

Case Study 2: Document Retrieval

Collapsed Gibbs and Variational Methods for LDA

Machine Learning/Statistics for Big Data
CSE599C1/STAT592, University of Washington

Emily Fox

February 7th, 2013

©Emily Fox 2013

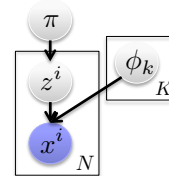
1

Example – Collapsed MoG Sampling

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad z^i \sim \pi$$

$$\{\mu_k, \Sigma_k\} \sim F(\phi) \quad x^i | z^i \sim N(x^i; \mu_{z^i}, \Sigma_{z^i})$$

- Collapsed sampler



For $i=1, \dots, N$

$$z^i(t) \sim p(z^i | z^{1(t)}, \dots, z^{i-1(t)}, z^{i+1(t)}, \dots, z^{N(t)}, x_{1:N}, \alpha, 1)$$

$z_{-i}^{(t)}$



©Emily Fox 2013

2

Example – Collapsed MoG Sampling

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad z^i \sim \pi$$

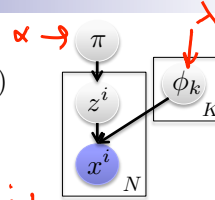
$$\{ \mu_k, \Sigma_k \} \sim F(\phi) \quad x^i | z^i \sim N(x^i; \mu_{z^i}, \Sigma_{z^i})$$

Derivation

$$p(z^i | z_{-i}, x_{1:N}, \alpha, \lambda) \propto p(z^i | z_{-i}, \alpha) p(x^i | z^i, z_{-i}, x_{-i}, \lambda)$$

$$p(z^i = k | z_{-i}, \alpha) = \int p(z^i = k | \pi) p(\pi | z_{-i}, \alpha) d\pi = \frac{n_k^i + \alpha_k}{N - 1 + \sum \alpha_k}$$

$$p(x^i | z_{1:N}, x_{-i}, \lambda) = \text{student-t} \quad \text{Dir post.} \quad \text{pred. likelihood}$$



Important facts:

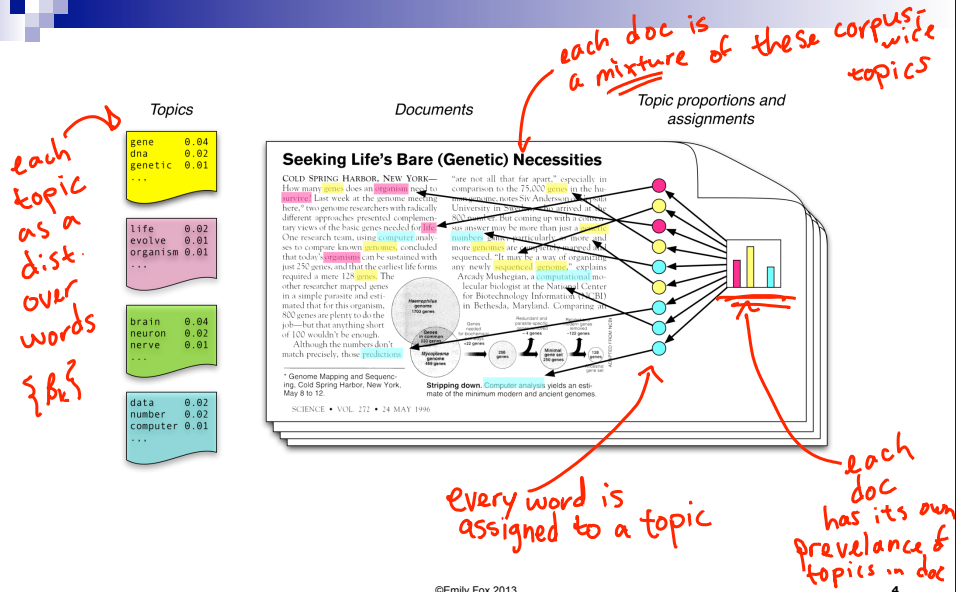
$$p(z_{1:N} | \alpha) = \frac{\Gamma(\sum_k \alpha_k) \prod_k \Gamma(n_k + \alpha_k)}{\prod_k \Gamma(\alpha_k) \Gamma(\sum_k n_k + \alpha_k)}$$

$$\frac{\Gamma(m+1)}{\Gamma(m)} = m$$

©Emily Fox 2013

3

Latent Dirichlet Allocation (LDA)



©Emily Fox 2013

4

LDA Generative Model

- Observations: $w_1^d, \dots, w_{N_d}^d$ $d=1, \dots, D$
- Associated topics: $z_1^d, \dots, z_{N_d}^d$ $d=1, \dots, D$ *corpus-wide topic*
- Parameters: $\theta = \{\{\pi^d\}, \{\beta_k\}\}$
- Generative model: *doc-specific preferences of topics*

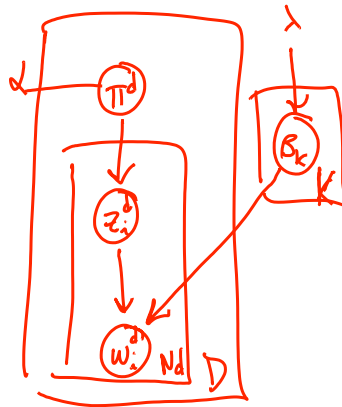
$$z_i^d \sim \pi^d \quad d=1, \dots, D$$

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

Priors:

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

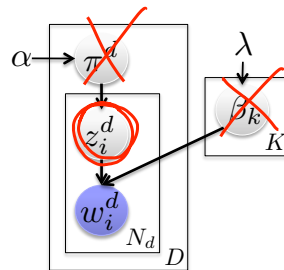
$$\beta_k \sim \text{Dir}(\lambda_1, \dots, \lambda_V) \quad k=1, \dots, K$$



©Emily Fox 2013

5

LDA Generative Model



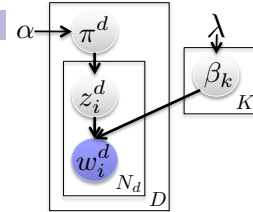
$$p(\cdot) = \prod_{k=1}^K p(\beta_k | \lambda) \prod_{d=1}^D p(\pi^d | \alpha) \left(\prod_{i=1}^{N_d} \underline{p(z_i^d | \pi^d)} \underline{p(w_i^d | z_i^d, \beta)} \right)$$

©Emily Fox 2013

6

Collapsed LDA Sampling

- Marginalize parameters
 - Document-specific topic weights
 - Corpus-wide topic-specific word distributions
- Sample topic indicators for each word
 - Derivation:



$$p(z_{1:N_d}^d | \alpha) = \frac{\Gamma(\sum_k \alpha_k) \prod_k \Gamma(n_k^d + \alpha_k)}{\prod_k \Gamma(\alpha_k) \Gamma(\sum_k n_k^d + \alpha_k)}$$

$$p(\{w_i^d | z_i^d = k\}, \lambda) = \frac{\Gamma(\sum_\nu \lambda_\nu) \prod_\nu \Gamma(v_\nu^k + \lambda_\nu)}{\prod_\nu \Gamma(\lambda_\nu) \Gamma(\sum_\nu v_\nu^k + \lambda_\nu)}$$

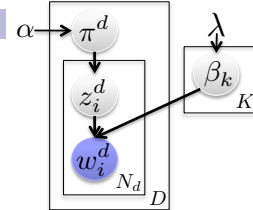
$$p(z | \alpha) = \prod_{d=1}^D p(z_{1:N_d}^d | \alpha) \quad p(w | z, \lambda) = \prod_{k=1}^K p(\{w_i^d | z_i^d = k\}, \lambda)$$

©Emily Fox 2013

7

Collapsed LDA Sampling

- Marginalize parameters
 - Document-specific topic weights
 - Corpus-wide topic-specific word distributions
- Sample topic indicators for each word
 - Algorithm:



©Emily Fox 2013

8

Sample Document

Etruscan	trade	price	temple	market

©Emily Fox 2013

9

Randomly Assign Topics

z_i^d	3	2	1	3	1
w_i^d	Etruscan	trade	price	temple	market

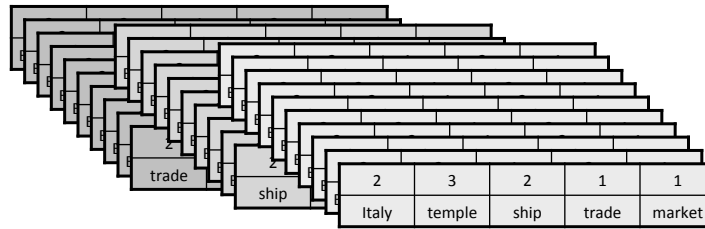
©Emily Fox 2013

10

Randomly Assign Topics

z_i^d
 w_i^d

3	2	1	3	1
Etruscan	trade	price	temple	market



©Emily Fox 2013

11

Maintain Global Statistics

z_i^d
 w_i^d

3	2	1	3	1
Etruscan	trade	price	temple	market

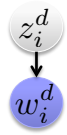
Total counts from all docs

	1	2	3
Etruscan	1	0	35
market	50	0	1
price	42	1	0
temple	0	0	20
trade	10	8	1
...			

©Emily Fox 2013

12

Resample Assignments



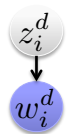
3	2	1	3	1
Etruscan	trade	price	temple	market

	1	2	3
Etruscan	1	0	35
market	50	0	1
price	42	1	0
temple	0	0	20
trade	10	8	1
...			

©Emily Fox 2013

13

What is the conditional distribution for this topic?



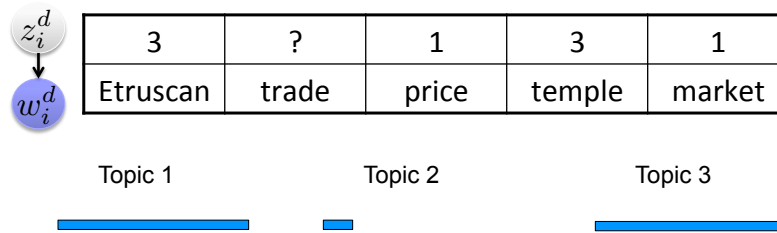
3	?	1	3	1
Etruscan	trade	price	temple	market

©Emily Fox 2013

14

What is the conditional distribution for this topic?

- Part I: How much does this document like each topic?

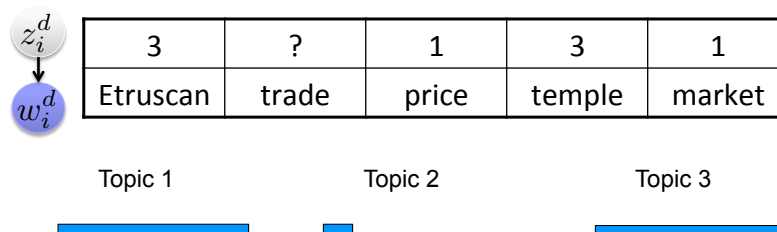


©Emily Fox 2013

15

What is the conditional distribution for this topic?

- Part I: How much does this document like each topic?
- Part II: How much does each topic like this word?



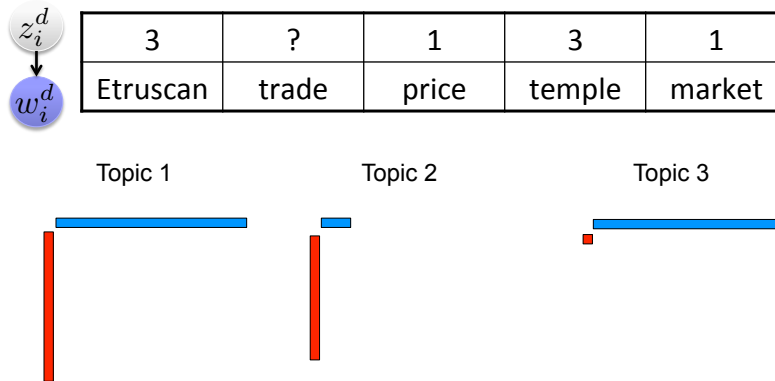
	1	2	3
trade	10	7	1

©Emily Fox 2013

16

What is the conditional distribution for this topic?

- Part I: How much does this document like each topic?
- Part II: How much does each topic like this word?

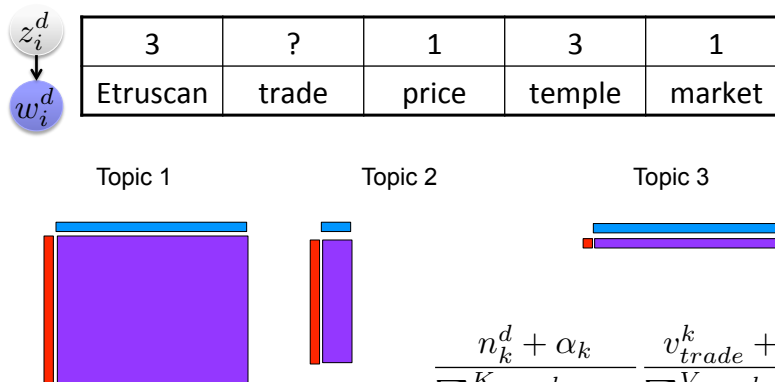


©Emily Fox 2013

17

What is the conditional distribution for this topic?

- Part I: How much does this document like each topic?
- Part II: How much does each topic like this word?

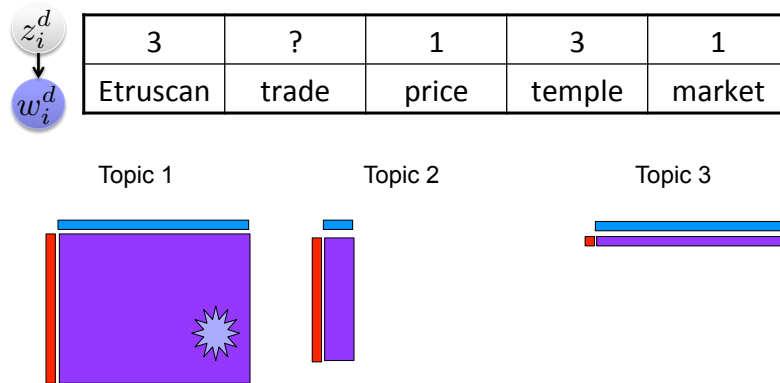


$$\frac{n_k^d + \alpha_k}{\sum_{j=1}^K n_j^d + \alpha_j} \frac{v_{trade}^k + \lambda_k}{\sum_{j=1}^V v_j^k + \lambda_j}$$

©Emily Fox 2013

18

Sample a New Topic Indicator



©Emily Fox 2013

19

Update Counts

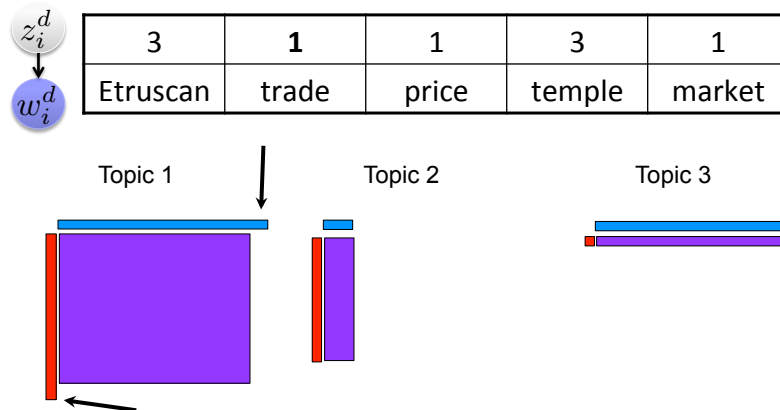
Diagram illustrating the update of counts for the topic indicator. The latent variable z_i^d is sampled from the distribution w_i^d . The resulting topic indicator is used to update the counts for each word in the topic-specific distribution. The diagram shows three topics: Topic 1 (Etruscan, trade, price, temple, market), Topic 2 (Etruscan, trade, price, temple, market), and Topic 3 (Etruscan, trade, price, temple, market).

	1	2	3
Etruscan	1	0	35
market	50	0	1
price	42	1	0
temple	0	0	20
trade	10	7	1
...			

©Emily Fox 2013

20

Geometrically...



©Emily Fox 2013

21

Issues with Generic LDA Sampling

- Slow mixing rates → Need many iterations
- Each iteration cycles through sampling topic assignments for *all* words in *all* documents
- Modern approaches:
 - Large-scale LDA. For example, [Mimno, David, Matthew D. Hoffman and David M. Blei. "Sparse stochastic inference for latent Dirichlet allocation." International Conference on Machine Learning. 2012.](#)
 - Distributed LDA. For example, [Ahmed, Amr, et al. "Scalable inference in latent variable models." Proceedings of the fifth ACM international conference on Web search and data mining \(2012\): 123-132](#)
- Alternative: Variational methods instead of sampling
 - Approximate posterior with an optimized variational distribution

©Emily Fox 2013

22

Variational Methods

- Recall task: Characterize the posterior
- Turn posterior inference into an optimization task
- Introduce a “tractable” family of distributions over parameters and latent variables
 - Family is indexed by a set of “free parameters”
 - Find member of the family closest to:
- Questions:
 - How do we measure “closeness”?
 - If the posterior is intractable, how can we approximate something we do not have to begin with?

©Emily Fox 2013

23

A Measure of Closeness

- Kullback-Leibler (KL) divergence
 - Measures “distance” between two distributions p and q
- Not symmetric
- p determines where the difference is important:
 - $p(x)=0$ and $q(x)\neq 0$
 - $p(x)\neq 0$ and $q(x)=0$
- Want
- Just as hard as the original problem!

©Emily Fox 2013

24

Reverse Divergence

- Divergence $D(q || p)$
 - true distribution p defines support of diff.
 - the “correct” direction
 - will be intractable to compute
- Reverse divergence $D(q || p)$
 - approximate distribution defines support
 - tends to give overconfident results
 - will be tractable

©Emily Fox 2013

25

Interpretations of Minimizing Reverse KL

- Similarity measure:
- Evidence lower bound (ELBO)
- Therefore, minimizing KL is equivalent to maximizing a lower bound on the marginal likelihood:
 - $\text{Max } \mathcal{L}(q) = \min D(q||p) = \text{max lower bound of } \log p(x)$

©Emily Fox 2013

26

Mean Field

- How do we choose a Q such that the following is tractable?
- Simplest case = mean field approximation
 - Assume each parameter and latent variable is conditionally independent given the set of free parameters
- Then, entropy term decomposes as

©Emily Fox 2013

27

Mean Field

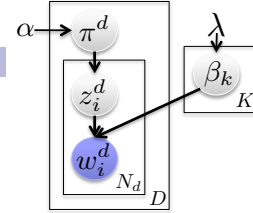
- Examine one free parameter, e.g., γ
 - Can rewrite joint as
$$E_q[\log p(\theta, z, x)] = E_q[\log p(\theta | z, x)] + E_q[\log p(z, x)]$$
 - Look at terms of ELBO just depending on γ
$$\mathcal{L}^\gamma =$$
- Likewise,
$$\mathcal{L}^{\phi^n} =$$
- This motivates using a coordinate ascent algorithm for optimization
 - Iteratively optimize each free parameter holding all others fixed

©Emily Fox 2013

28

Mean Field for LDA

- In LDA, our parameters are $\theta = \{\pi^d\}, \{\beta_k\}$
 $z = \{z_i^d\}$



- The variational distribution factorizes as

- The joint distribution factorizes as

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

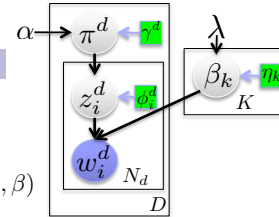
©Emily Fox 2013

29

Mean Field for LDA

$$q(\pi, \beta, z) = \prod_{k=1}^K q(\beta_k | \eta_k) \prod_{d=1}^D q(\pi^d | \gamma^d) \prod_{i=1}^{N_d} q(z_i^d | \phi_i^d)$$

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



- Examine the ELBO

$$\begin{aligned} \mathcal{L}(q) = & \sum_{k=1}^K E_q[\log p(\beta_k | \lambda)] + \sum_{d=1}^D E_q[\log p(\pi^d | \alpha)] \\ & + \sum_{d=1}^D \sum_{i=1}^{N_d} E_q[\log p(z_i^d | \pi^d)] + E_q[\log p(w_i^d | z_i^d, \beta)] \\ & - \sum_{k=1}^K E_q[\log q(\beta_k | \eta_k)] - \sum_{d=1}^D E_q[\log q(\pi^d | \gamma^d)] - \sum_{d=1}^D \sum_{i=1}^{N_d} E_q[\log q(z_i^d | \phi_i^d)] \end{aligned}$$

©Emily Fox 2013

30

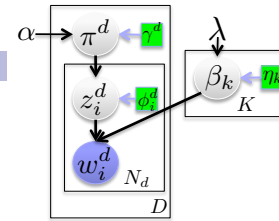
Mean Field for LDA

- Let's look at some of these terms

$$E_q[\log p(z_i^d \mid \pi^d)]$$

$$E_q[\log q(z_i^d \mid \phi_i^d)]$$

- Other terms follow similarly

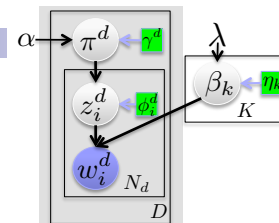


©Emily Fox 2013

31

Optimize via Coordinate Ascent

- Algorithm:

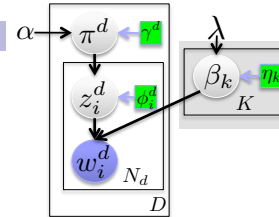


©Emily Fox 2013

32

Optimize via Coordinate Ascent

- Algorithm:



©Emily Fox 2013

33

Alternative Optimization Schemes

- Inefficient:

- Start from randomly initialized η_k (topics)
- Analyze whole corpus before updating η_k again
- If streaming data scenario, can't compute even one iteration!

- Didn't have to do coord. ascent. Could have used gradient ascent.

©Emily Fox 2013

34

Alternative Optimization Schemes

- Recall stochastic gradient ascent:

- Assume $M = 1$

- Unbiased, but noisy

- Here,

$$\mathcal{L} = E_q[\log p(\beta)] - E_q[\log q(\beta)] + \sum_{d=1}^D E_q[\log p(\pi^d)] - E_q[\log q(\pi^d)] \\ + \sum_{d=1}^D E_q[\log p(z^d, x^d | \pi^d, \beta)] - E_q[\log q(z^d)]$$

$$\mathcal{L}_t = E_q[\log p(\beta)] - E_q[\log q(\beta)] + D (E_q[\log p(\pi^t)] - E[\log q(\pi^t)]) \\ + D (E_q[\log p(z^t, x^t | \pi^t, \beta)] - E_q[\log q(z^t)])$$

©Emily Fox 2013

35

Stochastic Variational Inference for LDA

- Initialize $\eta^{(0)}$ randomly.
- Repeat (indefinitely):
 - Sample a document d uniformly from the data set.
 - For all k , initialize $\gamma_k^d = 1$
 - Repeat until converged

- For $i=1, \dots, N_d$

$$\phi_{ik}^d \propto \exp\{E[\log \pi_k^d] + E[\log \beta_{k,w_i^d}]\}$$

- Set $\gamma^d = \alpha + \sum_{i=1}^{N_d} \phi_i^d$

- Take a stochastic gradient step $\eta^{(t)} = \eta^{(t-1)} + \rho_t \nabla_{\eta} \mathcal{L}_d$

©Emily Fox 2013

36

Acknowledgements

- Thanks to Dave Blei, David Mimno, and Jordan Boyd-Graber for some material in this lecture relating to LDA