

Case Study 5: Mixed Membership Modeling

Variational Inference for LDA

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

Emily Fox

June 2nd, 2015

©Emily Fox 2015

1

Variational Methods Goal

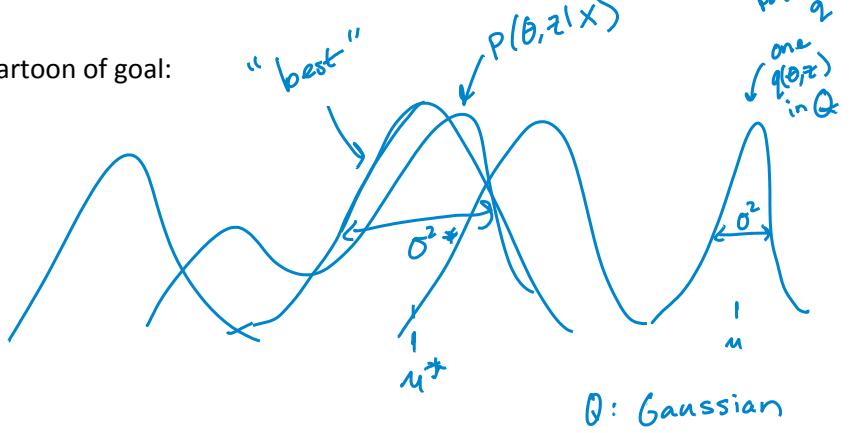
- Recall task: Characterize the posterior $p(\theta, z | x)$
 obs.
 params
 latent vars
 - ★ Turn posterior inference into an optimization task
 - Introduce a tractable family of distributions over parameters and latent variables
 - Family is indexed by a set of “free parameters”
 - Find member of the family closest to: $p(\theta, z | x)$
- Call family Q and want $q \in Q$
 that is closest to $p(\theta, z | x)$
- want q
 closest to p
 all 2D Gaussians

©Emily Fox 2015

2

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

3

Interpretations of Minimizing Reverse KL

- Evidence lower bound (ELBO)

$$\log p(x) = \underbrace{D(q(z, \theta) || p(z, \theta | x))}_{\text{KL div.} = \text{measure of "distance" between } p + q} + \mathcal{L}(q) \geq \mathcal{L}(q)$$

↑
ELBO
add to a const.

- Therefore,

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

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

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

©Emily Fox 2015

4

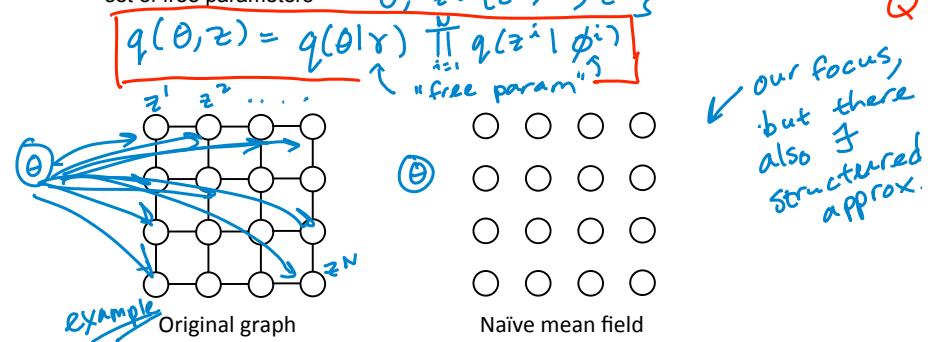
Mean Field

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

- 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 $\theta, z = \{z^1, \dots, z^N\}$

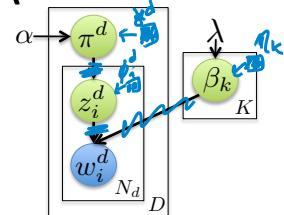


©Emily Fox 2015

5

Mean Field for LDA

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



- The variational distribution factorizes as

$$q(\pi, \beta, z) = \prod_{k=1}^K q(\beta_k | \eta_k) \prod_{d=1}^D \left[q(\pi^d | \alpha) \prod_{i=1}^{N_d} q(z_i^d | \phi_i^d) \right] \text{Dir}(\eta_{(1)}, \dots, \eta_{(K)}) \text{Dir}(\phi_1^d, \dots, \phi_K^d) \text{Mult}(\phi_i^d) \sum_{k=1}^K \phi_{ik}^d = 1$$

- The joint distribution factorizes as

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

optimize these free params

need to enforce this constraint

©Emily Fox 2015

6

Mean Field for LDA

mean field approx.

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

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

- Examine the ELBO

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

all terms from q

©Emily Fox 2015 7

Optimize via Coordinate Ascent

- Algorithm:

for d=1, ..., D

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

for i=1, ..., N_d

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

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

DATA PARALLEL across ϕ_i^d, γ^d

©Emily Fox 2015 8

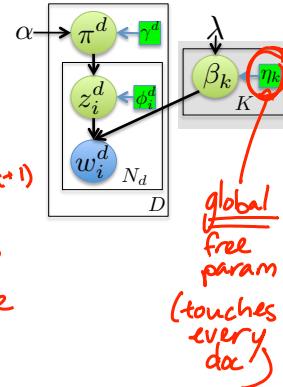
Optimize via Coordinate Ascent

- Algorithm:

for $k=1, \dots, K$

$$\frac{\partial L}{\partial \eta_k} = 0 \rightarrow \eta_k^{(t+1)} = \lambda + \sum_{d=1}^D \sum_{i=1}^{N_d} w_i^d f_i^{(t+1)}$$

aggregate



Map Reduce

©Emily Fox 2015

9

Generalizing

- Many Bayesian model have this form:

$$p(\theta, z^{1:N}, x^{1:N}) = p(\theta) \prod_{i=1}^N p(z^i | \theta) p(x^i | z^i, \theta)$$

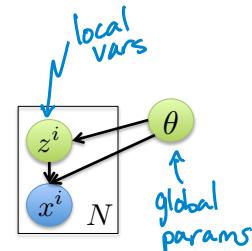
- Goal is to compute

$$p(\theta, z^i | x^i)$$

- Assume each complete conditional is in the exponential family

$$p(z^i | \theta, x^i) = h(z^i) \exp\{\eta_\ell(\theta, x^i)^T z^i - a(\eta_\ell(\theta, x^i))\}$$

↑ "local"



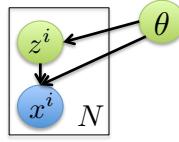
$$p(\theta | z, x) = h(\theta) \exp\{\eta_g(z, x)^T \theta - a(\eta_g(z, x))\}$$

↑ global

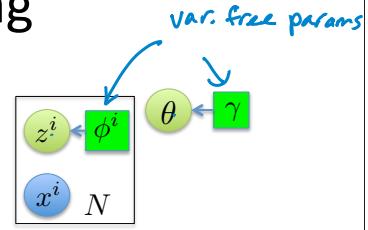
©Emily Fox 2015

10

Generalizing



mean field approx



var. free params

- Mean field variational approximation

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

- Match each component to have same family as model conditional

e.g. for global params
 $p(\theta | z, x) = h(\theta) \exp\{\eta_g(z, x)^T \theta - a(\eta_g(z, x))\}$

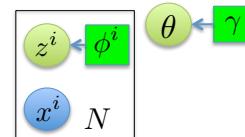
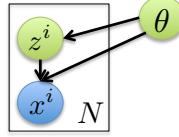
$$q(\theta | \gamma) = h(\theta) \exp\{\gamma^T \theta - a(\gamma)\}$$

- ✓ Same for local variational terms, too

©Emily Fox 2015

11

Generalizing



- Under these exponential family assumptions, the gradient is:

$$\nabla_\gamma \mathcal{L} = a''(\gamma)(E_\phi[\eta_g(z, x)] - \gamma)$$

- This leads to a simple coordinate update (Ghahramani and Beal, 2001)

$$\nabla_\gamma \mathcal{L} = 0 \rightarrow \gamma^* = E_\phi[\eta_g(z, x)]$$

©Emily Fox 2015

12

General Coord. Ascent Algorithm

Initialize γ randomly.

Repeat until the ELBO converges

- ① For each data point, update the local variational parameters:

$$\phi_i^{(t)} = \text{E}_{\gamma^{(t-1)}} [\ell(\theta, x_i)] \quad \text{for } i \in \{1, \dots, N\}.$$

- ② Update the global variational parameters:

$$\gamma^{(t)} = \text{E}_{\phi^{(t)}} [\ell_g(Z_{1:N}, x_{1:N})]$$

©Emily Fox 2015

13

Case Study 5: Mixed Membership Modeling

Stochastic Variational Inference

Machine Learning for Big Data
 CSE547/STAT548, University of Washington
 Emily Fox
 June 2nd, 2015

©Emily Fox 2015

14

Limitations of Batch Variational Methods

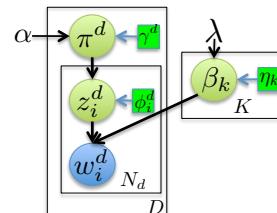


interested
in huge
datasets
(e.g. millions
of docs)

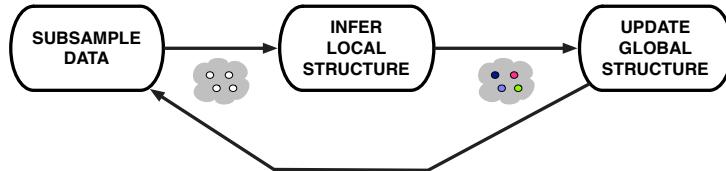
15

Limitations of Batch Variational Methods

- *Example = LDA*
 - Start from randomly initialized η_k (topics)
 - Analyze whole corpus before updating η_k again
- Streaming data: can't compute one iteration!
- *More generally...*
 - Do some local computation for each data point.
 - Aggregate these computations to re-estimate global structure.
 - Repeat.
- Inefficient, and cannot handle massive data sets.



Stochastic Variational Inference



- Stochastic variational inference harnesses:

- Idea #1: **Stochastic optimization** (Robbins and Monro, 1951)
 - Idea #2: **Natural gradients** (Amari, 1998)
- same ideas as behind SGD*

©Emily Fox 2015

17

Alternative Optimization Schemes

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

$$\nabla_t \mathcal{L} = 0$$

$$\gamma^{(t+1)} = \gamma^{(t)} + \epsilon_t \nabla_t \mathcal{L}$$

- Here,

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

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

touches all of the data (docs)!

- Use **stochastic gradient** step instead

- Consider one data point x^t sampled uniformly at random and define:

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

" t -ELBO"

$E[\nabla \mathcal{L}_t] = \nabla \mathcal{L}$ ← using all data

over subsampled data using subsample

©Emily Fox 2015

18

Alternative Optimization Schemes

- Recall the gradient of the ELBO for the global parameter:

$$\nabla_\gamma \mathcal{L} = a''(\gamma)(E_\phi[\eta_g(z, x)] - \gamma)$$

↑ global free param

- Even using just one data point, issue for scalability:

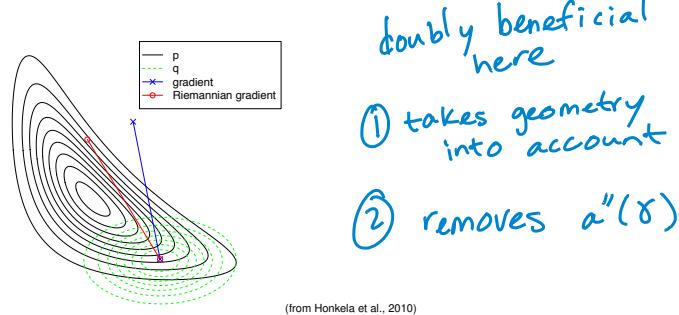
must compute $a''(\gamma)$...

computationally intensive

©Emily Fox 2015

19

Natural Gradient of the ELBO



- The **natural gradient** accounts for the geometry of parameter space
- Natural gradient of the ELBO:

$$\hat{\nabla}_\gamma \mathcal{L} = E_\phi[\eta_g(z, x)] - \gamma$$

↑ natural gradient "global"

©Emily Fox 2015

20

Noisy Natural Gradients

- Let $\eta_t(z^t, x^t)$ be the conditional distribution of the global parameters for the model where the observations are N replicates of x^t
- With this, the noisy natural gradient of the ELBO is

$$\hat{\nabla}_{\gamma} \mathcal{L}_t = E_{\phi_t}[\eta_t(z^t, x^t)] - \gamma$$

\uparrow natural grad of t -ELBO

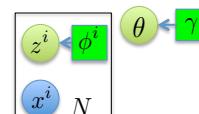
- Notes:
 - It only requires the local variational parameters of one data point.
 - In contrast, the full natural gradient requires all local parameters.
 - Thanks to conjugacy it has a simple form.

©Emily Fox 2015

21

SVI Algorithm Overview

Initialize global parameters γ randomly.
 Set the step-size schedule ϵ_t appropriately.
 Repeat forever



- Sample a data point uniformly,

$$x_t \sim \text{Uniform}(x_1, \dots, x_N).$$

- Compute its local variational parameter,

$$\phi_t = E_{\eta_t(\cdot|z^t)}[\eta_t(\theta, x_t)]$$

just as in
coord. asc.,
but only for
local var

- Pretend its the only data point in the data set,

$$\hat{\gamma} = E_{\phi_t}[\eta_t(z^t, x_t)]$$

- Update the current global variational parameter,

$$\gamma^{(t)} = (1 - \epsilon_t)\gamma^{(t-1)} + \epsilon_t \hat{\gamma}$$

$$\hat{\nabla}_{\gamma} \mathcal{L}_t = E_{\phi_t}[\cdot] - \gamma$$

$$\begin{aligned} \gamma^{(t)} &= \gamma^{(t-1)} + \epsilon_t \hat{\nabla}_{\gamma} \mathcal{L}_t \\ &= \gamma^{(t-1)} + \epsilon_t E_{\phi_t}[\cdot] - \epsilon_t \gamma^{(t-1)} \end{aligned}$$

©Emily Fox 2015

22

SVI for LDA

- In LDA, the full ELBO is given by

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

ELBO

all docs
- Assuming D documents, consider one sampled at random

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

t-ELBO

act as if doc t is repeated D times

just doc t

©Emily Fox 2015 23

SVI for LDA

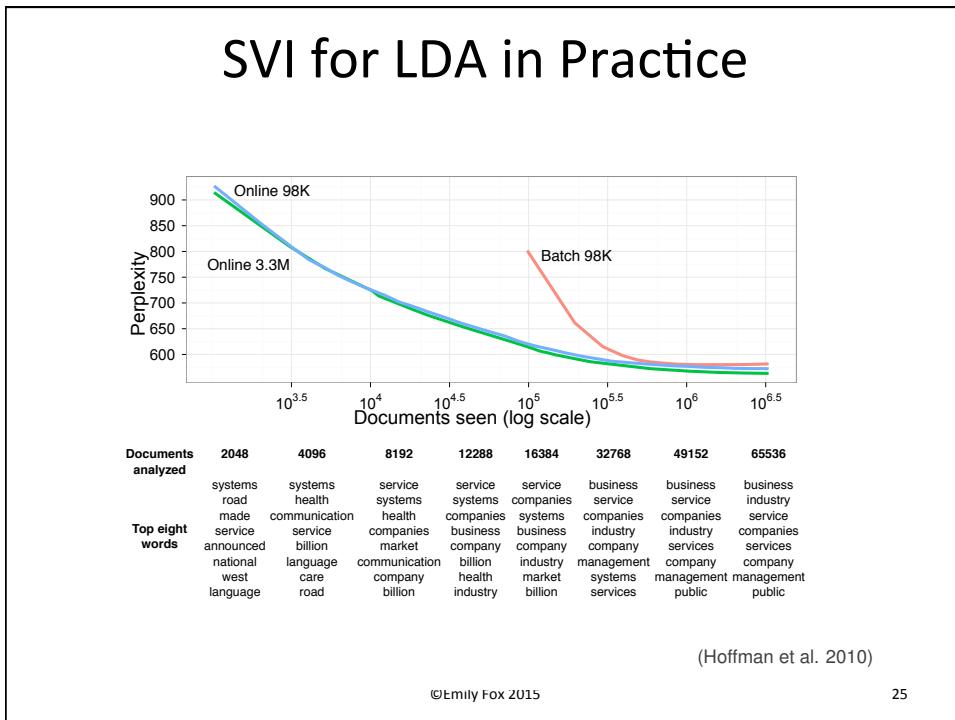
- Initialize $\eta^{(0)}$ randomly.
- Repeat (indefinitely):
 - Sample a document d uniformly from the data set.
 - For all k , initialize $\gamma_k^d = 1$
 - Repeat until converged
 - For $i=1, \dots, N_d$
 - Set $\phi_{ik}^d \propto \exp\{E[\log \pi_k^d] + E[\log \beta_k, w_i^d]\}$
 - Local free params for doc d
 - Take a stochastic gradient step

$$\eta^{(t)} = \eta^{(t-1)} + \epsilon_t \nabla_\eta \mathcal{L}_d$$
 - Global vars
 - $\gamma^d = \alpha + \sum_{i=1}^{N_d} \phi_i^d$
 - $\lambda + D \sum_{i=1}^{N_d} \phi_i^d w_i^d - \eta^{(t-1)}$
 - just as in coord. asc. for this doc

$\eta^{(t)} = (1-\epsilon_t) \eta^{(t-1)} + \epsilon_t \left(\lambda + D \sum_{i=1}^{N_d} \phi_i^d w_i^d - \eta^{(t-1)} \right)$

looks exactly like coord. asc. update for doc t replicated D times

©Emily Fox 2015 24



- ## What you need to know...
- Variational methods
 - Overall goal
 - Interpretation in terms of minimizing (reverse) KL
 - Mean field approximation
 - Mean field variational inference for LDA
 - Stochastic variational inference
 - General idea of using natural gradients + stochastic optimization
 - Resulting generic algorithm
 - SVI for LDA
- ©Emily Fox 2015 26

Reading

- **Inference in LDA:**

- Basic LDA and batch variational inference in LDA:
[Blei, David M., Andrew Y. Ng, and Michael I. Jordan. "Latent dirichlet allocation." the Journal of machine Learning research 3 \(2003\): 993-1022.](#)
- Stochastic variational inference:
[Hoffman, Matt, et al. "Stochastic Variational Inference." arXiv:1206.7051 \(2012\).](#)

©Emily Fox 2015

27

Acknowledgements

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

©Emily Fox 2015

28

Course Wrapup

Overview of CSE 547 / STAT 548 Topics Covered

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

Emily Fox

June 2nd, 2015

©Emily Fox 2015

29

What you need to know

- Case Study 1: Estimating Click Probabilities
 - Logistic regression
 - Regularization
 - Gradient descent, stochastic gradient decent
 - Hashing and sketching, random projections

©Emily Fox 2015

30

What you need to know

- Case Study 2: Document Retrieval and Clustering
 - Approach 1: **k-NN for nearest neighbor search**
 - *Algorithm:* Fast k-NN using KD-trees (exact)
 - *Algorithm:* Approximate k-NN using KD-trees and locality sensitive hashing
 - Approach 2: **k-means for clustering**
 - Data parallel problems
 - *Algorithm:* MapReduce framework and parallel k-means using MapReduce
 - Approach 3: **Gaussian mixture models (GMM) for clustering**
 - *Algorithm:* EM

©Emily Fox 2015

31

What you need to know

- Case Study 3: fMRI Prediction
 - Regularized linear models: Ridge regression and LASSO
 - Sparsistency
 - LASSO solvers:
 - LARS
 - Shotgun (stochastic coordinate descent)
 - Hogwild (stochastic gradient descent)
 - Averaging methods
 - ADMM
 - LASSO variants:
 - Fused LASSO
 - Graphical LASSO
 - Coping with large covariances using latent factor models

©Emily Fox 2015

32

What you need to know

- Case Study 4: Collaborative Filtering
 - Approach: Matrix factorization
 - *Algorithm*: Alternating least squares (ALS)
 - *Algorithm*: Stochastic gradient descent (SGD)
 - Cold-start problem and feature-based collaborative filtering
 - Model variants:
 - Non-negative matrix factorization
 - Probabilistic matrix factorization
 - *Algorithm*: Gibbs sampling
 - Probabilistic latent space models and network data
 - Graph parallel problems
 - GraphLab framework and application to distributed ALS and Gibbs for matrix factorization

©Emily Fox 2015

33

What you need to know

- Case Study 5: Document Mixed Membership Modeling
 - Approach 1: Bayesian document clustering model
 - Conditional independencies in directed graphical models
 - *Algorithm*: Gibbs sampling and collapsed Gibbs sampling
 - Approach 2: Latent Dirichlet allocation (LDA)
 - *Algorithm*: Collapsed Gibbs sampling
 - *Algorithm*: Variational methods and stochastic variational inference

©Emily Fox 2015

34

THANK YOU!!!

- You have been a great, interactive class!
...especially for a 9:30am lecture =)
- We're looking forward to the poster session
- Thanks to Marco and Alden, too!