

Case Study 5: Mixed Membership Modeling

LDA Collapsed Gibbs Sampler, Variational Inference

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

Emily Fox

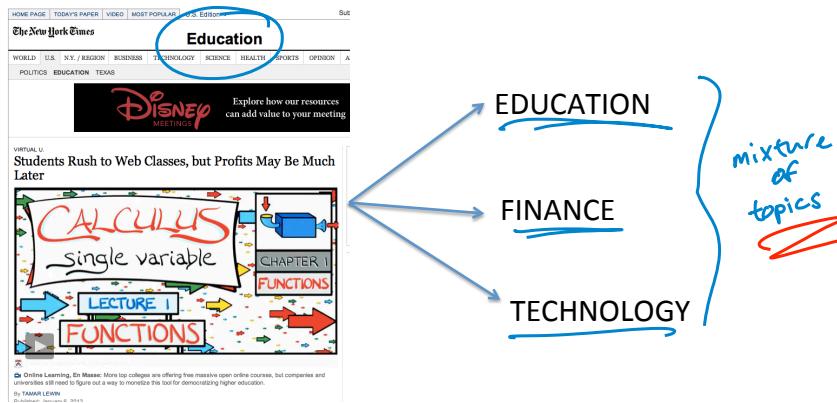
May 28th, 2015

©Emily Fox 2015

1

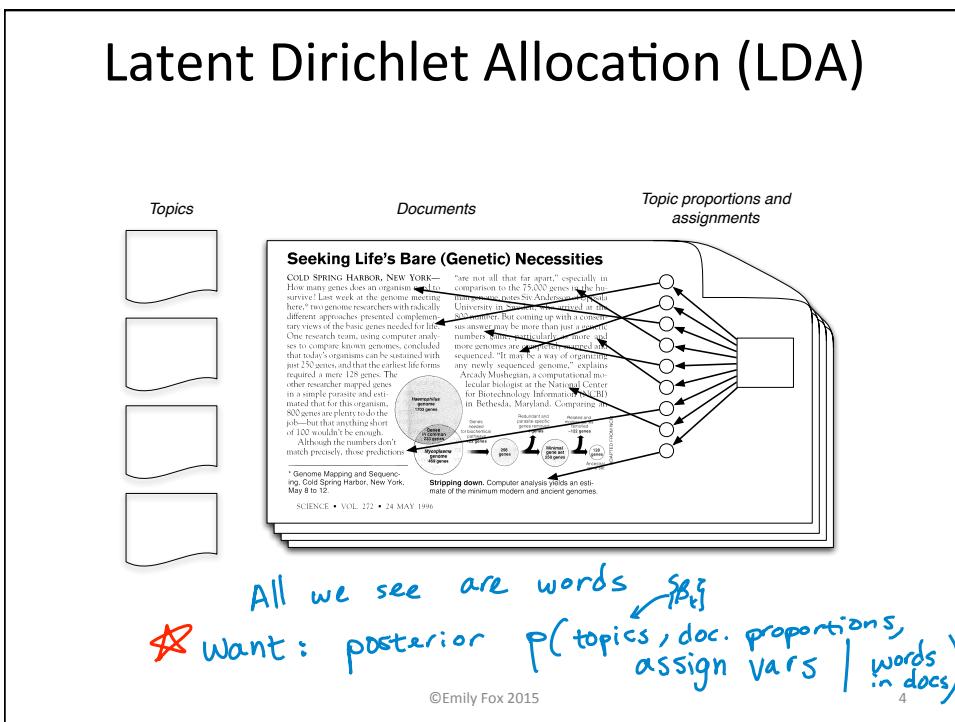
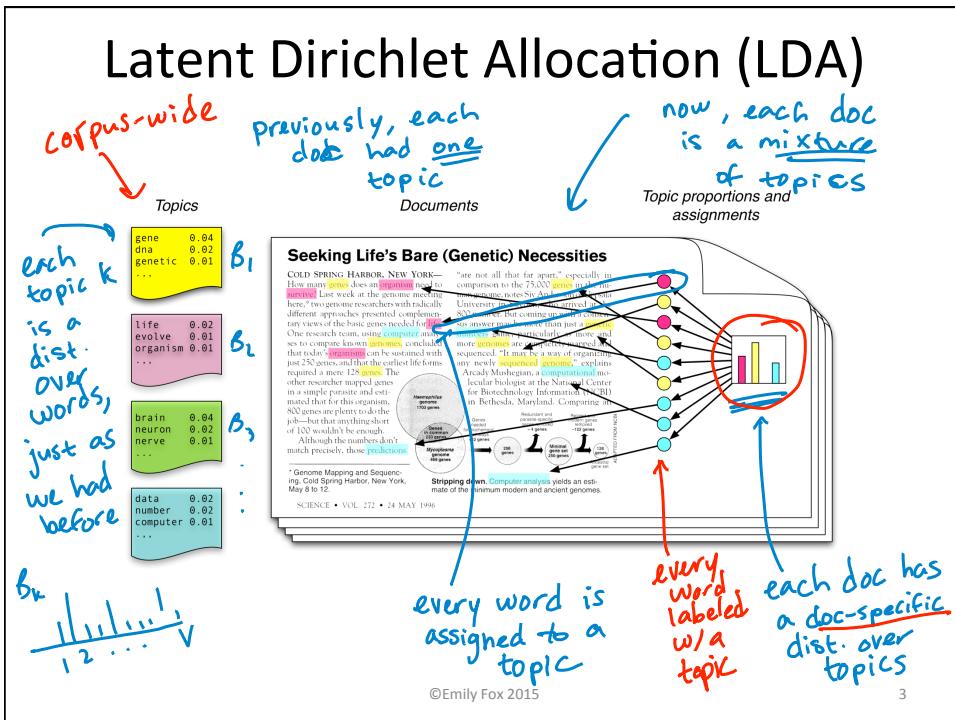
Task 3: Mixed Membership Models

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



©Emily Fox 2015

2



LDA Generative Model

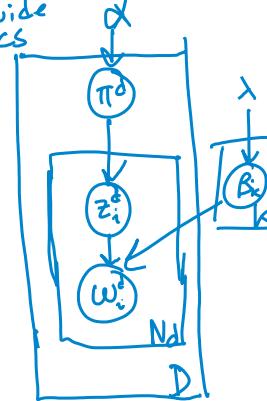
- Observations: $w_1^d, \dots, w_{N_d}^d \quad d=1, \dots, D$
- Associated topics: $z_1^d, \dots, z_{N_d}^d \quad \leftarrow \text{assign var. per word}$
- Parameters: $\theta = \{\{\pi^d\}, \{\beta_k\}\}$
- Generative model:
 - $z_i^d \sim \Pi^d \quad d=1, \dots, D$ (per word)
 - $w_i^d | z_i^d \sim \beta_{z_i^d} \quad i=1, \dots, N_d$ (doc-specific topic weights)
 - Π^d (corpus-wide topics) (topic weights)
 - β_k (topic words)

Priors:

$$\begin{aligned} \Pi^d &\sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad d=1, \dots, D \\ \beta_k &\sim \text{Dir}(\lambda_1, \dots, \lambda_V) \quad k=1, \dots, K \end{aligned}$$

©Emily Fox 2015

5



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"

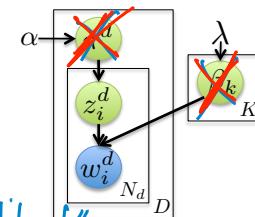
of words assigned to topic k in doc d (not counting i th word)

$\frac{n_{k,d} + \alpha_k}{N_d - 1 + \sum \alpha_k}$ • $\frac{m_{w_i^d, k} + \lambda_{w_i^d}}{\sum_\gamma m_{w_i^\gamma, k} + \lambda_\gamma}$

normalize within doc

examine entire corpus

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



©Emily Fox 2015

6

Select a Document

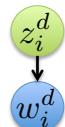
Etruscan	trade	price	temple	market

all words
in doc d

©Emily Fox 2015

7

Randomly Assign Topics



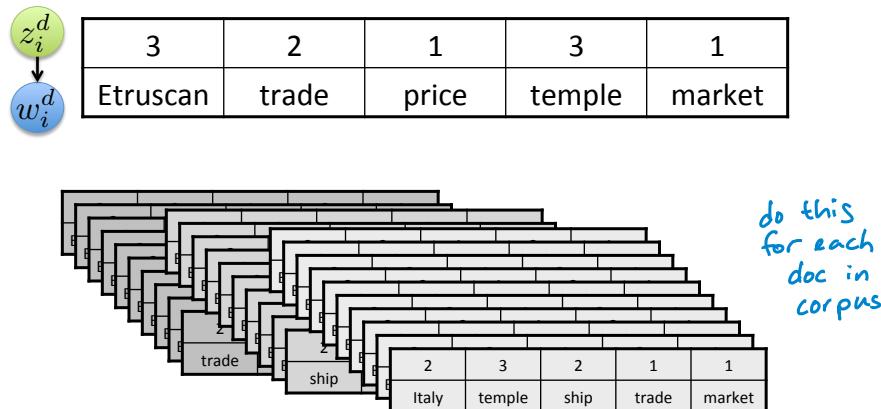
3	2	1	3	1
Etruscan	trade	price	temple	market

to initialize sampler
(one approach)

©Emily Fox 2015

8

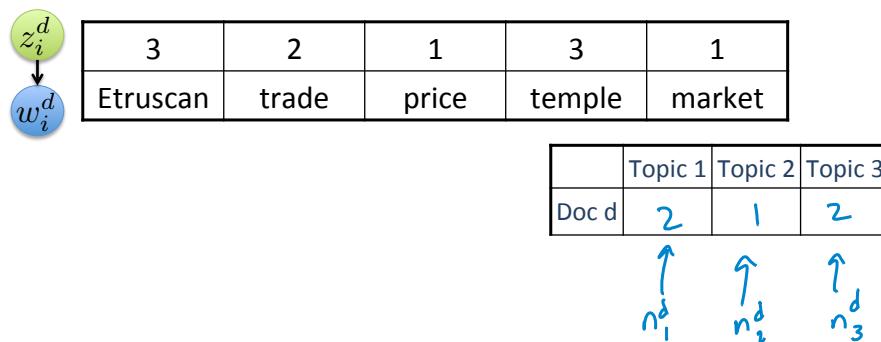
Randomly Assign Topics



©Emily Fox 2015

9

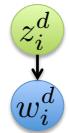
Maintain Local Statistics



©Emily Fox 2015

10

Maintain Global Statistics



3	2	1	3	1
Etruscan	trade	price	temple	market

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

	Topic 1	Topic 2	Topic 3
Doc d	2	1	2

©Emily Fox 2015

11

Resample Assignments



3	2	1	3	1
Etruscan	trade	price	temple	market

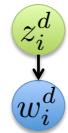
	Topic 1	Topic 2	Topic 3
Etruscan	1	0	35
market	50	0	1
price	42	1	0
temple	0	0	20
trade	10	7	1
...			

	Topic 1	Topic 2	Topic 3
Doc d	2	0	2

©Emily Fox 2015

12

What is the conditional distribution for this topic?



3	?	1	3	1
Etruscan	trade	price	temple	market

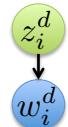
$$p(z_i^d \mid \text{everything else})$$

©Emily Fox 2015

13

What is the conditional distribution for this topic?

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

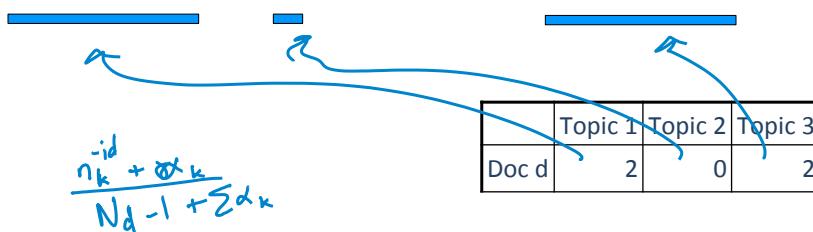


3	?	1	3	1
Etruscan	trade	price	temple	market

Topic 1

Topic 2

Topic 3

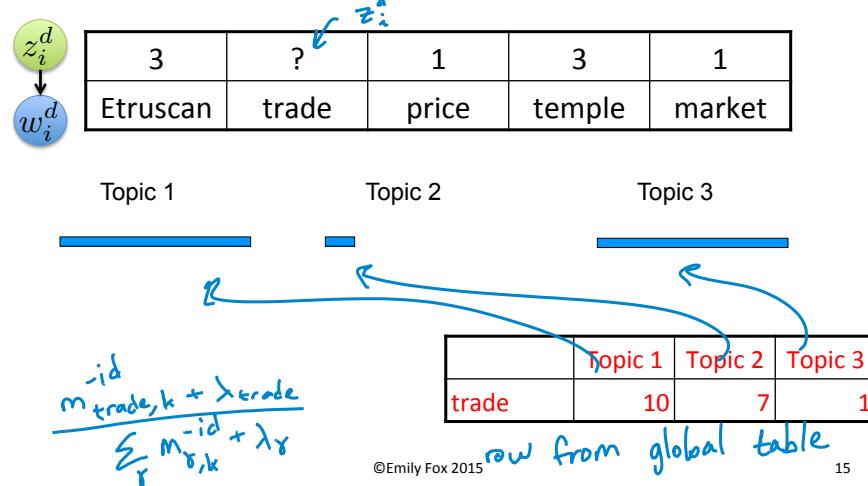


©Emily Fox 2015

14

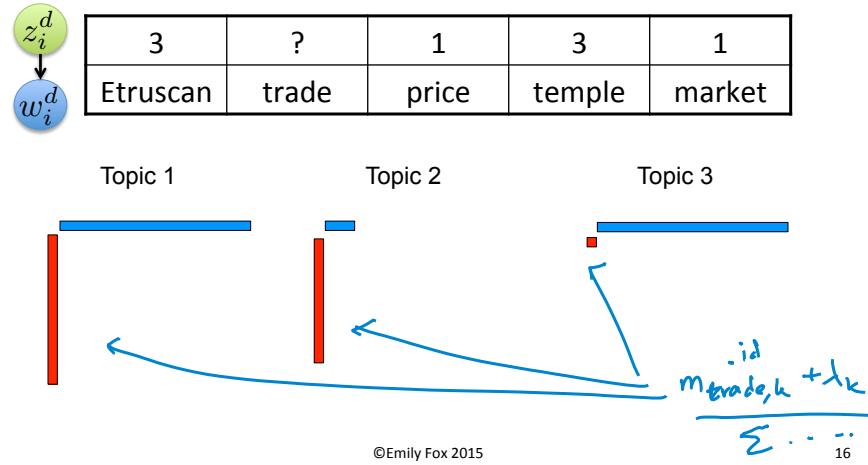
What is the conditional distribution for this topic?

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



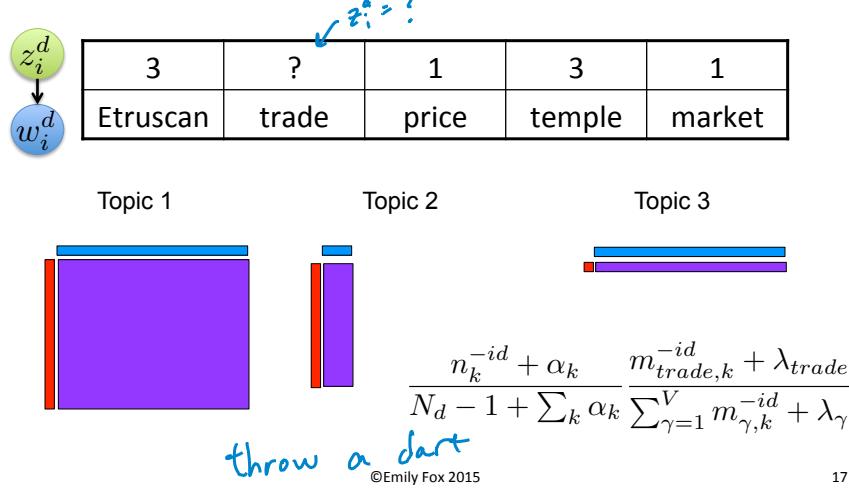
What is the conditional distribution for this topic?

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

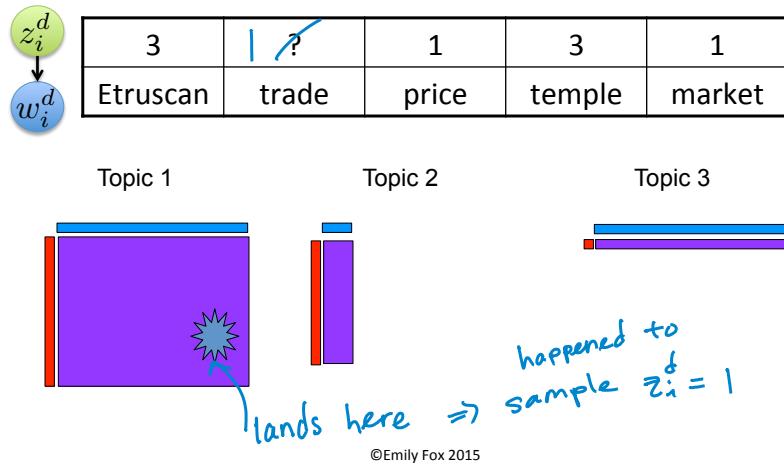


What is the conditional distribution for this topic?

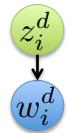
- Part I: How much does this document like each topic?
- Part II: How much does each topic like this word?



Sample a New Topic Indicator



Update Counts



3	1	1	3	1
Etruscan	trade	price	temple	market

	Topic 1	Topic 2	Topic 3
Etruscan	1	0	35
market	50	0	1
price	42	1	0
temple	0	0	20
trade	11	10	7
...			

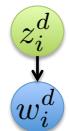
	Topic 1	Topic 2	Topic 3
Doc d	3	0	2

©Emily Fox 2015

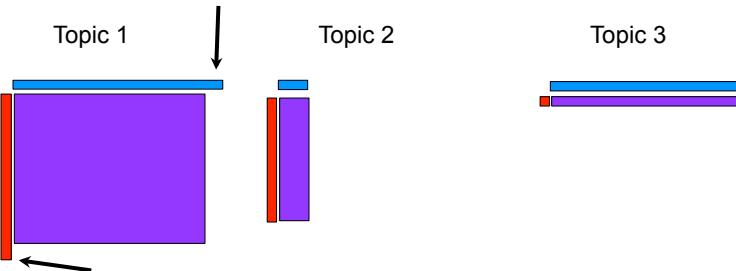
19

Geometrically...

inc. popularity of topic 1 in doc d
and prevalence of word "trade" in
topic 1 (corpus wide)



3	1	1	3	1
Etruscan	trade	price	temple	market



©Emily Fox 2015

20

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 include:
 - 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](#)
 - And many, many more!
- Alternative: Variational methods instead of sampling
 - Approximate posterior with an optimized variational distribution

©Emily Fox 2015

21

Case Study 5: Mixed Membership Modeling

Variational Methods

Machine Learning for Big Data
 CSE547/STAT548, University of Washington
 Emily Fox
 May 28th, 2015

©Emily Fox 2015

22

Variational Methods Goal

- Recall task: Characterize the posterior $p(\theta, z | x)$
 - params
 - latent vars
 - obs.
- 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: $p(\theta, z | x)$

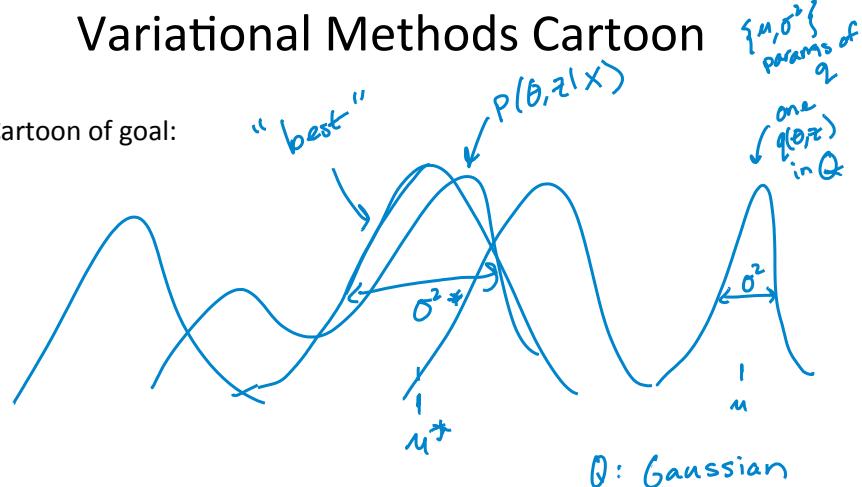
Call family Q and want $q \in Q$
 that is closest to $p(\theta, z | x)$

©Emily Fox 2015

23

Variational Methods Cartoon

- Cartoon of goal:



- Questions:

- (1) – How do we measure "closeness"?
- (2) – If the posterior is intractable, how can we approximate something we do not have to begin with?

©Emily Fox 2015

24

A Measure of Closeness

- Kullback-Leibler (KL) divergence
 - Measures “distance” between two distributions p and q

$$KL(p||q) \triangleq D(p||q) = E_p \left[\log \frac{p(\theta)}{q(\theta)} \right] = \int_{\theta} p(\theta) \log \frac{p(\theta)}{q(\theta)} d\theta$$

- If $p = q$ for all θ

$$D(p||q) = \int p(\theta) \log 1 d\theta = 0$$

- Otherwise, $D(p||q) > 0$

©Emily Fox 2015

25

A Measure of Closeness

$$KL(p||q) \triangleq D(p||q) = \int_{\theta} p(\theta) \log \frac{p(\theta)}{q(\theta)} d\theta$$

- Not symmetric $D(p||q) \neq D(q||p)$... not a true distance metric

- p determines where the difference is important:

$\exists \theta$ - $p(\theta)=0$ and $q(\theta)\neq 0$ $0 \log 0 = 0$

$\exists \theta$ - $p(\theta)\neq 0$ and $q(\theta)=0$ $e \log \frac{e}{0} = \infty$

If $D(p||q)$ Finite, $\text{supp}(q) \supseteq \text{supp}(p)$

- Want $\hat{q} = \arg \min_{q \in Q} D(p||q)$

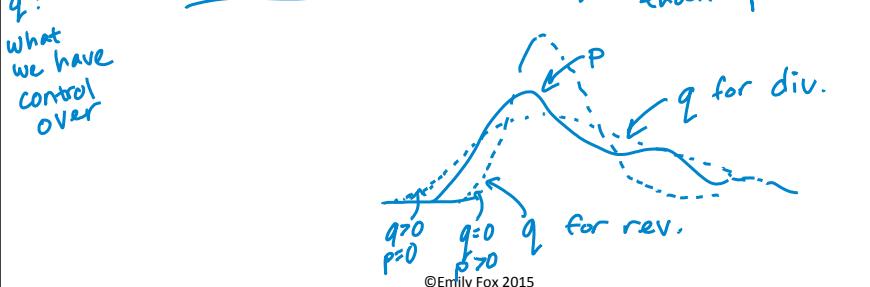
- Just as hard as the original problem! $E_p[\dots]$

©Emily Fox 2015

26

Reverse Divergence

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



27

Interpretations of Minimizing Reverse KL

$$D(q \parallel p) = E_q \left[\log \frac{q}{p} \right]$$

- Similarity measure:

$$\begin{aligned} D(q(\theta, z) \parallel p(\theta, z | x)) &= E_q [\log q(\theta, z)] - E_q [\log p(\theta, z | x)] \\ &= E_q [\log q(\theta, z)] - E_q [\log p(\theta, z, x)] + \log p(x) \end{aligned}$$

$\underbrace{- \mathcal{L}(q)}$

- Evidence lower bound (ELBO)

$$\log p(x) = D(q(z, \theta) \parallel p(\theta, z | x)) + \mathcal{L}(q) \geq \underline{\mathcal{L}(q)}$$

$\underbrace{\geq 0}$

↑
log marginal likelihood
= “evidence”

"ELBO"

©Emily Fox 2015

28

Interpretations of Minimizing Reverse KL

- Evidence lower bound (ELBO)

$$\log p(x) = \underbrace{D(q(z, \theta) || p(z, \theta | x))}_{\text{const.}} + \mathcal{L}(q) \geq \mathcal{L}(q)$$

↑
add to a const.
↑
ELBO

- Therefore,

- ELBO provides a lower bound on marginal likelihood
- Maximizing ELBO is equivalent to minimizing KL

$$\max \mathcal{L}(q) = \min D(q || p) = \max \text{lower bound on } \log p(x)$$

↑
what we can control
↑
depends on what we don't know
↑
ELBO

©Emily Fox 2015

29

Mean Field

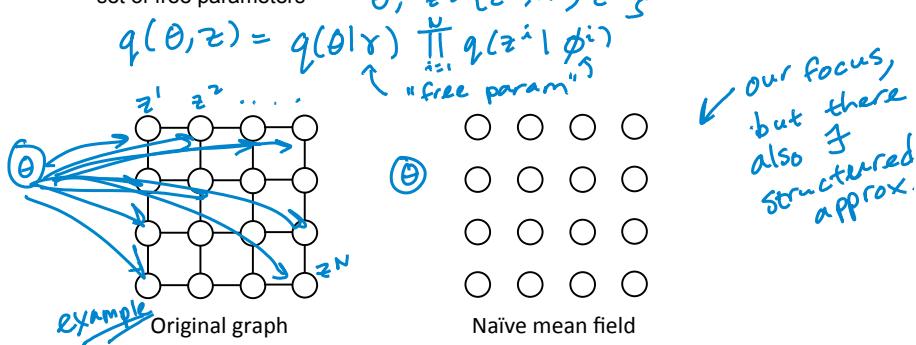
$$\mathcal{L}(q) = E_q[\log p(z, \theta, x)] - E_q[\log q(z, \theta)]$$

- How do we choose a Q such that the following is tractable?

$$\hat{q} = \arg \max_{q \in \mathcal{Q}} \mathcal{L}(q) \leftarrow \text{new objective}$$

- Simplest case = mean field approximation

- Assume each parameter and latent variable is conditionally independent given the set of free parameters $\theta, z = \{z^1, \dots, z^N\}$



Mean Field

$$\mathcal{L}(q) = E_q[\log p(z, \theta, x)] - E_q[\log q(z, \theta)]$$

Joint term entropy term Free params

- Naïve mean field decomposition:

$$q(z, \theta) = q(\theta | \gamma) \prod_{i=1}^N q(z^i | \phi^i)$$

- Under this approximation, entropy term decomposes as

$$-E_q[\log q(\theta, z)] = -E_q[\log q(\theta | \gamma)] - \sum_{i=1}^N E_q[\log q(z^i | \phi^i)]$$

decouples across γ, ϕ^i

- Can (always) rewrite joint term as

$$E_q[\log p(\theta, z, x)] = E_q[\log p(\theta | z, x)] + E_q[\log p(z, x)]$$

full conditional of θ full cond. of z^i

$$E_q[\log p(\theta, z, x)] = E_q[\log p(z^i | z_{\setminus i}, \theta, x)] + E_q[\log p(z_{\setminus i}, \theta, x)]$$

©Emily Fox 2015 31

Mean Field – Optimize γ

- Examine one free parameter, e.g., γ

$$\mathcal{L}(q) = E_q[\log p(\theta | z, x)] + E_q[\log p(z, x)] - E_q[\log q(\theta | \gamma)] - \sum_i E_q[\log q(z^i | \phi^i)]$$

consider the θ full cond. form of joint

- Look at terms of ELBO just depending on γ

$$\mathcal{L}^\gamma = E_q[\log p(\theta | z, x)] - E_q[\log q(\theta | \gamma)] + \text{const.}$$

really just $q_\theta = q(\theta | \gamma)$ needed here

doesn't depend on γ because under q_θ , $\theta \perp\!\!\!\perp z^i$!

©Emily Fox 2015 32

Mean Field – Optimize ϕ^i

- Examine another free parameter, e.g., ϕ^i

$$\mathcal{L}(q) = E_q[\log p(z^i | z_{\setminus i}, \theta, x)] + E_q[\log p(z_{\setminus i}, \theta, x)] - E_q[\log q(\theta | \gamma)] - \sum_i E_q[\log q(z^i | \phi^i)]$$

z^i full cond. form of joint

const. wrt ϕ^i

$$\mathcal{L}^{\phi^i} = E_q[\log p(z^i | z_{\setminus i}, \theta, x)] - E_q[\log q(z^i | \phi^i)]$$

really just $q_{\phi^i} = q(z^i | \phi^i)$

- This motivates using a coordinate ascent algorithm for optimization
 - Iteratively optimize each free parameter holding all others fixed

©Emily Fox 2015

33

Algorithm Outline

- Initialization:** Randomly select starting distribution $q_{\theta}^{(0)}$
- E-Step:** Given parameters, find posterior of hidden data

$$q_z^{(t)} = \arg \max_{q_z} \mathcal{L}(q_z, q_{\theta}^{(t-1)})$$

fixing prev. $\phi^{(t-1)}$
- M-Step:** Given posterior distributions, find likely parameters

$$q_{\theta}^{(t)} = \arg \max_{q_{\theta}} \mathcal{L}(q_z^{(t)}, q_{\theta})$$

fixing prev. $\phi^{(t)}$
- Iteration:** Alternate E-step & M-step until convergence

©Emily Fox 2015

34

Case Study 5: Mixed Membership Modeling

Variational Inference for LDA

Machine Learning for Big Data
 CSE547/STAT548, University of Washington
 Emily Fox
 May 28th, 2015

©Emily Fox 2015

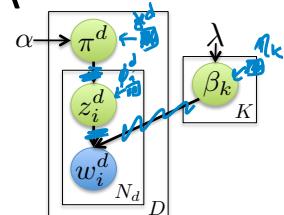
35

Mean Field for LDA

- In LDA, our parameters are $\theta = \{\pi^d\}, \{\beta_k\}$
 $z = \{z_i^d\}$
- The variational distribution factorizes as

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

$\underbrace{Dir(\eta_1, \dots, \eta_K)}_{\text{Dir}(\eta_1, \dots, \eta_K)}$
 $\underbrace{Dir(\alpha^1, \dots, \alpha^K)}_{\text{Dir}(\alpha^1, \dots, \alpha^K)}$
 $\underbrace{\text{Mult}(\phi_i^d)}_{\sum_{k=1}^K \phi_{ik}^d = 1 \quad \frac{1}{\sum_{k=1}^K \eta_k}}$
need to enforce this constraint



- 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 2015

36

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

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

all terms from q

from joint

©Emily Fox 2015 37

Mean Field for LDA

$$E_q(r, s) = \int q(r) q(s) f(r) f(s) dr ds = \int q(r) f(r) dr \int q(s) f(s) ds = E_q[f(r)] E_q[f(s)]$$

- Let's look at some of these terms

$$\begin{aligned} E_q[\log p(z_i^d | \pi^d)] &= E_q[\log \pi_{z_i^d}^d] = E_q\left[\sum_{k=1}^K (\log \pi_k^d) I(z_i^d=k)\right] \\ &= \sum_{k=1}^K E_q[I(z_i^d=k) \log \pi_k^d] = \sum_{k=1}^K E_q[I(z_i^d=k)] E_q[\log \pi_k^d] \end{aligned}$$

$z_i^d \perp\!\!\!\perp \pi_k^d$ under q

\Rightarrow why mean field approx. is so important

$\phi_{ik}^d \leftarrow \psi(\alpha_k^d) - \psi(\sum_k \alpha_k^d)$

$$E_q[\log q(z_i^d | \phi_i^d)] = \sum_k E_q[I(z_i^d=k) \log \phi_{ik}^d]$$

given

$$= \sum_k \phi_{ik}^d \log \phi_{ik}^d = \sum_k \phi_{ik}^d \log \phi_{ik}^d$$

- Other terms follow similarly

©Emily Fox 2015 38

Optimize via Coordinate Ascent

- Algorithm:

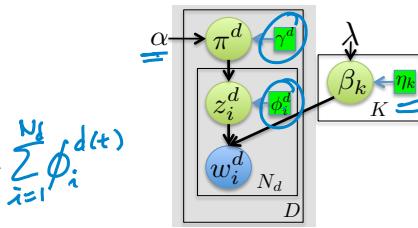
for $d=1, \dots, D$

$$\frac{\partial L}{\partial \gamma^d} = 0 \rightarrow \gamma^{d(t+1)} = \alpha + \sum_{i=1}^{N_d} \phi_i^{d(t)}$$

for $i=1, \dots, N_d$

$$\frac{\partial L}{\partial \phi_i^d} = 0 \rightarrow \phi_i^d \propto \exp \left\{ \Psi(\underline{\gamma}_{1:k}^{d(t+1)}) + \Psi(\underline{\eta}_{1:k, w_i^d}^{(t)}) - \Psi \left(\sum_v \underline{\eta}_{1:k, v}^{(t)} \right) \right\}$$

use Lagrange
multiplier to
enforce $\sum \phi_i^d = 1$



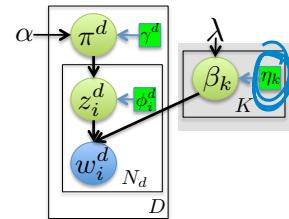
DATA PARALLEL
across ϕ_i^d

©Emily Fox 2015

39

Optimize via Coordinate Ascent

- Algorithm:



©Emily Fox 2015

40

What you need to know...

- Latent Dirichlet allocation (LDA)
 - Motivation and generative model specification
 - Collapsed Gibbs sampler
- Variational methods
 - Overall goal
 - Interpretation in terms of minimizing (reverse) KL
 - Mean field approximation

©Emily Fox 2015

41

Acknowledgements

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

©Emily Fox 2015

42