

Provable Gradient Variance Guarantees for Black-Box Variational Inference

Justin Domke, University of Massachusetts Amherst

1. Overview

Motivation:

- Recent black box variational inference methods used reparameterization gradient estimators.
- Their variance is poorly understood.
- This means we don't really understand when they work!

Contributions:

- If target distribution $\log p(\mathbf{z}, \mathbf{x})$ is M -smooth over \mathbf{z} , then

$$\mathbb{E} \|\mathbf{g}\|_2^2 \leq aM^2 \|\mathbf{w} - \bar{\mathbf{w}}\|_2^2$$

for some fixed $\bar{\mathbf{w}}$.

- This generalizes to the case where $\log p$ has different smoothness in different directions.
- This generalizes to consider data subsampling.
- All contributions unimprovable!

2. Background and Setup

2.1 Variational Inference

Goal is to maximize

$$\text{ELBO}(\mathbf{w}) = \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\mathbf{w}}(\mathbf{z})} \log p(\mathbf{z}, \mathbf{x})}_{l(\mathbf{w})} + \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\mathbf{w}}(\mathbf{z})} (-\log q_{\mathbf{w}}(\mathbf{z}))}_{h(\mathbf{w})}.$$

Equivalent to minimizing $KL(q_{\mathbf{w}}(\mathbf{z}) \| p(\mathbf{z} | \mathbf{x}))$.

2.2 Reparameterization Estimators

- Choose $s(\mathbf{u})$ and $\mathcal{T}_{\mathbf{w}}(\mathbf{u})$ such that if $\mathbf{u} \sim s$, then $\mathcal{T}_{\mathbf{w}}(\mathbf{u})$ has same distribution as $q_{\mathbf{w}}(\mathbf{z})$.
- Then,

$$l(\mathbf{w}) = \mathbb{E}_{\mathbf{u} \sim s} f(\mathcal{T}_{\mathbf{w}}(\mathbf{u}))$$

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

- Estimator:

$$\mathbf{g} = \nabla_{\mathbf{w}} f(\mathcal{T}_{\mathbf{w}}(\mathbf{u})).$$

- Goal of this paper: Bound $\mathbb{E} \|\mathbf{g}\|_2^2$

2.3 Location-Scale Families

- A location scale family $q_{\mathbf{w}}(\mathbf{z})$ is the distribution that results from drawing $\mathbf{u} \sim s$ and returning $\mathcal{T}(\mathbf{u})$, where

$$\mathcal{T}_{\mathbf{w}}(\mathbf{u}) = C\mathbf{u} + \mathbf{m}.$$

- For example, if $s = \mathcal{N}(0, I)$ then $q_{\mathbf{w}}(\mathbf{z}) = \mathcal{N}(\mathbf{m}, CC^T)$.
- s "standardized" if $(u_1, \dots, u_d) \sim s$ are i.i.d. with $\mathbb{E} u_1 = \mathbb{E} u_1^3 = 0$ and $\mathbb{V} u_1 = 1$.
- Bounds will depend on $\kappa = \mathbb{E}[u_1^4]$.

3. Main Results

3.1 Main Theorem

Suppose f is M -smooth, \mathbf{z}^* is a stationary point of f , and s is standardized. Let $\mathbf{g} = \nabla_{\mathbf{w}} f(\mathcal{T}_{\mathbf{w}}(\mathbf{u}))$ for $\mathbf{u} \sim s$. Then,

$$\mathbb{E} \|\mathbf{g}\|_2^2 \leq M^2 \left((d+1) \|\mathbf{m} - \mathbf{z}^*\|_2^2 + (d+\kappa) \|C\|_F^2 \right).$$

This result is unimprovable.

3.2 Main Proof

$$\begin{aligned} \mathbb{E} \|\mathbf{g}\|_2^2 &= \mathbb{E} \|\nabla_{\mathbf{w}} f(\mathcal{T}_{\mathbf{w}}(\mathbf{u}))\|_2^2 && \text{(Definition of } \mathbf{g} \text{)} \\ &= \mathbb{E} \|\nabla f(\mathcal{T}_{\mathbf{w}}(\mathbf{u}))\|_2^2 (1 + \|\mathbf{u}\|_2^2) && \text{(First Technical Lemma)} \\ &= \mathbb{E} \|\nabla f(\mathcal{T}_{\mathbf{w}}(\mathbf{u})) - \nabla f(\mathbf{z}^*)\|_2^2 (1 + \|\mathbf{u}\|_2^2) && (\nabla f(\mathbf{z}^*) = 0) \\ &\leq \mathbb{E} M^2 \|\mathcal{T}_{\mathbf{w}}(\mathbf{u}) - \mathbf{z}^*\|_2^2 (1 + \|\mathbf{u}\|_2^2) && (f \text{ is smooth)} \\ &= M^2 \left((d+1) \|\mathbf{m} - \mathbf{z}^*\|_2^2 + (d+\kappa) \|C\|_F^2 \right). && \text{(Second Technical Lemma)} \end{aligned}$$

To prove first technical lemma:

- Substitute definition of $\mathcal{T}_{\mathbf{w}}$
- Compute all components $\|\nabla_{w_i} f(\mathcal{T}_{\mathbf{w}}(\mathbf{u}))\|_2^2$
- Sum and simplify.

To prove second technical lemma:

- Substitute definition of $\mathcal{T}_{\mathbf{w}}$.
- Resulting expression has expectations between order 0 and 4 in \mathbf{u} .
- Compute each of these and simplify using that s is standardized.

3.3 Generalized Theorem

Definition: f is M -matrix-smooth if $\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{z})\|_2 \leq \|M(\mathbf{y} - \mathbf{z})\|_2$ (for symmetric M).

Suppose f is M -matrix smooth, \mathbf{z}^* is a stationary point of f , and s is standardized. Then,

$$\mathbb{E} \|\mathbf{g}\|_2^2 \leq (d+1) \|M(\mathbf{m} - \mathbf{z}^*)\|_2^2 + (d+\kappa) \|MC\|_F^2.$$

Unimprovable!

(Proof as above, with a trick of "absorbing" M into the parameters.)

3.4 Generalized Generalized Theorem

Suppose that $f(\mathbf{z}) = \sum_{n=1}^N f_n(\mathbf{z})$.

Suppose f_n is M_n -matrix-smooth, \mathbf{z}_n^* is a stationary point of f_n , and s is standardized.

Let $\mathbf{g} = \frac{1}{\pi(n)} \nabla f_n(\mathcal{T}_{\mathbf{w}}(\mathbf{u}))$ for $\mathbf{u} \sim s$ and $n \sim \pi$ independent. Then,

$$\mathbb{E} \|\mathbf{g}\|_2^2 \leq \sum_{n=1}^N \frac{1}{\pi(n)} \left((d+1) \|M_n(\mathbf{m} - \mathbf{z}_n^*)\|_2^2 + (d+\kappa) \|M_n C\|_F^2 \right).$$

This result is unimprovable.

(Proof uses previous result as a lemma, takes expectation over n .)

4. Experiments

Dataset	Type	# data	# dims
boston	r	506	13
fires	r	517	12
cpusmall	r	8192	13
ala	c	1695	124
ionosphere	c	351	35
australian	c	690	15
sonar	c	208	61
mushrooms	c	8124	113

Table 1: Regression (r) and classification (c) datasets

