# Reparameterizing Challenging Distributions

Francisco J. R. Ruiz

Joint work with Christian Naesseth, Scott Linderman, Michalis Titsias, David Blei

November 22nd, 2016



#### Overview

- Reparameterization allows for low-variance gradient estimators
- But it is available for some distributions only
- We show how to extend reparameterization to other distributions
- Allows 1-sample Monte Carlo estimation of gradient
- > Our goal: General variational inference for probabilistic models

#### Gradient of Expectations

Consider the gradient

 $\nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q(\mathbf{z};\boldsymbol{\lambda})}[f(\mathbf{z},\boldsymbol{\lambda})]$ 

w.r.t. some parameters  $\lambda$ 

#### Gradient of Expectations

Consider the gradient

$$abla_{\lambda} \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\lambda})} \left[ f(\mathbf{z}, \boldsymbol{\lambda}) \right]$$

w.r.t. some parameters  $\lambda$ 

▶ This is common in statistics/machine learning:

$$oldsymbol{\lambda}^{\star} = rg\max_{oldsymbol{\lambda}} \mathbb{E}_{q(\mathbf{z};oldsymbol{\lambda})}\left[f(\mathbf{z},oldsymbol{\lambda})
ight]$$

- Optimization of loss functions
- Reinforcement learning (policy gradient)
- Variational inference
- ▶ ...

Variational Inference (1/3)

$$egin{aligned} oldsymbol{\lambda}^{\star} &= rg\max_{oldsymbol{\lambda}} \mathbb{E}_{q(\mathbf{z};oldsymbol{\lambda})}\left[f(\mathbf{z},oldsymbol{\lambda})
ight] \end{aligned}$$

- > z: Latent variables in a probabilistic model
- 🕨 🗴: Data
- ▶ *p*(**x**, **z**): Probabilistic model
- $q(\mathbf{z}; \boldsymbol{\lambda})$ : Variational distribution

$$f(\mathsf{z}, \boldsymbol{\lambda}) = \log p(\mathsf{x}, \mathsf{z}) - \log q(\mathsf{z}; \boldsymbol{\lambda})$$

Variational Inference (2/3)

$$\lambda^{\star} = rg\max_{\lambda} \mathbb{E}_{q(\mathsf{z}; \lambda)} \left[\log p(\mathsf{x}, \mathsf{z}) - \log q(\mathsf{z}; \lambda)
ight]$$

- Variational inference approximates the posterior
- Minimizes the KL divergence between  $q(\mathbf{z}; \boldsymbol{\lambda})$  and the posterior

$$oldsymbol{\lambda}^{\star} = rgmin_{oldsymbol{\lambda}} D_{ ext{KL}}(q(\mathbf{z}; oldsymbol{\lambda}) || p(\mathbf{z} \,|\, \mathbf{x}))$$

# Variational Inference (3/3)

Optimization problem:

$$\lambda^{\star} = rg\max_{oldsymbol{\lambda}} \mathbb{E}_{q(\mathsf{z};oldsymbol{\lambda})} \left[\log p(\mathsf{x},\mathsf{z}) - \log q(\mathsf{z};oldsymbol{\lambda})
ight]$$

- Conditionally conjugate models: coordinate ascent
- Non-conjugate models: score function method, reparameterization

# Variational Inference (3/3)

Optimization problem:

$$\lambda^{\star} = rg\max_{oldsymbol{\lambda}} \mathbb{E}_{q(\mathbf{z};oldsymbol{\lambda})} \left[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; oldsymbol{\lambda})
ight]$$

- Conditionally conjugate models: coordinate ascent
- Non-conjugate models: score function method, reparameterization
- Score function and reparameterization are two ways to estimate

 $\nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q(\mathbf{z};\boldsymbol{\lambda})} [f(\mathbf{z},\boldsymbol{\lambda})]$ 

For simplicity, we focus on

$$f(\mathbf{z}) = \log p(\mathbf{x}, \mathbf{z})$$

(assume that the gradient of the entropy is tractable)

#### Score Function Method<sup>1</sup>

Monte Carlo estimator of the gradient

 $abla_{\lambda} \mathbb{E}_{q(\mathbf{z}; \lambda)} [f(\mathbf{z})]$ 

Score function method:

$$abla_{oldsymbol{\lambda}} \mathbb{E}_{q(\mathsf{z};oldsymbol{\lambda})} \left[ f(\mathsf{z}) 
ight] = \mathbb{E}_{q(\mathsf{z};oldsymbol{\lambda})} \left[ f(\mathsf{z}) 
abla_{oldsymbol{\lambda}} \log q(\mathsf{z};oldsymbol{\lambda}) 
ight]$$

Stochastic optimization

<sup>&</sup>lt;sup>1</sup>Paisley et al. (2012), Ranganath et al. (2014), Mnih and Gregor (2014)

# Score Function Method<sup>1</sup>

Monte Carlo estimator of the gradient

 $\nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q(\mathbf{z};\boldsymbol{\lambda})}[f(\mathbf{z})]$ 

Score function method:

$$abla_{oldsymbol{\lambda}} \mathbb{E}_{q(\mathsf{z};oldsymbol{\lambda})} \left[ f(\mathsf{z}) 
ight] = \mathbb{E}_{q(\mathsf{z};oldsymbol{\lambda})} \left[ f(\mathsf{z}) 
abla_{oldsymbol{\lambda}} \log q(\mathsf{z};oldsymbol{\lambda}) 
ight]$$

- Stochastic optimization
- Algorithm:
  - 1. Sample  $\mathbf{z}^{(s)} \stackrel{\text{iid}}{\sim} q(\mathbf{z}; \boldsymbol{\lambda})$
  - 2. Evaluate  $f(\mathbf{z}^{(s)})$  and  $\nabla_{\boldsymbol{\lambda}} \log q(\mathbf{z}^{(s)}; \boldsymbol{\lambda})$  for each sample s
  - 3. Obtain a Monte Carlo estimate of the gradient
  - 4. Take a gradient step for  $\lambda$

<sup>1</sup>Paisley et al. (2012), Ranganath et al. (2014), Mnih and Gregor (2014)

# Reparameterization Trick<sup>2</sup> (1/2)

Define an invertible transformation

$$\mathbf{z} = \mathcal{T}(\boldsymbol{\epsilon}; \boldsymbol{\lambda}), \qquad \boldsymbol{\epsilon} = \mathcal{T}^{-1}(\mathbf{z}; \boldsymbol{\lambda})$$

such that  $\pi(\epsilon)$  does not depend on  $\lambda$ .

<sup>&</sup>lt;sup>2</sup>Salimans and Knowles (2013), Kingma and Welling (2014), Rezende et al. (2014), Titsias & Lázaro-Gredilla (2014)

# Reparameterization Trick<sup>2</sup> (1/2)

Define an invertible transformation

$$\mathbf{z} = \mathcal{T}(\boldsymbol{\epsilon}; \boldsymbol{\lambda}), \qquad \boldsymbol{\epsilon} = \mathcal{T}^{-1}(\mathbf{z}; \boldsymbol{\lambda})$$

such that  $\pi(\epsilon)$  does not depend on  $\lambda$ .

Push the gradient inside the expectation

$$\nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q(\boldsymbol{z};\boldsymbol{\lambda})} \left[ f(\boldsymbol{z}) \right] = \mathbb{E}_{\pi(\boldsymbol{\epsilon})} \left[ \nabla_{\boldsymbol{z}} f(\boldsymbol{z}) \Big|_{\boldsymbol{z} = \mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda})} \nabla_{\boldsymbol{\lambda}} \mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda}) \right]$$

<sup>&</sup>lt;sup>2</sup>Salimans and Knowles (2013), Kingma and Welling (2014), Rezende et al. (2014), Titsias & Lázaro-Gredilla (2014)

# Reparameterization Trick<sup>2</sup> (1/2)

Define an invertible transformation

$$\mathbf{z} = \mathcal{T}(\boldsymbol{\epsilon}; \boldsymbol{\lambda}), \qquad \boldsymbol{\epsilon} = \mathcal{T}^{-1}(\mathbf{z}; \boldsymbol{\lambda})$$

such that  $\pi(\epsilon)$  does not depend on  $\lambda$ .

Push the gradient inside the expectation

$$\nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q(\boldsymbol{z};\boldsymbol{\lambda})} \left[ f(\boldsymbol{z}) \right] = \mathbb{E}_{\pi(\boldsymbol{\epsilon})} \left[ \nabla_{\boldsymbol{z}} f(\boldsymbol{z}) \Big|_{\boldsymbol{z} = \mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda})} \nabla_{\boldsymbol{\lambda}} \mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda}) \right]$$

Algorithm:

- 1. Sample  $\epsilon^{(s)}$  iid from  $\pi(\epsilon)$
- 2. Obtain  $\mathbf{z}^{(s)} = \mathcal{T}(\boldsymbol{\epsilon}^{(s)}; \boldsymbol{\lambda})$
- 3. Evaluate  $\nabla_{\mathbf{z}} f(\mathbf{z})$  and  $\nabla_{\boldsymbol{\lambda}} \mathcal{T}(\boldsymbol{\epsilon}; \boldsymbol{\lambda})$  for each sample s
- 4. Obtain a Monte Carlo estimate of the gradient
- 5. Take a gradient step for  $\lambda$

<sup>&</sup>lt;sup>2</sup>Salimans and Knowles (2013), Kingma and Welling (2014), Rezende et al. (2014), Titsias & Lázaro-Gredilla (2014)

# Reparameterization Trick (2/2)

• Simple example:  $q(\mathbf{z}; \boldsymbol{\lambda})$  is a Gaussian

$$oldsymbol{\epsilon} = \mathcal{T}^{-1}(\mathsf{z};oldsymbol{\lambda}) = \mathbf{\Sigma}^{-1/2}(\mathsf{z}-oldsymbol{\mu})$$

The transformed density is  $\pi(\epsilon) = \mathcal{N}(\mathbf{0}, \mathbf{1})$  (independent of  $\mu$ ,  $\Sigma$ )

<sup>&</sup>lt;sup>3</sup>Kucukelbir et al. (2015, 2016)

# Reparameterization Trick (2/2)

• Simple example:  $q(\mathbf{z}; \boldsymbol{\lambda})$  is a Gaussian

$$\epsilon = \mathcal{T}^{-1}(\mathsf{z}; oldsymbol{\lambda}) = \mathbf{\Sigma}^{-1/2}(\mathsf{z}-oldsymbol{\mu})$$

The transformed density is  $\pi(\epsilon) = \mathcal{N}(\mathbf{0}, \mathbf{1})$  (independent of  $\mu$ ,  $\Sigma$ )  $\blacktriangleright$  Used in ADVI<sup>3</sup>

- Non-linear transformation
- Gaussian distribution on the transformed space



<sup>3</sup>Kucukelbir et al. (2015, 2016)

#### Comparison

Score function

- + Any probabilistic model
- + Any variational distribution (as long as we can sample)
- High variance: Requires additional tricks and many samples
- Reparameterization
  - Differentiable probabilistic models (continuous z)<sup>4</sup>
  - Limited to some distributions (location-scale, inverse CDF)
  - + Low variance: Only 1 sample suffices in practice

 $<sup>^{4}\</sup>mbox{But}$  see Maddison et al. (2016), Jang et al. (2016), Kusner & Hernández-Lobato (2016)

#### Comparison

Score function

- + Any probabilistic model
- + Any variational distribution (as long as we can sample)
- High variance: Requires additional tricks and many samples
- Reparameterization
  - Differentiable probabilistic models (continuous z)<sup>4</sup>
  - Limited to some distributions (location-scale, inverse CDF)
  - + Low variance: Only 1 sample suffices in practice
- Our contribution
  - Extend reparameterization to other distributions

<sup>&</sup>lt;sup>4</sup>But see Maddison et al. (2016), Jang et al. (2016), Kusner & Hernández-Lobato (2016)

#### Motivation

• We may need more expressive variational families  $q(\mathbf{z}; \boldsymbol{\lambda})$ 



#### Contributions

#### The Generalized Reparameterization Gradient

Francisco J. R. Ruiz University of Cambridge Columbia University Michalis K. Titsias Athens University of Economics and Business David M. Blei Columbia University

#### **Rejection Sampling Variational Inference**

Christian A. Naesseth<sup>\*†‡</sup> Francisco J. R. Ruiz<sup>‡§</sup> Scott W. Linderman<sup>‡</sup> David M. Blei<sup>‡</sup> <sup>†</sup>Linköping University <sup>‡</sup>Columbia University <sup>§</sup>University of Cambridge

## The Generalized Reparameterization Gradient (1/3)

Define an invertible transformation

$$\mathbf{z} = \mathcal{T}(oldsymbol{\epsilon};oldsymbol{\lambda}), \qquad oldsymbol{\epsilon} = \mathcal{T}^{-1}(\mathbf{z};oldsymbol{\lambda})$$

but allow  $\pi(\epsilon; \lambda)$  to depend weakly on  $\lambda$ .

Gradient:

$$abla_{oldsymbol{\lambda}} \mathbb{E}_{q(\mathsf{z};oldsymbol{\lambda})} \left[ f(\mathsf{z}) 
ight] = \mathbf{g}^{ ext{rep}} + \mathbf{g}^{ ext{corr}}$$

### The Generalized Reparameterization Gradient (2/3)

$$abla_{\boldsymbol{\lambda}} \mathbb{E}_{q(\boldsymbol{\mathsf{z}}; \boldsymbol{\lambda})} \left[ f(\boldsymbol{\mathsf{z}}) \right] = \boldsymbol{\mathsf{g}}^{\mathrm{rep}} + \boldsymbol{\mathsf{g}}^{\mathrm{corr}}$$

 $\blacktriangleright~g^{\rm rep}$ : Reparameterization gradient

$$\mathbf{g}^{\mathrm{rep}} = \mathbb{E}_{\pi(\boldsymbol{\epsilon};\boldsymbol{\lambda})} \left[ \nabla_{\mathbf{z}} f(\mathbf{z}) \big|_{\mathbf{z} = \mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda})} \nabla_{\boldsymbol{\lambda}} \mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda}) \right]$$

▶ **g**<sup>corr</sup>: Correction term

$$\mathbf{g}^{ ext{corr}} = \mathbb{E}_{\pi(oldsymbol{\epsilon};oldsymbol{\lambda})} \left[ f(\mathcal{T}(oldsymbol{\epsilon};oldsymbol{\lambda})) 
abla_{oldsymbol{\lambda}} \log \pi(oldsymbol{\epsilon};oldsymbol{\lambda}) 
ight]$$

#### The Generalized Reparameterization Gradient (3/3)

$$\begin{aligned} \nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q(\boldsymbol{z};\boldsymbol{\lambda})} \left[ f(\boldsymbol{z}) \right] &= \boldsymbol{g}^{\text{rep}} + \boldsymbol{g}^{\text{corr}} \\ \boldsymbol{g}^{\text{rep}} &= \mathbb{E}_{\pi(\boldsymbol{\epsilon};\boldsymbol{\lambda})} \left[ \nabla_{\boldsymbol{z}} f(\boldsymbol{z}) \big|_{\boldsymbol{z} = \mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda})} \nabla_{\boldsymbol{\lambda}} \mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda}) \right] \\ \boldsymbol{g}^{\text{corr}} &= \mathbb{E}_{\pi(\boldsymbol{\epsilon};\boldsymbol{\lambda})} \left[ f(\mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda})) \nabla_{\boldsymbol{\lambda}} \log \pi(\boldsymbol{\epsilon};\boldsymbol{\lambda}) \right] \end{aligned}$$

• Under a transformation such that  $\pi(\epsilon; \lambda)$  does not depend on  $\lambda$ :

$$\label{eq:green} \begin{split} & \mathbf{g}^{\mathrm{rep}} = \mathrm{reparameterization \ gradient} \\ & \mathbf{g}^{\mathrm{corr}} = \mathbf{0} \end{split}$$

Under identity transformation:

$$\label{eq:grep} \begin{split} &g^{\rm rep} = 0 \\ &g^{\rm corr} = {\rm score\ function\ gradient} \end{split}$$

#### The Generalized Reparameterization Gradient (3/3)

$$\begin{aligned} \nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q(\boldsymbol{z};\boldsymbol{\lambda})} \left[ f(\boldsymbol{z}) \right] &= \boldsymbol{g}^{\text{rep}} + \boldsymbol{g}^{\text{corr}} \\ \boldsymbol{g}^{\text{rep}} &= \mathbb{E}_{\pi(\boldsymbol{\epsilon};\boldsymbol{\lambda})} \left[ \nabla_{\boldsymbol{z}} f(\boldsymbol{z}) \big|_{\boldsymbol{z} = \mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda})} \nabla_{\boldsymbol{\lambda}} \mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda}) \right] \\ \boldsymbol{g}^{\text{corr}} &= \mathbb{E}_{\pi(\boldsymbol{\epsilon};\boldsymbol{\lambda})} \left[ f(\mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda})) \nabla_{\boldsymbol{\lambda}} \log \pi(\boldsymbol{\epsilon};\boldsymbol{\lambda}) \right] \end{aligned}$$

• Under a transformation such that  $\pi(\epsilon; \lambda)$  does not depend on  $\lambda$ :

$$\label{eq:green} \begin{split} \boldsymbol{g}^{\mathrm{rep}} &= \mathrm{reparameterization\ gradient} \\ \boldsymbol{g}^{\mathrm{corr}} &= \boldsymbol{0} \end{split}$$

Under identity transformation:

$$\label{eq:grep} \begin{split} &g^{\rm rep} = 0 \\ &g^{\rm corr} = {\rm score\ function\ gradient} \end{split}$$

Goal: Find a transformation that makes  $\mathbf{g}^{\mathrm{corr}}$  small

#### Example: Gamma distribution

▶ We use transformations based on *standardization* 

• Example: 
$$q(z; \alpha, \beta) = \text{Gamma}(z; \alpha, \beta)$$

$$\epsilon = \mathcal{T}^{-1}(z; \alpha, \beta) = \frac{\log(z) - (\psi(\alpha) - \log(\beta))}{\sqrt{\psi_1(\alpha)}},$$

#### Example: Gamma distribution

▶ We use transformations based on *standardization* 

• Example:  $q(z; \alpha, \beta) = \text{Gamma}(z; \alpha, \beta)$ 

$$\epsilon = \mathcal{T}^{-1}(z; \alpha, \beta) = \frac{\log(z) - (\psi(\alpha) - \log(\beta))}{\sqrt{\psi_1(\alpha)}},$$



### G-REP: Full Algorithm

- 1. Draw a single sample  $\mathbf{z} \sim q(\mathbf{z}; \boldsymbol{\lambda})$
- 2. Obtain  $\epsilon = \mathcal{T}^{-1}(\mathsf{z}; \lambda)$
- 3. Estimate  $\mathbf{g}^{\mathrm{rep}}$  and  $\mathbf{g}^{\mathrm{corr}}$  (with 1 sample)
- 4. Take a gradient step for  $\lambda$

#### Results: MNIST

#### Model: Gamma-beta matrix factorization



#### Results: Olivetti Dataset

Model: Sparse gamma deep exponential family<sup>5</sup>



<sup>5</sup>Ranganath et al. (2015)

## Rejection Sampling Variational Inference

Every random variable that we can simulate on our computers is ultimately a transformation of elementary random variables

> In theory, this should allow for reparameterization of any distribution

# Rejection Sampling Variational Inference

Every random variable that we can simulate on our computers is ultimately a transformation of elementary random variables

- In theory, this should allow for reparameterization of any distribution
- Challenge: rejection sampling steps
  - We cannot push the gradient inside the integral

#### Reparameterized Rejection Sampling

In standard rejection sampling, we have:

- A target,  $q(z; \lambda)$
- A proposal, r(z; λ)
- A uniform random variable,  $u \sim \mathcal{U}(0,1)$
- An accept/reject step,  $u < \frac{q(z;\lambda)}{M_{\lambda}r(z;\lambda)}$

#### Reparameterized Rejection Sampling

In standard rejection sampling, we have:

- A target,  $q(z; \lambda)$
- A proposal, r(z; λ)
- A uniform random variable,  $u \sim \mathcal{U}(0,1)$
- An accept/reject step,  $u < \frac{q(z;\lambda)}{M_{\lambda}r(z;\lambda)}$
- ▶ In *reparameterized* rejection sampling, we have:
  - A target,  $q(z; \lambda)$
  - An elementary proposal,  $s(\varepsilon)$
  - A transformation,  $z = \mathcal{T}(\varepsilon; \lambda)$ , such that  $z \sim r(z; \lambda)$
  - A uniform random variable,  $u \sim \mathcal{U}(0, 1)$
  - An accept/reject step,  $u < \frac{q(\mathcal{T}(\varepsilon; \lambda); \lambda)}{M_{\lambda} r(\mathcal{T}(\varepsilon; \lambda); \lambda)}$

The accepted sample  $z = \mathcal{T}(arepsilon; oldsymbol{\lambda}) \sim q(z; oldsymbol{\lambda})$ 

# Rejection Sampling VI

Main idea:

- Integrate out the accept/reject variables u
- Consider the distribution of the *accepted* sample  $\varepsilon$

$$\pi(\varepsilon; \boldsymbol{\lambda}) = s(\varepsilon) \frac{q(\mathcal{T}(\varepsilon; \boldsymbol{\lambda}); \boldsymbol{\lambda})}{r(\mathcal{T}(\varepsilon; \boldsymbol{\lambda}); \boldsymbol{\lambda})}$$

• Use the transformation  $z = \mathcal{T}(\varepsilon; \lambda)$ 

# Example: Gamma Distribution (1/2)

Algorithm<sup>6</sup> to sample from a  $Gamma(\alpha, 1)$  (with  $\alpha \geq 1$ )

1. Generate  $\varepsilon \sim s(\varepsilon) = \mathcal{N}(0, 1)$ 

2. Transform as 
$$z = \mathcal{T}(\varepsilon; \alpha) = \left(\alpha - \frac{1}{3}\right) \left(1 + \frac{\varepsilon}{\sqrt{9\alpha - 3}}\right)^3$$

- 3. Generate  $u \sim \mathcal{U}(0, 1)$
- 4. Accept or reject  $\varepsilon$

4.1 If accept, return z and  $\varepsilon$ 

 $4.2\,$  If reject, go to Step 1 and repeat

<sup>6</sup>Marsaglia & Tsang (2000)

## Example: Gamma Distribution (2/2)



## Example: Gamma Distribution (2/2)



We can leverage 60+ years of research in rejection sampling to find good transformations!

#### Gradient of Rejection Sampling VI

The form of the gradient is similar to G-REP:

$$\begin{split} \nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q(\boldsymbol{z};\boldsymbol{\lambda})}\left[f(\boldsymbol{z})\right] &= \boldsymbol{g}^{\text{rep}} + \boldsymbol{g}^{\text{corr}} \\ \boldsymbol{g}^{\text{rep}} &= \mathbb{E}_{\pi(\boldsymbol{\epsilon};\boldsymbol{\lambda})}\left[\nabla_{\boldsymbol{z}} f(\boldsymbol{z})\big|_{\boldsymbol{z}=\mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda})} \nabla_{\boldsymbol{\lambda}} \mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda})\right] \\ \boldsymbol{g}^{\text{corr}} &= \mathbb{E}_{\pi(\boldsymbol{\epsilon};\boldsymbol{\lambda})}\left[f(\mathcal{T}(\boldsymbol{\epsilon};\boldsymbol{\lambda})) \nabla_{\boldsymbol{\lambda}} \log \pi(\boldsymbol{\epsilon};\boldsymbol{\lambda})\right] \end{split}$$

#### **RSVI: Full Algorithm**

- 1. Run the reparameterized rejection sampling to draw  $m{\epsilon} \sim \pi(m{\epsilon};m{\lambda})$
- 2. Transform  $\mathbf{z} = \mathcal{T}(\boldsymbol{\epsilon}; \boldsymbol{\lambda})$
- 3. Estimate  $\mathbf{g}^{\mathrm{rep}}$  and  $\mathbf{g}^{\mathrm{corr}}$  (with 1 sample)
- 4. Take a gradient step for  $\lambda$

#### Results: Variance of the Gradient



Model: Dirichlet/Multinomial with 100 components

(B denotes "shape augmentation")

#### Results: Olivetti Dataset



Model: Sparse gamma deep exponential family<sup>7</sup>

<sup>7</sup>Ranganath et al. (2015)

#### Summary: G-REP and RSVI

- Extend reparameterization trick to non-reparameterizable distributions (gamma, beta, Dirichlet, ...)
- Allow variational inference on continuous non-conjugate models
- ▶ Fast: Monte Carlo estimation with only 1 sample

# Thank you for your attention!

