

Case Study 5: Mixed Membership Modeling

Clustering Documents Revisited, Latent Dirichlet Allocation

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

Emily Fox

May 26th, 2015

©Emily Fox 2015

1

Task 2: Cluster Documents

■ Then examined:

- Cluster documents based on topic



sports



world news

©Emily Fox 2015

2

A Generative Model

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

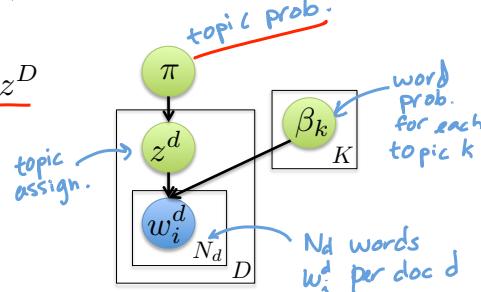
Sample topic:
 $z^d \sim \pi$

Sample words:

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

topic of doc d

Given topic $z^d=k$ for doc d,
draw each word from β_k



©Emily Fox 2015

3

Bayesian Document Model

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

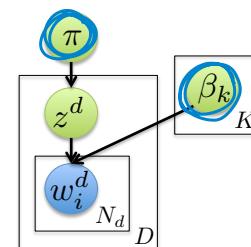
- Bayesian approach

*place priors on parameters

- Need distribution on pmf's

$$\sum_{k=1}^K \pi_k = 1$$

$$\sum_{v=1}^V \beta_{kv} = 1$$



What is a distribution on the simplex?
First, what is the simplex?

©Emily Fox 2015

4

Dirichlet Distributions

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

The Dirichlet distribution is defined on the simplex

$\alpha_k = 10 \forall k$

$\alpha_k = 0.1 \forall k$

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 2015

Model Summary

- Prior on model parameters
 - E.g., symmetric Dirichlet for π chose sym. Dir.

$\pi \sim \text{Dir}\left(\frac{\alpha}{K}, \dots, \frac{\alpha}{K}\right)$

Dirichlet prior for topic parameters $\beta_k \sim \text{Dir}\left(\frac{\lambda}{V}, \dots, \frac{\lambda}{V}\right)$

Sample observations as

$$z^d \sim \pi \quad d=1, \dots, D$$

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

©Emily Fox 2015

Posterior Inference via Sampling

- Iterate between sampling

$$\pi \sim p(\pi | \{z^d\}, \{\beta_k\}, \{w_i^d\})$$

"full conditionals"

For $k=1, \dots, K$

$$\beta_k \sim p(\beta_k | \pi, \{z^d\}, \{w_i^d\})$$

For $d=1, \dots, D$

$$z^d \sim p(z^d | \pi, \{\beta_k\}, \{w_i^d\})$$

- What form do these complete conditionals take?

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

©Emily Fox 2015

7

Markov Blanket

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

$$\begin{aligned} \text{Markov} \\ \text{blanket} &= -\text{all parents} \\ &\text{of} \\ &- \text{all children} \\ &- \text{all coparents} \end{aligned}$$

- Gibbs sampling iterates between full conditionals

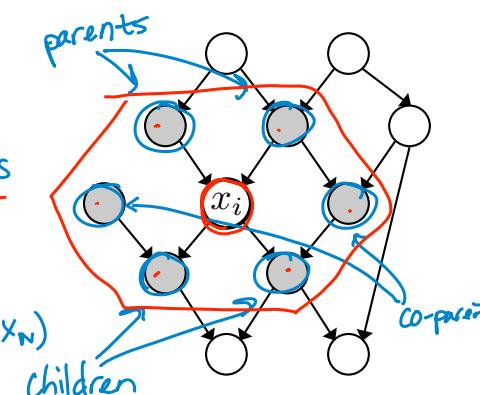
$$x_i \sim p(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_N)$$

→ simplify to

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

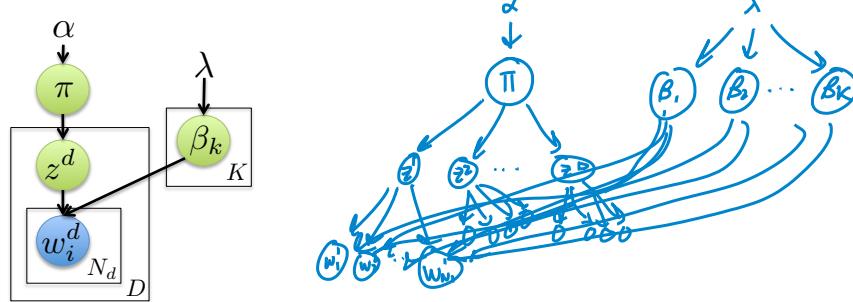
©Emily Fox 2015

8



Unplated Document Model

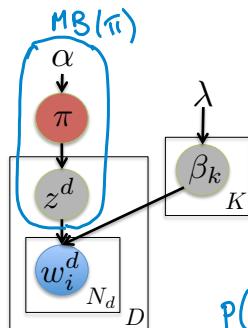
- Recall that the plate notation is really indicating



©Emily Fox 2015

9

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

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

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

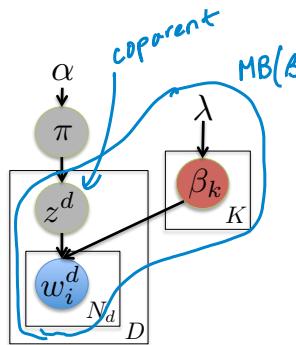
$$\propto \prod_k \pi_k^{N_k} \pi_k^{\alpha_k - 1} = \prod_k \pi_k^{N_k + \alpha_k - 1} = \text{Dir}(N_k + \alpha_k)$$

- Conjugacy: Posterior has same form as prior

©Emily Fox 2015

10

Complete Conditional for β_k



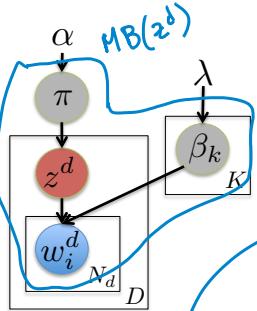
- Again, Dirichlet prior
 $\beta_k \sim \text{Dir}(\lambda_1, \dots, \lambda_V)$
- Consider docs d such that $z^d = k$
 - For these observations, $w_i^d \sim \beta_k$
 - Do any other docs depend on β_k ? No
- Then,
Count $m_{V,k} = |\{w_i^d : w_i^d = V, \forall d \text{ s.t. } z^d = k\}|$
 $\beta_k \sim \text{Dir}(m_{V,k} + \lambda_1, \dots, m_{V,k} + \lambda_V)$

□ Again, posterior has same form as prior

©Emily Fox 2015

11

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

r_{dk} is a discrete r.v.
- Sample each cluster indicator as
 $r_d = [r_{d1}, \dots, r_{dK}]$
 $z^d \sim \underline{r}_d$

©Emily Fox 2015

$$\frac{\prod_{i=1}^{N_d} \beta_{k_i, w_i}}{\prod_{j=1}^K \prod_{i=1}^{N_d} \beta_{j, w_i}}$$

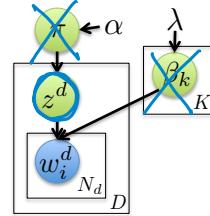
↑↑↑↑↑↑
1 2 ... K
responsibilities₁₂

Collapsed Gibbs Sampler

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

For $d=1, \dots, D$

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



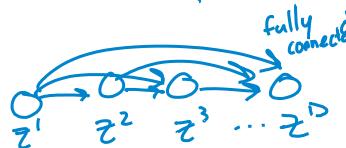
- Can improve efficiency if marginalized variables are high-dim

- Reduced dimension of search space
- But, often introduces dependences!

opportunities for parallelism
are lost



marg. PI



©Emily Fox 2015

13

Collapsed Sampler Full Conditional

full joint dist.:

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

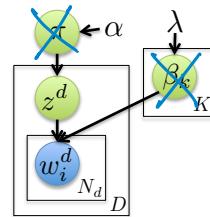
Derivation

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

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

depends on # of $z^c=k$ (N_k)

words in docs assigned to topic k not inc. d



©Emily Fox 2015

14

Collapsed Sampler Intuition (MoG)

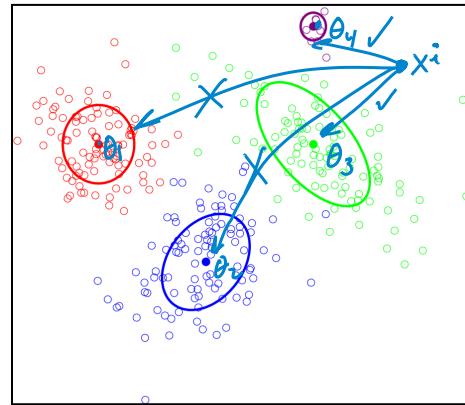
- Previously, $p(z^i = k | x^i, \pi, \theta) \propto \pi_k p(x^i | \theta_k)$

$\uparrow \{m_k, \Sigma_k\}$

- If you're not told π, θ_k

"prior" Approx. π by counts associated with each cluster

"likelihood" Approx. θ_k by obs. already assigned to those clusters

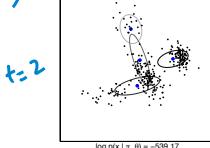


©Emily Fox 2015

16

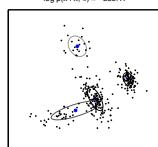
Example – Uncollapsed Results

one init. \rightarrow



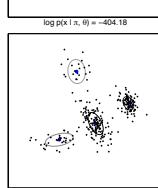
$\log p(x | z, \theta) = -539.17$

$t=10$



$\log p(x | z, \theta) = -404.18$

$t=50$



$\log p(x | z, \theta) = -397.40$

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

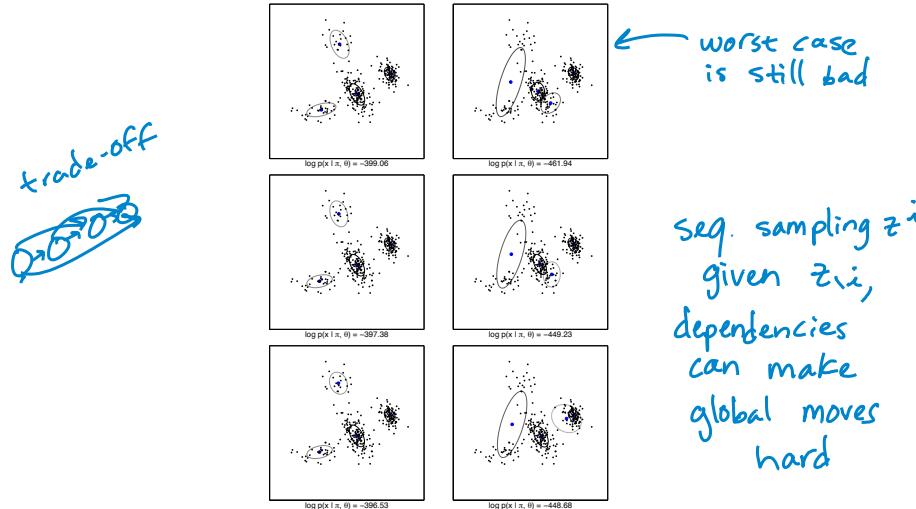
given z^1, \dots, z^n
low post. prob.
of drawing θ_k
here

will eventually
happen, but maybe
not in our lifetime

Figure courtesy of Erik Sudderth

©Emily Fox 2015

Example – Collapsed Results



Comparing Collapsed vs. Uncollapsed

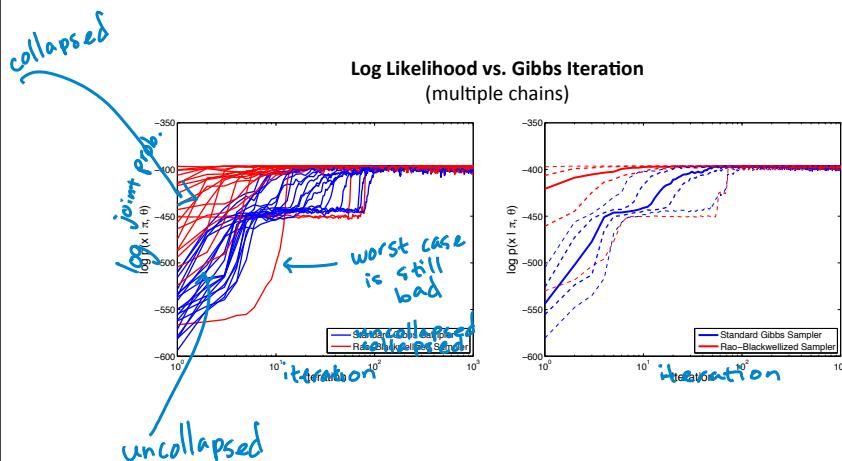


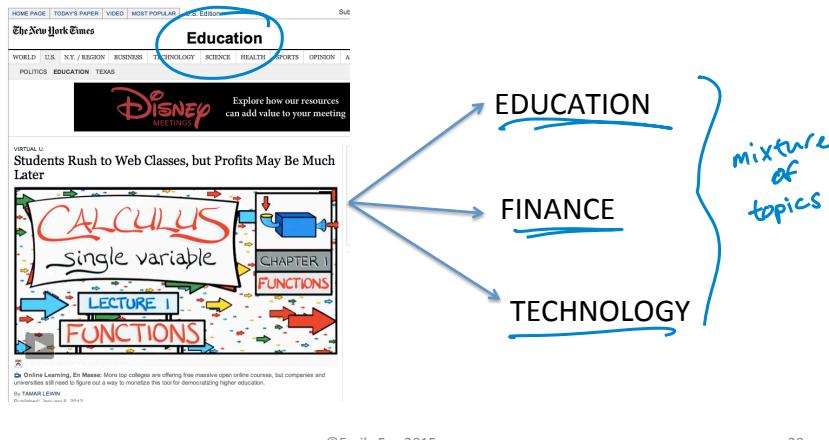
Figure courtesy of
Erik Sudderth

©Emily Fox 2015

19

Task 3: Mixed Membership Models

- Now: Document may belong to multiple clusters



©Emily Fox 2015

20

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.

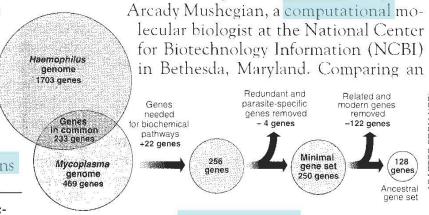
Although the numbers don't match precisely, those predictions

“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 newly sequenced genome,” explains

Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

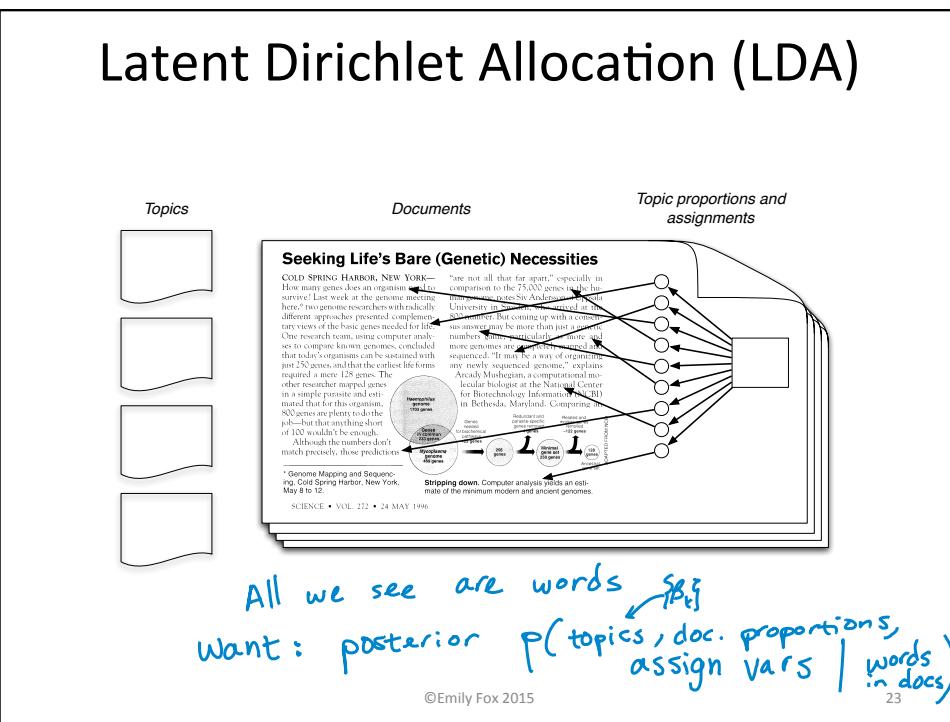
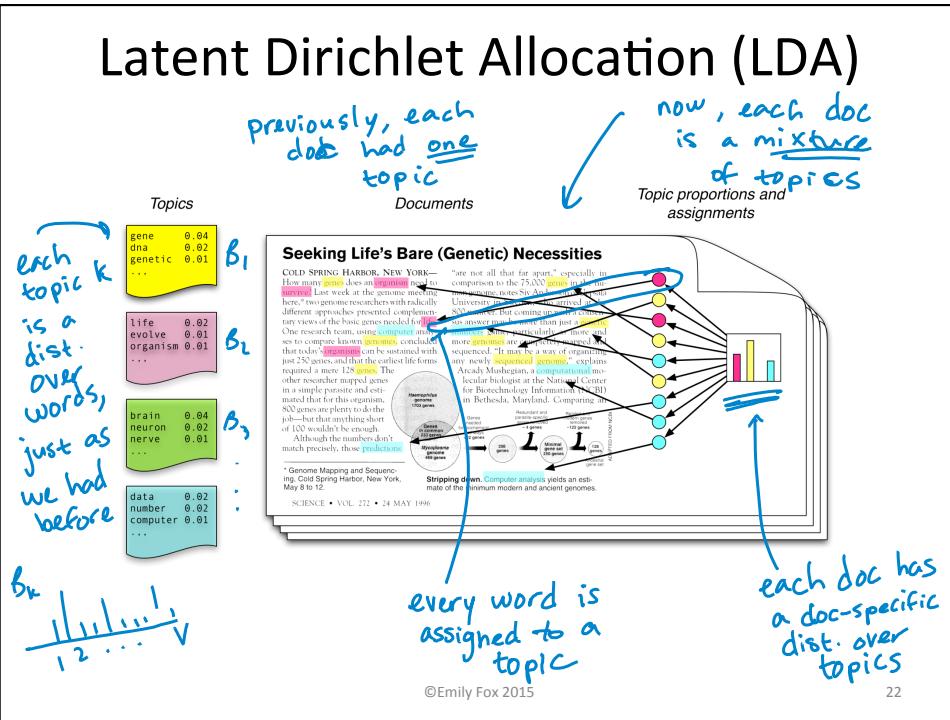
words from different topics

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



©Emily Fox 2015

21



LDA Generative Model

- Observations: $w_1^d, \dots, w_{N_d}^d \quad d=1, \dots, D$
- Associated topics: $z_1^d, \dots, z_{N_d}^d \quad \leftarrow \text{assign var. per word}$
- Parameters: $\theta = \{\{\pi^d\}, \{\beta_k\}\}$
- Generative model: $\begin{array}{l} \text{doc-specific topic weights} \\ \text{corpus-wide topic priors} \end{array}$

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

$$w_i^d | z_i^d \sim \beta_{z_i^d} \quad i=1, \dots, N_d$$

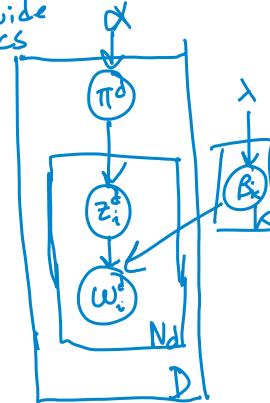
Priors:

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

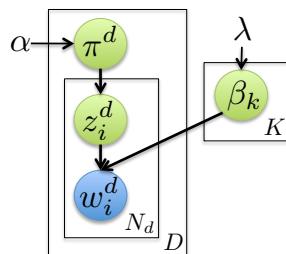
$$\beta_k \sim \text{Dir}(\lambda_1, \dots, \lambda_V) \quad k=1, \dots, K$$

©Emily Fox 2015

24



LDA Joint Probability



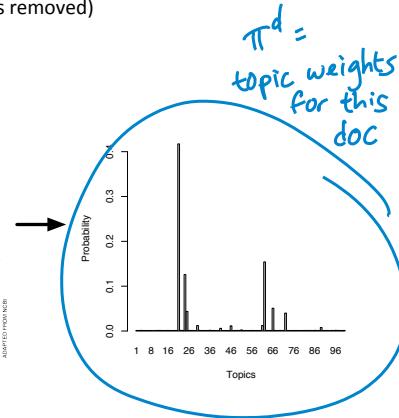
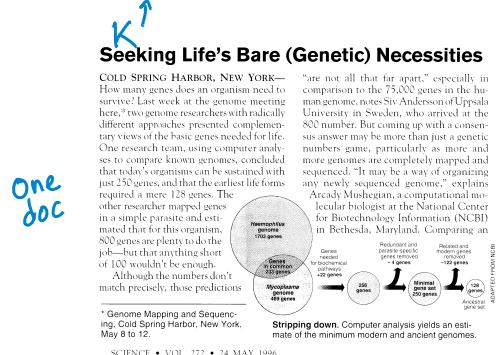
$$p(\cdot) = \prod_{k=1}^K p(\beta_k | \lambda) \prod_{d=1}^D p(\pi^d | \alpha) \left(\prod_{i=1}^{N_d} p(z_i^d | \pi^d) p(w_i^d | z_i^d, \beta) \right)$$

©Emily Fox 2015

25

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



©Emily Fox 2015

26

Example Inference – Topic Words

highest prob. words under each topic

topic 1	topic 2	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

β_1 β_2 β_3 β_4

annotated afterwards

©Emily Fox 2015

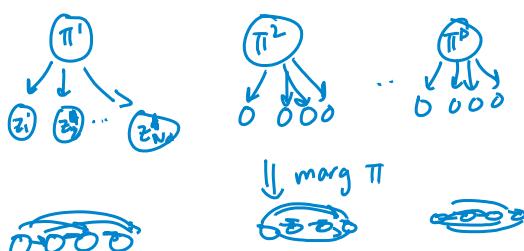
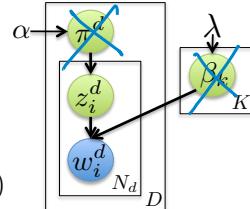
27

Collapsed LDA Sampling

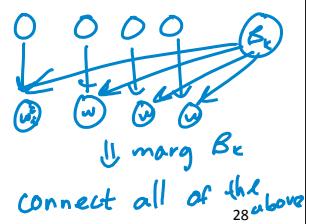
- Marginalize parameters
 - Document-specific topic weights
 - Corpus-wide topic-specific word distributions

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

- Unplate to see dependencies induced



All $z_i^d = k$: (in all of corpus)



©Emily Fox 2015

28

Collapsed LDA Sampling

- Sample topic indicators for each word
 - Algorithm:

$$p(z_i^d = k | z_{\setminus id}, \{w_i^d\}, \alpha, \lambda) \\ \propto p(z_i^d = k | \{z_j^d, j \neq i\}, \alpha) p(w_i^d | \{w_j^c : z_j^c = k, (j, c) \neq (i, d)\}, \lambda)$$

"prior"

"likelihood"

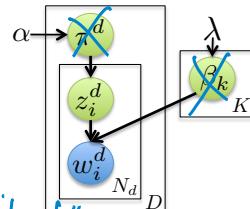
$$\times \frac{n_{k,d} + \delta_k}{N_d - 1 + \sum d_k} \cdot \frac{m_{w_i^d, k} + \lambda_{w_i^d}}{\sum_\gamma m_{w_i^\gamma, k} + \lambda_\gamma}$$

of words assigned to topic k in doc d (not counting i-th word)

normalize within doc

examine entire corpus

* of times word r appears in topic k (not w_i^d)



©Emily Fox 2015

29

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 2015

30

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](#)

©Emily Fox 2015

31

Acknowledgements

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

©Emily Fox 2015

32