# Automatic Differentiation Variational Inference

Philip Schulz and Wilker Aziz

https: //github.com/philschulz/VITutorial

# What we know so far

- DGMs:

# What we know so far

- ▶ DGMs: probabilistic models parameterised by neural networks

# What we know so far

- ▶ DGMs: probabilistic models parameterised by neural networks
- ▶ Objective:

# What we know so far

▶ DGMs: probabilistic models parameterised by neural networks
▶ Objective: lowerbound on likelihood (ELBO)

# What we know so far

▶ DGMs: probabilistic models parameterised by neural networks
▶ Objective: lowerbound on likelihood (ELBO)
  ▶ cannot be computed exactly

# What we know so far

▶ DGMs: probabilistic models parameterised by neural networks
▶ Objective: lowerbound on likelihood (ELBO)
  ▶ cannot be computed exactly
    we resort to Monte Carlo estimation

# What we know so far

- ▶ DGMs: probabilistic models parameterised by neural networks
- ▶ Objective: lowerbound on likelihood (ELBO)
  - ▶ cannot be computed exactly
    we resort to Monte Carlo estimation
- ▶ But the MC estimator is not differentiable

# What we know so far

- ▶ DGMs: probabilistic models parameterised by neural networks
- ▶ Objective: lowerbound on likelihood (ELBO)
  - ▶ cannot be computed exactly
    we resort to Monte Carlo estimation
- ▶ But the MC estimator is not differentiable
  - ▶ Score function estimator: applicable to any model

# What we know so far

- ▶ DGMs: probabilistic models parameterised by neural networks
- ▶ Objective: lowerbound on likelihood (ELBO)
  - ▶ cannot be computed exactly
    we resort to Monte Carlo estimation
- ▶ But the MC estimator is not differentiable
  - ▶ Score function estimator: applicable to any model
  - ▶ Reparameterised gradients
    so far seems applicable only to Gaussian variables

# Reparameterised gradients: Gaussian

We have seen one case, namely,
if $\epsilon \sim \mathcal{N}(0, I)$ and $Z \sim \mathcal{N}(\mu, \sigma^2)$

# Reparameterised gradients: Gaussian

We have seen one case, namely,
  if $\epsilon \sim \mathcal{N}(0, I)$ and $Z \sim \mathcal{N}(\mu, \sigma^2)$
Then

$$Z \sim \mu + \sigma\epsilon$$

and

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{\mathcal{N}(z|\mu, \sigma^2)} \left[ g(z) \right]$$

# Reparameterised gradients: Gaussian

We have seen one case, namely,
if $\epsilon \sim \mathcal{N}(0, I)$ and $Z \sim \mathcal{N}(\mu, \sigma^2)$
Then

$$Z \sim \mu + \sigma\epsilon$$

and

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{\mathcal{N}(z|\mu,\sigma^2)} [g(z)]$$

$$= \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial \lambda} g(z = \mu + \sigma\epsilon) \right]$$

# Reparameterised gradients: Gaussian

We have seen one case, namely,
if $\epsilon \sim \mathcal{N}(0, I)$ and $Z \sim \mathcal{N}(\mu, \sigma^2)$
Then

$$Z \sim \mu + \sigma\epsilon$$

and

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{\mathcal{N}(z|\mu,\sigma^2)}[g(z)]$$

$$= \mathbb{E}_{\mathcal{N}(0,I)}\left[\frac{\partial}{\partial \lambda} g(z = \mu + \sigma\epsilon)\right]$$

$$= \mathbb{E}_{\mathcal{N}(0,I)}\left[\frac{\partial}{\partial z} g(z = \mu + \sigma\epsilon)\frac{\partial z}{\partial \lambda}\right]$$

# Reparameterised gradients: Gaussian

Location

$$\frac{\partial}{\partial \mu} \mathbb{E}_{\mathcal{N}(z|\mu,\sigma^2)}\left[g(z)\right] = \mathbb{E}_{\mathcal{N}(0,I)}\left[\frac{\partial}{\partial z} g(z = \mu + \sigma\epsilon)\frac{\partial z}{\partial \mu}\right]$$

# Reparameterised gradients: Gaussian

Location

$$\frac{\partial}{\partial \mu} \mathbb{E}_{\mathcal{N}(z|\mu,\sigma^2)} \left[ g(z) \right] = \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma \epsilon) \frac{\partial z}{\partial \mu} \right]$$

$$= \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma \epsilon) \right]$$

# Reparameterised gradients: Gaussian

Location

$$\frac{\partial}{\partial \mu} \mathbb{E}_{\mathcal{N}(z|\mu,\sigma^2)} \left[ g(z) \right] = \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma \epsilon) \frac{\partial z}{\partial \mu} \right]$$

$$= \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma \epsilon) \right]$$

Scale

$$\frac{\partial}{\partial \sigma} \mathbb{E}_{\mathcal{N}(z|\mu,\sigma^2)} \left[ g(z) \right] = \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma \epsilon) \frac{\partial z}{\partial \sigma} \right]$$

# Reparameterised gradients: Gaussian

Location

$$\frac{\partial}{\partial \mu} \mathbb{E}_{\mathcal{N}(z|\mu,\sigma^2)} \left[ g(z) \right] = \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma \epsilon) \frac{\partial z}{\partial \mu} \right]$$

$$= \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma \epsilon) \right]$$

Scale

$$\frac{\partial}{\partial \sigma} \mathbb{E}_{\mathcal{N}(z|\mu,\sigma^2)} \left[ g(z) \right] = \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma \epsilon) \frac{\partial z}{\partial \sigma} \right]$$

$$= \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma \epsilon) \epsilon \right]$$

# Reparameterised gradients: Gaussian

Location

$$\frac{\partial}{\partial \mu} \mathbb{E}_{\mathcal{N}(z|\mu,\sigma^2)} \left[ g(z) \right] = \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma\epsilon) \frac{\partial z}{\partial \mu} \right]$$

$$= \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma\epsilon) \right]$$

Scale

$$\frac{\partial}{\partial \sigma} \mathbb{E}_{\mathcal{N}(z|\mu,\sigma^2)} \left[ g(z) \right] = \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma\epsilon) \frac{\partial z}{\partial \sigma} \right]$$

$$= \mathbb{E}_{\mathcal{N}(0,I)} \left[ \frac{\partial}{\partial z} g(z = \mu + \sigma\epsilon)\epsilon \right]$$

# Multivariate calculus recap

Let $x \in \mathbb{R}^K$ and let $\mathcal{T} : \mathbb{R}^K \to \mathbb{R}^K$ be differentiable and invertible

- $y = \mathcal{T}(x)$
- $x = \mathcal{T}^{-1}(y)$

# Jacobian

The Jacobian matrix $\mathbf{J} = J_{\mathcal{T}}(x)$ of $\mathcal{T}$
assessed at $x$ is the matrix of partial derivatives

$$J_{ij} = \frac{\partial y_i}{\partial x_j}$$

# Jacobian

The Jacobian matrix $\mathbf{J} = J_{\mathcal{T}}(x)$ of $\mathcal{T}$
assessed at $x$ is the matrix of partial derivatives

$$J_{ij} = \frac{\partial y_i}{\partial x_j}$$

Inverse function theorem

$$J_{\mathcal{T}^{-1}}(y) = (J_{\mathcal{T}}(x))^{-1}$$

# Differential (or inifinitesimal)

The **differential** $\mathrm{d}x$ of $x$
refers to an *infinitely small* change in $x$

# Differential (or inifinitesimal)

The **differential** $\mathrm{d}x$ of $x$
  refers to an *infinitely small* change in $x$

We can relate the differential $\mathrm{d}y$ of $y = \mathcal{T}(x)$ to $\mathrm{d}x$

# Differential (or inifinitesimal)

The **differential** $\mathrm{d}x$ of $x$
refers to an *infinitely small* change in $x$

We can relate the differential $\mathrm{d}y$ of $y = \mathcal{T}(x)$ to $\mathrm{d}x$

▶ Scalar case

$$\mathrm{d}y = \mathcal{T}'(x)\mathrm{d}x = \frac{\mathrm{d}y}{\mathrm{d}x}\mathrm{d}x = \frac{\mathrm{d}}{\mathrm{d}x}\mathcal{T}(x)\mathrm{d}x$$

where $\mathrm{d}y/\mathrm{d}x$ is the *derivative* of $y$ wrt $x$

# Differential (or inifinitesimal)

The **differential** $\mathrm{d}x$ of $x$
refers to an *infinitely small* change in $x$

We can relate the differential $\mathrm{d}y$ of $y = \mathcal{T}(x)$ to $\mathrm{d}x$

▶ Scalar case

$$\mathrm{d}y = \mathcal{T}'(x)\mathrm{d}x = \frac{\mathrm{d}y}{\mathrm{d}x}\mathrm{d}x = \frac{\mathrm{d}}{\mathrm{d}x}\mathcal{T}(x)\mathrm{d}x$$

where $\mathrm{d}y/\mathrm{d}x$ is the *derivative* of $y$ wrt $x$

▶ Multivariate case

$$\mathrm{d}y = |\det J_{\mathcal{T}}(x)| \, \mathrm{d}x$$

the absolute value absorbs the orientation

# Integration by substitution

We can integrate a function $g(x)$
by substituting $x = \mathcal{T}^{-1}(y)$

$$\int g(x)\mathrm{d}x$$

# Integration by substitution

We can integrate a function $g(x)$
by substituting $x = \mathcal{T}^{-1}(y)$

$$\int g(x)\mathrm{d}x = \int g(\underbrace{\mathcal{T}^{-1}(y)}_{x}) \underbrace{|\det J_{\mathcal{T}^{-1}}(y)| \, \mathrm{d}y}_{\mathrm{d}x}$$

# Integration by substitution

We can integrate a function $g(x)$
by substituting $x = \mathcal{T}^{-1}(y)$

$$\int g(x)\mathrm{d}x = \int g(\underbrace{\mathcal{T}^{-1}(y)}_{x}) \underbrace{|\det J_{\mathcal{T}^{-1}}(y)|\,\mathrm{d}y}_{\mathrm{d}x}$$

and similarly for a function $h(y)$

$$\int h(y)\mathrm{d}y$$

# Integration by substitution

We can integrate a function $g(x)$
by substituting $x = \mathcal{T}^{-1}(y)$

$$\int g(x)\mathrm{d}x = \int g(\underbrace{\mathcal{T}^{-1}(y)}_{x}) \underbrace{|\det J_{\mathcal{T}^{-1}}(y)|\,\mathrm{d}y}_{\mathrm{d}x}$$

and similarly for a function $h(y)$

$$\int h(y)\mathrm{d}y = \int h(\mathcal{T}(x))\,|\det J_{\mathcal{T}}(x)|\,\mathrm{d}x$$

# Change of density

Let $X$ take on values in $\mathbb{R}^K$ with density $f_X(x)$

# Change of density

Let $X$ take on values in $\mathbb{R}^K$ with density $f_X(x)$
and recall that $y = \mathcal{T}(x)$ and $x = \mathcal{T}^{-1}(y)$

# Change of density

Let $X$ take on values in $\mathbb{R}^K$ with density $f_X(x)$
and recall that $y = \mathcal{T}(x)$ and $x = \mathcal{T}^{-1}(y)$

Then $\mathcal{T}$ induces a density $f_Y(y)$ expressed as

$$f_Y(y) = f_X(x = \mathcal{T}^{-1}(y)) \left| \det J_{\mathcal{T}^{-1}}(y) \right|$$

# Change of density

Let $X$ take on values in $\mathbb{R}^K$ with density $f_X(x)$
and recall that $y = \mathcal{T}(x)$ and $x = \mathcal{T}^{-1}(y)$

Then $\mathcal{T}$ induces a density $f_Y(y)$ expressed as

$$f_Y(y) = f_X(x = \mathcal{T}^{-1}(y)) \left| \det J_{\mathcal{T}^{-1}}(y) \right|$$

and then it follows that

$$f_X(x) = f_Y(y = \mathcal{T}(x)) \left| \det J_{\mathcal{T}}(x) \right|$$

# Revisiting reparameterised gradients

Let $Z$ take on values in $\mathbb{R}^K$ with pdf $q(z|\lambda)$

# Revisiting reparameterised gradients

Let $Z$ take on values in $\mathbb{R}^K$ with pdf $q(z|\lambda)$

The idea is to count on a *standardisation* procedure

# Revisiting reparameterised gradients

Let $Z$ take on values in $\mathbb{R}^K$ with pdf $q(z|\lambda)$

The idea is to count on a *standardisation* procedure
a transformation $\mathcal{S}_\lambda : \mathbb{R}^K \to \mathbb{R}^K$ such that

# Revisiting reparameterised gradients

Let $Z$ take on values in $\mathbb{R}^K$ with pdf $q(z|\lambda)$

The idea is to count on a *standardisation* procedure
a transformation $\mathcal{S}_\lambda : \mathbb{R}^K \to \mathbb{R}^K$ such that

$$\mathcal{S}_\lambda(z) \sim \pi(\epsilon)$$
$$\mathcal{S}_\lambda^{-1}(\epsilon) \sim q(z|\lambda)$$

▶ $\pi(\epsilon)$ does not depend on parameters $\lambda$
we call it a *standard* density

# Revisiting reparameterised gradients

Let $Z$ take on values in $\mathbb{R}^K$ with pdf $q(z|\lambda)$

The idea is to count on a *standardisation* procedure
a transformation $\mathcal{S}_\lambda : \mathbb{R}^K \to \mathbb{R}^K$ such that

$$\mathcal{S}_\lambda(z) \sim \pi(\epsilon)$$
$$\mathcal{S}_\lambda^{-1}(\epsilon) \sim q(z|\lambda)$$

▶ $\pi(\epsilon)$ does not depend on parameters $\lambda$
  we call it a *standard* density
▶ $\mathcal{S}_\lambda(z)$ absorbs dependency on $\lambda$

# Reparameterised expectations

If we are interested in

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

# Reparameterised expectations

If we are interested in

$$\mathbb{E}_{q(z|\lambda)}\left[g(z)\right] = \int q(z|\lambda)g(z)\mathrm{d}z$$

# Reparameterised expectations

If we are interested in

$$\mathbb{E}_{q(z|\lambda)}[g(z)] = \int q(z|\lambda)g(z)\mathrm{d}z$$

$$= \int \underbrace{\pi(\mathcal{S}_\lambda(z))\,|\det J_{S_\lambda}(z)|}_{\text{change of density}}\,g(z)\mathrm{d}z$$

# Reparameterised expectations

If we are interested in

$$\mathbb{E}_{q(z|\lambda)}[g(z)] = \int q(z|\lambda)g(z)\mathrm{d}z$$

$$= \int \underbrace{\pi(\mathcal{S}_\lambda(z)) \left|\det J_{\mathcal{S}_\lambda}(z)\right|}_{\text{change of density}} g(z)\mathrm{d}z$$

$$= \int \pi(\epsilon)$$

# Reparameterised expectations

If we are interested in

$$\mathbb{E}_{q(z|\lambda)}\left[g(z)\right] = \int q(z|\lambda)g(z)\mathrm{d}z$$

$$= \int \underbrace{\pi(\mathcal{S}_\lambda(z))\left|\det J_{S_\lambda}(z)\right|}_{\text{change of density}} g(z)\mathrm{d}z$$

$$= \int \pi(\epsilon) \underbrace{\left|\det J_{\mathcal{S}_\lambda^{-1}}(\epsilon)\right|^{-1}}_{\text{inv func theorem}}$$

# Reparameterised expectations

If we are interested in

$$\mathbb{E}_{q(z|\lambda)}[g(z)] = \int q(z|\lambda)g(z)\mathrm{d}z$$

$$= \int \underbrace{\pi(\mathcal{S}_\lambda(z)) \left|\det J_{\mathcal{S}_\lambda}(z)\right|}_{\text{change of density}} g(z)\mathrm{d}z$$

$$= \int \pi(\epsilon) \underbrace{\left|\det J_{\mathcal{S}_\lambda^{-1}}(\epsilon)\right|^{-1}}_{\text{inv func theorem}} g(\underbrace{\mathcal{S}_\lambda^{-1}(\epsilon)}_{z})$$

# Reparameterised expectations

If we are interested in

$$\mathbb{E}_{q(z|\lambda)}\left[g(z)\right] = \int q(z|\lambda)g(z)\mathrm{d}z$$

$$= \int \underbrace{\pi(\mathcal{S}_\lambda(z))\left|\det J_{\mathcal{S}_\lambda}(z)\right|}_{\text{change of density}}g(z)\mathrm{d}z$$

$$= \int \pi(\epsilon)\underbrace{\left|\det J_{\mathcal{S}_\lambda^{-1}}(\epsilon)\right|^{-1}}_{\text{inv func theorem}}g(\underbrace{\mathcal{S}_\lambda^{-1}(\epsilon)}_{z})\underbrace{\left|\det J_{\mathcal{S}_\lambda^{-1}}(\epsilon)\right|}_{\text{change of var}}\mathrm{d}\epsilon$$

# Reparameterised expectations

If we are interested in

$$\mathbb{E}_{q(z|\lambda)}[g(z)] = \int q(z|\lambda)g(z)\mathrm{d}z$$

$$= \int \underbrace{\pi(\mathcal{S}_\lambda(z)) \left|\det J_{S_\lambda}(z)\right|}_{\text{change of density}} g(z)\mathrm{d}z$$

$$= \int \pi(\epsilon) \underbrace{\left|\det J_{\mathcal{S}_\lambda^{-1}}(\epsilon)\right|^{-1}}_{\text{inv func theorem}} g(\underbrace{\mathcal{S}_\lambda^{-1}(\epsilon)}_{z}) \underbrace{\left|\det J_{\mathcal{S}_\lambda^{-1}}(\epsilon)\right|}_{\text{change of var}} \mathrm{d}\epsilon$$

$$= \int \pi(\epsilon)g(\mathcal{S}_\lambda^{-1}(\epsilon))\mathrm{d}\epsilon$$

# Reparameterised expectations

If we are interested in

$$\mathbb{E}_{q(z|\lambda)}[g(z)] = \int q(z|\lambda)g(z)\mathrm{d}z$$

$$= \int \underbrace{\pi(\mathcal{S}_\lambda(z))\left|\det J_{S_\lambda}(z)\right|}_{\text{change of density}} g(z)\mathrm{d}z$$

$$= \int \pi(\epsilon) \underbrace{\left|\det J_{\mathcal{S}_\lambda^{-1}}(\epsilon)\right|^{-1}}_{\text{inv func theorem}} g(\underbrace{\mathcal{S}_\lambda^{-1}(\epsilon)}_{z}) \underbrace{\left|\det J_{\mathcal{S}_\lambda^{-1}}(\epsilon)\right|}_{\text{change of var}}\mathrm{d}\epsilon$$

$$= \int \pi(\epsilon)g(\mathcal{S}_\lambda^{-1}(\epsilon))\mathrm{d}\epsilon = \mathbb{E}_{\pi(\epsilon)}\left[g(\mathcal{S}_\lambda^{-1}(\epsilon))\right]$$

# Reparameterised gradients

For optimisation, we need tractable gradients

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{q(z|\lambda)}\left[g(z)\right] = \frac{\partial}{\partial \lambda} \mathbb{E}_{\pi(\epsilon)}\left[g(\mathcal{S}_\lambda^{-1}(\epsilon))\right]$$

# Reparameterised gradients

For optimisation, we need tractable gradients

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{q(z|\lambda)}\left[g(z)\right] = \frac{\partial}{\partial \lambda} \mathbb{E}_{\pi(\epsilon)}\left[g(\mathcal{S}_\lambda^{-1}(\epsilon))\right]$$

since now the measure of integration does not depend on $\lambda$, we can obtain a gradient estimate

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{q(z|\lambda)}\left[g(z)\right] = \mathbb{E}_{\pi(\epsilon)}\left[\frac{\partial}{\partial \lambda} g(\mathcal{S}_\lambda^{-1}(\epsilon))\right]$$

# Reparameterised gradients

For optimisation, we need tractable gradients

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{q(z|\lambda)} \left[ g(z) \right] = \frac{\partial}{\partial \lambda} \mathbb{E}_{\pi(\epsilon)} \left[ g(\mathcal{S}_\lambda^{-1}(\epsilon)) \right]$$

since now the measure of integration does not depend on $\lambda$, we can obtain a gradient estimate

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{q(z|\lambda)} \left[ g(z) \right] = \mathbb{E}_{\pi(\epsilon)} \left[ \frac{\partial}{\partial \lambda} g(\mathcal{S}_\lambda^{-1}(\epsilon)) \right]$$

$$\overset{\text{MC}}{\approx} \frac{1}{M} \sum_{\substack{i=1 \\ \epsilon_i \sim \pi(\epsilon)}}^{M} \frac{\partial}{\partial \lambda} g(\mathcal{S}_\lambda^{-1}(\epsilon_i))$$

# Standardisation functions

Location-scale family

▶ a family of distributions where for
$F_X(x) = \Pr\{X \leq x\}$
if $Y = a + bX$, then $F_Y(y|a, b) = F_X(\frac{z-a}{b})$

# Standardisation functions

Location-scale family

- ▶ a family of distributions where for
  $F_X(x) = \Pr\{X \leq x\}$
  if $Y = a + bX$, then $F_Y(y|a, b) = F_X(\frac{z-a}{b})$
- ▶ if we can draw from $f_X(x)$, we can draw from
  $f_Y(y|a, b)$

# Standardisation functions

Location-scale family

- ▶ a family of distributions where for
  $F_X(x) = \Pr\{X \leq x\}$
  if $Y = a + bX$, then $F_Y(y|a, b) = F_X(\frac{z-a}{b})$
- ▶ if we can draw from $f_X(x)$, we can draw from
  $f_Y(y|a, b)$
- ▶ the transformation absorbs the parameters $a, b$

# Standardisation functions

Location-scale family

- ▶ a family of distributions where for
  $F_X(x) = \Pr\{X \leq x\}$
  if $Y = a + bX$, then $F_Y(y|a, b) = F_X(\frac{z-a}{b})$

- ▶ if we can draw from $f_X(x)$, we can draw from
  $f_Y(y|a, b)$

- ▶ the transformation absorbs the parameters $a, b$

Examples: Gaussian, Laplace, Cauchy, Uniform

# Standardisation functions (cont.)

Inverse cdf

▶ for univariate $Z$ with pdf $f_Z(z)$ and cdf $F_Z(z)$

$$P \sim \mathcal{U}(0, 1) \qquad Z \sim F_Z^{-1}(P)$$

where $F_Z^{-1}(p)$ is the *quantile function*

# Standardisation functions (cont.)

Inverse cdf

▶ for univariate $Z$ with pdf $f_Z(z)$ and cdf $F_Z(z)$

$$P \sim \mathcal{U}(0,1) \qquad Z \sim F_Z^{-1}(P)$$

where $F_Z^{-1}(p)$ is the *quantile function*

Gumbel distribution

▶ $f_Z(z|\mu,\beta) = \beta^{-1}\exp(-z - \exp(-z))$

▶ $F_Z(z|\mu,\beta) = \exp\left(-\exp\left(-\frac{z-\mu}{\beta}\right)\right)$

▶ $F_Z^{-1}(p) = \mu - \beta\log(-\log p)$

# Beyond

Many interesting densities are not location-scale families

- ▶ e.g. Beta, Gamma

# Beyond

Many interesting densities are not location-scale families

▶ e.g. Beta, Gamma

The inverse cdf of a multivariate rv is seldom known in closed-form

▶ Dirichlet, von Mises-Fisher

# Automatic Differentiation VI

Motivation

▶ many models have intractable posteriors
their normalising constants (evidence) lacks
analytic solutions

# Automatic Differentiation VI

Motivation

- ▶ many models have intractable posteriors
  their normalising constants (evidence) lacks
  analytic solutions

- ▶ but many models are differentiable
  that's the main constraint for using NNs

# Automatic Differentiation VI

Motivation

- ▶ many models have intractable posteriors
  their normalising constants (evidence) lacks
  analytic solutions

- ▶ but many models are differentiable
  that's the main constraint for using NNs

Reparameterised gradients are a step towards
automatising VI for differentiable models

# Automatic Differentiation VI

Motivation

- ▶ many models have intractable posteriors
  their normalising constants (evidence) lacks
  analytic solutions

- ▶ but many models are differentiable
  that's the main constraint for using NNs

Reparameterised gradients are a step towards
automatising VI for differentiable models

- ▶ but not every model of interest employs rvs for
  which a standardisation function is known

# Example

Suppose we have some ordinal data which we assume to be Poisson-distributed

$$X|\lambda \sim \text{Poisson}(\lambda)$$

and suppose we want to impose

# Differentiable models

We focus on *differentiable probability models*

$$p(x, z) = p(x|z)p(z)$$

# Differentiable models

We focus on *differentiable probability models*

$$p(x, z) = p(x|z)p(z)$$

▶ members of this class have continuous latent variables $z$

# Differentiable models

We focus on *differentiable probability models*

$$p(x, z) = p(x|z)p(z)$$

▶ members of this class have continuous latent variables $z$

▶ and the gradient $\nabla_z \log p(x, z)$ is valid within the *support* of the prior
$\text{supp}(p(z)) = \{z \in \mathbb{R}^K : p(z) > 0\} \subseteq \mathbb{R}^K$

# VI optimisation problem

Let's focus on the design and optimisation of the variational approximation

$$\arg\min_{q(z)} \mathrm{KL}\left(q(z) \,\|\, p(z|x)\right)$$

# VI optimisation problem

Let's focus on the design and optimisation of the variational approximation

$$\arg\min_{q(z)} \text{KL}\left(q(z) \,||\, p(z|x)\right)$$

To automatise the search for a variational approximation $q(z)$ we must ensure that

$$\text{supp}(q(z)) \subseteq \text{supp}(p(z|x))$$

# VI optimisation problem

Let's focus on the design and optimisation of the variational approximation

$$\arg\min_{q(z)} \text{KL}\left(q(z) \,||\, p(z|x)\right)$$

To automatise the search for a variational approximation $q(z)$ we must ensure that

$$\text{supp}(q(z)) \subseteq \text{supp}(p(z|x))$$

▶ otherwise KL is not defined
$$\text{KL}\left(q \,||\, p\right) = \mathbb{E}_q\left[\log q\right] - \mathbb{E}_q\left[\log p\right] = \infty$$

# Support matching constraint

So let's constrain $q(z)$ to a family $\mathcal{Q}$ whose support is included in the support of the posterior

$$\arg\min_{q(z) \in \mathcal{Q}} \mathrm{KL}\left(q(z) \,\|\, p(z|x)\right)$$

where

$$\mathcal{Q} = \{q(z) : \mathrm{supp}(q(z)) \subseteq \mathrm{supp}(p(z|x))\}$$

# Support matching constraint

So let's constrain $q(z)$ to a family $\mathcal{Q}$ whose support is included in the support of the posterior

$$\underset{q(z) \in \mathcal{Q}}{\arg\min} \, \mathsf{KL}\left(q(z) \,||\, p(z|x)\right)$$

where

$$\mathcal{Q} = \{q(z) : \mathsf{supp}(q(z)) \subseteq \mathsf{supp}(p(z|x))\}$$

But what is the support of $p(z|x)$?

# Support matching constraint

So let's constrain $q(z)$ to a family $\mathcal{Q}$ whose support is included in the support of the posterior

$$\operatorname*{arg\,min}_{q(z) \in \mathcal{Q}} \operatorname{KL}\left(q(z) \,\|\, p(z|x)\right)$$

where

$$\mathcal{Q} = \{q(z) : \operatorname{supp}(q(z)) \subseteq \operatorname{supp}(p(z|x))\}$$

But what is the support of $p(z|x)$?

▶ typically the same as the support of $p(z)$

# Support matching constraint

So let's constrain $q(z)$ to a family $\mathcal{Q}$ whose support is included in the support of the posterior

$$\arg\min_{q(z)\in\mathcal{Q}} \text{KL}\left(q(z) \,||\, p(z|x)\right)$$

where

$$\mathcal{Q} = \{q(z) : \text{supp}(q(z)) \subseteq \text{supp}(p(z|x))\}$$

But what is the support of $p(z|x)$?

▶ typically the same as the support of $p(z)$
  as long as $p(x, z) > 0$ if $p(z) > 0$

# Parametric family

So let's constrain $q(z)$ to a family $\mathcal{Q}$ whose support is included in the support of the prior

$$\arg\min_{q(z)\in\mathcal{Q}} \mathrm{KL}\left(q(z) \parallel p(z|x)\right)$$

where

$$\mathcal{Q} = \{q(z;\phi) : \phi \in \Phi, \mathrm{supp}(q(z;\phi)) \subseteq \mathrm{supp}(p(z))\}$$

# Parametric family

So let's constrain $q(z)$ to a family $\mathcal{Q}$ whose support is included in the support of the prior

$$\arg\min_{q(z)\in\mathcal{Q}} \mathrm{KL}\left(q(z) \;\|\; p(z|x)\right)$$

where

$$\mathcal{Q} = \left\{q(z; \phi) : \phi \in \Phi, \mathrm{supp}(q(z; \phi)) \subseteq \mathrm{supp}(p(z))\right\}$$

▶ a parameter vector $\phi$ picks out a member of the family

# Constrained optimisation for the ELBO

We maximise the ELBO

$$\arg\max_{\phi \in \Phi} \mathbb{E}_{q(z;\phi)} \left[ \log p(x, z) \right] + \mathbb{H} \left( q(z|\phi) \right)$$

subject to

$$\mathcal{Q} = \left\{ q(z; \phi) : \phi \in \Phi, \text{supp}(q(z; \phi)) \subseteq \text{supp}(p(z)) \right\}$$

# Constrained optimisation for the ELBO

We maximise the ELBO

$$\arg\max_{\phi \in \Phi} \mathbb{E}_{q(z;\phi)}\left[\log p(x, z)\right] + \mathbb{H}\left(q(z|\phi)\right)$$

subject to

$$\mathcal{Q} = \{q(z; \phi) : \phi \in \Phi, \operatorname{supp}(q(z; \phi)) \subseteq \operatorname{supp}(p(z))\}$$

There are really two constraints here

# Constrained optimisation for the ELBO

We maximise the ELBO

$$\arg\max_{\phi \in \Phi} \mathbb{E}_{q(z;\phi)}\left[\log p(x, z)\right] + \mathbb{H}\left(q(z|\phi)\right)$$

subject to

$$\mathcal{Q} = \left\{q(z; \phi) : \phi \in \Phi, \operatorname{supp}(q(z; \phi)) \subseteq \operatorname{supp}(p(z))\right\}$$

There are really two constraints here

▶ support matching constraint

# Constrained optimisation for the ELBO

We maximise the ELBO

$$\arg\max_{\phi\in\Phi}\mathbb{E}_{q(z;\phi)}\left[\log p(x,z)\right]+\mathbb{H}\left(q(z|\phi)\right)$$

subject to

$$\mathcal{Q}=\left\{q(z;\phi):\phi\in\Phi,\mathrm{supp}(q(z;\phi))\subseteq\mathrm{supp}(p(z))\right\}$$

There are really two constraints here

▶ support matching constraint

▶ $\Phi$ can be an intricate subset of $\mathbb{R}^D$

# Constrained optimisation for the ELBO

We maximise the ELBO

$$\arg\max_{\phi \in \Phi} \mathbb{E}_{q(z;\phi)} \left[ \log p(x, z) \right] + \mathbb{H} \left( q(z|\phi) \right)$$

subject to

$$\mathcal{Q} = \{ q(z; \phi) : \phi \in \Phi, \operatorname{supp}(q(z; \phi)) \subseteq \operatorname{supp}(p(z)) \}$$

There are really two constraints here

- ▶ support matching constraint
- ▶ $\Phi$ can be an intricate subset of $\mathbb{R}^D$
  e.g. univariate Gaussian location lives in $\mathbb{R}$ but
  scale lives in $\mathbb{R}_{>0}$

# ADVI

From the point of view of a black-box procedure, this objective poses two problems

1. intractable expectations

# ADVI

From the point of view of a black-box procedure, this objective poses two problems

1. intractable expectations Reparameterised Gradients!
2. custom $\text{supp}(q(z; \phi))$

Idea

1. let's find a way to transform $\text{supp}(p(z))$ to the complete real coordinate space
2. then we pick a variational family over the complete real coordinate space for which a standardisation exists!