# DSC291: Machine Learning with Few Labels

## Unsupervised Learning

**Zhiting Hu**

Lecture 18, May 15, 2024

UC San Diego

**HALICIOĞLU DATA SCIENCE INSTITUTE**
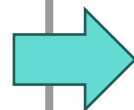
# Recap: EM and Variational Inference

- The EM algorithm:

  - E-step: $q^{t+1} = \arg\min_q F(q, \theta^t)$

    $= p(z|x, \theta^t) = \dfrac{p(z, x|\theta^t)}{\sum_z p(z, x|\theta^t)}$

**Intractable** when model $p(z, x|\theta)$ is complex

Approximate $p(z|x, \theta^t)$:
  - find a **tractable** $q(z|x, v^*)$ that is closest to $p(z|x, \theta^t)$

$q(z|x, v^*) = \min_v \mathrm{KL}\big(q(z|x, v) \,||\, p(z|x, \theta^t)\big)$

$= \min_v F(q(z|x, v), \theta^t) + const.$

$p(\mathbf{z}\,|\,\mathbf{x})$

$\mathrm{KL}(q(\mathbf{z}; v^*) \,||\, p(\mathbf{z}\,|\,\mathbf{x}))$

$q(\mathbf{z}; v)$   $v^*$

$v^{\text{init}}$
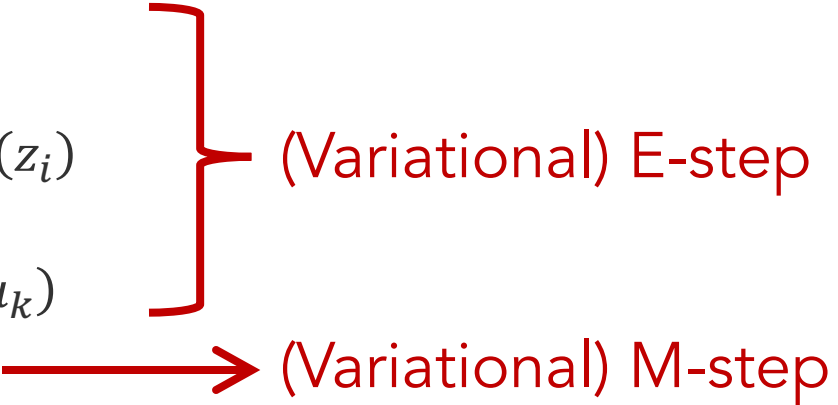
**Question:** What forms of $q(z|x, v)$ shall we choose?

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

# 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 $q(z_i)$    (Variational) E-step
    - End for
    - Update the global variational distributions $q(\mu_k)$
    - Update the parameters $\{\tau^2, \sigma^2, \pi\}$    (Variational) M-step
- Until ELBO converges

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

# Stochastic VI

Recap: Bayesian mixture of Gaussians

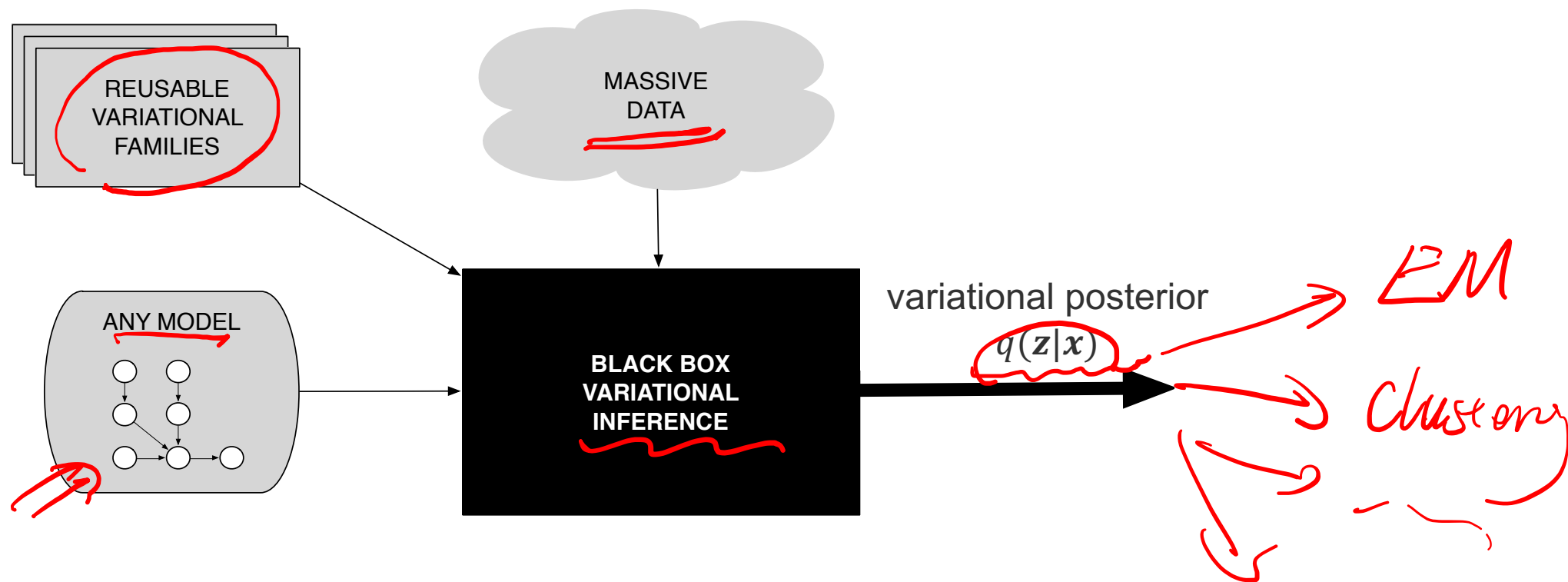Assume mean-field $q(\mu_{1:K}, z_{1:n}) = \prod_k q(\mu_k) \prod_i q(z_i)$

- Initialize the global variational distributions $q(\mu_k)$ and parameters $\{\tau^2, \sigma^2, \pi\}$
- Repeat:
  - Sample a data example $i \in \{1, 2, \ldots, D\}$
  - 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

[Hoffman et al., Stochastic Variational Inference, 2013]

# 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**
- **No mathematical work** beyond specifying the model
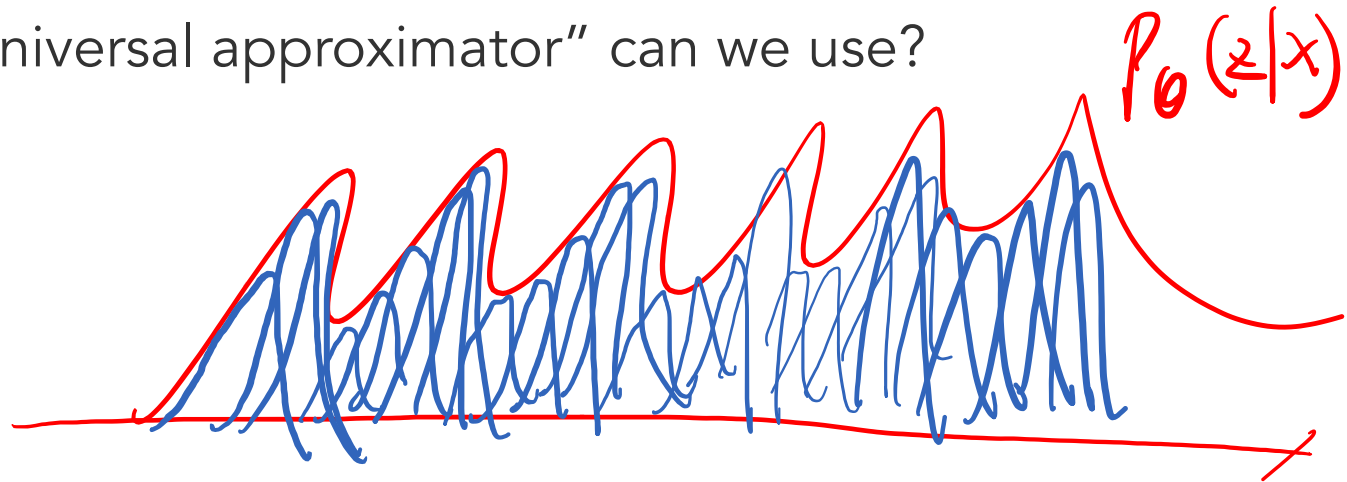- Perform inference with **massive data**

(Courtesy: Blei et al., 2018)

# Black-box Variational Inference (BBVI)

- Probabilistic model: $x$ -- observed variables, $z$ -- latent variables

- Variational distribution $q_\lambda(z|x)$ with parameters $\lambda$, 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)
  - **Question:** what other "universal approximator" can we use?

[Ranganath et al.,14]

# Black-box Variational Inference (BBVI)

- Probabilistic model: $\boldsymbol{x}$ -- observed variables, $\boldsymbol{z}$ -- latent variables
- Variational distribution $q_\lambda(\boldsymbol{z}|\boldsymbol{x})$ with parameters $\lambda$, 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)
  - **Question:** what other "universal approximator" can we use?
    Deep neural networks ↝ *VAE*

- ELBO to be maximized:

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

- Want to compute the gradient w.r.t variational parameters $\lambda$

[Ranganath et al.,14]

# The General Problem: Computing Gradients of Expectations

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

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

$$\mathbb{E}_{q_\lambda(z)}[g(z)]$$

- Computing exact gradients w.r.t. the parameters $\lambda$ is often infeasible

$\rightarrow KL$

- 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)

score-based GM
diffusion models

$\rightarrow$ VAE
$\rightarrow$ GAN

12

# Computing Gradients of Expectations w/ score function

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

- Log-derivative trick:  $\nabla_\lambda q_\lambda = q_\lambda \nabla_\lambda \log q_\lambda$

- **Question:** show that the gradient of $\mathcal{L}$ w.r.t. $\lambda$ is:

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

$$\nabla_\lambda \log q_\lambda = \frac{1}{q_\lambda} \nabla_\lambda q_\lambda$$

$$\nabla_\lambda \mathbb{E}_{q_\lambda}[f_\lambda] = \nabla_\lambda \int q_\lambda f_\lambda = \int \nabla_\lambda q_\lambda \cdot f_\lambda + q_\lambda \cdot \nabla_\lambda f_\lambda$$

$$= \int q_\lambda \cdot \nabla_\lambda \log q_\lambda f_\lambda + q_\lambda \cdot \nabla_\lambda f_\lambda$$

$$= \mathbb{E}_{q_\lambda}[\nabla_\lambda \log q_\lambda f_\lambda + \nabla_\lambda f_\lambda]$$

13

# Computing Gradients of Expectations w/ score function

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

- Log-derivative trick:  $\nabla_\lambda q_\lambda = q_\lambda \nabla_\lambda \log q_\lambda$

- Gradient of $\mathcal{L}$ w.r.t. $\lambda$:

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

  - score function: the gradient of the log of a probability distribution

- Monte Carlo estimation of the expectation:
  - Compute noisy unbiased gradients with Monte Carlo samples from $q_\lambda$

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

*(handwritten annotations)*

score- based GM

$a = \mathbb{E}_{p(y)}[g(y)]$

$y_i \sim p(y)$

$a \approx \frac{1}{n}\sum_i g(y_i)$

14

# Computing Gradients of Expectations w/ score function

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

- Log-derivative trick: $\nabla_\lambda q_\lambda = q_\lambda \nabla_\lambda \log q_\lambda$
- Gradient of $\mathcal{L}$ w.r.t. $\lambda$:

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

  - score function: the gradient of the log of a probability distribution
- Monte Carlo estimation of the expectation:
  - Compute noisy unbiased gradients with Monte Carlo samples from $q_\lambda$

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

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

# Computing Gradients of Expectations w/ reparametrization trick

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

- Assume that we can express the distribution $q_\lambda(z)$ with a transformation

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

  - E.g.,

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

- Reparameterization gradient:

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

  - **Question:** what's the gradient of $\mathcal{L}$ w.r.t. $\lambda$ ?

# 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

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

  - E.g.,

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

- Reparameterization gradient:

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

  - **Question:** what's the gradient of $\mathcal{L}$ w.r.t. $\lambda$ ?

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

$$\nabla_\lambda \mathcal{L} \approx \frac{1}{S} \sum_{s=1}^{S} \nabla_z f_\lambda(z) \nabla_\lambda t(\epsilon_i)$$

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}\left[\nabla_\lambda f_\lambda(z_\lambda)\right]$$

$$= \mathbb{E}_{\epsilon \sim s(\epsilon)}\left[\nabla_z f_\lambda \nabla_\lambda z_\lambda\right]$$

$$\epsilon_i \sim s(\epsilon)$$

17

# 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

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

  - E.g.,

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

- Reparameterization gradient

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

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

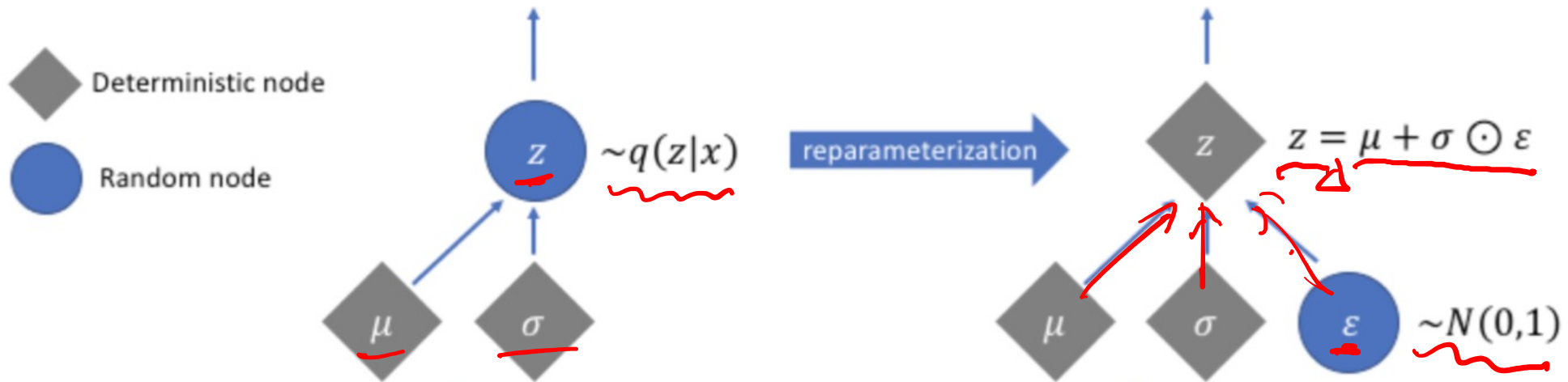*VAE 2014*

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

# Reparameterization trick

- Reparametrizing Gaussian distribution

$$\epsilon \sim Normal(0,1)$$
$$z = \epsilon\sigma + \mu$$

$$\Leftrightarrow \quad z \sim Normal(\mu, \sigma^2)$$



reparameterization

$$z = \mu + \sigma \odot \varepsilon$$

$$\sim N(0,1)$$

$$\sim q(z|x)$$

Deterministic node

Random node

19

# Reparameterization trick

- Reparametrizing Gaussian distribution

$$\epsilon \sim Normal(0,1)$$
$$z = \epsilon\sigma + \mu \qquad \Longleftrightarrow \qquad z \sim Normal(\mu, \sigma^2)$$

- Other reparameterizable distributions:
  $$\epsilon \sim Uniform(\epsilon)$$
  - Tractable inverse CDF $F^{-1}$:
    $$z = F^{-1}(\epsilon) \qquad \Longleftrightarrow \qquad z \sim q(z)$$
    - 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:  $\mathcal{L} = \mathbb{E}_{q_\lambda(\boldsymbol{z})}[f_\lambda(\boldsymbol{z})]$

- Score gradient

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

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

- Reparameterization gradient

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

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

# Recall: Black-box Variational Inference (BBVI)

- Probabilistic model: $x$ -- observed variables, $z$ -- latent variables
- Variational distribution $q_\lambda(z|x)$ with parameters $\lambda$, 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

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

- ELBO to be maximized:

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

- Want to compute the gradient w.r.t variational parameters $\lambda$

[Ranganath et al.,14]

# BBVI with the score gradient

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

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

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

- **Question:** what's the score gradient w.r.t. $\lambda$ ?

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

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

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

$$\text{where } z_s \sim q(z|\lambda).$$

[Ranganath et al.,14]

23

# BBVI with the reparameterization gradient

- ELBO:

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

- **Question:** what's the reparamerization gradient w.r.t. $\lambda$ ?

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

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

$$
\boxed{
\begin{aligned}
\mathcal{L} &= \mathbb{E}_{q_\lambda(\boldsymbol{z})}[f_\lambda(\boldsymbol{z})] \\
\nabla_\lambda \mathcal{L} &= \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_{\boldsymbol{z}} f_\lambda(\boldsymbol{z})\ \nabla_\lambda t(\epsilon, \lambda)]
\end{aligned}
}
$$

# Questions?