

Case Study 5: Mixed Membership Modeling

Variational Methods for LDA

Stochastic Variational Inference

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

Emily Fox

March 13th, 2014

©Emily Fox 2014

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: $p(\theta, z | x)$

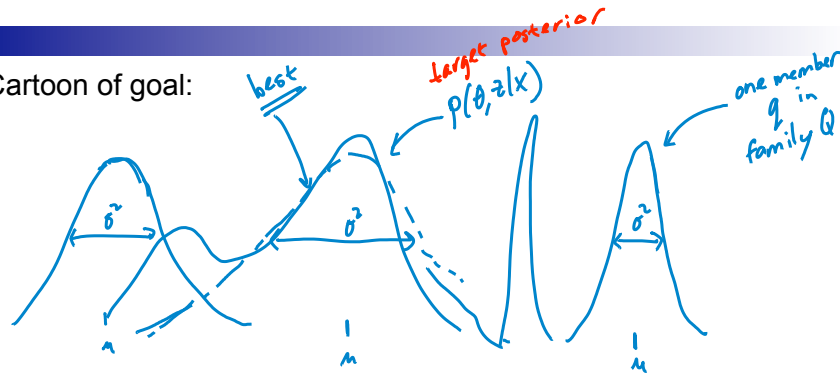
*Call the family Q and want $q \in Q$
that is closest to $p(\theta, z | x)$*

©Emily Fox 2014

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?

e.g., Q : all Gaussians

Interpretations of Minimizing Reverse KL

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

- Evidence lower bound (ELBO)

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

- ★ Therefore,

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

$$\max_{\text{what we can control}} \mathcal{L}(q) = \min_{\text{depends on what we don't know}} D(q \| p) = \max \text{ lower bound of } \log p(x)$$

Mean Field

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

joint *var. q*

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

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

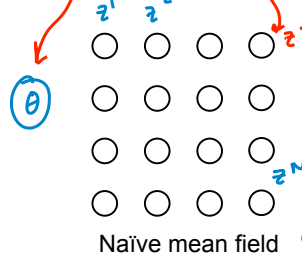
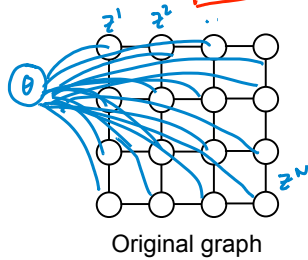
- Simplest case = mean field approximation

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

- Assume each parameter and latent variable is conditionally independent given the set of free parameters

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

$\gamma, \{\phi^i\}$ are "free params" = control knobs in getting q close to p



can also look at "structured" mean field approx (break by only some dependencies)

©Emily Fox 2014

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 θ -full-cond. form

- Look at terms of ELBO just depending on γ

don't depend on γ because under q,

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

$$z^i \perp \theta!$$

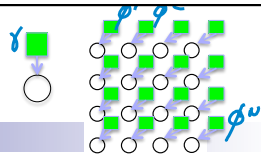
w.r.t. γ

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

©Emily Fox 2014

6

Mean Field – Optimize ϕ^i



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

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

consider the z^i -full-cond. form (under the first term)

const. wrt ϕ^i (under the second and third terms)

- Look at terms of ELBO just depending on ϕ^i

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

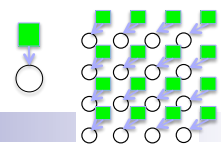
really just $q_{z^i} = q(z^i | \phi^i)$ here

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

©Emily Fox 2014

7

Algorithm Outline



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

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

latent vars z
- M-Step:** Given posterior distributions, find likely parameters θ

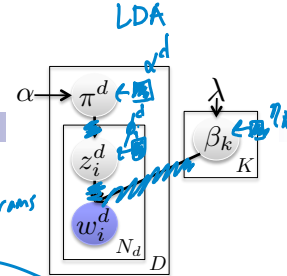
$$\text{optimize } \gamma \rightarrow q_{\theta}^{(t)} = \arg \max_{q_{\theta}} \mathcal{L}(q_z^{(t)}, q_{\theta})$$
- Iteration:** Alternate E-step & M-step until convergence

©Emily Fox 2014

8

Mean Field for LDA

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



- The variational distribution factorizes as

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

$\text{Dir}(\eta_1, \dots, \eta_K)$ $\text{Dir}(\alpha_1, \dots, \alpha_K)$ $\text{Mult}(\phi_i^d)$

$\sum_{k=1}^K \phi_{ik}^d = 1$ ← need to enforce this!

- The joint distribution factorizes as

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

$\text{Dir}(\lambda_1, \dots, \lambda_K)$ $\text{Dir}(\alpha_1, \dots, \alpha_K)$ $\text{Mult}(\pi^d)$

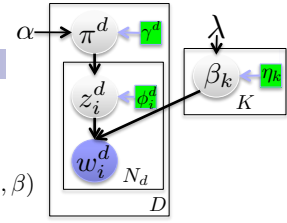
©Emily Fox 2014

9

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

$$\mathcal{L}(q) = \sum_{k=1}^K E_q[\log p(\beta_k | \lambda)] + \sum_{d=1}^D E_q[\log p(\pi^d | \alpha)]$$

$$+ \sum_{d=1}^D \sum_{i=1}^{N_d} E_q[\log p(z_i^d | \pi^d)] + E_q[\log p(w_i^d | z_i^d, \beta)]$$

$$- \sum_{k=1}^K E_q[\log q(\beta_k | \eta_k)] - \sum_{d=1}^D E_q[\log q(\pi^d | \gamma^d)] - \sum_{d=1}^D \sum_{i=1}^{N_d} E_q[\log q(z_i^d | \phi_i^d)]$$

From joint dist.

all terms from q

©Emily Fox 2014

10

Mean Field for LDA

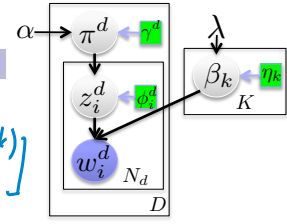
- Let's look at some of these terms

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

$$= \sum_{k=1}^K E_q[\mathbb{I}(z_i^d=k) \log \pi_k^d] = \sum_{k=1}^K E_q[\mathbb{I}(z_i^d=k)] E_q[\log \pi_k^d]$$

$z_i^d \perp \pi_k^d$ under q
 \Rightarrow why mean field is important

$\Pr(z_i^d=k) = \phi_{ik}^d$ $\Psi(\gamma_k^d) - \Psi(\sum_v \gamma_{k,v}^d)$



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

$\underbrace{\log \phi_{ik}^d}_{\text{given}}$

- Other terms follow similarly

©Emily Fox 2014

11

Optimize via Coordinate Ascent

- Algorithm:

for $d=1, \dots, D$

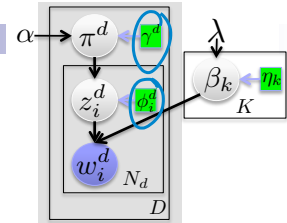
$$\frac{d\mathcal{L}}{d\gamma^d} = 0 \rightarrow \gamma^{d(t+1)} = \alpha + \sum_{i=1}^{N_d} \phi_i^d(t)$$

for $i=1, \dots, N_d$

$$\frac{d\mathcal{L}}{d\phi_i^d} = 0 \rightarrow \phi_i^d \propto \exp\left\{ \Psi(\gamma_{1:k}^{d(t+1)}) + \Psi(\eta_{k|k, w_i^d}^{(t)}) - \Psi\left(\sum_v \eta_{1:k, v}^{(t)}\right) \right\}$$

use Lagrange multipliers to enforce ϕ_i^d is a pmf

DATA PARALLEL



©Emily Fox 2014

12

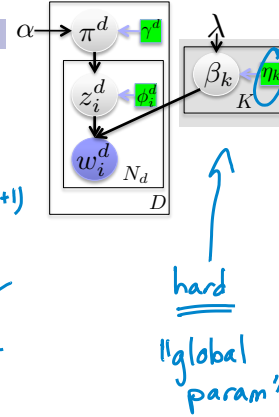
Optimize via Coordinate Ascent

- Algorithm:

for $k=1, \dots, K$

$$\frac{dL}{d\eta_k} = 0 \rightarrow \eta_k^{(t+1)} = \lambda + \underbrace{\sum_{d=1}^D \sum_{i=1}^{N_d} w_i^d \phi_i^d}_{\text{aggregate}} \eta_k^{(t)}$$

Map Reduce



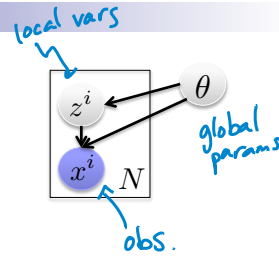
©Emily Fox 2014

13

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 $p(\theta, z | x)$

- 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"

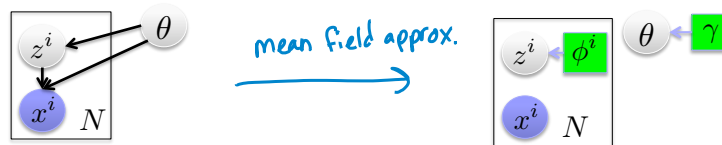
$$p(\theta | z, x) = h(\theta) \exp\{\eta_g(z, x)^T \theta - a(\eta_g(z, x))\}$$

"global"

©Emily Fox 2014

14

Generalizing



- 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

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

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

- Same for local variational terms, too

©Emily Fox 2014

15

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)

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

similarly for local vars

©Emily Fox 2014

16

General Coord. Ascent Algorithm

Initialize γ randomly.

Repeat until the ELBO converges

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

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

- 2 Update the global variational parameters:

$$\gamma^{(t)} = E_{\phi_i^{(t)}} [\mathbb{E}_g(Z_{1:N}, x_{1:N})]$$

Note: cycle through each obs. before updating global param.

Case Study 5: Mixed Membership Modeling

Stochastic Variational Inference

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

Emily Fox

March 13th, 2014

Limitations of Batch Variational Methods



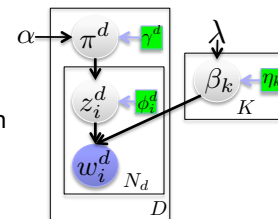
interested
in huge
datasets
(e.g. millions
of docs)

©Emily Fox 2014

20

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!

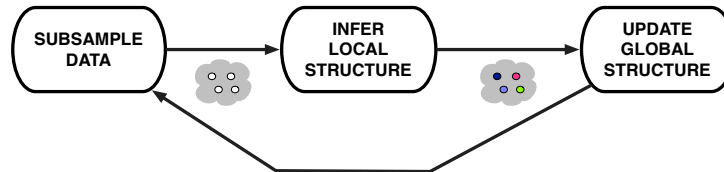


- **More generally...**
 - Do some local computation for each data point.
 - Aggregate these computations to re-estimate global structure.
 - Repeat.
- Inefficient, and cannot handle massive data sets.

©Emily Fox 2014

21

Stochastic Variational Inference



- Stochastic variational inference harnesses:

- Idea #1: **Stochastic optimization** (Robbins and Monro, 1951)
 - Idea #2: **Natural gradients** (Amari, 1998)
- just like in SGD* (handwritten note with an arrow pointing from Idea #1 to Idea #2)

Alternative Optimization Schemes

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

Handwritten notes:
 $\hookrightarrow \nabla_{\theta} \mathcal{L} = 0$
 $\gamma^{(t+1)} = \gamma^{(t)} + \epsilon_t \nabla_{\theta} \mathcal{L}$

- Here,

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

Handwritten note: touches all of the data! (does)

- Use **stochastic gradient** step instead

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

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

Handwritten note: "t-ELBO"

Handwritten notes:
 $E[\nabla \mathcal{L}_t] = \nabla \mathcal{L}$
 ← using all data
 ↑ using a subset
 ↑ over all subsets

Alternative Optimization Schemes

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

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

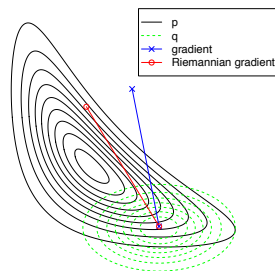
- Even using just one data point, issue for scalability:

*must compute $a''(\gamma)$...
computationally intensive*

©Emily Fox 2014

24

Natural Gradient of the ELBO



(from Honkela et al., 2010)

doubly beneficial:
① takes geometry into account
② removes $a''(\gamma)$ term

- 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 grad.*

©Emily Fox 2014

25

Noisy Natural Gradients

- Let $\eta_i(z^t, x^t)$ be the conditional distribution of the global variable for the model where the observations are N replicates of x^t

- With this, the noisy natural gradient of the ELBO is

$$\hat{\nabla}_{\gamma} \mathcal{L}_t = \underset{\substack{\uparrow \\ t\text{-ELBO}}}{\mathbb{E}_{\phi_t}} [\eta_t(z^t, x^t)] - \gamma$$

- 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.

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,

$$\phi_t = \mathbb{E}_{\gamma^{(t-1)}} [\eta(\theta, x_t)]$$

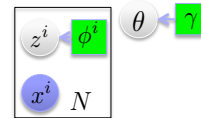
- Pretend its the only data point in the data set,

$$\hat{\gamma} = \mathbb{E}_{\phi_t} [\eta(z, x_t)]$$

- Update the current global variational parameter,

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

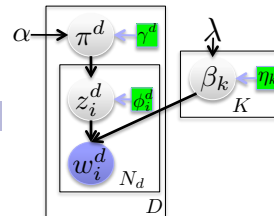
$$\gamma^{(t)} = \gamma^{(t-1)} + \epsilon_t [\mathbb{E}_{\phi_t} [\eta(z^t, x^t)] - \gamma^{(t-1)}]$$



← just as in coord. ascent, but just for local var only

← noisy natural gradient step

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, as before

- Assuming D documents, consider one sampled at random *as if we viewed it D times just one doc t*

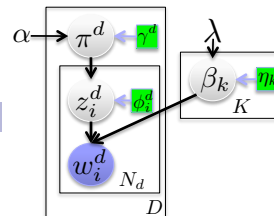
$$\mathcal{L}_t = E_q[\log p(\beta)] - E_q[\log q(\beta)] + D(E_q[\log p(\pi^t)] - E[\log q(\pi^t)]) + D(E_q[\log p(z^t, x^t | \pi^t, \beta)] - E_q[\log q(z^t)])$$

t-ELBO

©Emily Fox 2014

28

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 $i=1, \dots, N_d$

local free params for doc d

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

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

just like in coord. asc. for this doc.

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

global vars

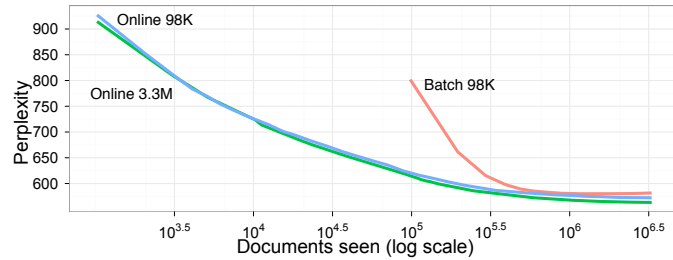
$$\eta^{(t)} = (1 - \epsilon_t) \eta^{(t-1)} + \epsilon_t \left(\lambda + D \sum_{i=1}^{N_d} \beta_i^d w_i^d - \eta^{(t-1)} \right)$$

looks exactly like what we had in coord. asc. update for doc d (here $D \sum_{i=1}^{N_d} \beta_i^d w_i^d$)

©Emily Fox 2014

29

SVI for LDA in Practice



Documents analyzed	2048	4096	8192	12288	16384	32768	49152	65536
Top eight words	systems road made service announced national west language	systems health communication service billion language care road	service systems health companies market communication company billion	service systems companies business company billion health industry	service companies systems business company industry market billion	business service companies industry company management systems services	business service companies industry services company management public	business industry service companies services company management public

(Hoffman et al. 2010)

©Emily Fox 2014

30

What you need to know...

- Variational methods
 - Mean field for LDA
- Stochastic variational inference
 - General idea of using natural gradients + stochastic optimization
 - Resulting generic algorithm
 - SVI for LDA

©Emily Fox 2014

31

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 \(2012\).](#)

Acknowledgements

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

Course Wrapup

Overview of CSE 547 / STAT 548 Topics Covered

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

Emily Fox

March 13th, 2014

©Emily Fox 2014

34

What you need to know

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

©Emily Fox 2014

35

What you need to know

■ Case Study 2: Document Retrieval and Clustering

- Approach 1: **k-NN**
- *Algorithm*: Fast k-NN using KD-trees (exact)
- *Algorithm*: Approximate k-NN using KD-trees and locality sensitive hashing

- Approach 2: **k-means**
- Data parallel problems
- *Algorithm*: MapReduce framework and parallel k-means using MapReduce

- Approach 3: **Gaussian mixture models (GMM)**
- *Algorithm*: EM

©Emily Fox 2014

36

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 2014

37

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

- Graph parallel problems
- GraphLab framework and application to distributed ALS and Gibbs for matrix factorization

©Emily Fox 2014

38

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
- *Algorithm*: Collapsed Gibbs sampling
- *Algorithm*: Variational methods and stochastic variational inference

©Emily Fox 2014

39

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 Alden and Chad, too!