

# Faster Stochastic Variational Inference using Proximal-Gradient Methods with General Divergence Functions

Mohammad Emtiyaz Khan<sup>1</sup>, Reza Babanezhad<sup>2</sup>, Wu Lin<sup>3</sup>, Mark Schmidt<sup>2</sup>, Masashi Sugiyama<sup>4</sup>

<sup>1</sup>Ecole Polytechnique Fédérale de Lausanne, <sup>2</sup>University of British Columbia, <sup>3</sup>University of Waterloo, <sup>4</sup>University of Tokyo

## Introduction

### Issues with existing methods:

- ▶ Most existing methods rely on “black-box” methods for optimization, and ignore the *geometry* of the variational-parameter space.
- ▶ Methods like stochastic variational inference (SVI) use *natural gradients* to exploit the geometry, but only apply to conjugate models.
- ▶ Theoretical convergence rate of natural-gradient methods for variational inference is unclear.

**Contributions:** We propose a **proximal-gradient** framework that unifies most of the existing approaches. Our method,

- ▶ can be stochastic to allow huge datasets,
- ▶ can exploit the geometry to improve performance,
- ▶ can yield a closed-form update even for non-conjugate models.

We also **analyze the convergence rate** of the proposed method, clearly showing the conditions under which a natural-gradient method can enable large steps than basic stochastic gradient methods.

## Variational Inference

In Bayesian inference, we need to marginalize the unknowns  $\mathbf{z}$  over the joint distribution  $p(\mathbf{y}, \mathbf{z})$  of the model, given data  $\mathbf{y}$ . Variational inference maximizes a lower bound to the integral w.r.t. a distribution  $q(\mathbf{z}|\lambda)$  where  $\lambda$  is the vector of variational parameters.

$$\log \int p(\mathbf{y}, \mathbf{z}) d\mathbf{z} = \log \int q(\mathbf{z}|\lambda) \frac{p(\mathbf{y}, \mathbf{z})}{q(\mathbf{z}|\lambda)} d\mathbf{z} \geq \max_{\lambda} \mathbb{E}_{q(\mathbf{z}|\lambda)} \left[ \log \frac{p(\mathbf{y}, \mathbf{z})}{q(\mathbf{z}|\lambda)} \right] := \underline{\mathcal{L}}(\lambda).$$

Existing methods for lower-bound optimization differ from each other in the choice of approximation and divergence functions (highlighted in red).

Gradient Descent:  $\lambda_{k+1} = \lambda_k + \beta_k \nabla \underline{\mathcal{L}}(\lambda_k)$

$$\iff \lambda_{k+1} = \operatorname{argmin}_{\lambda} -\lambda^T [\nabla \underline{\mathcal{L}}(\lambda_k)] + \frac{1}{2\beta_k} \|\lambda - \lambda_k\|_2^2,$$

Natural Gradient:  $\lambda_{k+1} = \operatorname{argmin}_{\lambda} -\lambda^T [\nabla \underline{\mathcal{L}}(\lambda_k)] + \frac{1}{\beta_k} \mathbb{D}_{KL}^{sym}[q(\mathbf{z}|\lambda) \| q(\mathbf{z}|\lambda_k)],$

Mirror Descent:  $\lambda_{k+1} = \operatorname{argmin}_{\lambda} -\lambda^T [\nabla \underline{\mathcal{L}}(\lambda_k)] + \frac{1}{\beta_k} \mathbb{D}_{Breg}(\lambda \| \lambda_k),$

Trust Region:  $\lambda_{k+1} = \operatorname{argmin}_{\lambda} -\underline{\mathcal{L}}(\lambda) + \frac{1}{\beta_k} \mathbb{D}_{KL}[q(\mathbf{z}|\lambda) \| q(\mathbf{z}|\lambda_k)],$

KL proximal:  $\lambda_{k+1} = \operatorname{argmin}_{\lambda} \lambda^T [\nabla f(\lambda_k)] + h(\lambda) + \frac{1}{\beta_k} \mathbb{D}_{KL}[q(\mathbf{z}|\lambda) \| q(\mathbf{z}|\lambda_k)].$

### Pros and cons:

- ▶ Gradient descent is generally applicable but ignores the geometry.
- ▶ Natural-gradient methods exploit the geometry but only apply to conjugate models. Mirror-descent does not cover all cases of natural-gradient and may not give a closed-form solution.
- ▶ Trust-region method may also lead to a difficult optimization problem.
- ▶ The KL-proximal method requires exact gradients.

## Proximal-Gradient SVI

We propose a proximal-gradient framework that unifies many existing approaches. Specifically, our method generalizes the KL-proximal method by allowing stochastic gradients and general divergence functions.

We split the ratio  $p(\mathbf{y}, \mathbf{z})/q(\mathbf{z}|\lambda) \equiv c \tilde{p}_d(\mathbf{z}|\lambda) \tilde{p}_e(\mathbf{z}|\lambda)$ , where  $\tilde{p}_d$  contains all factors that make the optimization difficult while  $\tilde{p}_e$  contains the rest.

$$\underline{\mathcal{L}}(\lambda) = \mathbb{E}_{q(\mathbf{z}|\lambda)} \left[ \log \frac{p(\mathbf{y}, \mathbf{z}|\theta)}{q(\mathbf{z}|\lambda)} \right] := \underbrace{\mathbb{E}_{q(\mathbf{z}|\lambda)} [\log \tilde{p}_d(\mathbf{z}|\lambda)]}_{-f(\lambda)} + \underbrace{\mathbb{E}_{q(\mathbf{z}|\lambda)} [\log \tilde{p}_e(\mathbf{z}|\lambda)]}_{-h(\lambda)},$$

We linearize the difficult terms and solve the following subproblem:

$$\lambda_{k+1} = \operatorname{argmin}_{\lambda \in \mathcal{S}} \left\{ \lambda^T [\hat{\nabla} f(\lambda_k)] + h(\lambda) + \frac{1}{\beta_k} \mathbb{D}(\lambda \| \lambda_k) \right\}. \quad (1)$$

where  $\hat{\nabla} f(\lambda_k)$  is the noisy gradient and  $\mathbb{D}$  is a divergence function. We make the following assumptions:

- ▶ (A1-A2)  $f$  is non-convex and  $L$ -smooth and  $h$  is convex.
- ▶ (A3-A4)  $\hat{\nabla} f(\lambda_k)$  is an unbiased estimate and its variance  $\leq \sigma^2$ .
- ▶ (A5-A6)  $\mathbb{D}(\lambda \| \lambda') > 0, \forall \lambda \neq \lambda'$ , and there exist an  $\alpha > 0$  such that for all  $\lambda, \lambda'$  generated by (1) we have:  $(\lambda - \lambda')^T [\nabla \lambda \mathbb{D}(\lambda \| \lambda')] \geq \alpha \|\lambda - \lambda'\|^2$ .

## Convergence

**Deterministic algorithm:** Suppose A1, A2, A5, and A6 be satisfied and we run  $t$  iterations of (1) with a fixed step-size  $\beta_k = \alpha/L$  for all  $k$  and an exact gradient  $\nabla f(\lambda)$ , then we have

$$\min_{k \in \{0, 1, \dots, t-1\}} \|\lambda_{k+1} - \lambda_k\|^2 \leq \frac{2C_0}{\alpha t}, \quad (2)$$

where  $C_0 = \underline{\mathcal{L}}^* - \underline{\mathcal{L}}(\lambda_0)$  is the initial (constant) sub-optimality.

**Stochastic algorithm:** Let A1-A6 be satisfied and we run  $t$  iterations of (1) for a fixed step-size  $\beta_k = \gamma \alpha_*/L$  (where  $0 < \gamma < 2$  is a scalar) and fixed batch-size  $M_k = M$  for all  $k$  with a stochastic gradient  $\hat{\nabla} f(\lambda)$ , then we have

$$\mathbb{E}_{R, \xi} (\|\lambda_{R+1} - \lambda_R\|^2) \leq \frac{1}{2-\gamma} \left[ \frac{2C_0}{\alpha_* t} + \frac{\gamma c \sigma^2}{ML} \right]. \quad (3)$$

where  $c$  is a constant such that  $c > 1/(2\alpha)$ ,  $\alpha_* := \alpha - 1/(2c)$ , and the expectation is taken with respect to the noise  $\xi := \{\xi_0, \xi_1, \dots, \xi_{t-1}\}$  and a random variable  $R$  which follows the uniform distribution.

## Example: Gaussian Process Models

We give an example for Gaussian-process models where we use  $q(\mathbf{z}|\lambda) := \mathcal{N}(\mathbf{z}|\mathbf{m}, \mathbf{V})$ , i.e.  $\lambda = \{\mathbf{m}, \mathbf{V}\}$  with mean  $\mathbf{m}$  and covariance  $\mathbf{V}$ . Consider  $N$  input-output pairs  $\{y_n, \mathbf{x}_n\}$  indexed by  $n$ . Let  $z_n := f(\mathbf{x}_n)$  be the latent function drawn from a GP with mean 0 and covariance  $\mathbf{K}$ . We use a non-Gaussian likelihood  $p(y_n|z_n)$  to model the output. We use the split,

$$\frac{p(\mathbf{y}, \mathbf{z})}{q(\mathbf{z}|\lambda)} = \underbrace{\prod_{n=1}^N p(y_n|z_n)}_{\tilde{p}_d(\mathbf{z}|\lambda)} \times \underbrace{\frac{\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{K})}{\mathcal{N}(\mathbf{z}|\mathbf{m}, \mathbf{V})}}_{\tilde{p}_e(\mathbf{z}|\lambda)}. \quad (4)$$

By substituting in the lower bound, we obtain the following:

$$-\underline{\mathcal{L}}(\lambda) := \underbrace{\sum_n \mathbb{E}_q[-\log p(y_n|z_n)]}_{f(\lambda)} + \underbrace{\mathbb{D}_{KL}[\mathcal{N}(\mathbf{z}|\mathbf{m}, \mathbf{V}) \| \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{K})]}_{h(\lambda)}. \quad (5)$$

We randomly pick a data-term  $\mathbb{E}_q[-\log p(y_n|z_n)]$  and denote it by  $f_n(m_n, v_n)$  where  $m_n$  and  $v_n$  are mean and variance of  $z_n$ . We compute its stochastic gradients and use a KL divergence function:

$$\lambda^T [\hat{\nabla} f(\lambda_k)] + h(\lambda) + \frac{1}{\beta_k} \mathbb{D}(\lambda \| \lambda_k) = \left[ m_n \left\{ \hat{\nabla}_{m_n} f_n(\lambda_k) \right\} + v_n \left\{ \hat{\nabla}_{v_n} f_n(\lambda_k) \right\} \right] + \mathbb{D}_{KL}[\mathcal{N}(\mathbf{m}, \mathbf{V}) \| \mathcal{N}(\mathbf{0}, \mathbf{K})] + \frac{1}{\beta_k} \mathbb{D}_{KL}[\mathcal{N}(\mathbf{m}, \mathbf{V}) \| \mathcal{N}(\mathbf{m}_k, \mathbf{V}_k)],$$

The minimum of this function can be obtained as shown below:

$$\mathcal{N}(\mathbf{z}|\mathbf{m}_{k+1}, \mathbf{V}_{k+1}) \propto \left[ e^{z_n - c_k z_n^2} \times \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{K}) \right]^{1-r_k} [\mathcal{N}(\mathbf{z}|\mathbf{m}_k, \mathbf{V}_k)]^{r_k}$$

where  $r_k := 1/(1 + \beta_k)$ . This equation can be implemented by solving two linear systems (see the paper for details).

## Experiments

Results for GP classification and Correlated topic model. We show that, compared to existing methods, PG-SVI requires less number of passes through the data to converge.

