

Case Study 5: Mixed Membership Modeling

Clustering Documents Revisited, Latent Dirichlet Allocation

Machine Learning for Big Data
CSE547/STAT548, University of Washington

Emily Fox
March 6th, 2014

©Emily Fox 2014

1

Task 2: Cluster Documents

- Then examined:
 - Cluster documents based on topic



©Emily Fox 2014

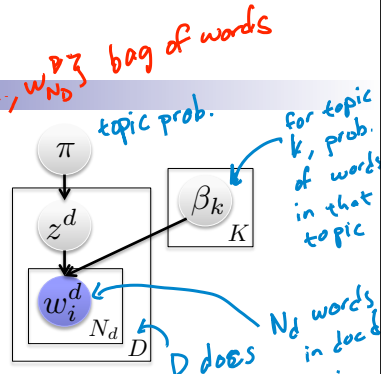
2

A Generative Model

- Documents: x^1, \dots, x^D
- Associated topics: z^1, \dots, z^D
- Parameters: $\theta = \{\pi, \beta\}$
- Generative model:

$z^d \sim \pi$ generate topic
 $w_i^d | z^d \sim \beta_{z^d} \quad i=1, \dots, N_d$

Given topic $z^d=k$ for doc d , draw each word from $\beta_k \leftarrow$ word prob. for topic k



©Emily Fox 2014

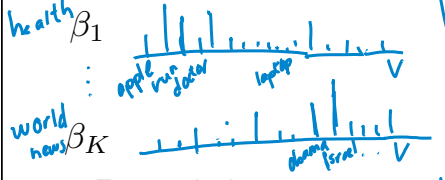
3

Model In Pictures

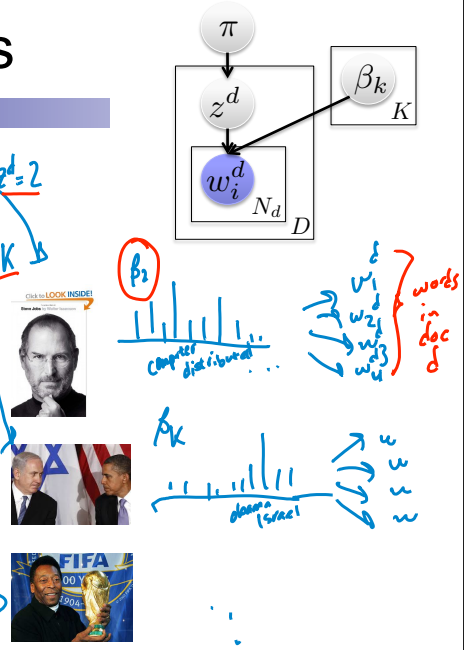
- Mixture weights (on topics)



- Topic distributions (on words)



- For each document,
 - $z^d \sim \pi$
 - $w_i^d | z^d \sim \beta_{z^d}$



©Emily Fox 2014

4

Bayesian Document Model

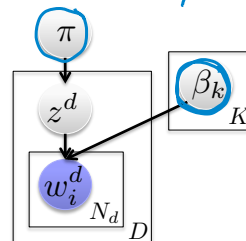
- Model parameters $\pi, \{\beta_k\}$ unknown

← can use EM as in case study 2

- Bayesian approach

place priors on parameters

- Need distribution on pmf's



$\sum_{k=1}^K \pi_k = 1$ $\sum_{k=1}^K \beta_k = 1$ ← π, β_k live on the simplex
 ① What is the simplex?
 ② What is a distribution on the simplex?

©Emily Fox 2014

5

Dirichlet Distributions

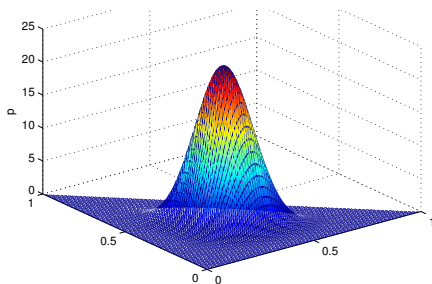
← distribution on the simplex

- The Dirichlet distribution is defined on the simplex

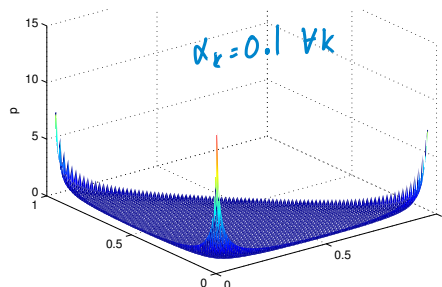
$\alpha_k = 10 \forall k$

$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$

$\Rightarrow \sum \pi_k = 1$ and $\pi_k \geq 0 \forall k$



$p(\pi | \alpha) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k - 1}$



Moments: $\mathbb{E}_\alpha[\pi_k] = \frac{\alpha_k}{\alpha_0}$

$\text{Var}_\alpha[\pi_k] = \frac{K-1}{K^2(\alpha_0+1)}$

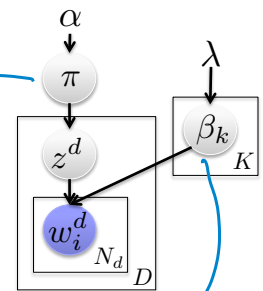
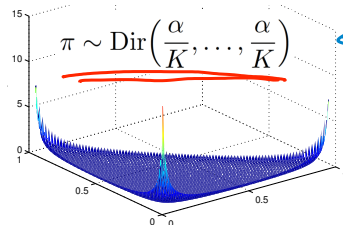
©Emily Fox 2014

6

Model Summary

- Prior on model parameters

- E.g., symmetric Dirichlet for π



- Dirichlet prior for topic parameters $\beta_k \sim \text{Dir}(\lambda_1, \dots, \lambda_V) \quad k=1, \dots, K$

- Sample observations as

$$\left. \begin{aligned} z^d &\sim \pi & d=1, \dots, D \\ w_i^d | z^d &\sim \beta_{z^d} & i=1, \dots, N_d \end{aligned} \right\}$$

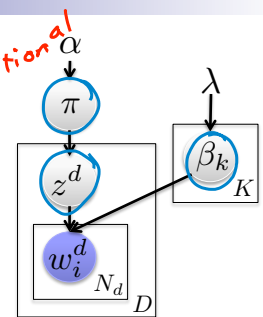
©Emily Fox 2014

7

Posterior Inference via Sampling

- Iterate between sampling

- ① $\pi \sim p(\pi | \{z^d\}, \{\beta_k\}, \{w_i^d\})$ *actual obs full conditional*
- ② For $k=1, \dots, K$
 $\beta_k \sim p(\beta_k | \pi, \{z^d\}, \{\beta_j, j \neq k\}, \{w_i^d\})$
- ③ For $d=1, \dots, D$
 $z^d \sim p(z^d | \pi, \{z^i, i \neq d\}, \{\beta_k\}, \{w_i^d\})$



- What form do these complete conditionals take?

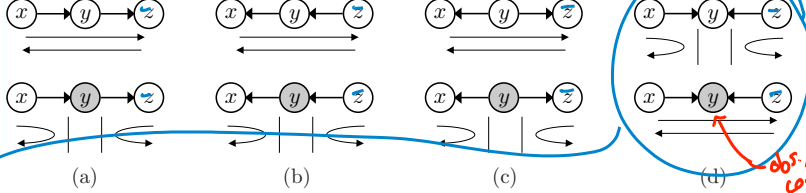
- First a look at statements of conditional independence in directed graphical models

©Emily Fox 2014

8

Conditional Independence in Bayes Nets

- Consider 4 different junction configurations



- Conditional versus unconditional independence:

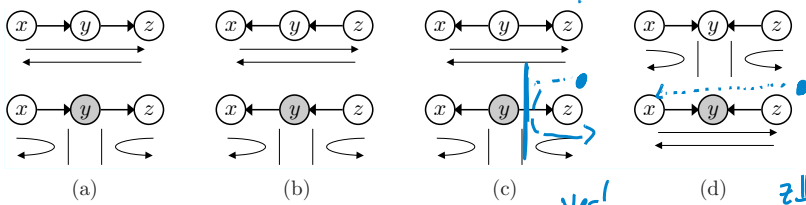
$P(x, y, z) = P(x)P(z)P(y|x, z)$ ^{int. over y} $\Rightarrow P(x, z) = P(x)P(z) \Rightarrow x \perp\!\!\!\perp z$
 $P(x, z|y) \propto P(x, y, z) = P(x)P(z)P(y|x, z) \neq P(x|y)P(z|y) \leftarrow x \not\perp\!\!\!\perp z | y$
 "explaining away": $x = \text{earthquake}$, $z = \text{burglar}$, $y = \text{car alarm}$
 If alarm ($y=1$), an increase in earthquake $P(x|y)$, means $P(z|y)$ lower
ind. a priori

©Emily Fox 2014

9

Bayes Ball Algorithm

- Consider 4 different junction configurations



- Bayes ball algorithm

start ball at one end or other.
 If ball passes to a node (straight arrows) then, not cond./marg. ind.
 If ball bounces back (walls + curved arrows), the nodes are cond./marg. ind.

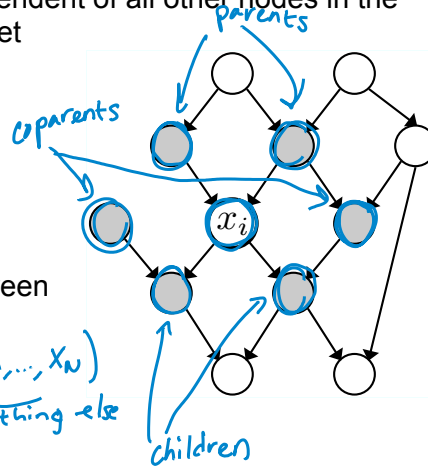
©Emily Fox 2014

10

Markov Blanket

- A node is conditionally independent of all other nodes in the graph given its Markov blanket

Markov blanket of x_i = - all parents
 - all children
 - all coparents



- Gibbs sampling iterates between full conditionals

$$x_i \sim p(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$$

everything else

→ simplify to

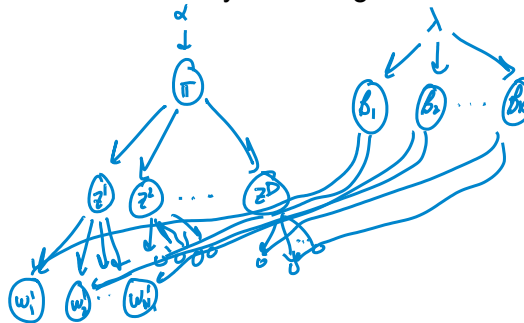
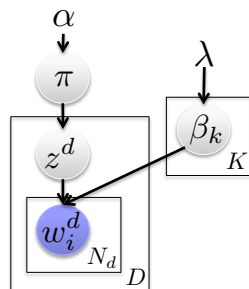
$$x_i \sim p(x_i | MB(x_i))$$

©Emily Fox 2014

11

Unplated Document Model

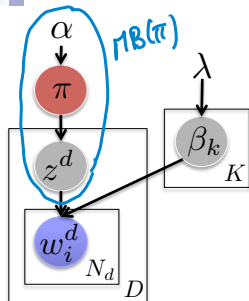
- Recall that the plate notation is really indicating



©Emily Fox 2014

12

Complete Conditional for π



- Recall conjugate Dirichlet prior

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad p(\pi | \alpha) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \pi_k^{\alpha_k - 1}$$

- Likelihood: $z^d \sim \pi \rightarrow \prod_{d=1}^D p(z^d | \pi)$

- Dirichlet posterior

- Count occurrences of $z^d = k$: $N_k = |\{z^d : z^d = k\}|$
- Then,

$$p(\pi | \{z^d\}, \alpha) \propto \prod_{d=1}^D p(z^d | \pi) p(\pi | \alpha)$$

Full cond. for π

$$\propto \prod_{k=1}^K \left(\prod_{d: z^d=k} \pi_k \right) \cdot \pi_k^{\alpha_k - 1}$$

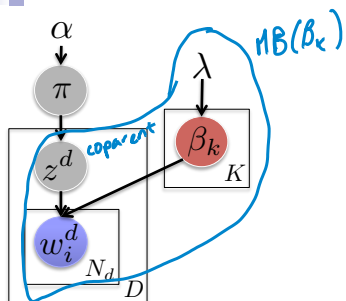
$$\propto \prod_{k=1}^K \pi_k^{N_k + \alpha_k - 1} = \text{Dir}(N_1 + \alpha_1, \dots, N_K + \alpha_K)$$

- Conjugacy: **Posterior** has same form as **prior**

©Emily Fox 2014

13

Complete Conditional for β_k



- Again, Dirichlet prior

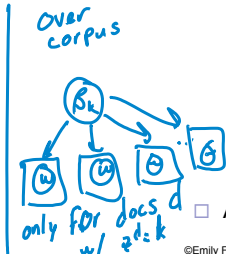
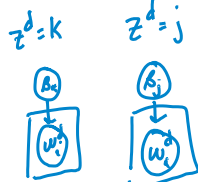
$$\beta_k \sim \text{Dir}(\lambda_1, \dots, \lambda_V)$$

- Consider docs d such that $z^d = k$

- For these observations, $w_i^d \sim \beta_k$
- Do any other docs depend on β_k ? **NO**

- Then, count $m_{v,k} = |\{w_i^d : w_i^d = v \ \forall d \text{ s.t. } z^d = k\}|$

$$\beta_k \sim \text{Dir}(m_{1,k} + \lambda_1, \dots, m_{V,k} + \lambda_V)$$

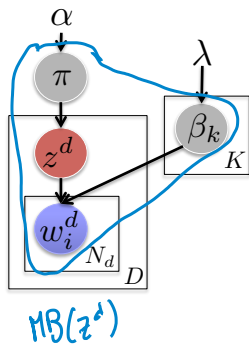


- Again, **posterior** has same form as **prior**

©Emily Fox 2014

14

Complete Conditional for z^d



- We have $z^d \sim \pi$ "prior"
- $w_i^d | z^d, \{\beta_k\} \sim \beta_{z^d}$ "likelihood"
- Calculate the posterior for each value of z^d ("responsibility" of each topic to the doc):

$$r_{dk} = p(z^d = k | \{w_i^d\}, \pi, \beta) = \frac{\pi_k p(\{w_i^d\} | \beta_k)}{\sum_j \pi_j p(\{w_i^d\} | \beta_j)}$$

- Sample each cluster indicator as $\prod_{i=1}^{N_d} \beta_{k, w_i^d}$

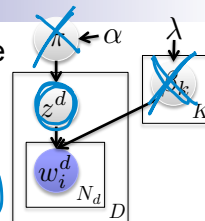


Collapsed Gibbs Sampler

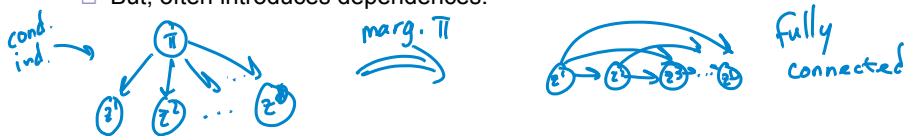
- In conjugate models, can analytically marginalize some variables and only sample remaining

for $d=1, \dots, D$

$$z^d \sim p(z^d | z^1, \dots, z^{d-1}, z^{d+1}, \dots, z^D, \{w_i^d\}, \lambda, \alpha)$$



- Can improve efficiency if marginalized variables are high-dim
 - Reduced dimension of search space
 - But, often introduces dependencies!



Collapsed Sampler Full Conditional

$$p(\cdot) = \int_{\pi} p(\pi | \alpha) \prod_{d=1}^D p(z^d | \pi) \left(\prod_{k=1}^K p(\beta_k | \lambda) \prod_{d=1}^D \prod_{i=1}^{N_d} p(w_i^d | z^d, \beta_k) \right) d\pi$$

$$p(z^d = k | z_{\setminus d}, \{w_i^d\}, \alpha, \lambda) \propto \int_{\pi} \int_{\beta_1} \dots \int_{\beta_K} p(\cdot)$$

■ Derivation

$$\propto p(z^d = k | z_{\setminus d}, \alpha) p(\{w_i^d\} | \{w_i^c : z^c = k, c \neq d\})$$

"prior" "likelihood"

©Emily Fox 2014

17

Collapsed Sampler Full Conditional

$$p(z^d = k | z_{\setminus d}, \{w_i^d\}, \alpha, \lambda) \propto p(z^d = k | z_{\setminus d}, \alpha) p(\{w_i^d\} | \{w_i^c : z^c = k, c \neq d\})$$

$$p(z^d = k | z_{\setminus d}, \alpha) = \frac{N_k^d + \alpha_k}{D - 1 + \sum \alpha_k}$$

$$p(w_i^d | z^d = k) = \frac{m_{i,k}^d + \lambda_{i,k}}{\sum_j (m_{i,j}^d + \lambda_{i,j})}$$

of docs assigned to topic k not counting doc d
 # of words assigned to topic k taking value of w_i (not inc doc d)

©Emily Fox 2014

18

Collapsed Sampler Intuition (MoG)

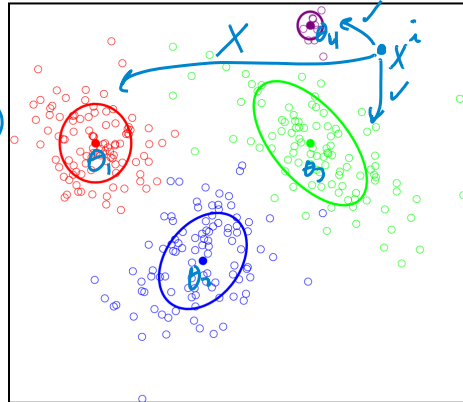
Previously, $p(z^i = k | x^i, \pi, \theta) \propto \pi_k p(x^i | \theta_k)$

If you're not told π, θ_k

$\{\mu_k, \Sigma_k\}$

"prior" Approx π by counts of occupancy of each cluster (plus prior)

"like" Approx θ_k based on obs. already assigned to cluster k



©Emily Fox 2014

19

Example – Uncollapsed Results

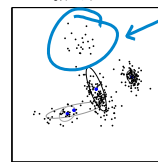
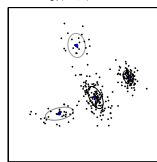
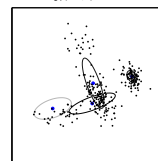
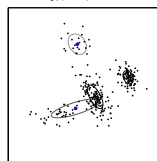
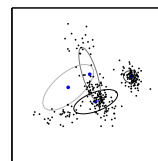
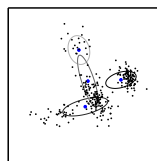
sampling $\pi, \{\theta_k\}$
 $\{z^i\}$

one init

t=2

t=10

t=50



given z^1, \dots, z^N
low post. prob. of drawing θ_k here

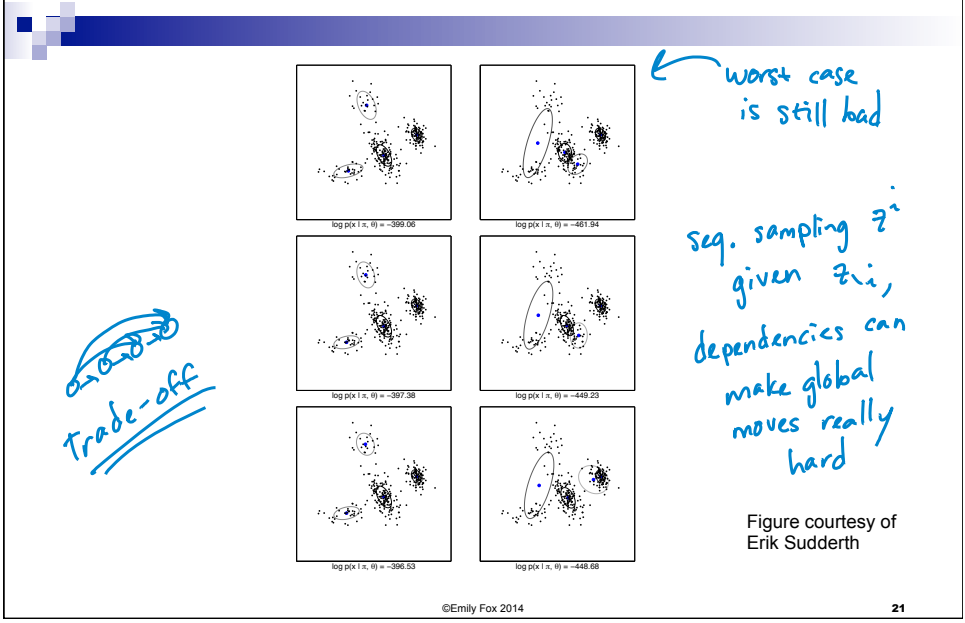
Figure courtesy of Erik Sudderth

will eventually happen, but maybe not in lifetime

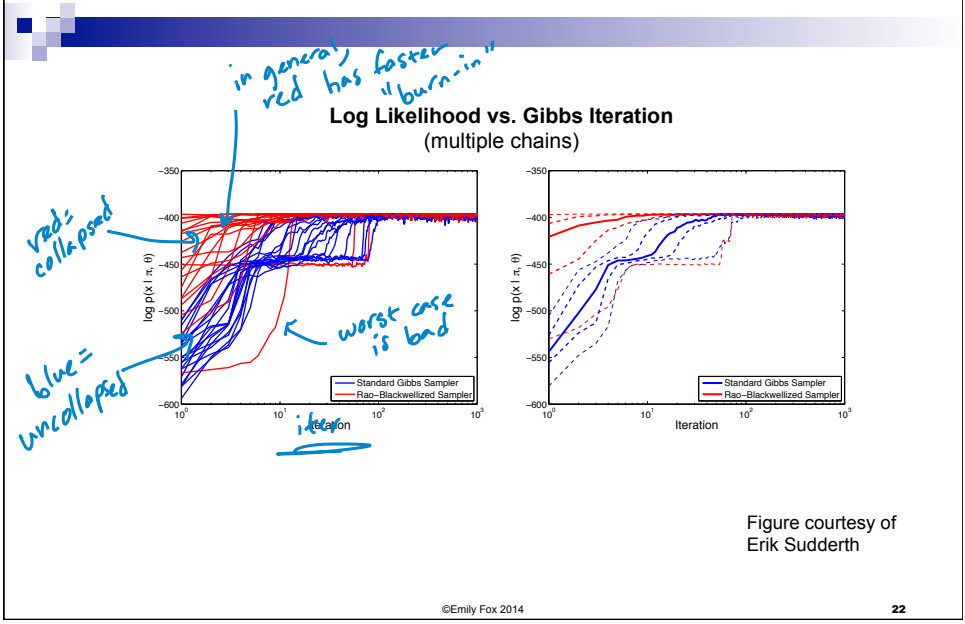
©Emily Fox 2014

20

Example – Collapsed Results

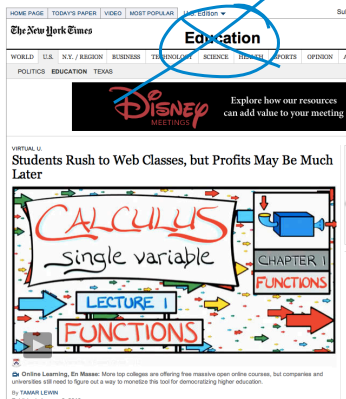


Comparing Collapsed vs. Uncollapsed



Task 3: Mixed Membership Models

- **Now:** Document may belong to multiple clusters



©Emily Fox 2014

23

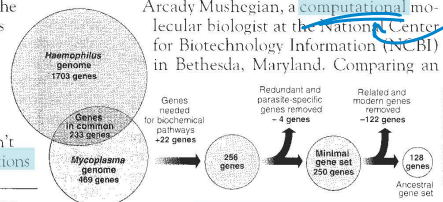
Latent Dirichlet Allocation (LDA)

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

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

"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



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

words from different topics

©Emily Fox 2014

24

Latent Dirichlet Allocation (LDA)

each topic k is a distribution over words in vocab, β_k , just as before β_k

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

Topic proportions and assignments

every word is assigned to a topic

each doc has its own prevalence of topics in that doc

©Emily Fox 2014 25

Latent Dirichlet Allocation (LDA)

Topics

Documents

Topic proportions and assignments

All we see are words β_k 's

Want: posterior $p(\text{topics}, \text{doc prop. of topics}, \text{assign. vars.} \mid \text{words in docs})$

©Emily Fox 2014 26

LDA Generative Model

- Observations: $w_1^d, \dots, w_{N_d}^d$
- Associated topics: $z_1^d, \dots, z_{N_d}^d$ ← *topic per word in doc d*
- Parameters: $\theta = \{\{\pi^d\}, \{\beta_k\}\}$
- Generative model:

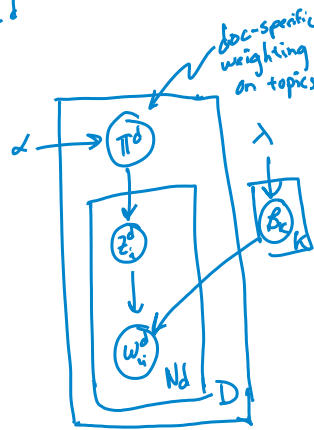
$$z_i^d \sim \pi^d \quad d=1, \dots, D \\ i=1, \dots, N_d$$

$$w_i^d | z_i^d \sim \beta_{z_i^d}$$

priors:

$$\pi^d \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad d=1, \dots, D$$

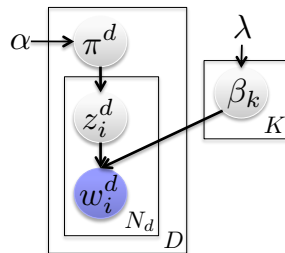
$$\beta_k \sim \text{Dir}(\lambda_1, \dots, \lambda_V) \quad k=1, \dots, K$$



©Emily Fox 2014

27

LDA Joint Probability



$$p(\cdot) = \prod_{k=1}^K p(\beta_k | \lambda) \prod_{d=1}^D p(\pi^d | \alpha) \left(\prod_{i=1}^{N_d} p(z_i^d | \pi^d) p(w_i^d | z_i^d, \beta) \right)$$

©Emily Fox 2014

28

Example Inference – Topic Weights

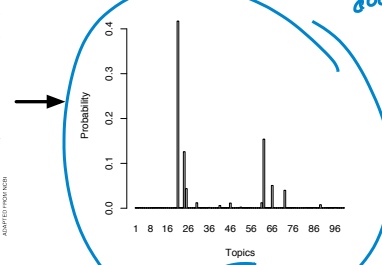
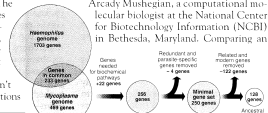
- **Data:** The OCR'ed collection of *Science* from 1990-2000
 - 17K documents
 - 11M words
 - 20K unique terms (stop words and rare words removed)
- **Model:** 100-topic LDA model

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 120 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

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

SCIENCE • VOL. 272 • 24 MAY 1996



one doc

$\pi d =$
topic weights
for this
doc

Example Inference – Topic Words

	topic1	topic2	...	
highest prob. words associated w/ each topic	human	evolution	disease	computer
	genome	evolutionary	host	models
	dna	species	bacteria	information
	genetic	organisms	diseases	data
	genes	life	resistance	computers
	sequence	origin	bacterial	system
	gene	biology	new	network
	molecular	groups	strains	systems
	sequencing	phylogenetic	control	model
	map	living	infectious	parallel
	information	diversity	malaria	methods
	genetics	group	parasite	networks
	mapping	new	parasites	software
	project	two	united	new
	sequences	common	tuberculosis	simulations

Shared throughout corpus

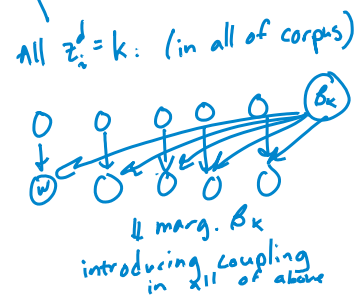
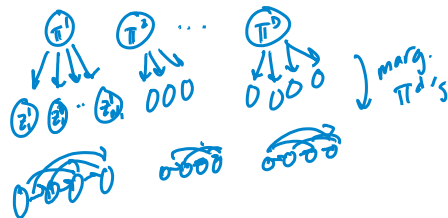
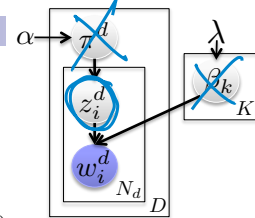
Collapsed LDA Sampling

- Marginalize parameters
 - Document-specific topic weights
 - Corpus-wide topic-specific word distributions

$$p(z_i^d = k \mid z_{\setminus id}, \{w_i^d\}, \alpha, \lambda)$$

$$\propto p(z_i^d = k \mid z_{\setminus id}, \alpha) p(w_i^d \mid z_i^d = k, z_{\setminus id}, w_{\setminus id}, \lambda)$$

- Unplate to see dependencies induced



©Emily Fox 2014

31

What you need to know...

- Bayesian specification of document clustering model
- Rules of conditional and unconditional independence in directed graphical models (Bayes nets)
 - Bayes' ball
 - Markov blanket
- Gibbs sampling for Bayesian document model
- Latent Dirichlet allocation (LDA)
 - Motivation and generative model specification
 - Collapsed Gibbs sampler

©Emily Fox 2014

32

Reading

■ Mixed Membership Models: KM Sec. 27.3

- Basic LDA:
[Blei, David M., Andrew Y. Ng, and Michael I. Jordan. "Latent dirichlet allocation." the Journal of machine Learning research 3 \(2003\): 993-1022.](#)
- Introduction:
[Blei, David M. "Probabilistic topic models." Communications of the ACM, vol. 55, no. 4 \(2012\): 77-84.](#)
- Sampling:
[Griffith, Thomas L. and Mark Steyvers. "Finding scientific topics." Proceedings of the National Academy of Sciences of the United States of America, Volume: 101, Supplement: 1 \(2004\): Pages: 5228-5235](#)

Acknowledgements

- Thanks to Dave Blei for some material in this lecture relating to LDA