

# Variational Inference (contd)

CS772A: Probabilistic Machine Learning

Piyush Rai

# Recap: Variational Inference (VI)

Variational  
distribution

Variational  
parameters

- Assuming  $p(\mathbf{Z}|\mathcal{D}, \Theta)$  is intractable, VI approximates it by a distr  $q(\mathbf{Z}|\phi)$  or  $q_\phi(\mathbf{Z})$

KL minimization

$$\phi^* = \operatorname{argmin}_\phi \text{KL}[q_\phi(\mathbf{Z}) || p(\mathbf{Z}|\mathcal{D}, \Theta)]$$

ELBO  
maximization

$$\begin{aligned} \phi^* &= \operatorname{argmax}_\phi \mathbb{E}_{q_\phi(\mathbf{Z})} [\log p(\mathcal{D}|\mathbf{Z}, \Theta)] - \text{KL}[q_\phi(\mathbf{Z}) || p(\mathbf{Z}|\Theta)] \\ &= \operatorname{argmax}_\phi \mathbb{E}_{q_\phi(\mathbf{Z})} [\log p(\mathcal{D}, \mathbf{Z}|\Theta) - \log q_\phi(\mathbf{Z})] = \operatorname{argmax}_\phi \mathcal{L}(\phi, \Theta) \end{aligned}$$

Can use gradient-based optimization to learn the parameters of the variational distribution

$$\phi_{t+1} = \phi_t + \eta_t \nabla_{\phi_t} \mathcal{L}(\phi, \Theta)$$

Mean-field  
assumption on  
the variational  
distribution

$$q(\mathbf{Z}|\phi) = \prod_{i=1}^M q(\mathbf{Z}_i|\phi_i)$$

$\mathbb{E}_{i \neq j}$  denotes expectations w.r.t.  $\prod_{i \neq j} q(\mathbf{Z}_i|\phi_i)$

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

Equivalent to writing  $\log q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\log p(\mathcal{D}, \mathbf{Z}|\Theta)] + \text{const}$



# Variational EM

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

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

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

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

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

This looks very similar to the expected CLL with the CP replaced by its variational approximation

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



# VI for models without “latent variables”

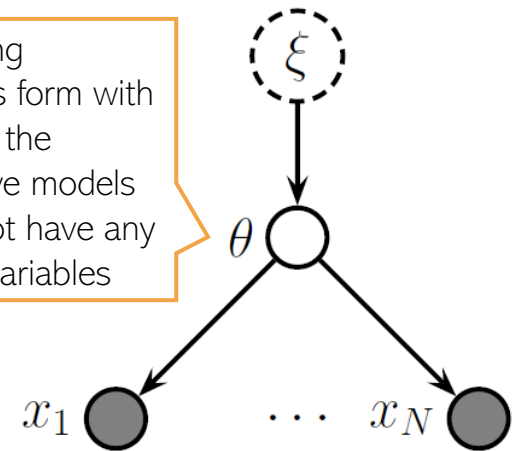
- Suppose we have a “fully observed” case (no missing data/latent variables but just some unknown global parameters  $\theta$  and known hyperparams  $\xi$ )
- A simple example of the model is shown in the figure below

$$p(\mathcal{D}, \theta | \xi) = p(\theta | \xi) \prod_{n=1}^N p(x_n | \theta)$$

If this CP is intractable, we can use VI to approximate this

$$p(\theta | \mathcal{D}, \xi) = \frac{p(\mathcal{D} | \theta) p(\theta | \xi)}{p(\mathcal{D} | \xi)}$$

Even supervised learning problems may have this form with  $\theta$  being the weights of the generative/discriminative models and the models may not have any missing data or latent variables



- 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  $\xi$



# Example: Mean-field VI without ELBO Derivatives

No “latent variables” here. Data  $\mathbf{X}$  is fully observed, and parameters  $\mu, \tau$  need to be estimated

- 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  $\tau$

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

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

- Posterior is also normal-gamma due to the jointly conjugate prior
- Let's still try mean-field VI for this model

Note that we aren't specifying the forms of these two distributions

- 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

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

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



# 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

$$\begin{aligned}\log q_\mu^*(\mu) &= \mathbb{E}_{q_\tau}[\log p(\mathbf{X}|\mu, \tau) + \log p(\mu|\tau)] + \text{const} \\ &= -\frac{\mathbb{E}_{q_\tau}[\tau]}{2} \left\{ \sum_{n=1}^N (x_n - \mu)^2 + \lambda_0(\mu - \mu_0)^2 \right\} + \text{const}\end{aligned}$$

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

$$\mu_N = \frac{\lambda_0 \mu_0 + N \bar{x}}{\lambda_0 + N} \quad \text{and} \quad \lambda_N = (\lambda_0 + N) \mathbb{E}_{q_\tau}[\tau]$$

This update depends on  $q_\tau$

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

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

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: A Closer Look

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

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

- 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, 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

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

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



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

- Many LVMs have local latent variables  $\mathbf{Z} = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N\}$  and global params  $\Theta$
- VI updates of local and global variables depend on each other (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  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$  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

Optimize this approximate ELBO w.r.t.  $\Theta$   
(note: this is an unbiased estimate\*)

$$\tilde{\mathcal{L}}(\phi, \Theta) = \frac{N}{B} \sum_{\mathbf{x}_i \in \mathcal{B}} \mathbb{E}_{q(\mathbf{z}_i|\phi_i)} [\log p(\mathbf{x}_i|\mathbf{z}_i, \Theta)] - \text{KL}[q_\phi(\mathbf{z}_i) || p(\mathbf{z}_i|\Theta)]$$



# Making VI Faster for LVMs: Amortized VI

- Instead of computing the optimal  $\phi_n$  for each  $q(\mathbf{z}_n|\phi_n)$ , learn a function to do so

$$q(\mathbf{z}_n|\phi_n) \approx q(\mathbf{z}_n|\hat{\phi}_n) \quad \text{where} \quad \hat{\phi}_n = \text{NN}_\phi(\mathbf{x}_n)$$

- 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  $q(\mathbf{z}_n|\phi_n)$  is replaced by  $q(\mathbf{z}_n|\text{NN}_\phi(\mathbf{x}_n))$
- 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  $\Theta$
- The ELBO's gradient w.r.t.  $\Theta$

$$\nabla_{\Theta} \mathcal{L}(\phi, \Theta) = \nabla_{\Theta} \mathbb{E}_{q_{\phi}(\mathbf{Z})} [\log p(\mathcal{D}, \mathbf{Z} | \Theta) - \log q_{\phi}(\mathbf{Z})]$$

Monte-Carlo approximation using samples of  $q_{\phi}(\mathbf{Z})$  is straightforward here

$$= \mathbb{E}_{q_{\phi}(\mathbf{Z})} [\nabla_{\Theta} \{\log p(\mathcal{D}, \mathbf{Z} | \Theta) - \log q_{\phi}(\mathbf{Z})\}]$$

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

- The ELBO's gradient w.r.t.  $\phi$

$$\nabla_{\phi} \mathcal{L}(\phi, \Theta) = \nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{Z})} [\log p(\mathcal{D}, \mathbf{Z} | \Theta) - \log q_{\phi}(\mathbf{Z})]$$

Monte-Carlo approximation using samples of  $q_{\phi}(\mathbf{Z})$  is NOT as straightforward

$$\neq \mathbb{E}_{q_{\phi}(\mathbf{Z})} [\nabla_{\phi} \{\log p(\mathcal{D}, \mathbf{Z} | \Theta) - \log q_{\phi}(\mathbf{Z})\}]$$

Gradient can't go inside expectation since  $q(\mathbf{Z})$  depends on  $\phi$



# 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  $\log q(\mathbf{Z}|\phi)$ 
  - Required gradients don't depend on the model; only on chosen **var. distribution** (hence “black-box”)
- Given  $S$  samples  $\{\mathbf{Z}_s\}_{s=1}^S$  from  $q(\mathbf{Z}|\phi)$ , we can get (noisy) gradient as follows

$$\nabla_{\phi} \mathcal{L}(q) \approx \frac{1}{S} \sum_{s=1}^S \nabla_{\phi} \log q(\mathbf{Z}_s|\phi)(\log p(\mathbf{X}, \mathbf{Z}_s) - \log q(\mathbf{Z}_s|\phi))$$

- 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

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

- 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  $q(\mathbf{Z} | \phi)$  (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  $\log p(\mathbf{X}, \mathbf{Z})$  and  $\log q(\mathbf{Z} | \phi)$  for any value of  $\mathbf{Z}$
- 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  $g$

$$\mathbf{Z} = g(\epsilon, \phi) \quad \text{where} \quad \epsilon \sim p(\epsilon)$$

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.

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

- Such gradients are called **pathwise gradients**\* (since we took a “path” from  $\epsilon$  to  $\mathbf{Z}$ )





# Reparametrization Trick: An Example

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

Or  $\phi = \{\mu, \mathbf{L}\}$   
where  $\mathbf{L} = \text{chol}(\Sigma)$

Note that we will still have  
 $q(\mathbf{w}|\phi) = \mathcal{N}(\mathbf{w}|\mu, \Sigma)$

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

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

Often even one or very few samples suffice

$$\begin{aligned} \nabla_{\mu} \mathbb{E}_{\mathcal{N}(\mathbf{w}|\mu, \Sigma)}[f(\mathbf{w})] &= \mathbb{E}_{\mathcal{N}(\mathbf{v}|\mathbf{0}, \mathbf{I})}[\nabla_{\mu} f(\mu + \mathbf{L}\mathbf{v})] \approx \nabla_{\mu} f(\mu + \mathbf{L}\mathbf{v}_s) \\ \nabla_{\mathbf{L}} \mathbb{E}_{\mathcal{N}(\mathbf{w}|\mu, \Sigma)}[f(\mathbf{w})] &= \mathbb{E}_{\mathcal{N}(\mathbf{v}|\mathbf{0}, \mathbf{I})}[\nabla_{\mathbf{L}} f(\mu + \mathbf{L}\mathbf{v})] \approx \nabla_{\mathbf{L}} f(\mu + \mathbf{L}\mathbf{v}_s) \end{aligned}$$

$$\frac{\partial f}{\partial \mathbf{w}} \frac{\partial \mathbf{w}}{\partial \mu}$$

Chain Rule

$$\frac{\partial f}{\partial \mathbf{w}} \frac{\partial \mathbf{w}}{\partial \mathbf{L}}$$

- Std. reparam. trick **assumes differentiability** (recent work on removing this req).

# Reparametrization Trick: Some Comments

- Standard Reparametrization Trick assumes the **model to be differentiable**

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

- In contrast, BBVI (score function gradients) only required  **$q(\mathbf{Z})$  to be differentiable**
- Thus rep. trick often isn't applicable, e.g., when  **$\mathbf{Z}$**  is discrete (e.g., binary /categorical)
  - Recent work on **continuous relaxation**<sup>†</sup> of discrete variables<sup>†</sup> (e.g., Gumbel Softmax for categorical)
- 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**#
- Assumes that we can **directly draw samples** from  **$p(\epsilon)$** . If not, then rep. trick isn't valid@

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

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

- Assuming  $q(\mathbf{u}|\psi) = \mathcal{N}(\mathbf{u}|\mu, \Sigma)$ , the ELBO becomes

Original ELBO for  $q(\mathbf{Z}|\phi)$

$$\mathcal{L}(\phi) = \mathbb{E}_{q(\mathbf{Z}|\phi)} [\log p(\mathcal{D}|\mathbf{Z}) + \log p(\mathbf{Z})] + H(q(\mathbf{Z}|\phi))$$

Transformed, equivalent ELBO

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

- We can optimize the above ELBO w.r.t.  $\psi$  to get  $q(\mathbf{u}|\psi)$  as a Gaussian
- The transformed density  $q(\mathbf{Z}|\phi)$  can be found using the transformation equation

