

DSC291: Machine Learning with Few Labels

Variational Inference

Zhiting Hu

Lecture 7, January 24, 2023

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

Recap: EM Algorithm

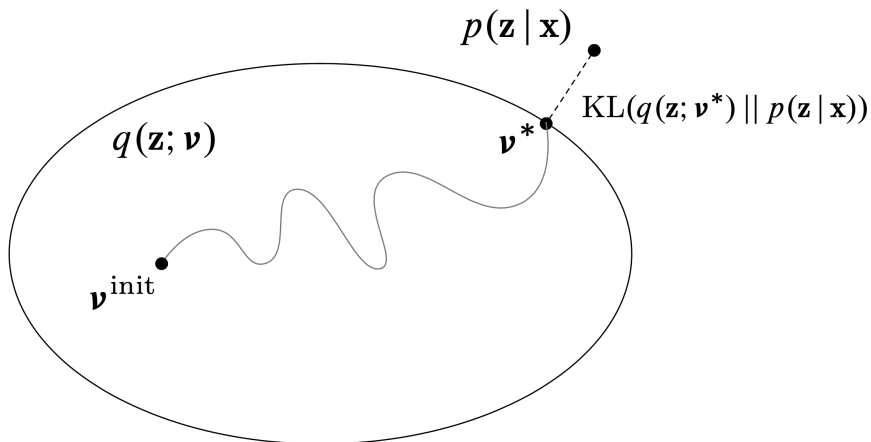
- The EM algorithm is coordinate-decent on $F(q, \theta)$
 - E-step: $q^{t+1} = \arg \min_q F(q, \theta^t) = p(\mathbf{z}|\mathbf{x}, \theta^t)$
 - M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta) = \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} q^{t+1}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$

$$\begin{aligned} \ell(\theta; \mathbf{x}) &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta)) \\ &= -F(q, \theta) + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta)) \end{aligned}$$

- Limitation: need to be able to compute $p(\mathbf{z}|\mathbf{x}, \theta)$, not possible for more complicated models --- solution: Variational inference

Recap: Variational Inference

- We often cannot compute posteriors $p(\mathbf{z}|\mathbf{x}, \theta)$, and so we need to approximate them, using variational methods.
- In variational Bayes, we'd like to find an approximation within some family that minimizes the KL divergence to the posterior, but we can't directly minimize this
- Therefore, we defined the ELBO, which we can maximize, and this is equivalent to minimizing the KL divergence.

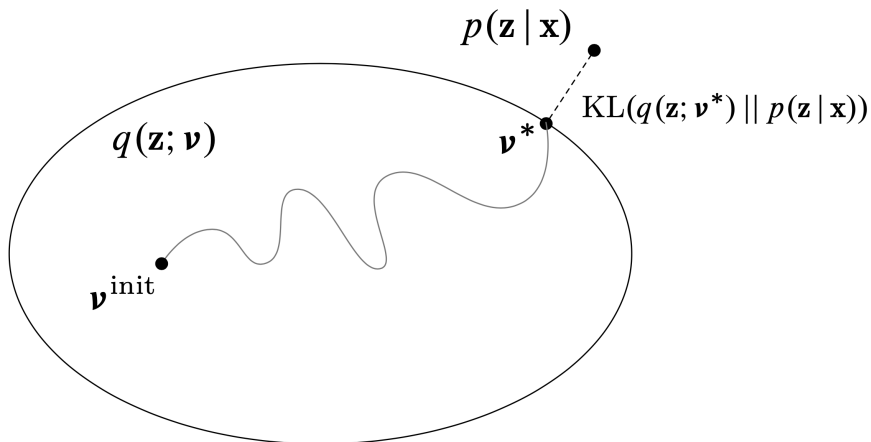


Evidence Lower Bound (ELBO)

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

Recap: Variational Inference

- We often cannot compute posteriors $p(\mathbf{z}|\mathbf{x}, \theta)$, and so we need to approximate them, using variational methods.
- In variational Bayes, we'd like to find an approximation within some family that minimizes the KL divergence to the posterior, but we can't directly minimize this
- Therefore, we defined the ELBO, which we can maximize, and this is equivalent to minimizing the KL divergence.



Evidence Lower Bound (ELBO)

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

- How do we choose the variational family $q(\mathbf{z}|\mathbf{x}, \mathbf{v})$?

Mean Field Variational Inference

- A popular family of variational approximations
- In this type of variational inference, we assume the variational distribution over the latent variables **factorizes** as

$$q(\mathbf{z}) = q(z_1, \dots, z_m) = \prod_{j=1}^m q(z_j)$$

- (where we omit variational parameters for ease of notation)
 - We refer to $q(z_j)$, the variational approximation for a single latent variable, as a “local variational approximation”
- In the above expression, the variational approximation $q(z_j)$ over each latent variable z_j is independent

Mean Field Variational Inference

- Note that this is a fairly general setup; we can also partition the latent variables z_1, \dots, z_m into R groups z_{G_1}, \dots, z_{G_R} , and use the approximation:

$$q(z_1, \dots, z_m) = q(z_{G_1}, \dots, z_{G_R}) = \prod_{r=1}^R q(z_{G_r})$$

- Often called “generalized mean field” versus (the above) “naïve mean field”.

Mean Field Variational Inference

- Note that this is a fairly general setup; we can also partition the latent variables z_1, \dots, z_m into R groups z_{G_1}, \dots, z_{G_R} , and use the approximation:

$$q(z_1, \dots, z_m) = q(z_{G_1}, \dots, z_{G_R}) = \prod_{r=1}^R q(z_{G_r})$$

- Often called “generalized mean field” versus (the above) “naïve mean field”.
- Typically, this approximation does not contain the true posterior (because the latent variables are dependent).
 - E.g.: in the (Bayesian) mixture of Gaussians model, all of the cluster assignments z_i for $i = 1, \dots, n$ are dependent on each other and on the cluster locations $\mu_{1:K}$ given data.

Optimizing the ELBO in Mean Field Variational Inference

How do we optimize the ELBO in mean field variational inference?

- Typically, we use coordinate ascent optimization.
- I.e. we optimize each latent variable's variational approximation $q(z_j)$ in turn while holding the others fixed.
 - At each iteration we get an updated “local” variational approximation.
 - And we iterate through each latent variable until convergence.

Optimizing the ELBO in Mean Field Variational Inference

- Recall that the ELBO is defined as:

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

- Note that we can decompose the entropy term of the ELBO (using the mean field variational approximation) as:

$$\mathbb{E}_q[\log q(z_{1:m})] = \sum_{j=1}^m \mathbb{E}_{q_j}[\log q(z_j)]$$

- Therefore, under the mean field approximation, the ELBO can be written:

$$\mathcal{L} = \mathbb{E}_{q_j} \mathbb{E}_{q_{-j}}[\log p(\mathbf{x}, \mathbf{z})] - \sum_{j=1}^m \mathbb{E}_{q_j}[\log q(z_j)]$$

Optimizing the ELBO in Mean Field Variational Inference

- Therefore, under the mean field approximation, the ELBO can be written:

$$\mathcal{L} = \mathbb{E}_{q_j} \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] - \sum_{j=1}^m \mathbb{E}_{q_j} [\log q(z_j)]$$

- Next, we want to derive the coordinate ascent update for a latent variable z_j , keeping all other latent variables fixed.
 - i.e. we want the $\operatorname{argmax}_{q_j} \mathcal{L}$.
- Removing the parts that do not depend on $q(z_j)$, we can write:

$$\mathcal{L} = \mathbb{E}_{q_j} \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q_j} [\log q(z_j)] + \text{const.}$$

Optimizing the ELBO in Mean Field Variational Inference

- To find this argmax, we take the derivative of \mathcal{L} w.r.t $q(z_j)$ and set the derivative to zero :

$$\frac{d\mathcal{L}}{dq(z_j)} = \mathbb{E}_{q_j} \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] - \log q(z_j) - 1 = 0$$

- From this, we arrive at the coordinate ascent update:

$$q^*(z_j) \propto \exp \left\{ \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] \right\}$$

Optimizing the ELBO in Mean Field Variational Inference

- The coordinate ascent update:

$$q^*(z_j) \propto \exp \left\{ \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] \right\}$$

- The optimal solution for factor $q(z_j)$ is obtained simply by considering the log of the joint distribution over all observed and latent variables and then taking the expectation with respect to all of the other factors $q(z_k)$, $k \neq j$, then taking exponential and normalizing
- Note that the only assumption we made so far is the mean-field factorization:
$$q(\mathbf{z}) = q(z_1, \dots, z_m) = \prod_{j=1}^m q(z_j)$$
- We haven't yet made any assumptions on the form of $q(z_j)$

Simple example:

- Consider a univariate Gaussian distribution $p(x) = \mathcal{N}(x|\mu, \tau^{-2})$, given a dataset $\mathcal{D} = \{x_1, \dots, x_N\}$:

$$p(\mathcal{D}|\mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{N/2} \exp\left\{-\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2\right\}$$

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

$$p(\tau) = \text{Gam}(\tau|a_0, b_0)$$

- $\text{Gam}(\tau|a_0, b_0) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$: gamma distribution
- For this simple problem the posterior distribution can be found exactly. But we use it as an example for tutorial anyway

$$q^*(z_j) \propto \exp \left\{ \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] \right\}$$

Simple example:

$$p(\mathcal{D}|\mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{N/2} \exp \left\{ -\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\} \quad \begin{aligned} p(\mu|\tau) &= \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1}) \\ p(\tau) &= \text{Gam}(\tau|a_0, b_0) \end{aligned}$$

- Introduce the factorized variational approximation: $q(\mu, \tau) = q_\mu(\mu)q_\tau(\tau)$
- Solution to q_μ :

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

- We can see q_μ^* is a Gaussian $\mathcal{N}(x|\mu_N, \lambda_N^{-1})$:

$$\begin{aligned} \mu_N &= \frac{\lambda_0\mu_0 + N\bar{x}}{\lambda_0 + N} \\ \lambda_N &= (\lambda_0 + N)\mathbb{E}[\tau] \end{aligned}$$

$$q^*(z_j) \propto \exp \left\{ \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] \right\}$$

Simple example:

$$p(\mathcal{D}|\mu, \tau) = \left(\frac{\tau}{2\pi} \right)^{N/2} \exp \left\{ -\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\} \quad \begin{array}{l} p(\mu|\tau) = \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1}) \\ p(\tau) = \text{Gam}(\tau|a_0, b_0) \end{array}$$

- Introduce the factorized variational approximation: $q(\mu, \tau) = q_\mu(\mu)q_\tau(\tau)$

- Solution to q_τ :
$$\begin{aligned} \ln q_\tau^*(\tau) &= \mathbb{E}_\mu [\ln p(\mathcal{D}|\mu, \tau) + \ln p(\mu|\tau)] + \ln p(\tau) + \text{const} \\ &= (a_0 - 1) \ln \tau - b_0\tau + \frac{N}{2} \ln \tau \\ &\quad - \frac{\tau}{2} \mathbb{E}_\mu \left[\sum_{n=1}^N (x_n - \mu)^2 + \lambda_0(\mu - \mu_0)^2 \right] + \text{const} \end{aligned}$$

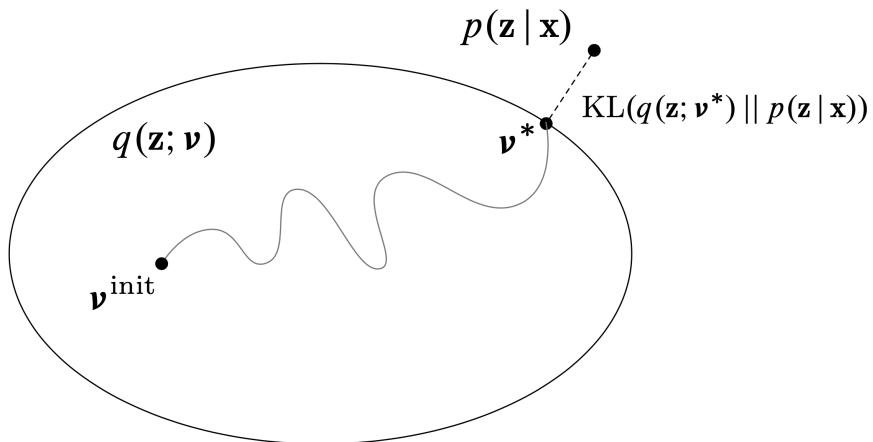
- We can see q_τ^* is a gamma distribution $\text{Gam}(\tau|a_N, b_N)$:

$$a_N = a_0 + \frac{N}{2}$$

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

Quick Recap

- We often cannot compute posteriors, and so we need to approximate them, using variational methods.
- In variational Bayes, we'd like to find an approximation within some family that minimizes the KL divergence to the posterior, but we can't directly minimize this
- Therefore, we defined the ELBO, which we can maximize, and this is equivalent to minimizing the KL divergence.



Evidence Lower Bound (ELBO)

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z} | \mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q(\mathbf{z} | \mathbf{x})} \right] + \text{KL}(q(\mathbf{z} | \mathbf{x}) || p(\mathbf{z} | \mathbf{x}, \theta))$$

Quick Recap

- We defined a family of approximations called “mean field” approximations, in which there are no dependencies between latent variables

$$q(\mathbf{z}) = q(z_1, \dots, z_m) = \prod_{j=1}^m q(z_j)$$

- We optimize the ELBO with coordinate ascent updates to iteratively optimize each local variational approximation under mean field assumptions

$$q^*(z_j) \propto \exp \left\{ \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] \right\}$$

Key Takeaways

- KL Divergence $\text{KL}(q(\mathbf{x}) \parallel p(\mathbf{x})) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$

- The key equation of EM and VI

Evidence Lower Bound (ELBO)

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta))$$

- Free energy $F(q, \theta)$
- EM: E-step and M-step optimizing ELBO w.r.t q and θ
- Mean-field VI: optimizing ELBO w.r.t factorized q components

Stochastic VI; Black-box VI

VI with coordinate ascent

Example: Bayesian mixture of Gaussians

- Treat the mean μ_k and cluster proportion π as latent variables

$$\mu_k \sim \mathcal{N}(0, \tau^2) \text{ for } k = 1, \dots, K$$

$$\pi \sim \text{Dirichlet}(\alpha)$$

- For each data $i = 1, \dots, n$

$$z_i \sim \text{Cat}(\pi).$$

$$x_i \sim \mathcal{N}(\mu_{z_i}, \sigma^2).$$

- We have
 - observed variables $x_{1:n}$
 - latent variables $\mu_{1:k}$, π and $z_{1:n}$
 - Hyper-parameters $\{\tau^2, \sigma^2\}$

VI with coordinate ascent

Example: Bayesian mixture of Gaussians

Assume mean-field $q(\mu_{1:K}, \pi, z_{1:n}) = \prod_k q(\mu_k) q(\pi) \prod_i q(z_i)$

- Initialize the global variational distributions $q(\mu_k)$ and $q(\pi)$
- **Repeat:**
 - **For** each data example $i \in \{1, 2, \dots, D\}$
 - Update the local variational distribution $q(z_i)$
 - **End for**
 - Update the global variational distributions $q(\mu_k)$ and $q(\pi)$
- **Until** ELBO converges

- What if we have millions of data examples? This could be very slow.

Stochastic VI

Example: Bayesian mixture of Gaussians

Assume mean-field $q(\mu_{1:K}, \pi, z_{1:n}) = \prod_k q(\mu_k) q(\pi) \prod_i q(z_i)$

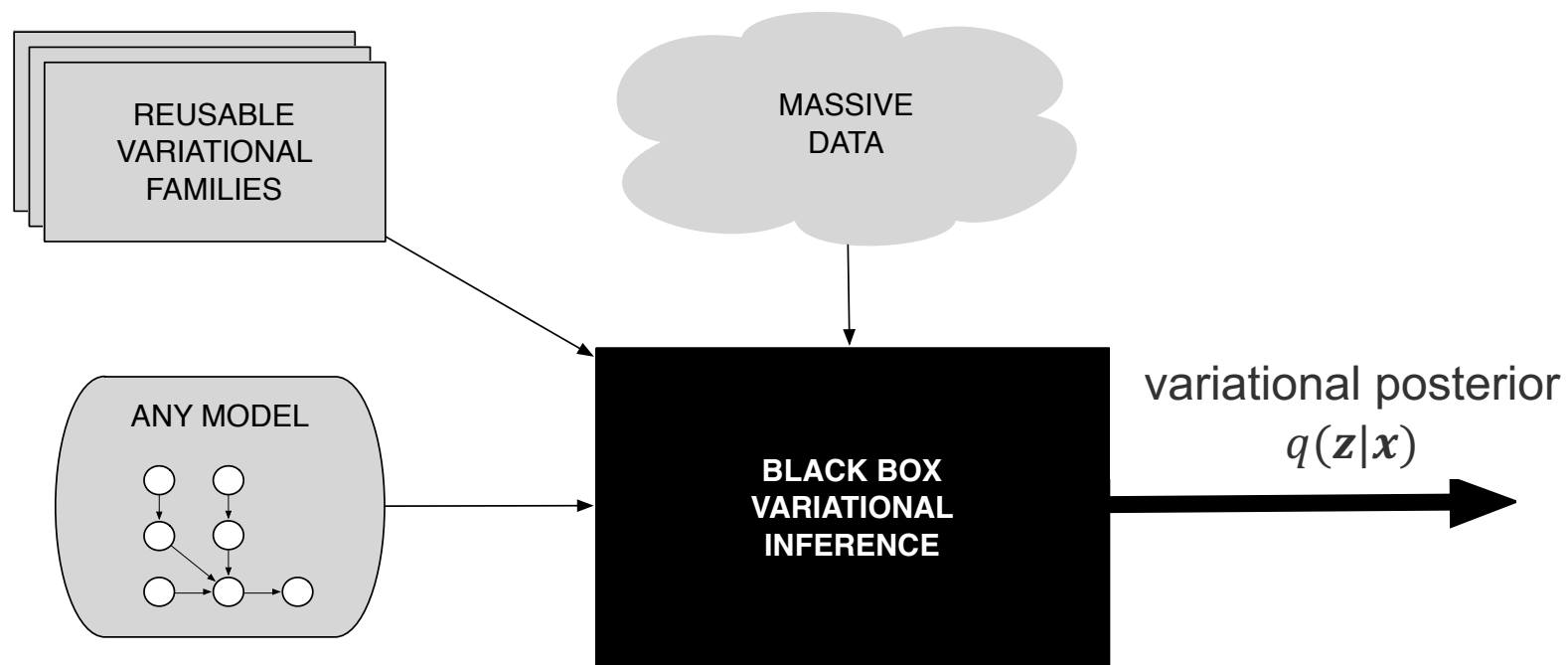
- Initialize the global variational distributions $q(\mu_k)$ and $q(\pi)$
- **Repeat:**
 - Sample a data example $i \in \{1, 2, \dots, D\}$
 - Update the local variational distribution $q(z_i)$
 - Update the global variational distributions $q(\mu_k)$ and $q(\pi)$ with **natural gradient ascent**
- **Until** ELBO converges

- (Setting natural gradient = 0 gives the traditional mean-field update)

Black-box Variational Inference (BBVI)

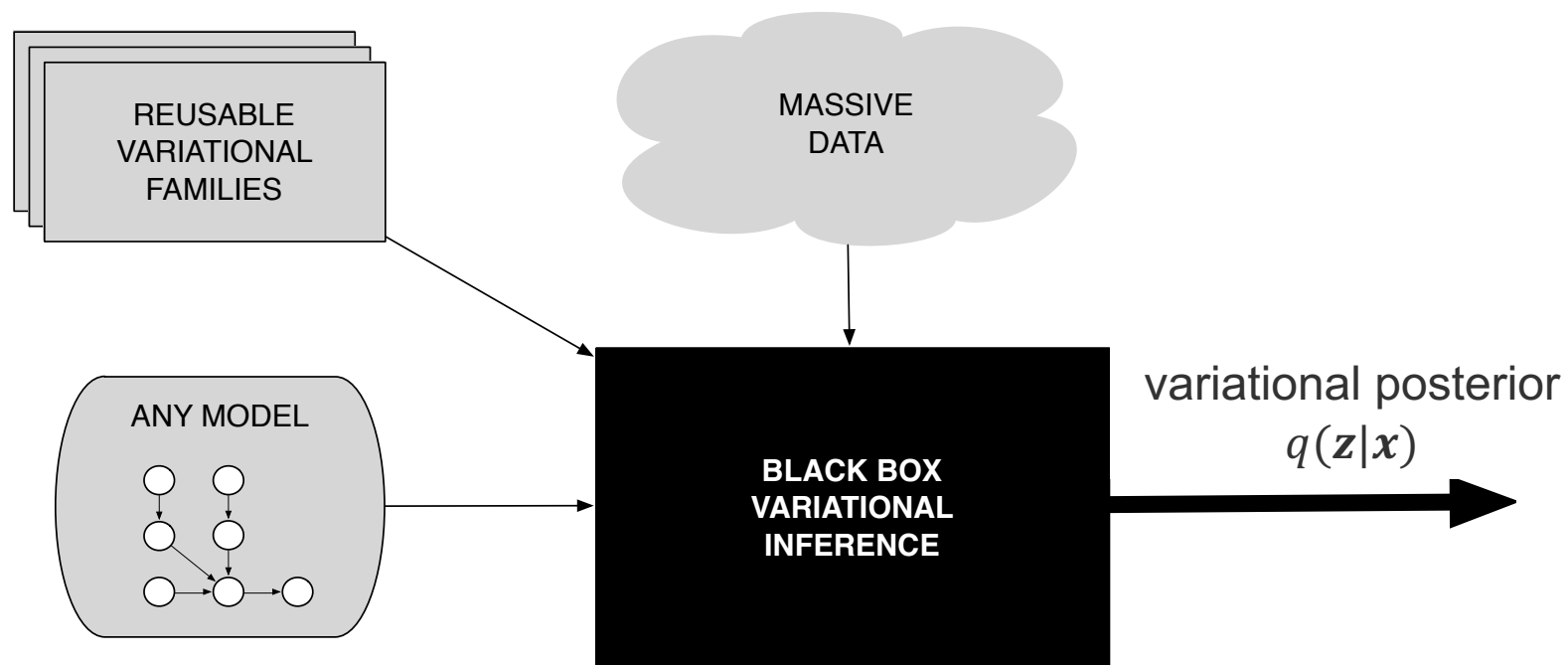
- We have derived variational inference specific for Bayesian Gaussian (mixture) models
- There are innumerable models
- Can we have a solution that does not entail model-specific work?

Black-box Variational Inference (BBVI)



- Easily use variational inference with **any model**
- Perform inference with **massive data**
- **No mathematical work** beyond specifying the model

Black-box Variational Inference (BBVI)



- Sample from $q(\cdot)$
- Form noisy gradients (without model-specific computation)
- Use stochastic optimization

Black-box Variational Inference (BBVI)

- Probabilistic model: \mathbf{x} -- observed variables, \mathbf{z} -- latent variables
- Variational distribution $q_{\lambda}(\mathbf{z}|\mathbf{x})$ with parameters λ , e.g.,
 - Gaussian mixture distribution:
 - "A Gaussian mixture model is a universal approximator of densities, in the sense that any smooth density can be approximated with any specific nonzero amount of error by a Gaussian mixture model with enough components." (Deep Learning book, pp.65)
 - Deep neural networks
- ELBO:

$$\mathcal{L}(\lambda) = \mathbb{E}_{q(\mathbf{z}|\lambda)}[\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q(\mathbf{z}|\lambda)}[\log q(\mathbf{z}|\lambda)]$$

- Want to compute the gradient w.r.t variational parameters λ

The General Problem: Computing Gradients of Expectations

- When the objective function \mathcal{L} is defined as an expectation of a (differentiable) test function $f_\lambda(\mathbf{z})$ w.r.t. a probability distribution $q_\lambda(\mathbf{z})$

$$\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$$

- Computing exact gradients w.r.t. the parameters λ is often unfeasible
- Need stochastic gradient estimates
 - The score function estimator (a.k.a log-derivative trick, REINFORCE)
 - The reparameterization trick (a.k.a the pathwise gradient estimator)

Computing Gradients of Expectations w/ score function

- Loss: $\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$
- Log-derivative trick: $\nabla_\lambda q_\lambda = q_\lambda \nabla_\lambda \log q_\lambda$
- Gradient w.r.t. λ :

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z}) \nabla_\lambda \log q_\lambda(\mathbf{z}) + \nabla_\lambda f_\lambda(\mathbf{z})]$$

- **score function**: the gradient of the log of a probability distribution
- Compute noisy unbiased gradients with Monte Carlo samples from q_λ

$$\nabla_\lambda \mathcal{L} \approx \frac{1}{S} \sum_{s=1}^S f_\lambda(\mathbf{z}_s) \nabla_\lambda \log q_\lambda(\mathbf{z}_s) + \nabla_\lambda f_\lambda(\mathbf{z}_s) \quad \text{where } \mathbf{z}_s \sim q_\lambda(\mathbf{z})$$

- Pros: generally applicable to any distribution $q(\mathbf{z}|\lambda)$
- Cons: empirically has high variance \rightarrow slow convergence
 - To reduce variance: Rao-Blackwellization, control variates, importance sampling, ...

Computing Gradients of Expectations w/ reparametrization trick

- Loss: $\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$
- Assume that we can express the distribution $q_\lambda(\mathbf{z})$ with a transformation

$$\begin{aligned} \epsilon &\sim s(\epsilon) \\ \mathbf{z} &= t(\epsilon, \lambda) \end{aligned} \iff \mathbf{z} \sim q(\mathbf{z}|\lambda)$$

- E.g.,

$$\begin{aligned} \epsilon &\sim \text{Normal}(0, 1) \\ \mathbf{z} &= \epsilon\sigma + \mu \end{aligned} \iff \mathbf{z} \sim \text{Normal}(\mu, \sigma^2)$$

- Reparameterization gradient

$$\mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[f_\lambda(\mathbf{z}(\epsilon, \lambda))]$$

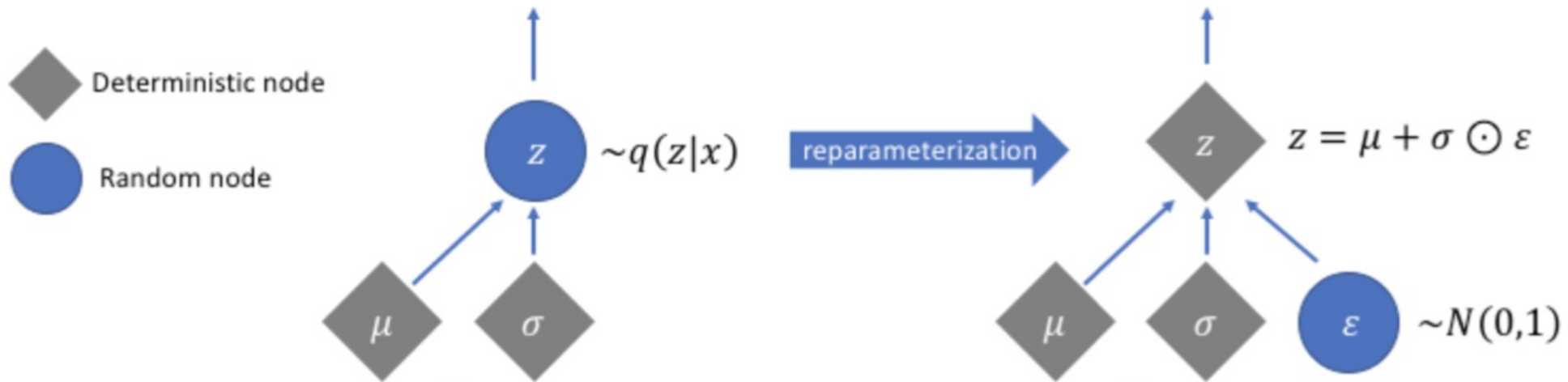
$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_{\mathbf{z}} f_\lambda(\mathbf{z}) \nabla_\lambda t(\epsilon, \lambda)]$$

- Pros: empirically, lower variance of the gradient estimate
- Cons: Not all distributions can be reparameterized

Reparameterization trick

- Reparameterizing Gaussian distribution

$$\begin{aligned} \epsilon &\sim \text{Normal}(0, 1) \\ z &= \epsilon\sigma + \mu \end{aligned} \iff z \sim \text{Normal}(\mu, \sigma^2)$$



Reparameterization trick

- Reparametrizing Gaussian distribution

$$\begin{aligned} \epsilon &\sim \text{Normal}(0, 1) \\ z &= \epsilon\sigma + \mu \end{aligned} \iff z \sim \text{Normal}(\mu, \sigma^2)$$

- Other reparameterizable distributions: $\epsilon \sim \text{Uniform}(\epsilon) \iff z \sim q(z)$
 - Tractable inverse CDF F^{-1} : $z = F^{-1}(\epsilon)$
 - Exponential, Cauchy, Logistic, Rayleigh, Pareto, Weibull, Reciprocal, Gompertz, Gumbel, Erlang
 - Location-scale:
 - Laplace, Elliptical, Student's t, Logistic, Uniform, Triangular, Gaussian
 - Composition:
 - Log-Normal (exponentiated normal) Gamma (sum of exponentials) Dirichlet (sum of Gammas) Beta, Chi-Squared, F

Computing Gradients of Expectations: Summary

- Loss: $\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$

- Score gradient

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z}) \nabla_\lambda \log q_\lambda(\mathbf{z}) + \nabla_\lambda f_\lambda(\mathbf{z})]$$

- Pros: generally applicable to any distribution $q(\mathbf{z}|\lambda)$
- Cons: empirically has high variance \rightarrow slow convergence

- Reparameterization gradient

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_{\mathbf{z}} f_\lambda(\mathbf{z}) \nabla_\lambda t(\epsilon, \lambda)]$$

- Pros: empirically, lower variance of the gradient estimate
- Cons: Not all distributions can be reparameterized

Recall: Black-box Variational Inference (BBVI)

- Probabilistic model: \mathbf{x} -- observed variables, \mathbf{z} -- latent variables
- Variational distribution $q_{\lambda}(\mathbf{z}|\mathbf{x})$ with parameters λ , e.g.,
 - Gaussian mixture distribution:
 - "A Gaussian mixture model is a universal approximator of densities, in the sense that any smooth density can be approximated with any specific nonzero amount of error by a Gaussian mixture model with enough components." (Deep Learning book, pp.65)

- Deep neural networks

$$\mathcal{L}(\lambda) \triangleq \mathbb{E}_{q_{\lambda}(\mathbf{z})}[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})]$$

- ELBO:

$$\mathcal{L}(\lambda) = \mathbb{E}_{q(\mathbf{z}|\lambda)}[\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q(\mathbf{z}|\lambda)}[\log q(\mathbf{z}|\lambda)]$$

- Want to compute the gradient w.r.t variational parameters λ

BBVI with the score gradient

- ELBO:

$$\mathcal{L}(\lambda) = \mathbb{E}_{q(\mathbf{z}|\lambda)}[\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q(\mathbf{z}|\lambda)}[\log q(\mathbf{z}|\lambda)]$$

- Gradient w.r.t. λ (using the log-derivative trick)

$$\nabla_{\lambda} \mathcal{L} = \mathbb{E}_q[\nabla_{\lambda} \log q(\mathbf{z}|\lambda)(\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\lambda))]$$

- Compute noisy unbiased gradients of the ELBO with Monte Carlo samples from the variational distribution

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

where $\mathbf{z}_s \sim q(\mathbf{z}|\lambda)$.

BBVI with the reparameterization gradient

- ELBO:

$$\mathcal{L}(\lambda) = \mathbb{E}_{q(\mathbf{z}|\lambda)}[\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q(\mathbf{z}|\lambda)}[\log q(\mathbf{z}|\lambda)]$$

- Gradient w.r.t. λ

$$\begin{aligned} \epsilon &\sim s(\epsilon) \\ z &= t(\epsilon, \lambda) \end{aligned} \iff z \sim q(z|\lambda)$$

$$\nabla_{\lambda} \mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)} [\nabla_{\mathbf{z}} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})] \nabla_{\lambda} t(\epsilon, \lambda)]$$

Questions?

Questions?