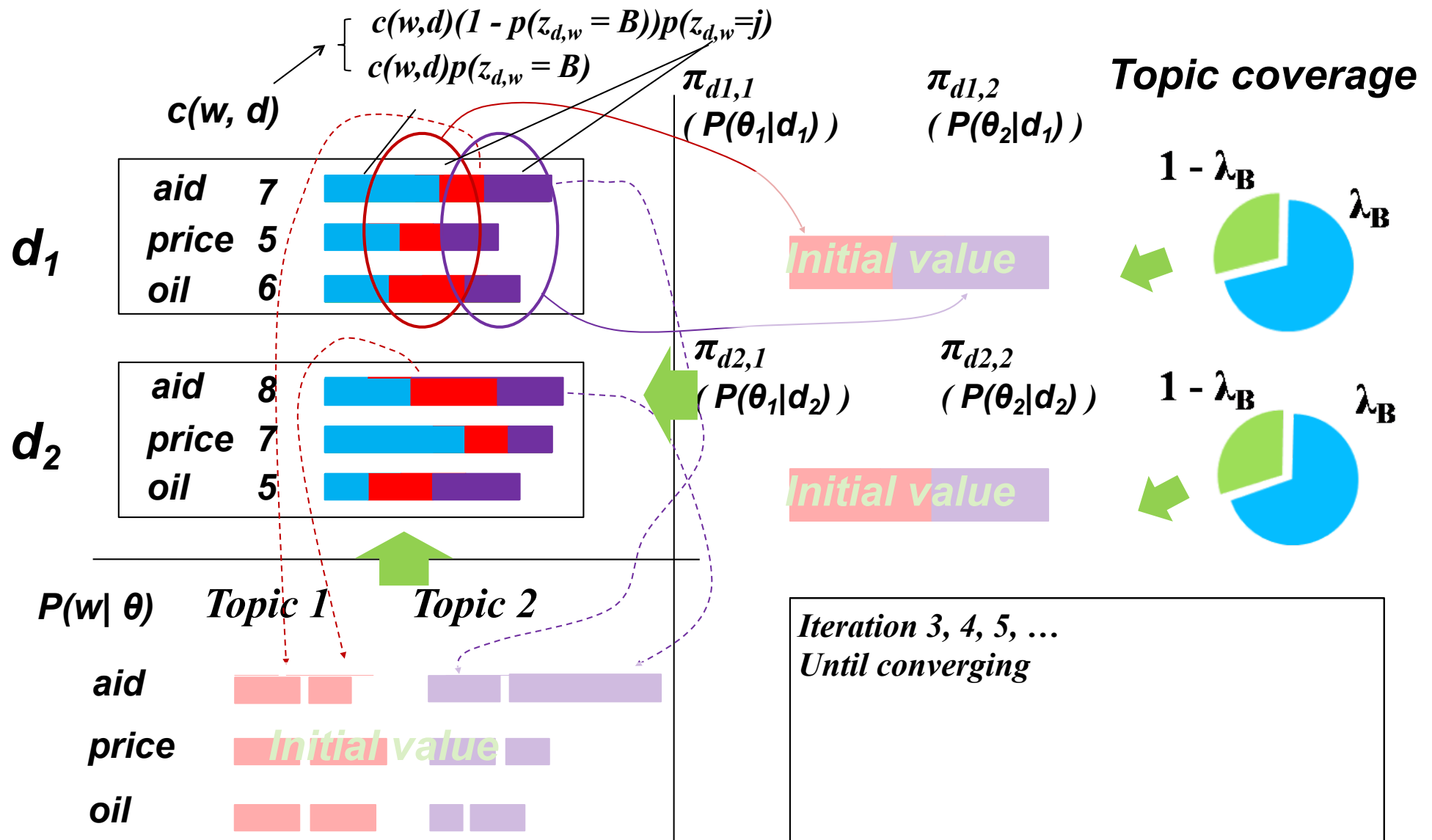
$$p^{(n+1)}(w | \theta_j) = \frac{\sum_{d \in C} c(w, d) (1 - p(z_{d,w} = B)) p(z_{d,w} = j)}{\sum_{w' \in V} \sum_{d \in C} c(w', d) (1 - p(z_{d,w'} = B)) p(z_{d,w'} = j)}$$

Sum over all docs
in the collection

Fractional counts contributing to

- using topic j in generating d
- generating w from topic j

How the algorithm works



Sample pLSA topics from TDT Corpus [Hofmann 99b]

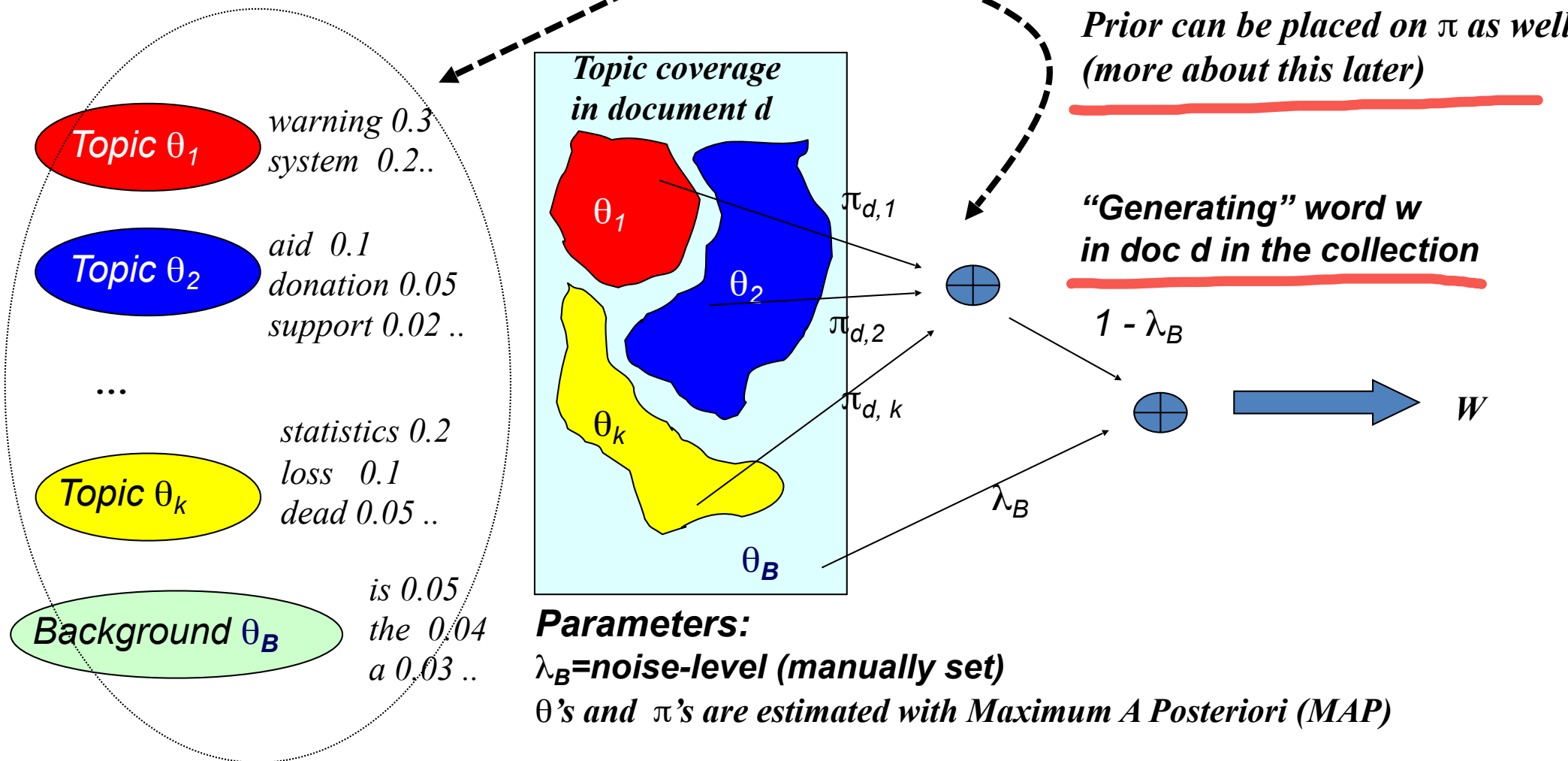
“plane”	“space shuttle”	“family”	“Hollywood”
plane	space	home	film
airport	shuttle	family	movie
crash	mission	like	music
flight	astronauts	love	new
safety	launch	kids	best
aircraft	station	mother	hollywood
air	crew	life	love
passenger	nasa	happy	actor
board	satellite	friends	entertainment
airline	earth	cnn	star

pLSA with prior knowledge

- What if we have some domain knowledge in mind
 - We want to see topics such as “battery” and “memory” for opinions about a laptop
 - We want words like “apple” and “orange” co-occur in a topic
 - One topic should be fixed to model background words (infinitely strong prior!)
- We can easily incorporate such knowledge as priors of pLSA model

Maximum a Posteriori (MAP) estimation

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



MAP estimation

- Choosing conjugate priors *Pseudo counts of w from prior θ'*
 - Dirichlet prior for multinomial distribution

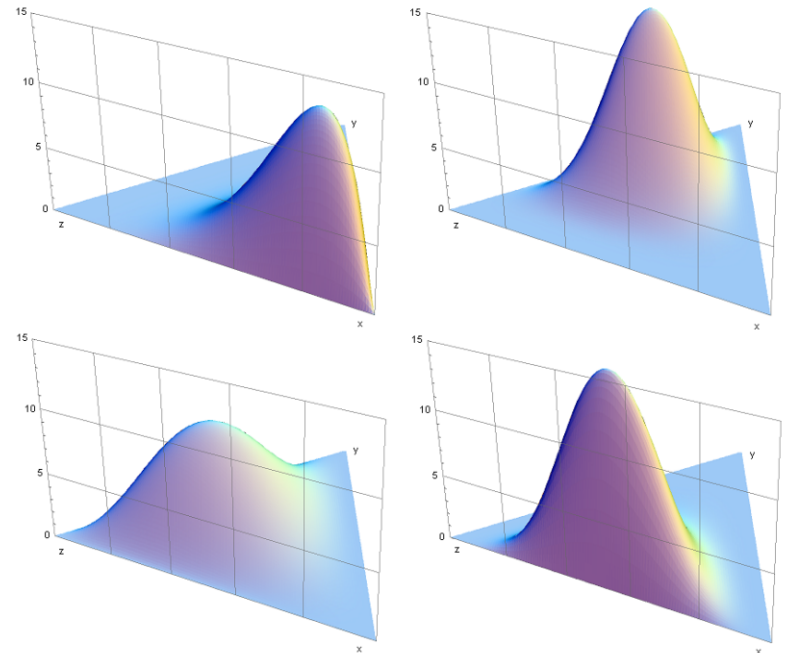
$$p^{(n+1)}(w|\theta_j) = \frac{\sum_{d \in C} c(w, d)(1 - p(z_{d,w} = B))p(z_{d,w} = j) + \mu p(w|\theta'_j)}{\sum_{w' \in V} \sum_{d \in C} c(w', d)(1 - p(z_{d,w'} = B))p(z_{d,w'} = j) + \mu}$$

- What if $\mu=0$? What if $\mu=+\infty$? *Sum of all pseudo counts*
- A consequence of using conjugate prior is that the prior can be converted into “pseudo data” which can then be “merged” with the actual data for parameter estimation

Some background knowledge

- Conjugate prior
 - Posterior distribution 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



Prior as pseudo counts



**Known
Background**
 $p(w | B)$

the 0.2
a 0.1
we 0.01
to 0.02
...

**Unknown
topic model**
 $p(w|\theta_1)=?$

“Text mining”

...
text =?
mining =?
association =?
word =?
...

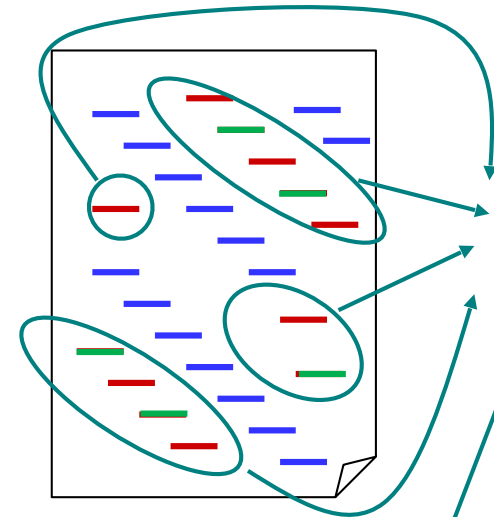
**Unknown
topic model**
 $p(w|\theta_2)=?$

**“information
retrieval”**

...
information =?
retrieval =?
query =?
document =?
...

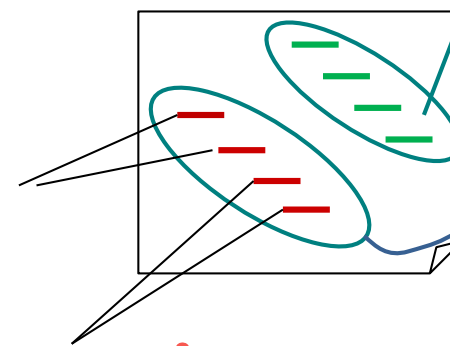
**Suppose,
we know
the identity
of each
word ...**

Observed Doc(s)



**MAP
Estimator**

Pseudo Doc



text

mining

Size = μ

Deficiency of pLSA

Then what is fully
generative model???

- Not a fully generative model
 - Can't compute probability of a new document
 - Topic coverage $p(\pi | d)$ is per-document estimated
 - Heuristic workaround is possible
- Many parameters → high complexity of models
 - Many local maxima
 - Prone to overfitting