

Reparameterization Gradient for Non-differentiable Models

Wonyeol Lee Hangeol Yu Hongseok Yang
KAIST

Published at NeurIPS 2018

Reparameterization Gradient for Non-differentiable Models

Wonyeol Lee Hangeol Yu Hongseok Yang
KAIST

Published at NeurIPS 2018

Reparameterization Gradient for Non-differentiable Models

Wonyeol Lee Hangeol Yu Hongseok Yang
KAIST

Published at NeurIPS 2018

Reparameterization Gradient for Non-differentiable Models

Wonyeol Lee Hangeol Yu Hongseok Yang
KAIST

Published at NeurIPS 2018

Backgrounds

Posterior inference

- Latent variable $z \in \mathbb{R}^n$.
- Observed variable $x \in \mathbb{R}^m$.
- Joint density $p(x,z)$.
- Want to infer posterior $p(z|x^0)$ given a particular value x^0 of x .

Variational inference

1. Fix a family of variational distr. $\{q_{\theta}(z)\}_{\theta}$.
2. Find $q_{\theta}(z)$ that approximates $p(z|x^0)$ well.

Variational inference

differentiable & easy-to-sample



1. Fix a family of variational distr. $\{q_{\theta}(z)\}_{\theta}$.
2. Find $q_{\theta}(z)$ that approximates $p(z|x^0)$ well.

Variational inference

differentiable & easy-to-sample

1. Fix a family of variational distr. $\{q_{\theta}(z)\}_{\theta}$.
2. Find $q_{\theta}(z)$ that approximates $p(z|x^0)$ well.

Typically, by solving

$$\operatorname{argmax}_{\theta}(\operatorname{ELBO}_{\theta})$$

where $\operatorname{ELBO}_{\theta} = \mathbb{E}_{q_{\theta}(z)}[\log(p(x^0, z)/q_{\theta}(z))]$.

Variational inference

differentiable & easy-to-sample

1. Fix a family of variational distr. $\{q_{\theta}(z)\}_{\theta}$.
2. Find $q_{\theta}(z)$ that approximates $p(z|x^0)$ well.

Typically, by solving

$$\operatorname{argmax}_{\theta}(\operatorname{ELBO}_{\theta})$$

where $\operatorname{ELBO}_{\theta} = \mathbb{E}_{q_{\theta}(z)}[\dots z \dots z \dots]$.

Gradient ascent

$$\theta_{n+1} = \theta_n + \eta \times \nabla_{\theta} \text{ELBO}_{\theta=\theta_n}$$

Gradient ascent

$$\theta_{n+1} = \theta_n + \eta \times \nabla_{\theta} \text{ELBO}_{\theta=\theta_n}$$

- Difficult to compute $\nabla_{\theta} \text{ELBO}_{\theta}$.

Gradient ascent

$$\theta_{n+1} = \theta_n + \eta \times \overline{\nabla_{\theta} \text{ELBO}_{\theta=\theta_n}}$$

- Difficult to compute $\nabla_{\theta} \text{ELBO}_{\theta}$.
- Use an estimated gradient instead.

Reparameterization estimator

- Works if $p(x^0, z)$ is differentiable wrt. z .
- Need distr. $q(\varepsilon)$ & smooth function $f_\theta(\varepsilon)$ s.t.

$f_\theta(\varepsilon)$ for $\varepsilon \sim q(\varepsilon)$ has the distr. $q_\theta(z)$.

- Derived from the equation:

$$\nabla_\theta \text{ELBO}_\theta = \mathbb{E}_{q(\varepsilon)} [\nabla_\theta (.. f_\theta(\varepsilon) .. f_\theta(\varepsilon) ..)]$$

Reparameterization estimator

- Works if $p(x^0, z)$ is differentiable wrt. z .
- Need distr. $q(\varepsilon)$ & smooth function $f_\theta(\varepsilon)$ s.t.
 $f_\theta(\varepsilon)$ for $\varepsilon \sim q(\varepsilon)$ has the distr. $q_\theta(z)$.
- Derived from the equation:

$$\nabla_\theta \text{ELBO}_\theta = \mathbb{E}_{q(\varepsilon)} [\nabla_\theta (.. f_\theta(\varepsilon) .. f_\theta(\varepsilon) ..)]$$

$$\nabla_{\theta} \text{ELBO}_{\theta} =$$

imator

- Works if $p(x^0, z)$ is differentiable wrt. z .
- Need distr. $q(\varepsilon)$ & smooth function $f_{\theta}(\varepsilon)$ s.t.
 $f_{\theta}(\varepsilon)$ for $\varepsilon \sim q(\varepsilon)$ has the distr. $q_{\theta}(z)$.
- Derived from the equation:

$$\nabla_{\theta} \text{ELBO}_{\theta} = \mathbb{E}_{q(\varepsilon)} [\nabla_{\theta} (.. f_{\theta}(\varepsilon) .. f_{\theta}(\varepsilon) ..)]$$

$$\nabla_{\theta} \text{ELBO}_{\theta} = \nabla_{\theta} \mathbb{E}_{q_{\theta}(z)}[\dots z \dots z \dots]$$

imator

- Works if $p(x^0, z)$ is differentiable wrt. z .
- Need distr. $q(\varepsilon)$ & smooth function $f_{\theta}(\varepsilon)$ s.t.
 $f_{\theta}(\varepsilon)$ for $\varepsilon \sim q(\varepsilon)$ has the distr. $q_{\theta}(z)$.
- Derived from the equation:

$$\nabla_{\theta} \text{ELBO}_{\theta} = \mathbb{E}_{q(\varepsilon)}[\nabla_{\theta}(\dots f_{\theta}(\varepsilon) \dots f_{\theta}(\varepsilon) \dots)]$$

$$\begin{aligned}\nabla_{\theta} \text{ELBO}_{\theta} &= \nabla_{\theta} \mathbb{E}_{q_{\theta}(z)}[\dots z \dots z \dots] \\ &= \nabla_{\theta} \mathbb{E}_{q(\varepsilon)}[\dots f_{\theta}(\varepsilon) \dots f_{\theta}(\varepsilon) \dots]\end{aligned}$$

imator

- Works if $p(x^0, z)$ is differentiable wrt. z .
- Need distr. $q(\varepsilon)$ & smooth function $f_{\theta}(\varepsilon)$ s.t.

$f_{\theta}(\varepsilon)$ for $\varepsilon \sim q(\varepsilon)$ has the distr. $q_{\theta}(z)$.

- Derived from the equation:

$$\nabla_{\theta} \text{ELBO}_{\theta} = \mathbb{E}_{q(\varepsilon)}[\nabla_{\theta}(\dots f_{\theta}(\varepsilon) \dots f_{\theta}(\varepsilon) \dots)]$$

$$\begin{aligned}
\nabla_{\theta} \text{ELBO}_{\theta} &= \nabla_{\theta} \mathbb{E}_{q_{\theta}(z)}[\dots z \dots z \dots] \\
&= \nabla_{\theta} \mathbb{E}_{q(\varepsilon)}[\dots f_{\theta}(\varepsilon) \dots f_{\theta}(\varepsilon) \dots] \\
&= \mathbb{E}_{q(\varepsilon)}[\nabla_{\theta}(\dots f_{\theta}(\varepsilon) \dots f_{\theta}(\varepsilon) \dots)]
\end{aligned}$$

imator

- Works if $p(x^0, z)$ is differentiable wrt. z .
- Need distr. $q(\varepsilon)$ & smooth function $f_{\theta}(\varepsilon)$ s.t.
 $f_{\theta}(\varepsilon)$ for $\varepsilon \sim q(\varepsilon)$ has the distr. $q_{\theta}(z)$.
- Derived from the equation:

$$\nabla_{\theta} \text{ELBO}_{\theta} = \mathbb{E}_{q(\varepsilon)}[\nabla_{\theta}(\dots f_{\theta}(\varepsilon) \dots f_{\theta}(\varepsilon) \dots)]$$

Non-differentiable models from probabilistic programming

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```



```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

$$p(z,x=0) = [z>0]r_1(z) + [z\leq 0]r_2(z)$$

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

$$p(z,x=0) = [z>0]r_1(z) + [z\leq 0]r_2(z)$$

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

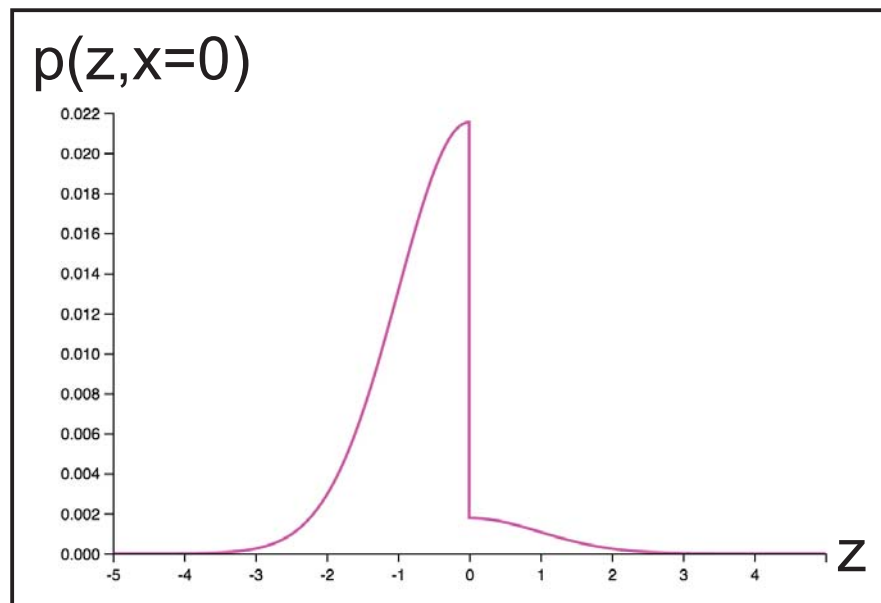
$$p(z,x=0) = [z>0]r_1(z) + [z\leq 0]r_2(z)$$

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

$$p(z,x=0) = [z>0]r_1(z) + [z\leq 0]r_2(z)$$



```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

≈

```
(let
  [ε (sample (normal 0 1))
   z (+ ε θ)]
  z)
```

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

$$p(z,x=0) = [z>0]r_1(z) + [z\leq 0]r_2(z)$$

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

≈

```
(let
  [ε (sample (normal 0 1))
   z (+ ε θ)]
  z)
```

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

$$p(z,x=0) = [z>0]r_1(z) + [z\leq 0]r_2(z)$$

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

≈

```
(let
  [ε (sample (normal 0 1))
   z (+ ε θ)]
  z)
```

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

$$p(z,x=0) = [z>0]r_1(z) + [z\leq 0]r_2(z)$$


```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```



```
(let
  [ε (sample (normal 0 1))
   z (+ ε θ)]
  z)
```

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

$$p(z,x=0) = [z>0]r_1(z) + [z\leq 0]r_2(z)$$

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

≈

```
(let
  [ε (sample (normal 0 1))
   z (+ ε θ)]
  z)
```

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

$$p(z,x=0) = [z>0]r_1(z) + [z\leq 0]r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0,1)$$

$$z = \varepsilon + \theta$$

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

≈

```
(let
  [ε (sample (normal 0 1))
   z (+ ε θ)]
  z)
```

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

$$p(z,x=0) = [z>0]r_1(z) + [z\leq 0]r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0,1)$$

$$z = \varepsilon + \theta$$

How to find a good θ ?

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

≈

```
(let
  [ε (sample (normal 0 1))
   z (+ ε θ)]
  z)
```

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

$$p(z,x=0) = [z>0]r_1(z) + [z\leq 0]r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0,1)$$

$$z = \varepsilon + \theta$$

How to find a good θ ?

By gradient ascent on $ELBO_{\theta}$.

$$\theta_{n+1} \leftarrow \theta_n + \eta \times \nabla_{\theta} ELBO_{\theta=\theta_n}$$

```
(let
  [z (sample (normal 0 1))]
  (if (> z 0)
      (observe (normal 3 1) 0)
      (observe (normal -2 1) 0))
  z)
```

≈

```
(let
  [ε (sample (normal 0 1))
   z (+ ε θ)]
  z)
```

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

$$p(z,x=0) = [z>0]r_1(z) + [z\leq 0]r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0,1)$$

$$z = \varepsilon + \theta$$

How to find a good θ ?

By gradient ascent on $ELBO_{\theta}$.

$$\theta_{n+1} \leftarrow \theta_n + \eta \times \nabla_{\theta} ELBO_{\theta=\theta_n}$$

$$\nabla_{\theta} \text{ELBO}_{\theta}$$

$$r_1(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1)\mathcal{N}(x=0|-2,1)$$

$$p(z, x=0) = [z > 0]r_1(z) + [z \leq 0]r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0,1)$$

$$z = \varepsilon + \theta$$

How to find a good θ ?

By gradient ascent on ELBO_{θ} .

$$\theta_{n+1} \leftarrow \theta_n + \eta \times \nabla_{\theta} \text{ELBO}_{\theta=\theta_n}$$

$$\nabla_{\theta} \text{ELBO}_{\theta}$$

$$= \nabla_{\theta} \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \log(r_2(\varepsilon + \theta))]$$

$$r_1(z) = \mathcal{N}(z|0, 1) \mathcal{N}(x=0|3, 1)$$

$$r_2(z) = \mathcal{N}(z|0, 1) \mathcal{N}(x=0|-2, 1)$$

$$p(z, x=0) = [z > 0] r_1(z) + [z \leq 0] r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0, 1)$$

$$z = \varepsilon + \theta$$

How to find a good θ ?

By gradient ascent on ELBO_{θ} .

$$\theta_{n+1} \leftarrow \theta_n + \eta \times \nabla_{\theta} \text{ELBO}_{\theta=\theta_n}$$

$$\nabla_{\theta} \text{ELBO}_{\theta}$$

$$= \nabla_{\theta} \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \log(r_2(\varepsilon + \theta))]$$

$$r_1(z) = \mathcal{N}(z|0, 1) \mathcal{N}(x=0|3, 1)$$

$$r_2(z) = \mathcal{N}(z|0, 1) \mathcal{N}(x=0|-2, 1)$$

$$p(z, x=0) = [z > 0] r_1(z) + [z \leq 0] r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0, 1)$$

$$z = \varepsilon + \theta$$

How to find a good θ ?

By gradient ascent on ELBO_{θ} .

$$\theta_{n+1} \leftarrow \theta_n + \eta \times \nabla_{\theta} \text{ELBO}_{\theta=\theta_n}$$

$$\nabla_{\theta} \text{ELBO}_{\theta}$$

$$= \nabla_{\theta} \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \log(r_2(\varepsilon + \theta))]$$

$$r_1(z) = \mathcal{N}(z|0, 1) \mathcal{N}(x=0|3, 1)$$

$$r_2(z) = \mathcal{N}(z|0, 1) \mathcal{N}(x=0|-2, 1)$$

$$p(z, x=0) = [z > 0] r_1(z) + [z \leq 0] r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0, 1)$$

$$z = \varepsilon + \theta$$

How to find a good θ ?

By gradient ascent on ELBO_{θ} .

$$\theta_{n+1} \leftarrow \theta_n + \eta \times \nabla_{\theta} \text{ELBO}_{\theta=\theta_n}$$

$$\nabla_{\theta} \text{ELBO}_{\theta}$$

$$= \nabla_{\theta} \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \log(r_2(\varepsilon + \theta))]$$

$$= \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \nabla_{\theta} \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \nabla_{\theta} \log(r_2(\varepsilon + \theta))]$$

$$r_1(z) = \mathcal{N}(z|0, 1) \mathcal{N}(x=0|3, 1)$$

$$r_2(z) = \mathcal{N}(z|0, 1) \mathcal{N}(x=0|-2, 1)$$

$$p(z, x=0) = [z > 0] r_1(z) + [z \leq 0] r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0, 1)$$

$$z = \varepsilon + \theta$$

How to find a good θ ?

By gradient ascent on ELBO_{θ} .

$$\theta_{n+1} \leftarrow \theta_n + \eta \times \nabla_{\theta} \text{ELBO}_{\theta=\theta_n}$$

$$\nabla_{\theta} \text{ELBO}_{\theta}$$

$$= \nabla_{\theta} \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \log(r_2(\varepsilon + \theta))]$$

$$= \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \nabla_{\theta} \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \nabla_{\theta} \log(r_2(\varepsilon + \theta))]$$

$$= \mathbb{E}_{q(\varepsilon)} [-\theta - \varepsilon] = -\theta$$

$$r_1(z) = \mathcal{N}(z|0,1) \mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1) \mathcal{N}(x=0|-2,1)$$

$$p(z, x=0) = [z > 0] r_1(z) + [z \leq 0] r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0,1)$$

$$z = \varepsilon + \theta$$

How to find a good θ ?

By gradient ascent on ELBO_{θ} .

$$\theta_{n+1} \leftarrow \theta_n + \eta \times \nabla_{\theta} \text{ELBO}_{\theta=\theta_n}$$

$$\nabla_{\theta} \text{ELBO}_{\theta}$$

$$= \nabla_{\theta} \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \log(r_2(\varepsilon + \theta))]$$

$$= \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \nabla_{\theta} \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \nabla_{\theta} \log(r_2(\varepsilon + \theta))]$$

$$= \mathbb{E}_{q(\varepsilon)} [-\theta - \varepsilon] = -\theta$$

$$r_1(z) = \mathcal{N}(z|0,1) \mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1) \mathcal{N}(x=0|-2,1)$$

$$p(z, x=0) = [z > 0] r_1(z) + [z \leq 0] r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0,1)$$

$$z = \varepsilon + \theta$$

How to find a good θ ?

By gradient ascent on ELBO_{θ} .

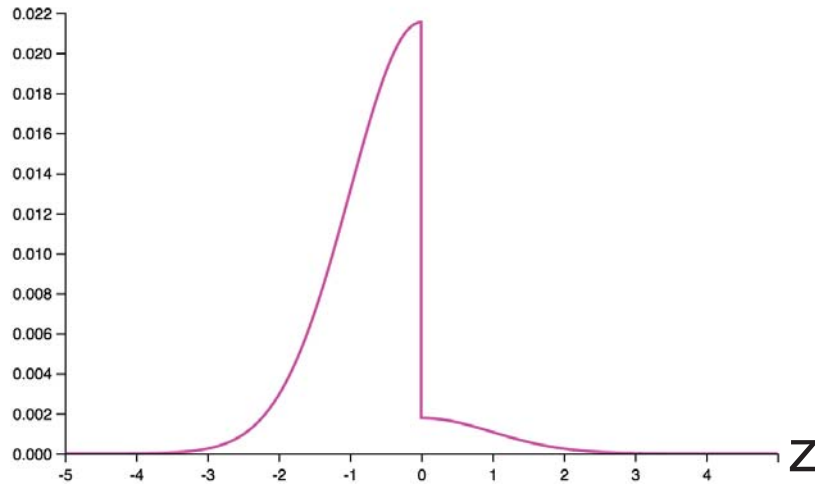
$$\theta_{n+1} \leftarrow \theta_n + \eta \times \nabla_{\theta} \text{ELBO}_{\theta=\theta_n}$$

$$\nabla_{\theta} \text{ELBO}_{\theta}$$

$$= \nabla_{\theta} \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \log(r_2(\varepsilon + \theta))]$$

$$= \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \nabla_{\theta} \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \nabla_{\theta} \log(r_2(\varepsilon + \theta))]$$

$$= \mathbb{E}_{q(\varepsilon)} [-\theta - \varepsilon] = -\theta$$



$$q(\varepsilon) = \mathcal{N}(\varepsilon | 0, 1)$$

$$z = \varepsilon + \theta$$

BO_{θ} .

$$\theta_{n+1} \leftarrow \theta_n + \eta \times \nabla_{\theta} \text{ELBO}_{\theta=\theta_n}$$

$$\nabla_{\theta} \text{ELBO}_{\theta}$$

$$= \nabla_{\theta} \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \log(r_2(\varepsilon + \theta))]$$

$$\neq \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \nabla_{\theta} \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \nabla_{\theta} \log(r_2(\varepsilon + \theta))]$$

$$= \mathbb{E}_{q(\varepsilon)} [-\theta - \varepsilon] = -\theta$$

$$r_1(z) = \mathcal{N}(z|0,1) \mathcal{N}(x=0|3,1)$$

$$r_2(z) = \mathcal{N}(z|0,1) \mathcal{N}(x=0|-2,1)$$

$$p(z, x=0) = [z > 0] r_1(z) + [z \leq 0] r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0,1)$$

$$z = \varepsilon + \theta$$

How to find a good θ ?

By gradient ascent on ELBO_{θ} .

$$\theta_{n+1} \leftarrow \theta_n + \eta \times \nabla_{\theta} \text{ELBO}_{\theta=\theta_n}$$

$$\nabla_{\theta} \text{ELBO}_{\theta}$$

$$= \nabla_{\theta} \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \log(r_2(\varepsilon + \theta))]$$

$$= \mathbb{E}_{q(\varepsilon)} [[\varepsilon > -\theta] \nabla_{\theta} \log(r_1(\varepsilon + \theta)) + [\varepsilon \leq -\theta] \nabla_{\theta} \log(r_2(\varepsilon + \theta))]$$

+ Correction Term

$$r_1(z) = \mathcal{N}(z|0, 1) \mathcal{N}(x=0|3, 1)$$

$$r_2(z) = \mathcal{N}(z|0, 1) \mathcal{N}(x=0|-2, 1)$$

$$p(z, x=0) = [z > 0] r_1(z) + [z \leq 0] r_2(z)$$

$$q(\varepsilon) = \mathcal{N}(\varepsilon|0, 1)$$

$$z = \varepsilon + \theta$$

How to find a good θ ?

By gradient ascent on ELBO_{θ} .

$$\theta_{n+1} \leftarrow \theta_n + \eta \times \nabla_{\theta} \text{ELBO}_{\theta=\theta_n}$$

Why doesn't it work?

$$\nabla_{\theta} \int H(\theta, x) dx = \int \nabla_{\theta} H(\theta, x) dx$$

- Careful when exchanging gradient and integration.

Why doesn't it work?

$$\nabla_{\theta} \int H(\theta, x) dx \neq \int \nabla_{\theta} H(\theta, x) dx$$

- Careful when exchanging gradient and integration.
- May fail unexpectedly.

Why doesn't it work?

$$\nabla_{\theta} \int H(\theta, x) dx = \int \nabla_{\theta} H(\theta, x) dx$$

+ Correction Term

- Careful when exchanging gradient and integration.
- May fail unexpectedly.
- May hold unexpectedly, but with correction.

Our results formally

Non-differentiable models

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

Non-differentiable models

$$p(x^0, \mathbf{z}) = \sum_k [z \in A_k] \cdot r_k(\mathbf{z})$$

- $\mathbf{z} \in \mathbb{R}^n$.

Non-differentiable models

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

- $z \in \mathbb{R}^n$.
- $\{A_1, \dots, A_K\}$ forms a **partition** of \mathbb{R}^n .
- r_k is **differentiable**.

Non-differentiable models

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

- $z \in \mathbb{R}^n$.
- $\{A_1, \dots, A_K\}$ forms a **partition** of \mathbb{R}^n .
- r_k is **differentiable**.
- ∂A_k has Lebesgue **measure zero**.

Wishful thinking

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

Wishful thinking

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} [\sum_k [f_\theta(\epsilon) \in A_k] \cdot H_k(\epsilon, \theta)]$$

Wishful thinking

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} [\sum_k [f_\theta(\epsilon) \in A_k] \cdot H_k(\epsilon, \theta)]$$

Wishful thinking

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

Wishful thinking

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

Wishful thinking

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$\nabla_\theta \text{ELBO}_\theta = \nabla_\theta \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

Wishful thinking

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$\begin{aligned} \nabla_\theta \text{ELBO}_\theta &= \nabla_\theta \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)] \\ &= \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot \nabla_\theta H_k(\epsilon, \theta)] \end{aligned}$$

Wishful thinking

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$\nabla_\theta \text{ELBO}_\theta = \nabla_\theta \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$\neq \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot \nabla_\theta H_k(\epsilon, \theta)]$$

Correction

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta) \right]$$

$$\nabla_\theta \text{ELBO}_\theta = \nabla_\theta \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta) \right]$$

$$= \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot \nabla_\theta H_k(\epsilon, \theta) \right]$$

$$+ \sum_k \text{surface integral over } \partial f_\theta^{-1}(A_k)$$

Correction

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$\nabla_\theta \text{ELBO}_\theta = \nabla_\theta \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$= \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot \nabla_\theta H_k(\epsilon, \theta)]$$

$$+ \sum_k \text{surface integral over } \partial f_\theta^{-1}(A_k)$$

Correction

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$\begin{aligned} \nabla_\theta \text{ELBO}_\theta &= \nabla_\theta \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)] \\ &= \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot \nabla_\theta H_k(\epsilon, \theta)] \\ &\quad + \sum_k \text{ surface integral over } \partial f_\theta^{-1}(A_k) \end{aligned}$$

Correction

$$\int_{B_{k,\theta}} \dots q(\epsilon) d\epsilon \quad [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$\nabla_\theta \text{ELBO}_\theta = \nabla_\theta \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$= \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot \nabla_\theta H_k(\epsilon, \theta)]$$

$$+ \sum_k \text{surface integral over } \partial f_\theta^{-1}(A_k)$$

Correction

$$\nabla_{\theta} \int_{B_{k,\theta}} \cdots q(\epsilon) d\epsilon \quad [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_{\theta} = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$\nabla_{\theta} \text{ELBO}_{\theta} = \nabla_{\theta} \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$= \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot \nabla_{\theta} H_k(\epsilon, \theta)]$$

$$+ \sum_k \text{surface integral over } \partial f_{\theta}^{-1}(A_k)$$

Correction

$$\begin{aligned} & \nabla_{\theta} \int_{B_{k,\theta}} \dots q(\epsilon) d\epsilon \\ &= \int_{B_{k,\theta}} \nabla_{\theta} (\dots q(\epsilon)) d\epsilon + \int_{\partial B_{k,\theta}} (\dots) \cdot d\mathbf{\Sigma} \end{aligned}$$

$$\text{ELBO}_{\theta} = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$\nabla_{\theta} \text{ELBO}_{\theta} = \nabla_{\theta} \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$= \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot \nabla_{\theta} H_k(\epsilon, \theta)]$$

$$+ \sum_k \text{surface integral over } \partial f_{\theta}^{-1}(A_k)$$

Correction

$$\begin{aligned} & \nabla_{\theta} \int_{B_{k,\theta}} \dots q(\epsilon) d\epsilon \\ &= \int_{B_{k,\theta}} \nabla_{\theta} (\dots q(\epsilon)) d\epsilon + \int_{\partial B_{k,\theta}} (\dots) \cdot d\Sigma \end{aligned}$$

$$\text{ELBO}_{\theta} = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$\nabla_{\theta} \text{ELBO}_{\theta} = \nabla_{\theta} \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$= \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot \nabla_{\theta} H_k(\epsilon, \theta)]$$

$$+ \sum_k \text{ surface integral over } \partial f_{\theta}^{-1}(A_k)$$

Correction

$$\begin{aligned} & \nabla_{\theta} \int_{B_{k,\theta}} \dots q(\epsilon) d\epsilon \\ &= \int_{B_{k,\theta}} \nabla_{\theta} (\dots q(\epsilon)) d\epsilon + \int_{\partial B_{k,\theta}} (\dots) \cdot d\Sigma \end{aligned}$$

$$\text{ELBO}_{\theta} = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$\nabla_{\theta} \text{ELBO}_{\theta} = \nabla_{\theta} \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$= \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot \nabla_{\theta} H_k(\epsilon, \theta)]$$

$$+ \sum_k \text{ surface integral over } \partial f_{\theta}^{-1}(A_k)$$

Correction

$$\begin{aligned} & \nabla_{\theta} \int_{B_{k,\theta}} \dots q(\epsilon) d\epsilon \\ &= \int_{B_{k,\theta}} \nabla_{\theta} (\dots q(\epsilon)) d\epsilon + \int_{\partial B_{k,\theta}} (\dots) \cdot d\Sigma \end{aligned}$$

$$\text{ELBO}_{\theta} = \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$\nabla_{\theta} \text{ELBO}_{\theta} = \nabla_{\theta} \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot H_k(\epsilon, \theta)]$$

$$= \mathbb{E}_{q(\epsilon)} [\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot \nabla_{\theta} H_k(\epsilon, \theta)]$$

$$+ \sum_k \text{surface integral over } \partial f_{\theta}^{-1}(A_k)$$

Correction

$$p(x^0, z) = \sum_k [z \in A_k] \cdot r_k(z)$$

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta) \right]$$

$$\nabla_\theta \text{ELBO}_\theta = \nabla_\theta \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta) \right]$$

$$= \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot \nabla_\theta H_k(\epsilon, \theta) \right]$$

$$+ \sum_k \text{surface integral over } \partial f_\theta^{-1}(A_k)$$

Correction

Accounts for the impact of **moving** the boundaries.

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta) \right]$$

$$\nabla_\theta \text{ELBO}_\theta = \nabla_\theta \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta) \right]$$

$$= \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot \nabla_\theta H_k(\epsilon, \theta) \right]$$

$$+ \sum_k \text{surface integral over } \partial f_\theta^{-1}(A_k)$$

Correction

Accounts for the impact of **moving** the boundaries.

Can be estimated by **manifold sampling**.

$$\text{ELBO}_\theta = \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta) \right]$$

$$\nabla_\theta \text{ELBO}_\theta = \nabla_\theta \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot H_k(\epsilon, \theta) \right]$$

$$= \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_\theta^{-1}(A_k)] \cdot \nabla_\theta H_k(\epsilon, \theta) \right]$$

$$+ \sum_k \text{surface integral over } \partial f_\theta^{-1}(A_k)$$

Two ingredients

- Differentiation under moving domain:

$$\nabla_{\theta} \int_{B_{\theta}} g(\epsilon, \theta) d\epsilon = \int_{B_{\theta}} (\nabla_{\theta} g + \nabla_{\epsilon} \cdot (g\mathbf{V}))(\epsilon, \theta) d\epsilon$$

Two ingredients

- Differentiation under moving domain:

$$\nabla_{\theta} \int_{B_{\theta}} g(\epsilon, \theta) d\epsilon = \int_{B_{\theta}} (\nabla_{\theta} g + \nabla_{\epsilon} \cdot (g\mathbf{V}))(\epsilon, \theta) d\epsilon$$

- Divergence theorem:

$$\int_B (\nabla \cdot \mathbf{G}) dV = \int_{\partial B} \mathbf{G} \cdot d\boldsymbol{\Sigma}$$

Two ingredients

$$\begin{aligned}\nabla_{\theta} \int_{B_{\theta}} g(\epsilon, \theta) d\epsilon &= \int_{B_{\theta}} (\nabla_{\theta} g + \nabla_{\epsilon} \cdot (g\mathbf{V}))(\epsilon, \theta) d\epsilon \\ &= \int_{B_{\theta}} \nabla_{\theta} g(\epsilon, \theta) d\epsilon + \int_{\partial B_{\theta}} (g\mathbf{V})(\epsilon, \theta) \cdot d\boldsymbol{\Sigma}\end{aligned}$$

Correction term

Surface integral over $\partial f_{\theta}^{-1}(A_k)$

$$= \int_{\partial f_{\theta}^{-1}(A_k)} (q(\epsilon)H_k(\epsilon, \theta) \mathbf{V}(\epsilon, \theta)) \cdot d\boldsymbol{\Sigma}$$

- $\mathbf{V}(\epsilon, \theta)_{ij} = \left(\frac{\partial f_{\theta}}{\partial \theta_i} \right)_j$
- $\boldsymbol{\Sigma}$ is a normal vector of $\partial f_{\theta}^{-1}(A_k)$.

Correction term

Surface integral over $\partial f_{\theta}^{-1}(A_k)$

$$= \int_{\partial f_{\theta}^{-1}(A_k)} (q(\epsilon)H_k(\epsilon, \theta)\mathbf{V}(\epsilon, \theta)) \cdot d\mathbf{\Sigma}$$

Requires manifold sampling

Hard to estimate in general cases

Correction term

Surface integral over $\partial f_\theta^{-1}(A_k)$

$$= \int_{\partial f_\theta^{-1}(A_k)} (q(\epsilon)H_k(\epsilon, \theta)\mathbf{V}(\epsilon, \theta)) \cdot d\mathbf{\Sigma}$$

- Easy to estimate if $\partial f_\theta^{-1}(A_k)$ is a **hyperplane**.

Correction term

Surface integral over $\partial f_\theta^{-1}(A_k)$

$$= \int_{\partial f_\theta^{-1}(A_k)} (q(\epsilon)H_k(\epsilon, \theta)\mathbf{V}(\epsilon, \theta)) \cdot d\mathbf{\Sigma}$$

- Easy to estimate if $\partial f_\theta^{-1}(A_k)$ is a **hyperplane**.
- Assume the branch condition of each if-statement is **linear** in z .

Subsampling k

$$\begin{aligned} \nabla_{\theta} \text{ELBO}_{\theta} &= \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot \nabla_{\theta} H_k(\epsilon, \theta) \right] \\ &\quad + \sum_k \text{ surface integral over } \partial f_{\theta}^{-1}(A_k) \end{aligned}$$

Subsampling k

$$\begin{aligned} \nabla_{\theta} \text{ELBO}_{\theta} &= \mathbb{E}_{q(\epsilon)} \left[\sum_k [\epsilon \in f_{\theta}^{-1}(A_k)] \cdot \nabla_{\theta} H_k(\epsilon, \theta) \right] \\ &\quad + \sum_k \text{ surface integral over } \partial f_{\theta}^{-1}(A_k) \end{aligned}$$

- For computational efficiency, we **subsample** k surface integrals.

Experiments

Implementation

- Implemented a **black-box variational inference** engine for a simple probabilistic programming language
- Supports `sample`, `observe`, `if`, ...
- Written in Python, using **autograd** package.

Benchmarks

textmsg

- Models #'s of per-day **SNS msg's**, where SNS-usage pattern changes on some day.
- **Non-differentiable** part: the day of change in SNS-usage pattern.
- Given #'s of per-day SNS msg's over 2 months, infer **the day** when the pattern changes.

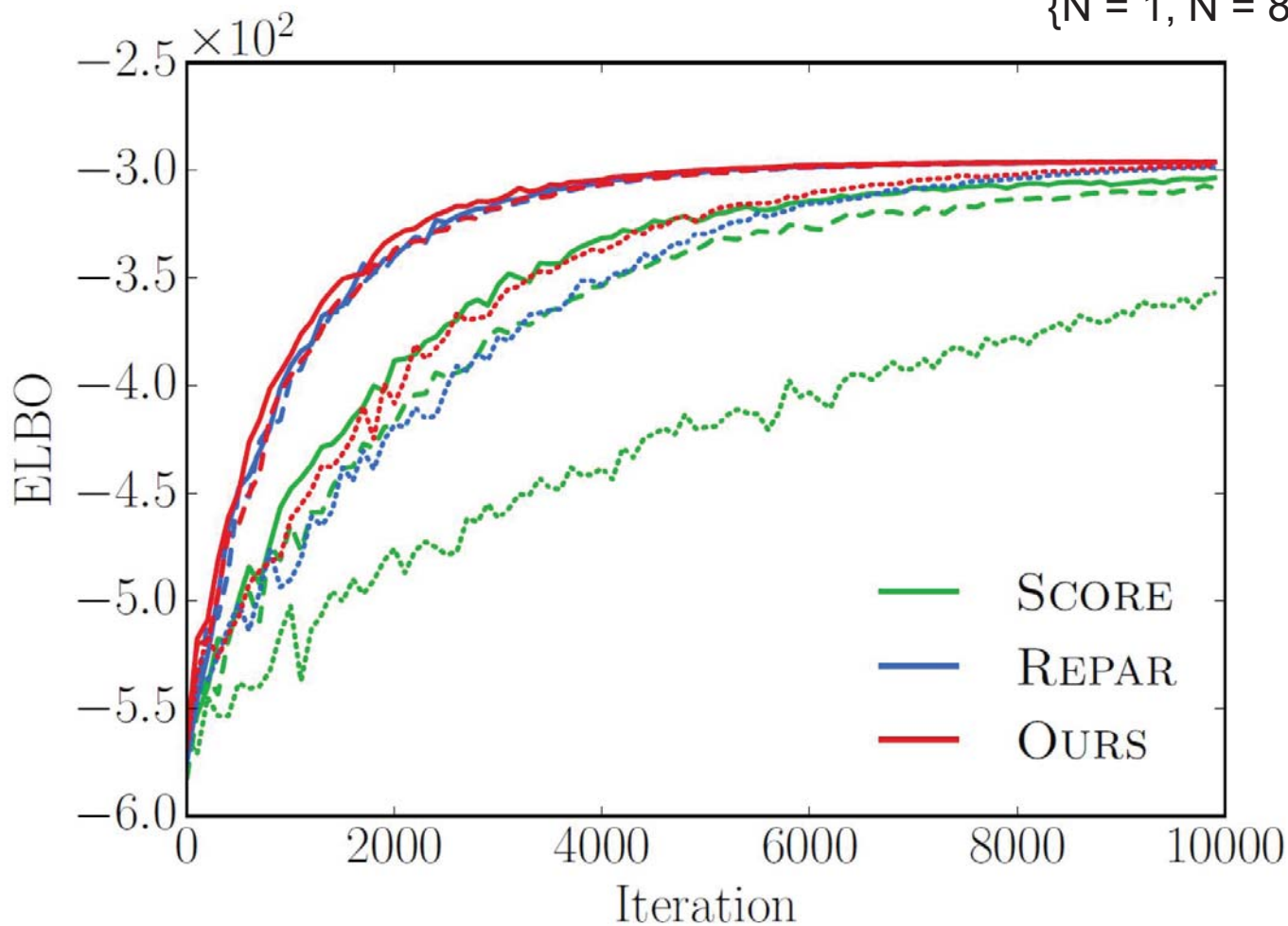
Benchmarks

temperature

- Models random dynamics of a **controller** that tries to keep room temp. stable.
- **Non-differentiable** part: on/off of air conditioner, on which evolution of room temp. depends.
- Given noisy observations of temp. at each step, infer **on/off status** of the controller at each step.

ELBO

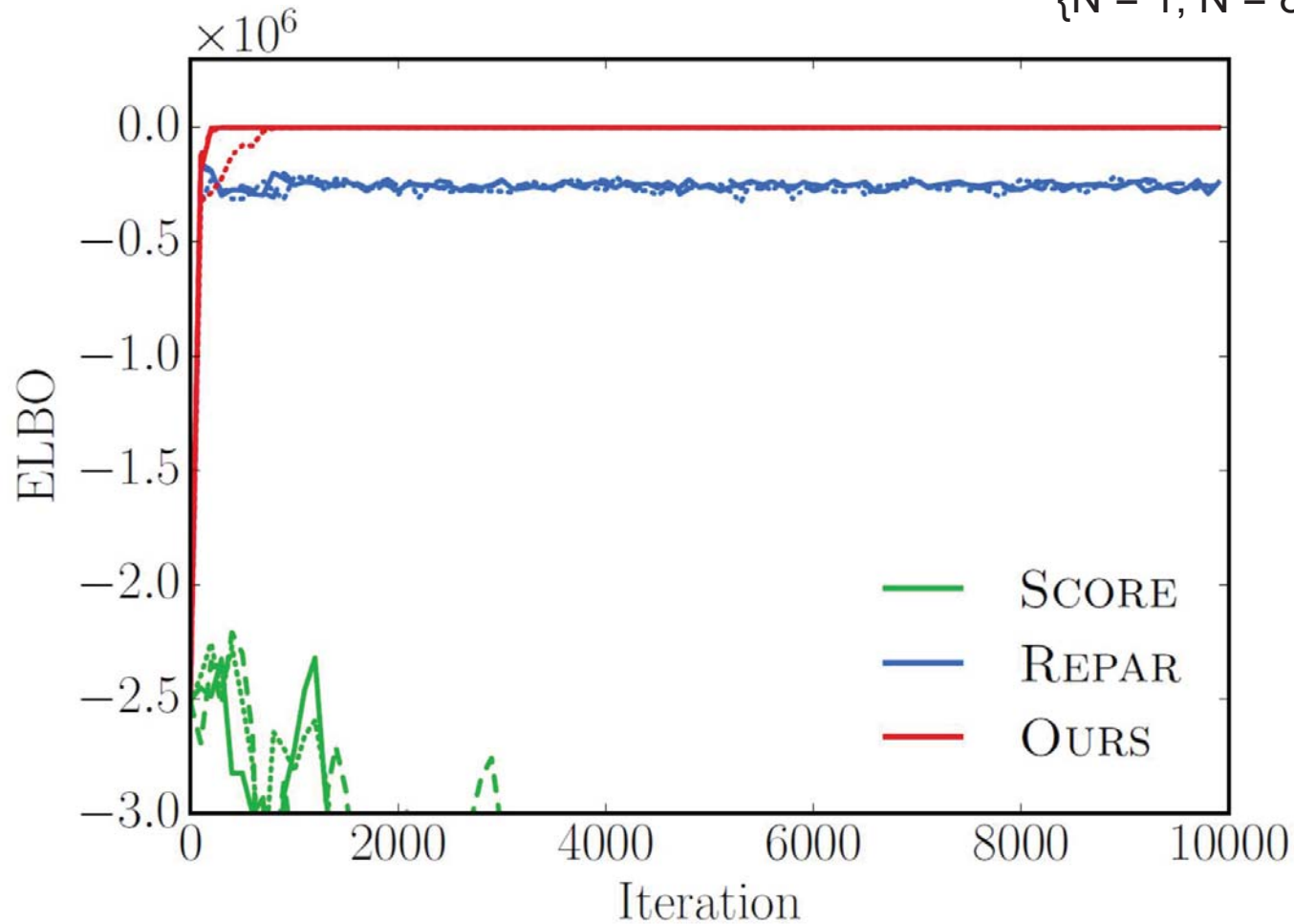
{dotted, dashed, solid} lines:
{N = 1, N = 8, N = 16}



textmsg (stepsize = 0.001)

ELBO

{dotted, dashed, solid} lines:
{N = 1, N = 8, N = 16}



temperature (stepsize = 0.01)

Computation time

Estimator	temperature	textmsg	influenza
SCORE	21.7	4.9	18.7
REPARAM	46.1	15.4	251.4
OURS	79.2	24.9	269.8

High-level message

$$\nabla_{\theta} \int H(\theta, x) dx = \int \nabla_{\theta} H(\theta, x) dx$$

- Careful when exchanging gradient and integration.

High-level message

$$\nabla_{\theta} \int H(\theta, x) dx = \int \nabla_{\theta} H(\theta, x) dx$$

- Careful when exchanging gradient and integration.
- May **fail** unexpectedly.

High-level message

$$\nabla_{\theta} \int H(\theta, x) dx = \int \nabla_{\theta} H(\theta, x) dx$$

- Careful when exchanging gradient and integration.
- May **fail** unexpectedly.
- May **hold** unexpectedly, but **with correction**.

Any questions?