# Variational Inference (wrap-up), Approximate Inference via Sampling

CS772A: Probabilistic Machine Learning

Piyush Rai

### 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  $\Theta$ 

CS772A: PML

• The ELBO's gradient w.r.t.  $\phi$ 

Gradient can't go inside  $\nabla_{\phi} \mathcal{L}(\phi, \Theta) = \nabla_{\phi} \mathbb{E}_{q_{\phi}(Z)}[\log p(\mathcal{D}, Z | \Theta) - \log q_{\phi}(Z)]$ expectation since q(Z)depends on  $\phi$ 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

#### 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  $v \sim \mathcal{N}(0, I) \frac{|v|}{|v|} \frac{|v|$

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

- 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

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

#### 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

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

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

**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  $\mathbf{z}_{K} = f_{K} \circ \cdots \circ f_{1}(\mathbf{z}_{0}), \quad \mathbf{z}_{0} \sim q_{0}(\mathbf{z}_{0}), \quad \text{A simple unimodal variational distribution (e.g. <math>\mathcal{N}(0, I)$  $\mathbf{z}_{K} \sim q_{K}(\mathbf{z}_{K}) = q_{0}(\mathbf{z}_{0}) \prod_{k=1}^{K} \left| \det \frac{\partial f_{k}}{\partial \mathbf{z}_{k-1}} \right|^{-1}$ 



#### 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







### Sampling for Approximate Inference

Some typical tasks that we have to solve in probabilistic/fully-Bayesian inference

Posterior  
distribution 
$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)} = \frac{p(D|\theta)p(\theta)}{\int p(D|\theta)p(\theta)d\theta}$$
  
Posterior  
predictive  
distribution  $p(D^{new}|D) = \int p(D^{new}|\theta)p(\theta|D)d\theta = \mathbb{E}_{p(\theta|D)}[p(D^{new}|\theta)]$   
Needed for model  
selection (and in  
computing  
posterior too) Marginal  
marginal  
p(D|m) =  $\int p(D|\theta)p(\theta|m)d\theta = \mathbb{E}_{p(\theta|m)}[p(D|\theta)]$   
Needed in EM  
Expected  
complete data  
log-likelihood  $Exp-CLL = \int p(z|\theta, x)p(x, z|\theta)dz = \mathbb{E}_{p(z|\theta, x)}[p(x, z|\theta)]$   
Needed in V  
Evidence lower  
Needed in V  
Evidence lower  
Determine  $\mathcal{L}(q) = \mathbb{E}_{q}[\log p(x, z)] - \mathbb{E}_{q}[\log p(z)]$ 

Sampling methods provide a general way to (approximately) solve these problems

**CS772A: PML** 

More general than VI methods which only approximate the posterior distribution

# Approximating a Prob. Distribution using Samples<sup>12</sup>

Can approximate any distribution using a set of randomly drawn samples from it



- The samples can also be used for computing expectations (Monte-Carlo averaging)
- Usually straightforward to generate samples if it is a simple/standard distribution
- The interesting bit: Even if the distribution is "difficult" (e.g., an intractable posterior), it is often possible to generate random samples from such a distribution, as we will see.

### The Empirical Distribution

- Sampling based approx. can be formally represented using an empirical distribution
- Given L points/samples  $z^{(1)}, z^{(2)}, \dots, z^{(L)}$ , empirical distr. defined by these is





13

# Sampling: Some Basic Methods

- Most of these basic methods are based on the idea of transformation
  - Generate a random sample x from a distribution q(x) which is easy to sample from
  - Apply a transformation on x to make it random sample z from a complex distr p(z)
- Some popular examples of transformation methods
  - Inverse CDF method
    - $x \sim \text{Unif}(0, 1) \Rightarrow z = \text{Inv-CDF}_{p(z)}(x) \sim p(z)$
  - Reparametrization method

$$x \sim \mathcal{N}(0, 1) \Rightarrow z = \mu + \sigma x \sim \mathcal{N}(\mu, \sigma^2)$$

• Box-Mueller method: Given  $(x_1, x_2)$  from Unif(0, 1), generate  $(z_1, z_2)$  from  $\mathcal{N}(0, \mathbf{I}_2)$ 

$$z_1 = \sqrt{-2 \ln x_1} \cos(2\pi x_2), \ z_1 = \sqrt{-2 \ln x_1} \sin(2\pi x_2)$$

- Transformation Methods are simple but have limitations
  - Mostly limited to standard distributions and/or distributions with very few variables





## **Rejection Sampling**

- Goal: Generate a random sample from a distribution of the form  $p(z) = \frac{p(z)}{Z_p}$ , assuming
  - We can only <u>evaluate</u> the value of numerator  $\widetilde{p}(z)$  for any z
  - The denominator (normalization constant)  $Z_p$  is intractable and we don't know its value Should have the same support as p(z)
- Assume a proposal distribution q(z) we can generate samples from, and

 $Mq(z) \geq \tilde{p}(z)$   $\forall z$  (where M > 0 is some const.)

- Rejection Sampling then works as follows
  - Sample a random variable  $z_*$  from q(z)
  - Sampling a uniform r.v.  $u \sim \text{Unif}[0, Mq(z_*)]$
  - If  $u \leq \widetilde{p}(z_*)$  then accept  $z_*$ , otherwise reject it
- All accepted  $z_*$ 's will be random samples from p(z). Proof on next slide



### **Rejection Sampling**

- Why  $z \sim q(z)$  + accept/reject rule is equivalent to  $z \sim p(z)$ ?
- Let's look at the pdf of the z's that were accepted, i.e., p(z|accept)

$$p(\operatorname{accept}|z) = \int_{0}^{\tilde{p}(z)} \frac{1}{Mq(z)} du = \frac{\tilde{p}(z)}{Mq(z)}$$

$$p(z, \operatorname{accept}) = q(z)p(\operatorname{accept}|z) = \frac{\tilde{p}(z)}{M}$$

$$p(\operatorname{accept}) = \int \frac{\tilde{p}(z)}{M} dz = \frac{Z_{p}}{M}$$

$$p(z|\operatorname{accept}) = \frac{p(z, \operatorname{accept})}{p(\operatorname{accept})} = \frac{\tilde{p}(z)}{Z_{p}} = p(z)$$



# Computing Expectations via Monte Carlo Sampling<sup>1</sup>

Often we are interested in computing expectations of the form

$$\mathbb{E}[f] = \int f(z)p(z)dz$$

where f(z) is some function of the random variable  $z \sim p(z)$ 

- A simple approx. scheme to compute the above expectation: Monte Carlo integration
  - Generate L independent samples from  $p(z): \{z^{(\ell)}\}_{\ell=1}^{L} \sim p(z) \prec$  Assuming we know how to sample from p(z)
  - Approximate the expectation by the following empirical average

$$\mathbb{E}[f] \approx \hat{f} = \frac{1}{L} \sum_{\ell=1}^{L} f(z^{(\ell)})$$

Since the samples are independent of each other, we can show the following (exercise)

Unbiased  
expectation 
$$\mathbb{E}[\hat{f}] = \mathbb{E}[f]$$
 and  $\operatorname{var}[\hat{f}] = \frac{1}{L}\operatorname{var}[f] = \frac{1}{L}\mathbb{E}[(f - \mathbb{E}[f])^2]$  as *L* increases  
as *L* increases

Variance in our

# Computing Expectations via Importance Sampling<sup>18</sup>

- How to compute Monte Carlo expec. if we don't know how to sample from p(z)?
- One way is to use transformation methods or rejection sampling
- Another way is to use Importance Sampling (assuming p(z) can be <u>evaluated</u> at least)
  - Generate L indep samples from a proposal q(z) we know how sample from:  $\{z^{(\ell)}\}_{\ell=1}^{L} \sim q(z)$
  - Now approximate the expectation as follows

$$\mathbb{E}[f] = \int f(z)p(z)dz = \int f(z)\frac{p(z)}{q(z)}q(z)dz \approx \frac{1}{L}\sum_{\ell=1}^{L} f(z^{(\ell)})\frac{p(z^{(\ell)})}{q(z^{(\ell)})}$$

See PRML 11.1.4

- This is basically "weighted" Monte Carlo integration
  - $w^{(\ell)} = \frac{p(z^{(\ell)})}{q(z^{(\ell)})}$  denotes the importance weight of each sample  $z^{(\ell)}$
- IS works even when we can only evaluate  $p(z) = \frac{\tilde{p}(z)}{Z_n}$  up to a prop. constant
- Note: Monte Carlo and Importance Sampling are NOT sampling methods!
  - These are only uses for computing expectations (approximately)

### Limitations of the Basic Methods

- Transformation based methods: Usually limited to drawing from standard distributions
- Rejection Sampling and Importance Sampling: Require good proposal distributions



- In general, difficult to find good prop. distr. especially when z is high-dim
- More sophisticated sampling methods like MCMC work well in such high-dim spaces



# MCMC: The Basic Scheme

- The chain run infinitely long (i.e., upon convergence) will give ONE sample from  $p({m z})$
- But we usually require several samples to approximate p(z)
- This is done as follows
  - Start the chain at an initial  $m{z}^{(0)}$
  - Using the proposal  $q(\mathbf{z}|\mathbf{z}^{(\ell)})$ , run the chain long enough, say  $T_1$  steps
  - Discard the first  $T_1 1$  samples (called "burn-in" samples) and take last sample  $\mathbf{z}^{(T_1)}$
  - Continue from  $\mathbf{z}^{(T_1)}$  up to  $T_2$  steps, discard intermediate samples, take last sample  $\mathbf{z}^{(T_2)}$ 
    - This discarding (called "thinning") helps ensure that  $z^{(T_1)}$  and  $z^{(T_2)}$  are uncorrelated
  - Repeat the same for a total of S times
  - In the end, we now have S approximately independent samples from p(z)
- Note: Good choices for  $T_1$  and  $T_i T_{i-1}$  (thinning gap) are usually based on heuristics





MCMC is exact in theory but approximate in practice since

we can't run the chain for

infinitely long in practice

Requirement for Monte Carlo approximation



CS772A: PML

Will treat it as our first sample from p(z)