# Variational Inference (contd)

CS772A: Probabilistic Machine Learning Piyush Rai

Secap: Variational Inference (VI)  
Variational distribution  
Variational parameters  
Assuming 
$$p(Z|\mathcal{D}, \Theta)$$
 is intractable, VI approximates it by a distr  $q(Z|\phi)$  or  $q_{\phi}(Z)$   
L minimization  
 $\phi^* = \operatorname{argmin}_{\phi} \operatorname{KL}[q_{\phi}(Z)||p(Z|\mathcal{D}, \Theta)]$   
ELO  
maximization  
 $\phi^* = \operatorname{argmax}_{\phi} \mathbb{E}_{q_{\phi}(Z)}[\log p(\mathcal{D}|Z, \Theta)] - \operatorname{KL}[q_{\phi}(Z)||p(Z|\Theta)]$   
 $= \operatorname{argmax}_{\phi} \mathbb{E}_{q_{\phi}(Z)}[\log p(\mathcal{D}, Z|\Theta) - \log q_{\phi}(Z)] = \operatorname{argmax}_{\phi} \mathcal{L}(\phi, \Theta)$   
Can use gradient-based optimization  
to learn the parameters of the  
variational distribution  
 $\phi_{t+1} = \phi_t + \eta_t \nabla_{\phi = \phi_t} \mathcal{L}(\phi, \Theta)$   
Mean-field  
assumption on  
the variational  
 $drifticty$   
 $q(Z|\phi) = \prod_{i=1}^{M} q(Z_i|\phi_i)$   
This for simple enough model,  
we using mean-field Vi, we can  
get optimal  $q$  "directly" withous  
 $expectations and expectations are the cases where these
expectations and by a distribution
 $f_i(Z_j) = \frac{\exp(\mathbb{E}_{i\neq j}[\log p(\mathcal{D}, Z|\Theta)])}{\int \exp(\mathbb{E}_{i\neq j}[\log p(\mathcal{D}, Z|\Theta)] + \operatorname{const}}$   
 $f_i(Z_j) = \frac{\exp(\mathbb{E}_{i\neq j}[\log p(\mathcal{D}, Z|\Theta)]}{\int \exp(\mathbb{E}_{i\neq j}[\log p(\mathcal{D}, Z|\Theta)] + \operatorname{const}}$   
 $f_i(z_j) = \frac{\exp(\mathbb{E}_{i\neq j}[\log p(\mathcal{D}, Z|\Theta)] + \operatorname{const}}{\int \exp(\mathbb{E}_{i\neq j}[\log p(\mathcal{D}, Z|\Theta)] + \operatorname{const}}$$ 

#### Example: Mean-field VI without ELBO Derivatives

- Consider data  $\mathbf{X} = \{x_1, x_2, \dots, x_N\}$  from a one-dim Gaussian  $\mathcal{N}(\mu, \tau^{-1})$
- Assume the following normal-gamma prior on  $\mu$  and au

 $p(\mu|\tau) = \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1})$   $p(\tau) = \text{Gamma}(\tau|a_0, b_0)$ 

- Posterior is also normal-gamma due to the jointly conjugate prior
- Let's still try mean-field VI for this model
- With mean-field assumption on the variational posterior  $q(\mu, \tau) = q_{\mu}(\mu)q_{\tau}(\tau)$

$$\log q_{\mu}^{*}(\mu) = \mathbb{E}_{q_{\tau}}[\log p(\mathbf{X}, \mu, \tau)] + \text{const}$$
  
$$\log q_{\tau}^{*}(\tau) = \mathbb{E}_{q_{\mu}}[\log p(\mathbf{X}, \mu, \tau)] + \text{const}$$

• In this example, the log-joint  $\log p(\mathbf{X}, \mu, \tau) = \log p(\mathbf{X}|\mu, \tau) + \log p(\mu|\tau) + \log p(\tau)$ . Thus

 $\log q_{\mu}^{*}(\mu) = \mathbb{E}_{q_{\tau}}[\log p(\mathbf{X}|\mu,\tau) + \log p(\mu|\tau)] + \text{const} \quad (\text{only keeping terms that involve } \mu)$  $\log q_{\tau}^{*}(\tau) = \mathbb{E}_{q_{\mu}}[\log p(\mathbf{X}|\mu,\tau) + \log p(\mu|\tau) + \log p(\tau)] + \text{const}$ 

No "latent variables" here. Data **X** is fully observed, and parameters  $\mu$ ,  $\tau$  need to be estimated

Assume the hyperparameters  $\mu_0, \lambda_0, a_0, b_0$  are known

Note that we aren't even specifying the forms of these two distributions! We'll be able <u>identify the forms</u> in a few steps after working with the expectations

Example: Mean-field VI without ELBO Derivatives

• Substituting  $p(\mathbf{X}|\mu,\tau) = \prod_{n=1}^{N} p(x_n|\mu,\tau)$  and  $p(\mu|\tau)$ , we get

$$\log q_{\mu}^{*}(\mu) = \mathbb{E}_{q_{\tau}}[\log p(\mathbf{X}|\mu,\tau) + \log p(\mu|\tau)] + \text{const}$$
$$= -\frac{\mathbb{E}_{q_{\tau}}[\tau]}{2} \left\{ \sum_{n=1}^{N} (x_{n}-\mu)^{2} + \lambda_{0}(\mu-\mu_{0})^{2} \right\} + \text{const}$$

• (Verify) The above is log of a Gaussian. This  $q_{\mu}^* = \mathcal{N}(\mu | \mu_N, \lambda_N^{-1})$  with

$$\mu_N = rac{\lambda_0 \mu_0 + N ar{x}}{\lambda_0 + N}$$
 and  $\lambda_N = (\lambda_0 + N) \mathbb{E}_{q_\tau}[\tau]^2$  This update depends on  $q_\tau$ 

CS772A: PML

• Proceeding in a similar way (verify), we can show that  $q_{\tau}^* = \text{Gamma}(\tau | a_N, b_N)$ 

$$a_N = a_0 + rac{N+1}{2}$$
 and  $b_N = b_0 + rac{1}{2}\mathbb{E}_{q_\mu}\left[\sum_{n=1}^N (x_n - \mu)^2 + \lambda_0(\mu - \mu_0)^2\right]^2$  This update depends on  $q_\mu$ 

• Note: Updates of  $q_{\mu}^{*}$  and  $q_{\tau}^{*}$  depend on each other (hence alternating updates needed)

Mean-Field VI for Locally Conjugate Models

- Since  $\log q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j}[\log p(\mathbf{X}, \mathbf{Z})] + \text{const} = \mathbb{E}_{i \neq j}[\log p(\mathbf{X}, \mathbf{Z}_j, \mathbf{Z}_{-j})] + \text{const}$  $\log q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j}[\log p(\mathbf{Z}_j | \mathbf{X}, \mathbf{Z}_{-j})] + \text{const}$
- Thus finding optimal  $q_j^*(Z_j)$  only requires expectations of params of CP  $p(Z_j|X,Z_{-j})$
- For locally conjugate models, we know CP is easy and is an exp-fam distr of the form

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

• 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  $Z_j$ 



#### Variational EM

- In LVMs, latent vars Z and parameters  $\Theta$  <u>both</u> may be unknown. In such cases, we can use variational EM (VEM). Same as EM except VEM uses VI to approx. CP of Z
- VEM alternates between the following two steps
  - Maximize the ELBO w.r.t.  $\phi$  (gives the variational approximation q(Z) of CP of Z)

 $\phi^{(t)} = \operatorname{argmax}_{\phi} \mathbb{E}_{q_{\phi}(Z)} \left[ \log p(\mathcal{D}, Z | \Theta^{(t-1)}) - \log q_{\phi}(Z) \right]$ 

• Maximize the ELBO w.r.t.  $\Theta$  (gives us point estimate of  $\Theta$ )

 $\Theta^{(t)} = \operatorname{argmax}_{\Theta} \mathbb{E}_{q_{\phi^{(t)}}(Z)} \left[ \log p(\mathcal{D}, Z | \Theta) - \log q_{\phi^{(t)}}(Z) \right]$ 

 $= \operatorname{argmax}_{\Theta} \mathbb{E}_{q_{\phi}(t)(Z)}[\log p(\mathcal{D}, Z|\Theta)] \stackrel{\text{This looks very similar to the}}{=} \operatorname{expected CLL with the CP rep}$ 

expected CLL with the CP replaced by its variational approximation

• Note: If we want posterior for  $\Theta$  as well, treat it similar to Z and apply variational approximation (instead of using VEM) if the posterior isn't tractable

# VI for models <u>without</u> "latent variables"

Recall the Gaussian mean and variance estimation problem

- Suppose we have a "fully observed" case (no missing data/latent variables but just some unknown global parameters  $\theta$  and known hyperparams  $\xi$ )
- A simple example of the model is shown in the figure below



- If  $\xi$  are also unknown then one way would be to alternate like Variational EM
  - Approximating the CP  $p(\theta | \mathcal{D}, \xi)$  using VI
  - Using MLE-II to get point estimates of the hyperparameters  $\pmb{\xi}$

# Making VI Faster for LVMs: Stochastic VI (SVI)

- Many LVMs have local latent variables  $Z = \{z_1, z_2, ..., z_N\}$  and global params  $\Theta$
- VI updates of local and global variables <u>depend on each other</u> (similar to EM)
- This makes things slow (for VI and also for EM) especially when N is large
  - We must update  $q(\mathbf{z}_n | \phi_n)$ , i.e., compute  $\phi_n$ , for each latent variable before updating  $\Theta$
- Also need all the data  $X = \{x_1, x_2, \dots, x_N\}$  in memory to do these updates
- Stochastic VI\* is an efficient way using minibatches of data
- In each iteration, SVI takes a minibatch  $\mathcal{B}$  of  $|\mathcal{B}| \ll N$  data points, updates  $q(\mathbf{z}_n | \phi_n)$  examples in that minibatches and approximates the ELBO as follows



#### Making VI Faster for LVMs: Amortized VI

Instead of computing the optimal  $\phi_n$  for each  $q(z_n | \phi_n)$ , learn a function to do so

 $q(z_n|\phi_n) \approx q(z_n|\hat{\phi}_n)$  where  $\hat{\phi}_n = NN_{\phi}(x_n)$ 

- ullet Function is usually a neural network with weights  $\phi$ 
  - Usually referred to as "inference network" or "recognition model"
- Amortization: We are shifting the cost of finding  $\phi_n$  for each data point to finding the weights  $\phi$  of the neural network shared by all data points
- Can also combine amortized VI with stochastic VI
  - Each iteration only uses a minibatch to optimize NN weights  $\phi$  and global params  $\Theta$
- ELBO expression remains the same but  $q(z_n | \phi_n)$  is replaced by  $q(z_n | NN_{\phi}(x_n))$
- Amortized VI quality can be poor but it is fast and can give a quick solution
  - We can refine this solution other methods (e.g., using sampling; will see later)
  - This refinement based approach is called "semi-amortized VI"

# VI using ELBO's gradients

- For simple locally conjugate models, VI updates are usually easy
  - Sometimes, can find the optimal q even without taking the ELBO's gradients
- For complex models, we have to use the more general gradient-based approach.
- Consider the setting when we have latent variables  $\mathbf{Z}$  and parameters  $\boldsymbol{\Theta}$
- The ELBO's gradient w.r.t.  $\Theta$

 $\nabla_{\Theta} \mathcal{L}(\phi, \Theta) = \nabla_{\Theta} \mathbb{E}_{q_{\phi}(Z)}[\log p(\mathcal{D}, Z | \Theta) - \log q_{\phi}(Z)]$ Monte-Carlo approximation using samples of  $q_{\phi}(z)$  is  $= \mathbb{E}_{q_{\phi}(z)} \left[ \nabla_{\Theta} \left\{ \log p(\mathcal{D}, Z | \Theta) - \log q_{\phi}(Z) \right\} \right]$ straightforward here

Gradient can go inside expectation since q(Z)doesn't depend on  $\Theta$ 

• The ELBO's gradient w.r.t.  $\phi$ 

Gradient can't go inside  $\nabla_{\phi} \mathcal{L}(\phi, \Theta) = \nabla_{\phi} \mathbb{E}_{q_{\phi}(Z)}[\log p(\mathcal{D}, Z | \Theta) - \log q_{\phi}(Z)]$ Monte-Carlo approximation  $\neq \mathbb{E}_{q_{\phi}(\mathbf{Z})} \Big[ \nabla_{\phi} \{ \log p(\mathbf{\mathcal{D}}, \mathbf{Z} | \Theta) - \log q_{\phi}(\mathbf{Z}) \} \Big]$ using samples of  $q_{\phi}(\mathbf{Z})$  is NOT as straightforward



# Black-Box Variational Inference (BBVI)

- Black-box Var. Inference\* (BBVI) approximates ELBO derivatives using Monte-Carlo
- Uses the following identity for the ELBO's derivative

$$\begin{aligned} \nabla_{\phi} \mathcal{L}(q) &= \nabla_{\phi} \mathbb{E}_{q}[\log p(\mathbf{X}, \mathbf{Z}) - \log q(\mathbf{Z}|\phi)] \\ &= \mathbb{E}_{q}[\nabla_{\phi} \log q(\mathbf{Z}|\phi)(\log p(\mathbf{X}, \mathbf{Z}) - \log q(\mathbf{Z}|\phi))] \quad (\text{proof on next slide}) \end{aligned}$$

- Thus ELBO gradient can be written solely in terms of expec. of gradient of  $\log q(Z|\phi)$ 
  - Required gradients don't depend on the model; only on chosen var. distribution (hence "black-box")
- Given S samples  $\{Z_s\}_{s=1}^S$  from  $q(Z|\phi)$ , we can get (noisy) gradient as follows

$$abla_{\phi}\mathcal{L}(q) pprox rac{1}{S} \sum_{s=1}^{S} 
abla_{\phi} \log q(\mathbf{Z}_{s}|\phi) (\log p(\mathbf{X}, \mathbf{Z}_{s}) - \log q(\mathbf{Z}_{s}|\phi))$$

Above is also called the "score function" based gradient (also REINFORCE method)

Gradient of a log-likelihood or log-probability function w.r.t. its params is called score function; hence the name

# Proof of BBVI Identity

The ELBO gradient can be written as

$$\begin{aligned} \nabla_{\phi} \mathcal{L}(q) &= \nabla_{\phi} \int (\log p(\mathbf{X}, \mathbf{Z}) - \log q(\mathbf{Z}|\phi)) q(\mathbf{Z}|\phi) d\mathbf{Z} \\ &= \int \nabla_{\phi} [(\log p(\mathbf{X}, \mathbf{Z}) - \log q(\mathbf{Z}|\phi)) q(\mathbf{Z}|\phi)] d\mathbf{Z} \quad (\nabla \text{ and } \int \text{ interchangeable; dominated convergence theorem}) \\ &= \int \nabla_{\phi} [(\log p(\mathbf{X}, \mathbf{Z}) - \log q(\mathbf{Z}|\phi))] q(\mathbf{Z}|\phi) + \nabla_{\phi} q(\mathbf{Z}|\phi) [(\log p(\mathbf{X}, \mathbf{Z}) - \log q(\mathbf{Z}|\phi))] d\mathbf{Z} \\ &= \mathbb{E}_{q} [-\nabla_{\phi} \log q(\mathbf{Z}|\phi)] + \int \nabla_{\phi} q(\mathbf{Z}|\phi) [(\log p(\mathbf{X}, \mathbf{Z}) - \log q(\mathbf{Z}|\phi))] d\mathbf{Z} \end{aligned}$$

- Note that  $\mathbb{E}_{q}[\nabla_{\phi} \log q(\mathbf{Z}|\phi)] = \mathbb{E}_{q}\left[\frac{\nabla_{\phi}q(\mathbf{Z}|\phi)}{q(\mathbf{Z}|\phi)}\right] = \int \nabla_{\phi}q(\mathbf{Z}|\phi)d\mathbf{Z} = \nabla_{\phi}\int q(\mathbf{Z}|\phi)d\mathbf{Z} = \nabla_{\phi}1 = 0$
- Also note that  $\nabla_{\phi} q(\mathbf{Z}|\phi) = \nabla_{\phi} [\log q(\mathbf{Z}|\phi)] q(\mathbf{Z}|\phi)$ , using which

$$\int \nabla_{\phi} q(\mathbf{Z}|\phi) [(\log p(\mathbf{X}, \mathbf{Z}) - \log q(\mathbf{Z}|\phi))] d\mathbf{Z} = \int \nabla_{\phi} \log q(\mathbf{Z}|\phi) [(\log p(\mathbf{X}, \mathbf{Z}) - \log q(\mathbf{Z}|\phi))] q(\mathbf{Z}|\phi) d\mathbf{Z}$$
$$= \mathbb{E}_{q} [\nabla_{\phi} \log q(\mathbf{Z}|\phi) (\log p(\mathbf{X}, \mathbf{Z}) - \log q(\mathbf{Z}|\phi))]$$

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



# Benefits of BBVI

- Recall that BBVI approximates the ELBO gradients by the Monte Carlo expectations  $\nabla_{\phi} \mathcal{L}(q) \approx \frac{1}{S} \sum_{s=1}^{S} \nabla_{\phi} \log q(\mathbf{Z}_{s} | \phi) (\log p(\mathbf{X}, \mathbf{Z}_{s}) - \log q(\mathbf{Z}_{s} | \phi))$
- Enables applying VI for a wide variety of probabilistic models
- Can also work with small minibatches of data rather than full data
- BBVI has very few requirements
  - Should be able to sample from  $q(\mathbf{Z}|\boldsymbol{\phi})$  (usually sampling routines exists!)
  - Should be able to compute  $\nabla_{\phi} \log q(\mathbf{Z}|\phi)$  (automatic differentiation methods exist!)
  - Should be able to evaluate  $\log p(X, Z)$  and  $\log q(Z|\phi)$  for any value of Z
- Some tricks needed to control the variance in the Monte Carlo estimate of the ELBO gradient (if interested in the details, please refer to the BBVI paper)



#### Reparametrization Trick

- Another Monte-Carlo approx. of ELBO grad (with often lower var than BBVI gradient)
- Suppose we want to compute ELBO's gradient  $\nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{Z})}[\log p(\mathbf{X}, \mathbf{Z}) \log q_{\phi}(\mathbf{Z})]$
- Assume a deterministic transformation g

 $Z = g(\epsilon, \phi)$  where  $\epsilon \sim p(\epsilon)$  Assumed to not depend on  $\phi$ 

With this reparametrization, and using LOTUS rule, the ELBO's gradient would be

 $\nabla_{\phi} \mathbb{E}_{p(\epsilon)}[\log p(\mathbf{X}, g(\epsilon, \phi)) - \log q_{\phi}(g(\epsilon, \phi))] = \mathbb{E}_{p(\epsilon)} \nabla_{\phi}[\log p(\mathbf{X}, g(\epsilon, \phi)) - \log q_{\phi}(g(\epsilon, \phi))]$ 

• Given S i.i.d. random samples  $\{\epsilon_s\}_{s=1}^S$  from  $p(\epsilon)$ , we can get a Monte-Carlo approx.

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathsf{Z})}[\log p(\mathsf{X}, \mathsf{Z}) - \log q_{\phi}(\mathsf{Z})] \approx \frac{1}{S} \sum_{s=1}^{S} [\nabla_{\phi} \log p(\mathsf{X}, g(\epsilon_{s}, \phi)) - \nabla_{\phi} \log q_{\phi}(g(\epsilon_{s}, \phi))]$$

**CS772A: PML** 

• Such gradients are called pathwise gradients\* (since we took a "path" from  $\epsilon$  to Z)

#### **Reparametrization Trick: An Example**

- Suppose our variational distribution is  $q(w|\phi) = \mathcal{N}(w|\mu, \Sigma)$ , so  $\phi = \{\mu, \Sigma\}$
- Suppose our ELBO has a difficult expectation term  $\mathbb{E}_{a}[f(w)]$
- However, note that we need ELBO gradient, not ELBO itself. Let's use the trick
- Reparametrize w as  $w = \mu + Lv$  where  $v \sim \mathcal{N}(0, I) \frac{|v|}{|v|} \frac{|v|$

 $\nabla_{\mu,\mathsf{L}}\mathbb{E}_{\mathcal{N}(\boldsymbol{w}|\mu,\boldsymbol{\Sigma})}[f(\boldsymbol{w})] = \nabla_{\mu,\mathsf{L}}\mathbb{E}_{\mathcal{N}(\boldsymbol{v}|0,\mathsf{I})}[f(\mu+\mathsf{L}\boldsymbol{v})] = \mathbb{E}_{\mathcal{N}(\boldsymbol{v}|0,\mathsf{I})}[\nabla_{\mu,\mathsf{L}}f(\mu+\mathsf{L}\boldsymbol{v})]$ 

Or  $\boldsymbol{\phi} = \{\boldsymbol{\mu}, \mathbf{L}\}$ 

Often even one or very

 $\partial f \partial w$ 

 $\partial w \partial \mu$ 

∂f ∂w

**CS772A: PML** 

few samples suffice

where  $\mathbf{L} = \operatorname{chol}(\Sigma)$ 

- The above is now straightforward
  - Easily take derivatives of f(w) w.r.t. variational params  $\mu$ , L
  - Replace exp. by Monte-Carlo averaging using samples of **v** from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$

$$\nabla_{\mu} \mathbb{E}_{\mathcal{N}(\boldsymbol{w}|\mu,\Sigma)}[f(\boldsymbol{w})] = \mathbb{E}_{\mathcal{N}(\boldsymbol{v}|0,\mathbf{I})}[\nabla_{\mu}f(\mu + \mathbf{L}\boldsymbol{v})] \approx \nabla_{\mu}f(\mu + \mathbf{L}\boldsymbol{v}_{s})$$
Chain Rule
$$\nabla_{\mathbf{L}} \mathbb{E}_{\mathcal{N}(\boldsymbol{w}|\mu,\Sigma)}[f(\boldsymbol{w})] = \mathbb{E}_{\mathcal{N}(\boldsymbol{v}|0,\mathbf{I})}[\nabla_{\mathbf{L}}f(\mu + \mathbf{L}\boldsymbol{v})] \approx \nabla_{\mathbf{L}}f(\mu + \mathbf{L}\boldsymbol{v}_{s})$$
Chain Rule
$$\frac{\partial f}{\partial \boldsymbol{w}} \frac{\partial f}{\partial \boldsymbol{$$

∂w ∂L Std. reparam. trick assumes differentiability (recent work on removing this

#### Reparametrization Trick: Some Comments

- Standard Reparametrization Trick assumes the model to be differentiable  $\nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathsf{Z})}[\log p(\mathsf{X}, \mathsf{Z}) - \log q_{\phi}(\mathsf{Z})] = \mathbb{E}_{p(\epsilon)}[\nabla_{\phi} \log p(\mathsf{X}, g(\epsilon, \phi)) - \nabla_{\phi} \log q_{\phi}(g(\epsilon, \phi))]$
- In contrast, BBVI (score function gradients) only required q(Z) to be differentiable
- Thus rep. trick often isn't applicable, e.g., when Z is discrete (e.g., binary /categorical)
  - Recent work on continuous relaxation<sup>†</sup> of discrete variables<sup>†</sup>(e.g., Gumbel Softmax for categorical)
- ullet The transformation function g may be difficult to find for general distributions
  - Recent work on generalized reparametrizations\*
- Also, the transformation function g needs to be invertible (difficult/expensive)
  - Recent work on implicit reparametrized gradients<sup>#</sup>

• Assumes that we can directly draw samples from  $p(\epsilon)$ . If not, then rep. trick isn't valid<sup>@</sup>

<sup>†</sup>Categorical Reparameterization with Gumbel-Softmax (Jang et al, 2017), \* The Generalized Reparameterization Gradient (Ruiz et al, 2016), # Implicit Reparameterization Gradients (Figurnov et al, 2018), @ Reparameterization Gradients through Acceptance-Rejection Sampling Algorithms (Naesseth et al, 2016)

# Automatic Differentiation Variational Inference

- Suppose Z is D-dim r.v. with constraints (e.g., non-negativity) and distribution  $q(Z|\phi)$
- Assume a transformation T such that u = T(Z) s.t.  $u \in \mathbb{R}^D$  (unconstrained) then

$$q(\boldsymbol{u}) = q(\boldsymbol{Z}) \left| \det\left(\frac{\partial \boldsymbol{Z}}{\partial \boldsymbol{u}}\right) \right|$$

- Recall the original ELBO expression  $\mathcal{L}(\phi) = \mathbb{E}_{q(Z|\phi)} \left[ \log \frac{p(D,Z)}{q(Z|\phi)} \right]$
- Assuming  $q(\boldsymbol{u}|\boldsymbol{\psi}) = \mathcal{N}(\boldsymbol{u}|\boldsymbol{\mu},\boldsymbol{\Sigma})$  and using  $\boldsymbol{Z} = T^{-1}(\boldsymbol{u})$ , the ELBO becomes

$$\mathcal{L}(\psi) = \mathbb{E}_{q(\boldsymbol{u}|\psi)} \left[ \log \frac{p(\boldsymbol{D}, T^{-1}(\boldsymbol{u})) |\det(\partial \boldsymbol{Z}/\partial \boldsymbol{u})|}{q(\boldsymbol{u}|\psi)} \right]$$

 $= \mathbb{E}_{q(\boldsymbol{u}|\boldsymbol{\psi})}[\log p(\boldsymbol{D}, T^{-1}(\boldsymbol{u})) + \log|\det(\partial \boldsymbol{Z}/\partial \boldsymbol{u})|] + H(q(\boldsymbol{u}|\boldsymbol{\psi}))$ 

17

CS772A: PML

• Optimize  $\mathcal{L}(\psi)$  w.r.t.  $\psi$  to get  $q(u|\psi)$  as a Gaussian and use distribution transformation equation to get  $q(Z|\phi)$ 

#### Structured Variational Inference

- Here "structured" may refer to anything that makes VI approx. more expressive, e.g.,
  - Removing the independence assumption of mean-field VI
  - In general, learning more complex forms for the variational approximation family  $q(\pmb{Z}|\pmb{\phi})$
- To remove the mean-field assumption in VI, various approaches exist
  - Structured mean-field (Saul et al, 1996)
  - Hierarchical VI (Ranganath et al, 2016): Variational params  $\phi_1, \phi_2, \dots, \phi_M$  "tied" via a shared prior

$$q(\mathbf{z}_1,\ldots,\mathbf{z}_M|\theta) = \int \left[\prod_{m=1}^M q(\mathbf{z}_m|\phi_m)\right] p(\phi|\theta) d\phi$$

- Recent work on learning more expressive variational approx. for general VI
  - Boosting or mixture of simpler distributions, e.g.,  $q(Z) = \sum_{c=1}^{C} \rho_c q_c(Z)$  Even simple unimodal components will give a multimodal q(Z)
  - Normalizing flows\*: Turn a simple var. distr. into a complex one via series of invertible transfor.

A much more complex (e.g., multimodal) variational distribution obtained via the flow idea

#### Other Divergence Measures

- VI minimizes KL(q||p) but other divergences can be minimized as well
  - Recall that VI with minimization of KL(q||p) leads to underestimated variances
- A general form of divergence is Renyi's  $\alpha$ -divergence defined as

$$D_{\alpha}^{R}(p(\boldsymbol{Z})||q(\boldsymbol{Z})) = \frac{1}{\alpha - 1} \log \int p(\boldsymbol{Z})^{\alpha} q(\boldsymbol{Z})^{1 - \alpha} d\boldsymbol{Z}$$

- KL(p||q) is a special case with  $\alpha \to 1$  (can verify using L'Hopital rule of taking limits)
- An even more general form of divergence is f-Divergence

$$D_f(p(\mathbf{Z})||q(\mathbf{Z})) = \int q(\mathbf{Z}) f\left(\frac{p(\mathbf{Z})}{q(\mathbf{Z})}\right) d\mathbf{Z}$$

Many recent variational inference algorithms are based on minimizing such divergences

# Variational Inference: Some Comments

- Many probabilistic models nowadays rely on VI to do approx. inference
- Even mean-field with locally-conjugacy used in lots of models
  - This + SVI gives excellent scalability as well on large datasets
- Progress in various areas has made VI very popular and widely applicable
  - Stochastic Optimization (e.g., SGD)
  - Automatic Differentiation
  - Monte-Carlo gradient of ELBO
- Note: Most of these ideas apply also to Variational EM
- Many VI and advanced VI algos are implemented in probabilistic prog. packages (e.g., Tensorflow Probability, PyTorch, etc), making VI easy even for complex models
- Still a very active area of research, especially for doing VI in complex models
  - Models with discrete latent variables
  - Reducing the variance in Monte-Carlo estimate of ELBO gradients
  - More expressive variational distribution for better approximation





