

DSC291: Machine Learning with Few Labels

Unsupervised Learning

Zhiting Hu

Lecture 16, May 8, 2024

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: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta) = \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} q^{t+1}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$

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

Recap: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:

- Z is a latent class indicator vector:

$$p(z_n) = \text{multi}(z_n : \pi) = \prod_k (\pi_k)^{z_n^k}$$

- X is a conditional Gaussian variable with a class-specific mean/covariance

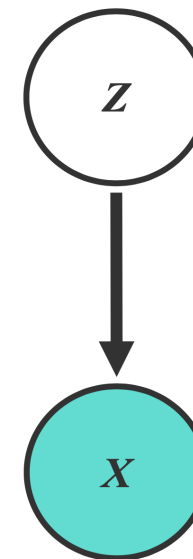
$$p(x_n | z_n^k = \mathbf{1}, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$

- The likelihood of a sample:

$$\begin{aligned} p(x_n | \mu, \Sigma) &= \sum_k p(z^k = \mathbf{1} | \pi) p(x, | z^k = \mathbf{1}, \mu, \Sigma) \\ &= \sum_{z_n} \prod_k \left((\pi_k)^{z_n^k} N(x_n : \mu_k, \Sigma_k)^{z_n^k} \right) = \sum_k \pi_k N(x, | \mu_k, \Sigma_k) \end{aligned}$$

mixture proportion

mixture component



Recap: Gaussian Mixture Models (GMMs)

- E-step: computing the posterior of z_n given the current estimate of the parameters (i.e., π, μ, Σ)

$$\begin{aligned} p(z^k = 1 \mid \mathbf{x}) &= \frac{p(z^k = 1)p(\mathbf{x} \mid z^k = 1)}{p(\mathbf{x})} \\ &= \frac{p(z^k = 1)p(\mathbf{x} \mid z^k = 1)}{\sum_{j=1}^K p(z^j = 1)p(\mathbf{x} \mid z^j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x} \mid \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} \mid \mu_j, \Sigma_j)} \\ &:= \gamma_k \end{aligned}$$

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

Variational Inference

Inference

- Given a model, the goals of inference can include:
 - Computing the likelihood of observed data $p(\mathbf{x}^*)$
 - Computing the marginal distribution over a given subset of variables in the model $p(\mathbf{x}_A)$
 - Computing the conditional distribution over a subsets of nodes given a disjoint subset of nodes $p(\mathbf{x}_A|\mathbf{x}_B)$
 - Computing a mode of the density (for the above distributions) $\operatorname{argmax}_{\mathbf{x}} p(\mathbf{x})$
 -

Variational Inference

- Observed variables \mathbf{x} , latent variables \mathbf{z}
- Variational (Bayesian) inference, a.k.a. **variational Bayes**, is most often used to **approximately** infer the conditional distribution over the latent variables given the observations (and parameters)
 - i.e., the **posterior distribution** over the latent variables

$$p(\mathbf{z}|\mathbf{x}, \theta) = \frac{p(\mathbf{z}, \mathbf{x}|\theta)}{\sum_{\mathbf{z}} p(\mathbf{z}, \mathbf{x}|\theta)}$$

Motivating Example

- Why do we often need to use an approximate inference methods (such as variational Bayes) to compute the posterior distribution?
- It's because we cannot directly compute the posterior distribution for many interesting models
 - I.e. the posterior density is in an intractable form (often involving integrals) which cannot be easily analytically solved.
- As a motivating example, we will try to compute the posterior for a (Bayesian) mixture of Gaussians.

Bayesian mixture of Gaussians

- The mean μ_k is treated as a (latent) random variable

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

- 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 $\mu_{1:k}$ and $z_{1:n}$
- parameters $\{\tau^2, \pi, \sigma^2\}$

- $p(x_{1:n}, z_{1:n}, \mu_{1:k} | \tau^2, \pi, \sigma^2) = \prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i) p(x_i | z_i, \mu_{1:K})$

Bayesian mixture of Gaussians

- We can write the posterior distribution as

$$p(\mu_{1:K}, z_{1:n} | x_{1:n}) = \frac{\prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i) p(x_i | z_i, \mu_{1:K})}{\int_{\mu_{1:K}} \sum_{z_{1:n}} \prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i) p(x_i | z_i, \mu_{1:K})}$$

- The numerator can be computed for any choice of the latent variables
- The problem is the denominator (the marginal probability of the observations)
 - This integral cannot easily be computed analytically
- We need some approximation..

Variational Inference

The main idea behind variational inference:

- Choose a family of distributions over the latent variables $z_{1:m}$ with its own set of variational parameters ν , i.e.

$$q(z_{1:m}|\nu)$$

- Then, we find the setting of the parameters that makes our approximation q closest to the posterior distribution.
 - This is where optimization algorithms come in.
- Then we can use q with the fitted parameters in place of the posterior.
 - E.g. to form predictions about future data, or to investigate the posterior distribution over the hidden variables, find modes, etc.

Variational Inference

- We want to minimize the KL divergence between our approximation $q(\mathbf{z}|\mathbf{x})$ and our posterior $p(\mathbf{z}|\mathbf{x})$

$$\text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}))$$

- But we can't actually minimize this quantity w.r.t q because $p(\mathbf{z}|\mathbf{x})$ is unknown
- **Question:** how can we minimize the KL divergence?
 - Hint: recall what we did in EM:

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

Variational Inference

- We want to minimize the KL divergence between our approximation $q(\mathbf{z}|\mathbf{x})$ and our posterior $p(\mathbf{z}|\mathbf{x})$

$$\text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}))$$

- But we can't actually minimize this quantity w.r.t q because $p(\mathbf{z}|\mathbf{x})$ is unknown
- **Question:** how can we minimize the KL divergence?

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

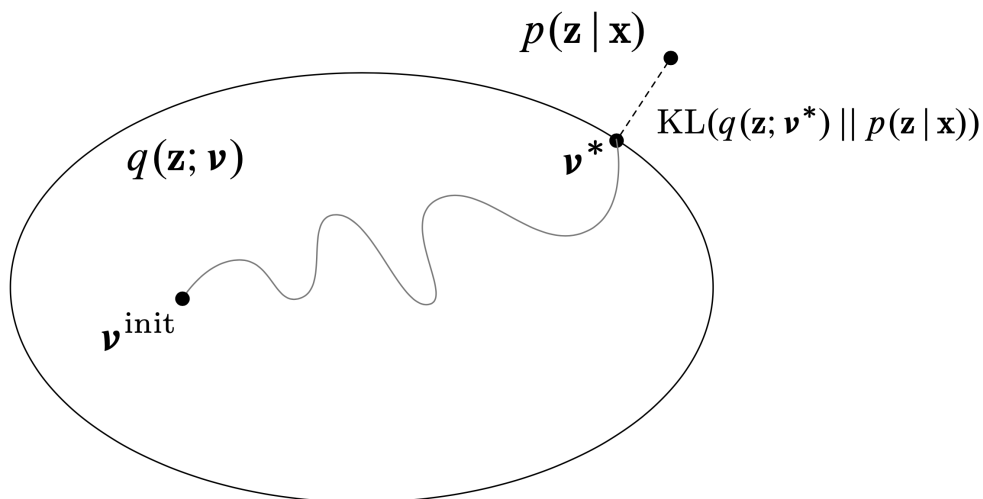
Evidence Lower Bound (ELBO)

- The ELBO is equal to the negative KL divergence up to a constant $\ell(\theta; \mathbf{x})$
- We maximize the ELBO over q to find an "optimal approximation" to $p(\mathbf{z}|\mathbf{x})$

Variational Inference

- Choose a family of distributions over the latent variables \mathbf{z} with its own set of variational parameters ν , i.e. $q(\mathbf{z}|\mathbf{x}, \nu)$
- We maximize the ELBO over q to find an “optimal approximation” to $p(\mathbf{z}|\mathbf{x})$

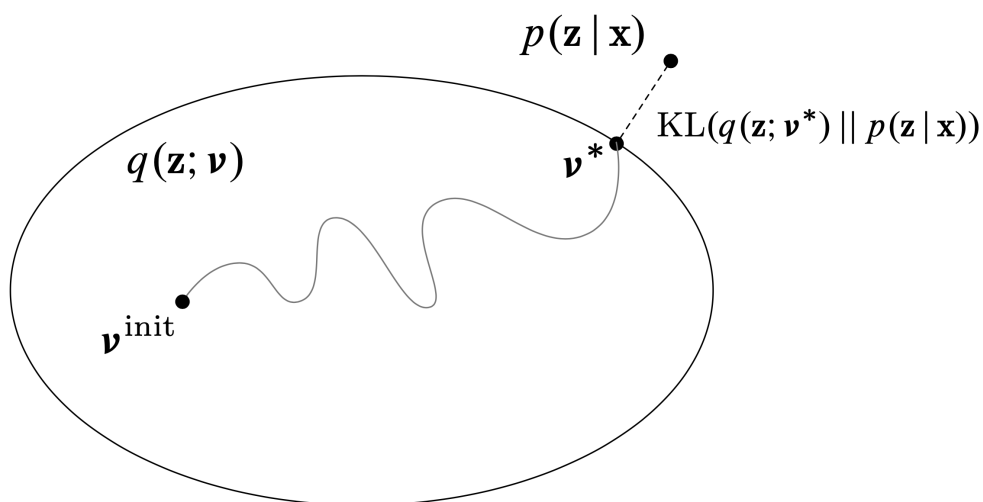
$$\begin{aligned} & \operatorname{argmax}_{\nu} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x}, \nu)} \right] \\ & = \operatorname{argmax}_{\nu} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} [\log p(\mathbf{x}, \mathbf{z}|\theta)] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} [\log q(\mathbf{z}|\mathbf{x}, \nu)] \end{aligned}$$



Variational Inference

- Choose a family of distributions over the latent variables \mathbf{z} with its own set of variational parameters \mathbf{v} , i.e. $q(\mathbf{z}|\mathbf{x}, \mathbf{v})$
- We maximize the ELBO over q to find an “optimal approximation” to $p(\mathbf{z}|\mathbf{x})$

$$\begin{aligned} & \operatorname{argmax}_{\mathbf{v}} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} \right] \\ & = \operatorname{argmax}_{\mathbf{v}} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} [\log p(\mathbf{x}, \mathbf{z}|\theta)] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} [\log q(\mathbf{z}|\mathbf{x}, \mathbf{v})] \end{aligned}$$

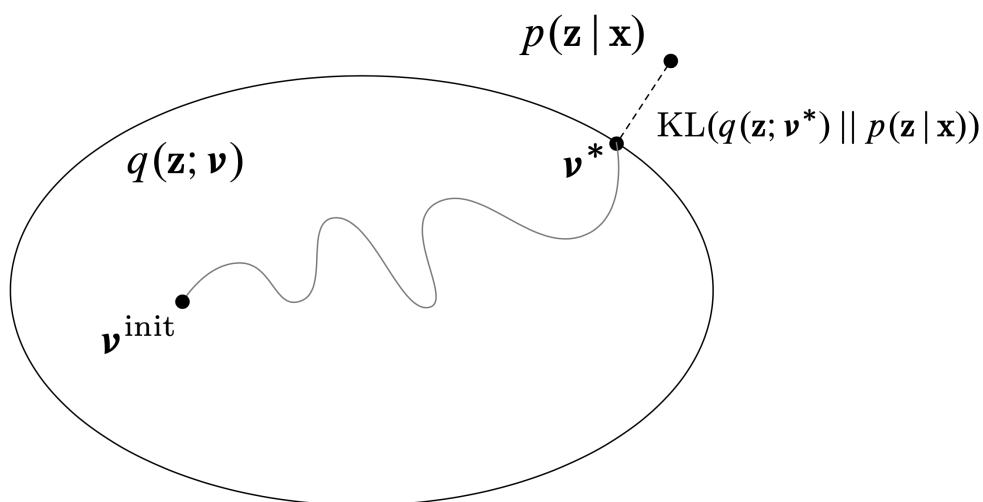


Question: How do we choose the variational family $q(\mathbf{z}|\mathbf{x}, \mathbf{v})$?

Variational Inference

- Choose a family of distributions over the latent variables \mathbf{z} with its own set of variational parameters \mathbf{v} , i.e. $q(\mathbf{z}|\mathbf{x}, \mathbf{v})$
- We maximize the ELBO over q to find an “optimal approximation” to $p(\mathbf{z}|\mathbf{x})$

$$\begin{aligned} & \operatorname{argmax}_{\mathbf{v}} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} \right] \\ &= \operatorname{argmax}_{\mathbf{v}} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} [\log p(\mathbf{x}, \mathbf{z}|\theta)] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} [\log q(\mathbf{z}|\mathbf{x}, \mathbf{v})] \end{aligned}$$



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

- Factorized distribution -> mean field VI
- Mixture of Gaussian distribution -> black-box VI
- Neural-based distribution -> Variational Autoencoders (VAEs)

Example: 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

Example: Mean Field Variational Inference

- 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.

Example: 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.

Mean Field Variational Inference with Coordinate Ascent

Recap: Bayesian mixture of Gaussians

- Treat the mean μ_k as latent variables

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

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

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

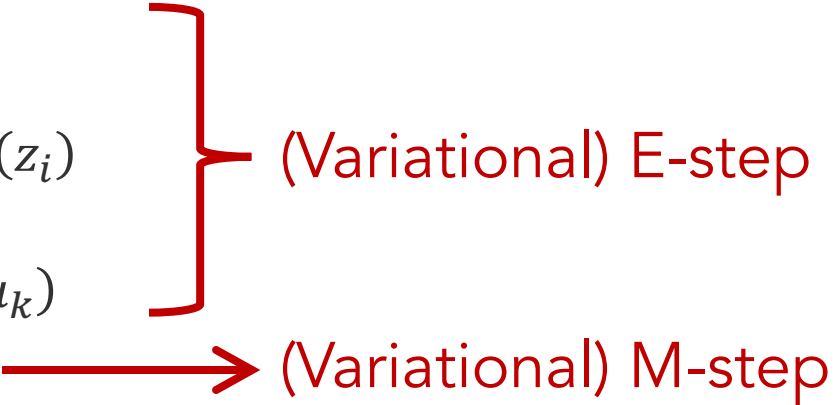
$$x_i \sim \mathcal{N}(\mu_{z_i}, \sigma^2).$$

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

Mean Field Variational Inference with Coordinate Ascent

Recap: Bayesian mixture of Gaussians

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

- Initialize the global variational distributions $q(\mu_k)$ and parameters $\{\tau^2, \sigma^2, \pi\}$
 - **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 $q(\mu_k)$
 - Update the parameters $\{\tau^2, \sigma^2, \pi\}$
 - **Until** ELBO converges
- 

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

Stochastic VI

Recap: Bayesian mixture of Gaussians

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

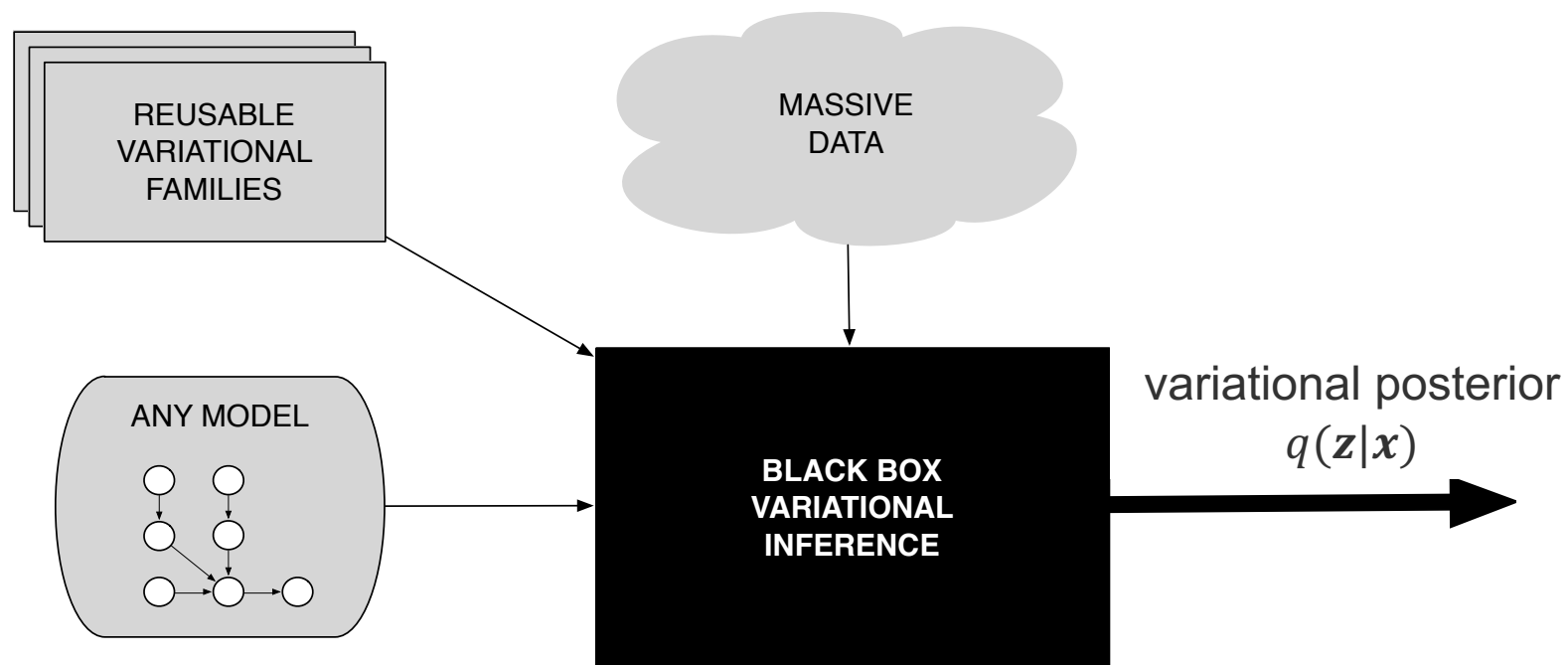
- Initialize the global variational distributions $q(\mu_k)$ and parameters $\{\tau^2, \sigma^2, \pi\}$
- **Repeat:**
 - **Sample** a data example $i \in \{1, 2, \dots, D\}$
 - Update the local variational distribution $q(z_i)$
 - Update the global variational distributions $q(\mu_k)$ with **natural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

Black-box Variational Inference

Black-box Variational Inference (BBVI)

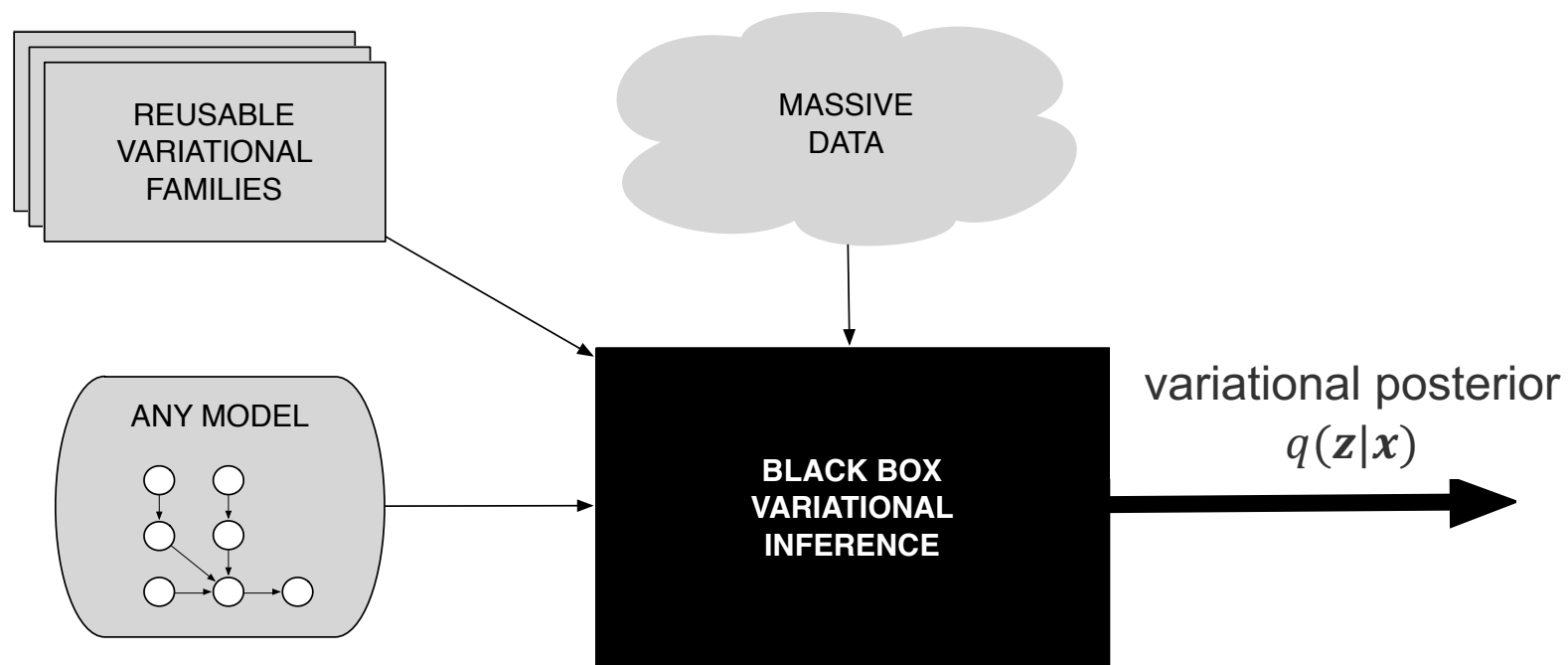
- We have derived variational inference specific for Bayesian Gaussian (mixture) models
- There are innumerable models
- Can we have a solution that does not entail model-specific work?

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

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

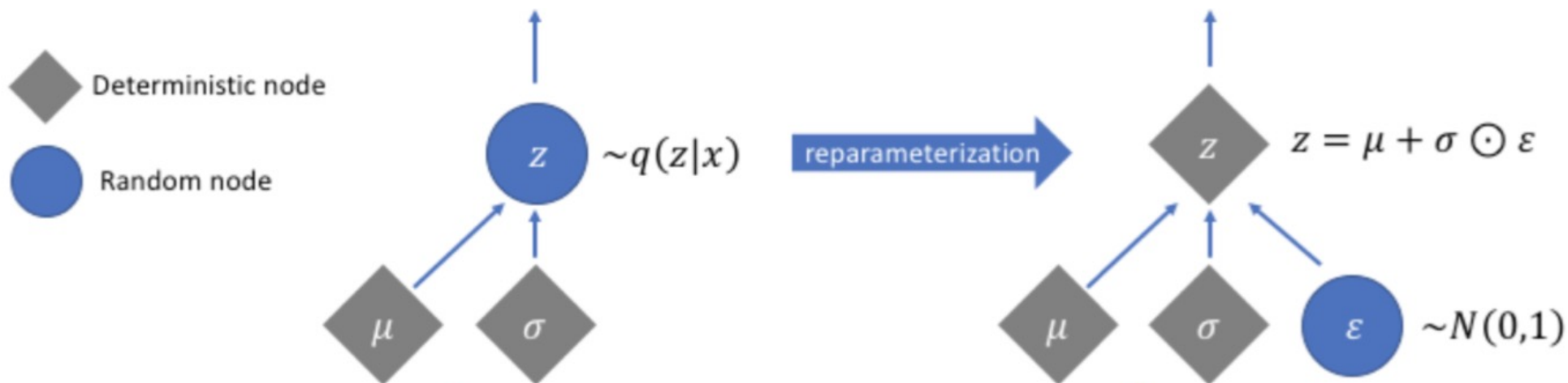
$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_{\mathbf{z}} f_\lambda(\mathbf{z}) \nabla_\lambda t(\epsilon, \lambda)]$$

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

Reparameterization trick

- Reparametrizing Gaussian distribution

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



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?