# Provable Smoothness Guarantees for Black-Box Variational Inference

## Justin Domke

College of Information and Computer Sciences, UMass Amherst

## This paper in one slide

Variational Inference (VI): Approximate $p(\boldsymbol{z}|\boldsymbol{x})$ with $q_{\boldsymbol{w}}(\boldsymbol{z})$ by solving

$$\max_{w} \text{ELBO}(\boldsymbol{w}), \quad -\text{ELBO}(\boldsymbol{w}) = \underbrace{- \mathop{\mathbb{E}}_{z \sim q_{\boldsymbol{w}}} \log p(z, \boldsymbol{x})}_{\text{Energy term } l(\boldsymbol{w})} + \underbrace{\mathop{\mathbb{E}}_{z \sim q_{\boldsymbol{w}}} \log q_{\boldsymbol{w}}(z)}_{\text{Neg-Entropy term } h(\boldsymbol{w})} \quad .$$

**This paper**: If $p(\boldsymbol{z}, \boldsymbol{x})$ is *nice* then $l(\boldsymbol{w})$ is also *nice* (for Gaussian $q_{\boldsymbol{w}}$)

- $\log p(\boldsymbol{z}, \boldsymbol{x})$ smooth over $\boldsymbol{z}$ $\Rightarrow l(\boldsymbol{w})$ smooth
- $\log p(\boldsymbol{z}, \boldsymbol{x})$ strongly concave over $\boldsymbol{z} \Rightarrow l(\boldsymbol{w})$ strongly convex

**Implications**: If you can do MAP inference, then you can do VI, *as long as you're careful.*

## Motivation

Black-Box VI. Do SGD on $\mathrm{ELBO}(\boldsymbol{w})$.

**Example Problem**: Three different initializations, three different step sizes. (Exact gradients)

# Goals

Black-Box VI often works, but also often fails!

To give a convergence guarantee for SGD you need two things:

- A bound on the gradient estimator's variance.
- A proof that the objective is smooth or (strongly) convex (or both).

## Main Result: Smoothness

- $\phi(x)$ is $M$-smooth if $\|\nabla\phi(x) - \nabla\phi(x')\|_2 \le M\|x - x'\|_2$.

Theorem: Say $q_w$ is a **location-scale family** with a **standardized base distribution** (e.g. a Gaussian) and $f(z)$ is $M$-smooth. Then,

$$l(w) = \mathop{\mathbb{E}}_{z \sim q_w} f(z)$$

is also $M$-smooth.

Proof: Define inner-product space + Bessel's inequality + several laborious exact calculations for location-scale families.

# Secondary Result: Strong Convexity

- $\phi(x)$ $c$-strongly convex if $\phi(y) \geq \phi(x) + \nabla\phi(x)^\top (y - x) + \frac{c}{2} \|y - x\|_2^2$

Theorem: Say $q_w$ is a **location-scale family** with a **standardized base distribution** (e.g. a Gaussian) and $f(z)$ is $c$-strongly convex. Then,

$$l(w) = \mathop{\mathbb{E}}_{z \sim q_w} f(z)$$

is also $c$-strongly convex.

Proof: Comparatively easy.

## Convergence Considerations

Say $\log p(z, x)$ is $M$-smooth. Want to opt. $-\text{ELBO}(w) = l(w) + h(w)$.

Main result: $l(w)$ is $M$-smooth.
Problem: $h(w)$ is *not* smooth.

One solution:
- Define $\mathscr{W}_M = \left\{ w \,\middle|\, \text{Cov of } q_w \succeq \frac{1}{M} \right\}$.
- Result: Optimum of ELBO is in $\mathscr{W}_M$.
- Result: $h(w)$ is $M$-smooth over $\mathscr{W}_M$ (so $l + h$ is $2M$-smooth)
- So projected gradient descent works.

Another solution: Do proximal gradient descent.

# Demonstration

Compare three algorithms:

- Projected optimization (step $1/(2M)$)
- Proximal optimization (step $1/M$)
- Naive optimization (step $1/M$)

Initialize $q_{\boldsymbol{w}}$ with mean 0 and covariance $\rho^2 I$ where $\rho$ is a scaling factor.