

# Variational Inference (Contd)

CS772A: Probabilistic Machine Learning

Piyush Rai

# Quick Recap: Variational Inference (VI)

Defines a class of distributions parametrized by  $\phi$

- Approximate the true posterior  $p(\mathbf{Z}|\mathbf{X})$  by an approx. distribution  $q(\mathbf{Z}|\phi)$  or  $q_\phi(\mathbf{Z})$

$$\phi^* = \operatorname{argmin}_\phi \operatorname{KL}[q_\phi(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X})]$$

Often, we will simply write it as  $\operatorname{argmin}_q \operatorname{KL}[q||p_z]$

- Due to the below identity, minimizing the KL is equivalent to maximizing the **ELBO**

Log-evidence of model  $m$

Evidence lower bound (ELBO)

Non-negative

$$\log p(\mathbf{X}|m) = \mathcal{L}(q) + \operatorname{KL}(q||p_z)$$

- The ELBO is defined as  $\mathcal{L}(q) = \int q(\mathbf{Z}) \log \left[ \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \right] d\mathbf{Z}$

$$\mathcal{L}(q) = \mathcal{L}(\phi) = \mathbb{E}_q[\log p(\mathbf{X}, \mathbf{Z})] - \mathbb{E}_q[\log q(\mathbf{Z})]$$

Find  $q$  such that  $\mathbf{z}$  explains data well

$$= \mathbb{E}_q[\log p(\mathbf{X}|\mathbf{Z})] - \operatorname{KL}[q(\mathbf{Z})||p(\mathbf{Z})]$$

Find  $q$  that is "simple", i.e., is close to the prior

- VI optimizes (maximizes) the above w.r.t.  $q$ , i.e., w.r.t. its variational parameters  $\phi$



# Quick Recap: Mean-Field VI

- Assume  $q(\mathbf{Z}|\phi) = \prod_{i=1}^M q(\mathbf{Z}_i|\phi_i)$ . Simplifies ELBO expression/its maximization
- Learning the optimal  $q$  then reduces to learning the optimal  $q_1, q_2, \dots, q_M$
- For mean-field VI, each optimal factor  $q_j$  is given by

For locally conjugate models,  $q_j^*(\mathbf{z}_j)$  will have the same form as prior  $p(\mathbf{z}_j)$

$$q_j^* = \frac{\exp(\mathbb{E}_{i \neq j}[\log p(\mathbf{X}, \mathbf{Z})])}{\int \exp(\mathbb{E}_{i \neq j}[\log p(\mathbf{X}, \mathbf{Z})]) d\mathbf{Z}_j}$$

This general expression holds even if there is no local conjugacy

$\mathbb{E}_{i \neq j}$  denotes the expectation w.r.t. the distribution  $\prod_{i \neq j} q_i$

- Updates of optimal  $q_1, q_2, \dots, q_M$  depend on each other because of the expectations
- Therefore, MFVI works by updating the  $q_j$ 's in a cyclic fashion
  - Leads to the coordinate ascent VI (CAVI) algorithm



# Mean-Field VI: A Closer Look

- Since  $\log q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j}[\log p(\mathbf{X}, \mathbf{Z})] + \text{const} = \mathbb{E}_{i \neq j}[\log p(\mathbf{X}, \mathbf{Z}_j, \mathbf{Z}_{-j})] + \text{const}$

$$\log q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j}[\log p(\mathbf{Z}_j | \mathbf{X}, \mathbf{Z}_{-j})] + \text{const}$$

For any model

- Thus opt variational distr  $q_j^*(\mathbf{Z}_j)$  basically requires expectations of CP  $p(\mathbf{Z}_j | \mathbf{X}, \mathbf{Z}_{-j})$
- For *locally conjugate* models, CP can be easily found and is an *exp-fam distr* of the form

$$p(\mathbf{Z}_j | \mathbf{X}, \mathbf{Z}_{-j}) = h(\mathbf{Z}_j) \exp \left[ \eta(\mathbf{X}, \mathbf{Z}_{-j})^\top \mathbf{Z}_j - A(\eta(\mathbf{X}, \mathbf{Z}_{-j})) \right]$$

Gibbs sampling samples from each CP. MFVI uses each CP to compute the corresponding  $q_j$

- Using the above, we can rewrite the optimal variational distribution as follows

$$\log q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} \left[ \log \left( h(\mathbf{Z}_j) \exp \left[ \eta(\mathbf{X}, \mathbf{Z}_{-j})^\top \mathbf{Z}_j - A(\eta(\mathbf{X}, \mathbf{Z}_{-j})) \right] \right) \right] + \text{const}$$

$$\implies q_j^*(\mathbf{Z}_j) \propto h(\mathbf{Z}_j) \exp \left[ \mathbb{E}_{i \neq j}[\eta(\mathbf{X}, \mathbf{Z}_{-j})]^\top \mathbf{Z}_j \right] \quad (\text{verify})$$

- Thus, with local conj, we just require expectation of nat. params. of CP of  $\mathbf{Z}_j$



# VI by Computing ELBO Gradients

- Can also do VI by computing ELBO's gradient and doing gradient based optimization
- Gradient based approach is broadly applicable, not just for mean-field VI
  1. Assume  $q(\mathbf{Z})$  to be from some family of distributions with variational parameters  $\phi$
  2. Write down the full ELBO expression (will give us a function of var. parameters  $\phi$ )

$$\begin{aligned}\mathcal{L}(q) = \mathcal{L}(\phi) &= \mathbb{E}_q[\log p(\mathbf{X}, \mathbf{Z})] - \mathbb{E}_q[\log q(\mathbf{Z})] \\ &= \int q(\mathbf{Z}) \log p(\mathbf{X}|\mathbf{Z}) d\mathbf{Z} + \int q(\mathbf{Z}) \log p(\mathbf{Z}) d\mathbf{Z} - \int q(\mathbf{Z}) \log q(\mathbf{Z}) d\mathbf{Z}\end{aligned}$$

3. Compute ELBO gradients, i.e.,  $\nabla_{\phi} \mathcal{L}(\phi)$  and use gradient methods to find optimal  $\phi$
- Step 2 may be simplified due to the problem structure or the form of  $q(\mathbf{Z})$ 
    - **i.i.d. observations** simplify  $\log p(\mathbf{X}|\mathbf{Z})$ ; **conditionally independent priors** simplify  $\log p(\mathbf{Z})$
    - Locally-conjugate models
    - The mean-field assumption simplifies  $q(\mathbf{Z})$  as  $q = \prod_{i=1}^M q_i$ 
      - Moreover, the last term reduces to sum of entropies of  $q_i$ 's (which usually has known forms)



# Mean-Field VI by Taking ELBO's Gradients

- Mean-field assumption  $q(\mathbf{Z}|\phi) = \prod_{i=1}^M q(\mathbf{Z}_i|\phi_i)$  results in following optimal distribution

$$q_j^*(\mathbf{Z}_j) = \frac{\exp(\mathbb{E}_{i \neq j}[\log p(\mathbf{X}, \mathbf{Z})])}{\int \exp(\mathbb{E}_{i \neq j}[\log p(\mathbf{X}, \mathbf{Z})]) d\mathbf{Z}_j}$$

This approach is applicable even if we don't have mean-field assumption

Note that here we do not have to assume the form of this variational distribution. We simply compute the RHS and find what it is (in the locally-conjugate case, it will be the same distribution as the prior)

- Alternatively, we can take ELBO's partial deriv w.r.t.  $\phi_1, \phi_2, \dots, \phi_M$  to find their optimal values
- Consider a Bayesian linear regression model

Likelihood

$$y_i \sim \text{Normal}(x_i^T w, \alpha^{-1}), \quad w \sim \text{Normal}(0, \lambda^{-1} I), \quad \alpha \sim \text{Gamma}(a, b)$$

Prior on  $w$

$\lambda$  assumed fixed

Prior on variance of Gaussian likelihood

Needed in ELBO

Joint distribution on data and unknowns

$$p(y, w, \alpha | x) = p(\alpha) p(w) \prod_{i=1}^N p(y_i | x_i, w, \alpha)$$

Assumed variational posterior with mean-field assumption

$$q(w, \alpha) = q(\alpha) q(w) = \text{Gamma}(\alpha | a', b') \text{Normal}(w | \mu', \Sigma')$$

Note that in this approach, we have to assume a form for each variational distribution. It is common to assume them to have the same form as the respective priors

- Now doing VI amounts to maximizing ELBO to find the optimal variational params  $a', b', \mu', \Sigma'$



# Mean-Field VI by Taking ELBO's Gradients

- The ELBO is

For the Bayesian linear regression model, instead of  $p(\mathbf{X}, \mathbf{Z})$ , it will be of the form  $p(\mathbf{y}, \mathbf{Z}|\mathbf{X})$

$$\begin{aligned}\mathcal{L}(q) = \mathcal{L}(\phi) &= \mathbb{E}_q[\log p(\mathbf{X}, \mathbf{Z})] - \mathbb{E}_q[\log q(\mathbf{Z})] = \mathbb{E}_q[\log p(\mathbf{Z})] + \mathbb{E}_q[\log p(\mathbf{X}|\mathbf{Z})] - \mathbb{E}_q[\log q(\mathbf{Z})] \\ &= \int q(\mathbf{Z}) \log p(\mathbf{Z}) d\mathbf{Z} + \int q(\mathbf{Z}) \log p(\mathbf{X}|\mathbf{Z}) d\mathbf{Z} + \int q(\mathbf{Z}) \log q(\mathbf{Z}) d\mathbf{Z}\end{aligned}$$

- Thus the ELBO in the Bayesian linear regression model will be (assuming i.i.d. obs)

$$\begin{aligned}\mathcal{L}(a', b', \mu', \Sigma') &= \int q(\alpha) \ln p(\alpha) d\alpha + \int q(w) \ln p(w) dw \\ &\quad + \sum_{i=1}^N \int \int q(\alpha) q(w) \ln p(y_i | x_i, w, \alpha) dw d\alpha - \int q(\alpha) \ln q(\alpha) d\alpha - \int q(w) \ln q(w) dw\end{aligned}$$

Expectations of the log of the prior

Expectations of the log of the likelihood

Expectations of the log of the var. distributions (= their entropies)

- Substituting the priors, likelihoods, and variational distributions

$$\begin{aligned}\mathcal{L}(a', b', \mu', \Sigma') &= (a' - 1)(\psi(a') - \ln b') - b' \frac{a'}{b'} + \text{constant} - \frac{\lambda}{2}(\mu'^T \mu' + \text{tr}(\Sigma')) + \text{constant} + \frac{N}{2}(\psi(a') - \ln b') - \sum_{i=1}^N \frac{1}{2} \frac{a'}{b'} \left( (y_i - x_i^T \mu')^2 + x_i^T \Sigma' x_i \right) + \text{constant} \\ &\quad + a' - \ln b' + \ln \Gamma(a') + (1 - a')\psi(a') + \frac{1}{2} \ln |\Sigma'| + \text{constant}\end{aligned}$$

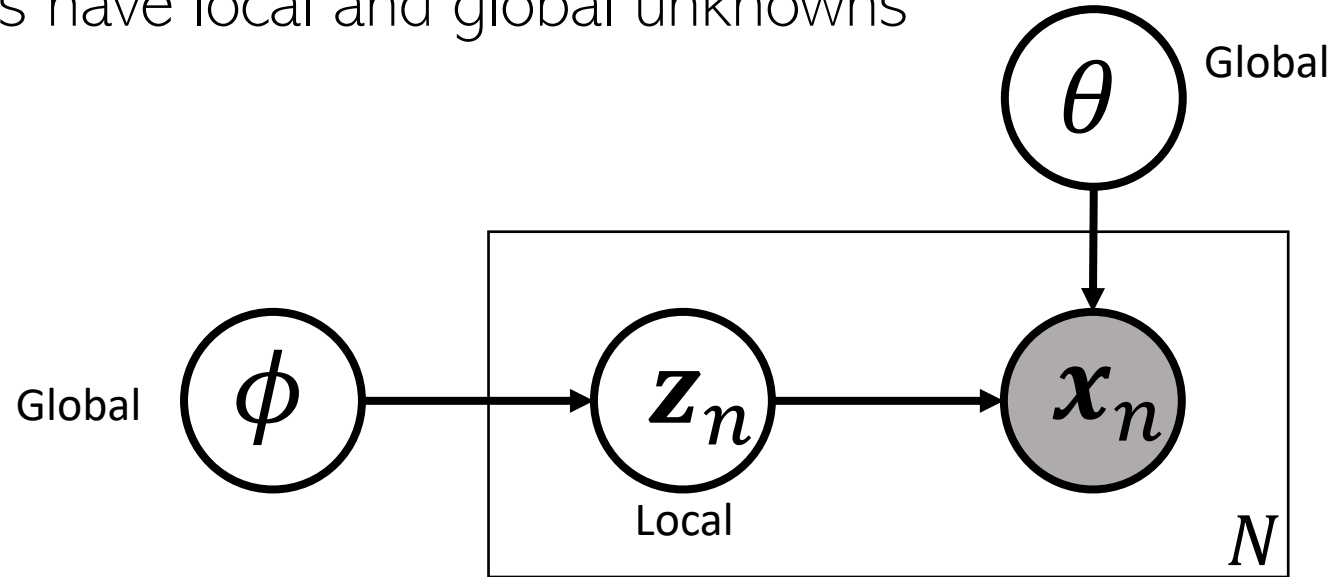
Digamma function (log of gamma function)

- Can now maximize the above ELBO w.r.t.  $a', b', \mu', \Sigma'$  in an alternating fashion
- For most models, ELBO or its gradients won't have a simple form (methods like "black-box" variational inference, reparametrization trick etc will be needed in those cases)

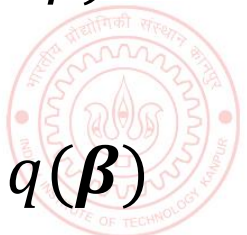


# MFVI for LVMs with Local and Global Unknowns

- Many LVMs have local and global unknowns



- Examples: Gaussian Mixture Model, Prob. PCA, Variational Autoencoder (VAE), etc
- Denote all local unknowns  $\{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N\}$  as  $\mathbf{Z}$  and global unknown as  $\boldsymbol{\beta} = (\theta, \phi)$
- The goal is to infer the posterior  $p(\mathbf{Z}, \boldsymbol{\beta} | \mathbf{X})$  which is intractable in general
- Mean-field VI will approximate this posterior as  $p(\mathbf{Z}, \boldsymbol{\beta} | \mathbf{X}) \approx q(\mathbf{Z}, \boldsymbol{\beta}) \approx q(\mathbf{Z})q(\boldsymbol{\beta})$

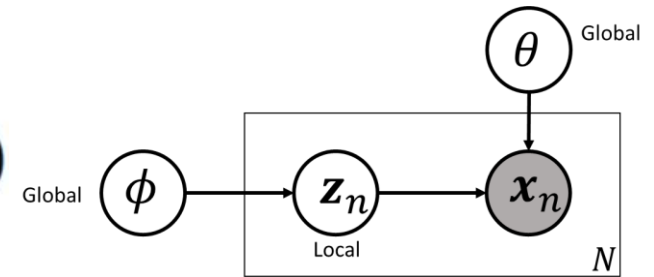




# MFVI for LVMs with Local and Global Unknowns

- Assuming independence, the joint distribution of data  $\mathbf{X}$  and unknowns  $\boldsymbol{\beta} = (\theta, \phi)$

$$p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\beta}) = p(\boldsymbol{\beta}) \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\beta}) p(\mathbf{z}_n | \boldsymbol{\beta}) = p(\boldsymbol{\beta}) \prod_{n=1}^N p(\mathbf{x}_n, \mathbf{z}_n | \boldsymbol{\beta})$$



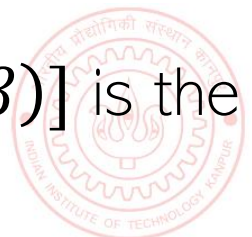
- Assume the joint dist. of  $\mathbf{x}_n$  and  $\mathbf{z}_n$  to be an exp-fam dist with natural params  $\boldsymbol{\beta}$

$$p(\mathbf{x}_n, \mathbf{z}_n | \boldsymbol{\beta}) = h(\mathbf{x}_n, \mathbf{z}_n) \exp \left[ \boldsymbol{\beta}^\top \overset{\text{Sufficient statistics}}{t(\mathbf{x}_n, \mathbf{z}_n)} - A(\boldsymbol{\beta}) \right]$$

- Assume a prior on  $\boldsymbol{\beta}$ , that is conjugate to the above exp-fam dist

$$p(\boldsymbol{\beta} | \boldsymbol{\alpha}) = h(\boldsymbol{\beta}) \exp \left[ \boldsymbol{\alpha}^\top [\boldsymbol{\beta}, -A(\boldsymbol{\beta})] - A(\boldsymbol{\alpha}) \right]$$

where  $\boldsymbol{\alpha} = [\alpha_1, \alpha_2]^\top$  are the hyperparameters of the prior  $p(\boldsymbol{\beta})$  and  $[\boldsymbol{\beta}, -A(\boldsymbol{\beta})]$  is the sufficient statistics vector for this exp-family distribution



# MFVI for LVMs with Local and Global Vars

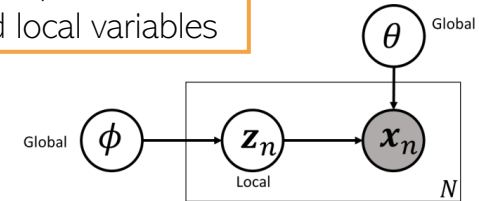
- Recall that mean-field VB can be obtained using CP of each unknown
- Optimal var. distribution for each unknown requires *exp. of nat. params of its CP*

$$p(\beta|\alpha) = h(\beta) \exp [\alpha^\top [\beta, -A(\beta)] - A(\alpha)]$$

- Due to conj, CP of global vars  $\beta = (\theta, \phi)$ , will have the same form as prior  $p(\beta|\alpha)$

$$p(\beta|\mathbf{X}, \mathbf{Z}) = p(\beta|\hat{\alpha}) \quad \text{where} \quad \hat{\alpha} = \left[ \alpha_1 + \sum_{n=1}^N t(\mathbf{x}_n, \mathbf{z}_n), \alpha_2 + N \right]$$

Updates to the natural parameters requires a summing suff-stats over all data and local variables



- Likewise, CP of each local variable  $\mathbf{z}_n$

Due to the independence structure

Assuming CP is an exp-fam distribution (will be the case if the prior  $p(\mathbf{z}_n|\phi)$  and likelihood  $p(\mathbf{x}_n|\mathbf{z}_n, \theta)$  are exp-family and conjugate to each other)

$$p(\mathbf{z}_n|\mathbf{Z}_{-n}, \mathbf{X}, \beta) = p(\mathbf{z}_n|\mathbf{x}_n, \beta) = h(\mathbf{z}_n) \exp [\eta(\mathbf{x}_n, \beta)^\top \mathbf{z}_n - A(\eta(\mathbf{x}_n, \beta))]$$

Nat. params depends on data  $\mathbf{x}_n$  and global var  $\beta$

- Having these CPs, we can compute the mean-field updates for  $q(\beta)$  and  $q(\mathbf{z}_n)$



# MFVI for LVMs with Local and Global Vars

- Let's assume our mean-field approximation to be of the form

$$q(\boldsymbol{\beta}, \mathbf{Z}) = q(\boldsymbol{\beta}|\boldsymbol{\lambda}) \prod_{n=1}^N q(\mathbf{z}_n|\phi_n)$$

- CPs are exp-fam, so optimal  $q$ 's depend on expected suff-stats of CP's nat. params

- The optimal variational dist. for local vars  $\mathbf{z}_n$  will be  $q(\mathbf{z}_n|\phi_n)$  with

Basically requires expectation over the  $q(\boldsymbol{\beta}|\boldsymbol{\lambda})$  distribution

$$\phi_n = \mathbb{E}_{\boldsymbol{\lambda}} [\boldsymbol{\eta}(\mathbf{x}_n, \boldsymbol{\beta})] \quad \forall n$$

- The optimal variational dist. for global vars  $\boldsymbol{\beta}$  will be  $q(\boldsymbol{\beta}|\boldsymbol{\lambda})$  with

Basically requires expectation over the  $q(\mathbf{z}_n|\phi_n)$  distribution

$$\boldsymbol{\lambda} = \left[ \alpha_1 + \sum_{n=1}^N \mathbb{E}_{\phi_n} [t(\mathbf{x}_n, \mathbf{z}_n)], \alpha_2 + N \right]^T$$

- Mean-Field updates alternate between estimating  $\phi_n$ 's and  $\boldsymbol{\lambda}$  until convergence
- Potential bottleneck: Updating  $\boldsymbol{\lambda}$  requires waiting for all  $\phi_n$ 's to be updated (thus slow)
  - Can be handled using [online VI \(stochastic VI\)](#)



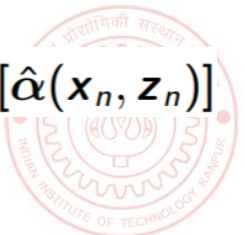
# Stochastic Variational Inference (SVI)

- An “online” algorithm<sup>†</sup> to speed-up VI for LVMs with local and global variables
- Recall the mean-field VI updates ( $q(\boldsymbol{\beta}, \mathbf{Z}) = q(\boldsymbol{\beta}|\boldsymbol{\lambda}) \prod_{n=1}^N q(\mathbf{z}_n|\phi_n)$ ) for such models

$$\phi_n = \mathbb{E}_{\boldsymbol{\lambda}} [\eta(\mathbf{x}_n, \boldsymbol{\beta})] \quad \forall n \quad \text{and} \quad \boldsymbol{\lambda} = \left[ \alpha_1 + \sum_{n=1}^N \mathbb{E}_{\phi_n} [t(\mathbf{x}_n, \mathbf{z}_n)], \alpha_2 + N \right]^{\top} = \mathbb{E}_{\phi} [\hat{\boldsymbol{\alpha}}(\mathbf{X}, \mathbf{Z})]$$

Local var. params:  $\phi_n$   
 Nat. param of CP of  $\mathbf{z}_n$ :  $\eta(\mathbf{x}_n, \boldsymbol{\beta})$   
 Global var. params:  $\boldsymbol{\lambda}$   
 Slow; requires all local var params  $\phi_n$ 's to be computed already:  $\mathbb{E}_{\phi_n} [t(\mathbf{x}_n, \mathbf{z}_n)]$   
 Nat. param of CP of  $\boldsymbol{\beta}$ :  $\hat{\boldsymbol{\alpha}}(\mathbf{X}, \mathbf{Z})$

- SVI uses minibatches to make the global param  $\boldsymbol{\lambda}$  updates more efficient
  1. Initialize  $\boldsymbol{\lambda}$  randomly as  $\boldsymbol{\lambda}^{(0)}$  and set current iteration number as  $i = 1$
  2. Set the learning rate (decaying as) as  $\epsilon_i = (i + 1)^{-\kappa}$  where  $\kappa \in (0.5, 1]$
  3. Choose a data point  $n$  uniformly randomly, i.e.,  $n \sim \text{Uniform}(1, 2, \dots, N)$  Assuming minibatch size = 1
  4. Compute local var. param  $\phi_n$  for data point  $\mathbf{x}_n$  as  $\phi_n = \mathbb{E}_{\boldsymbol{\lambda}^{(i-1)}} [\eta(\mathbf{x}_n, \boldsymbol{\beta})]$
  5. Update  $\boldsymbol{\lambda}$  as  $\boldsymbol{\lambda}^{(i)} = (1 - \epsilon_i)\boldsymbol{\lambda}^{(i-1)} + \epsilon_i \boldsymbol{\lambda}_n$  where  $\boldsymbol{\lambda}_n = [\alpha_1 + \mathbb{E}_{\phi_n} [t(\mathbf{x}_n, \mathbf{z}_n)], \alpha_2 + 1]^{\top} = \mathbb{E}_{\phi_n} [\hat{\boldsymbol{\alpha}}(\mathbf{x}_n, \mathbf{z}_n)]$
  6. Set  $i = i + 1$ . If ELBO not converged, go to Step 2



# What is SVI Doing?

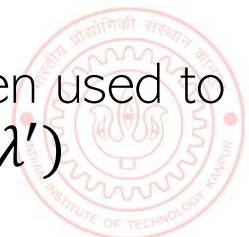
- SVI updates the global var params  $\lambda$  using **stochastic optimization**<sup>†</sup> of the ELBO
- However, instead of usual gradient of ELBO w.r.t.  $\lambda$ , SVI uses the **natural gradient**
- Denoting the double derivative of the log-partition function of CP of  $\beta$  as  $A''$

Usual gradient:  $\nabla_{\lambda} \text{ELBO} = A''(\lambda)(\mathbb{E}_{\phi}[\hat{\alpha}(\mathbf{X}, \mathbf{Z})] - \lambda)$  If interested in the proof, can see the derivation in the SVI paper

Natural gradient:  $g(\lambda) = A''(\lambda)^{-1} \times \nabla_{\lambda} \text{ELBO} = \mathbb{E}_{\phi}[\hat{\alpha}(\mathbf{X}, \mathbf{Z})] - \lambda$

Note:  $A''(\lambda)$  is cov. of suff-stats of CP of  $\beta$  and  $A''(\lambda)^{-1}$  is the Fisher information matrix

- Using the natural gradient has some nice advantages
  - Nat. grad. based updates of  $\lambda$  have simple form + easy to compute (no need to compute  $A''$ )
  - $\lambda^{(i)} = \lambda^{(i-1)} + \epsilon_i g(\lambda)|_{\lambda^{(i-1)}} = (1 - \epsilon_i)\lambda^{(i-1)} + \epsilon_i \mathbb{E}_{\phi}[\hat{\alpha}(\mathbf{X}, \mathbf{Z})]$  (assuming full batch)
  - Natural grad. are more intuitive/meaningful: Euclidean distance isn't often meaningful when used to compute distance between parameters of probability distributions, e.g.,  $q(\beta|\lambda)$  and  $q(\beta|\lambda')$



# SVI: Some Comments

- Often operates on minibatches: For iteration  $i$  minibatch  $\mathcal{B}_i$ , update  $\lambda$  as follows

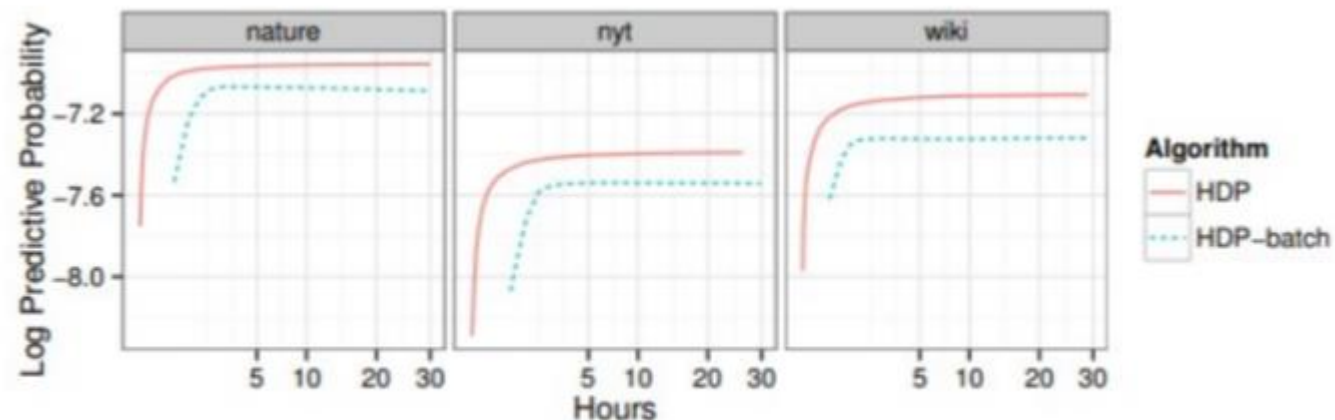
Global var. param computed on this minibatch

$$\hat{\lambda} = \frac{1}{|\mathcal{B}_i|} \sum_{n \in \mathcal{B}_i} \lambda_n$$

Now blending with the older estimate of  $\lambda$  from iteration  $i - 1$

$$\lambda^{(i)} = (1 - \epsilon_i) \lambda^{(i-1)} + \epsilon_i \hat{\lambda}$$

- Decaying learning rate  $\epsilon_i$  is necessary for convergence (need  $\sum_i \epsilon_i = \infty$  and  $\sum_i \epsilon_i^2 < \infty$ )
- SVI successfully used on many large-scale problems (topic modeling, citation network analysis, etc). Much faster convergence (and better results) compared to batch VI



SVI vs Batch VI on a nonparametric Bayesian Topic Model  
(Hierarchical Dirichlet Process)





# Coming Up Next

- VI for non-conjugate models
  - Black-box VI (BBVI) for general purpose VI
  - Reparametrization Trick for general purpose VI
  - Some model-specific tricks
- Other recent advances in VI
  - Amortized VI
  - Structured VI

