## On the Convergence of Black-Box Variational Inference

**Kyurae Kim** 

Jisu Oh

Kaiwen Wu

University of Pennsylvania kyrkim@seas.upenn.edu

North Carolina State University joh26@ncsu.edu

University of Pennsylvania kaiwenwu@seas.upenn.edu

Yi-An Ma

University of California, San Diego yianma@ucsd.edu

Jacob R. Gardner University of Pennsylvania jacobrg@seas.upenn.edu

#### **Abstract**

We provide the first convergence guarantee for black-box variational inference (BBVI) with the reparameterization gradient. While preliminary investigations worked on simplified versions of BBVI (e.g., bounded domain, bounded support, only optimizing for the scale, and such), our setup does not need any such algorithmic modifications. Our results hold for log-smooth posterior densities with and without strong log-concavity and the location-scale variational family. Notably, our analysis reveals that certain algorithm design choices commonly employed in practice, such as nonlinear parameterizations of the scale matrix, can result in suboptimal convergence rates. Fortunately, running BBVI with proximal stochastic gradient descent fixes these limitations and thus achieves the strongest known convergence guarantees. We evaluate this theoretical insight by comparing proximal SGD against other standard implementations of BBVI on large-scale Bayesian inference problems.

#### 1 Introduction

Despite the practical success of black-box variational inference (BBVI; Kucukelbir *et al.*, 2017; Ranganath *et al.*, 2014; Titsias & Lázaro-Gredilla, 2014), also known as stochastic gradient variational Bayes and Monte Carlo variational inference, whether it converges under appropriate assumptions on the target problem have been an open problem for a decade. While our understanding of BBVI has been advancing (Bhatia *et al.*, 2022; Challis & Barber, 2013; Domke, 2019, 2020; Hoffman & Ma, 2020), a full convergence guarantee that extends to the practical implementations as used in probabilistic programming languages (PPL) such as Stan (Carpenter *et al.*, 2017), Turing (Ge *et al.*, 2018), Tensorflow Probability (Dillon *et al.*, 2017), Pyro (Bingham *et al.*, 2019), and PyMC (Patil *et al.*, 2010) has yet to be demonstrated.

Due to our lack of understanding, a consensus on how we should implement our BBVI algorithms has yet to be achieved. For example, when the variational family is chosen to be the location-scale family, the "scale" matrix can be parameterized linearly or nonlinearly, and both parameterizations are used by default in popular software packages. (See Table 1 in Kim *et al.* 2023.) Surprisingly, as we will show, seemingly innocuous design choices like these can substantially impact the convergence of BBVI. This is critical as BBVI has been shown to be less robust (*e.g.*, sensitive to initial points, stepsizes, and such) than competing inference methods such as Markov chain Monte Carlo (MCMC). (See Dhaka *et al.*, 2020; Domke, 2020; Welandawe *et al.*, 2022; Yao *et al.*, 2018.) Instead, the evaluation of BBVI algorithms has been relying on expensive empirical evaluations (Agrawal *et al.*, 2020; Dhaka *et al.*, 2021; Giordano *et al.*, 2018; Yao *et al.*, 2018).

To rigorously analyze the design of BBVI algorithms, we establish the first convergence guarantee for the implementations *precisely* as used in practice. We provide results for BBVI with the reparameterization gradient (RP; Kingma & Welling, 2014; Titsias & Lázaro-Gredilla, 2014) and the location-scale variational family, arguably the most widely used combination in practice. Our results apply to log-smooth posteriors, which is a routine assumption for analyzing the convergence of stochastic optimization (Garrigos & Gower, 2023) and sampling algorithms (Dwivedi *et al.*, 2019, §2.3). The key is to show that evidence lower bound (ELBO; Jordan *et al.*, 1999) satisfies regularity conditions required by convergence proofs of stochastic gradient descent (SGD; Bottou, 1999; Nemirovski *et al.*, 2009; Robbins & Monro, 1951), the workhorse underlying BBVI.

Our analysis reveals that nonlinear scale matrix parameterizations used in practice are suboptimal: they provably break strong convexity and sometimes even convexity. Even if the posterior is strongly log-concave, the ELBO is not strongly convex anymore. This contrasts with linear parameterizations, which guarantee the ELBO to be strongly convex if the posterior is strongly log-concave (Domke, 2020). Under linear parameterizations, however, the ELBO is no longer smooth, making optimization challenging. Because of this, Domke (2020) proposed to use proximal SGD, which Agrawal & Domke (2021, Appendix A) report to have better performance than vanilla SGD with nonlinear parameterizations. Indeed, we show that BBVI with proximal SGD achieves the *fastest* known converges rates of SGD, unlike vanilla BBVI. Thus, we provide a concrete reason for employing proximal SGD. We evaluate this insight on large-scale Bayesian inference problems by implementing an Adam-like (Kingma & Ba, 2015) variant of proximal SGD proposed by Yun *et al.* (2021).

Concurrently to this work, convergence guarantees on BBVI with the RP and the sticking-the-landing estimator (STL; Roeder *et al.*, 2017) under the linear parameterization were published by Domke *et al.* (2023). To achieve this, they show that a quadratic bound on the gradient variance is sufficient to guarantee the convergence of projected and proximal SGD. In contrast, we focus on analyzing the ELBO under nonlinear parameterizations and connect it to existing analysis strategies. A more in-depth comparison of the two works is provided in Appendix E.

- Convergence Guarantee for BBVI: Theorem 3 establishes a convergence guarantee for BBVI with assumptions matching the implementations used in practice. That is, without algorithmic simplifications and unrealistic assumptions such as bounded domain or bounded support.
- **Optimality of Linear Parameterizations:** Theorem 2 shows that, for location-scale variational families, nonlinear scale parameterizations prevent the ELBO from being strongly-convex even when the target posterior is strongly log-concave.
- **3** Convergence Guarantee for Proximal BBVI: Theorem 4 guarantees that, if proximal SGD is used, BBVI on  $\mu$ -strongly log-concave posteriors can obtain a solution  $\epsilon$ -close to the global optimum with  $\mathcal{O}(1/\epsilon)$  iterations.
- **Q** Evaluation of Proximal BBVI in Practice: In Section 5, we evaluate the utility of proximal SGD on large-scale Bayesian inference problems.

#### 2 Background

**Notation** Random variables are denoted in serif (e.g., x, x), vectors are in bold (e.g., x, x), and matrices are in bold capitals (e.g. A). For a vector  $\mathbf{x} \in \mathbb{R}^d$ , we denote the inner product as  $\mathbf{x}^T \mathbf{x}$  and  $\langle \mathbf{x}, \mathbf{x} \rangle$ , the  $\ell_2$ -norm as  $\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^T \mathbf{x}}$ . For a matrix  $\mathbf{A}$ ,  $\|\mathbf{A}\|_F = \sqrt{\operatorname{tr}(\mathbf{A}^T \mathbf{A})}$  denotes the Frobenius norm.  $\mathbb{S}_{++}^d$  is the set of positive definite matrices. For some function f,  $\mathbf{D}_i f$  denotes the ith coordinate of  $\nabla f$ , and  $\mathbf{C}^k(\mathcal{X}, \mathcal{Y})$  is the set of k-time differentiable continuous functions mapping from  $\mathcal{X}$  to  $\mathcal{Y}$ .

#### 2.1 Black-Box Variational Inference

Variational inference (VI, Blei *et al.*, 2017; Jordan *et al.*, 1999; Zhang *et al.*, 2019) aims to minimize the exclusive (or backward/reverse) Kullback-Leibler (KL) divergence as:

```
 \begin{array}{ll} \text{minimize} \ \ D_{\text{KL}}\left(q_{\lambda},\pi\right)\triangleq\mathbb{E}_{\mathbf{z}\sim q_{\lambda}}-\log\pi\left(\mathbf{z}\right)-\mathbb{H}\left(q_{\lambda}\right), \\ \text{where} \quad D_{\text{KL}}\left(q_{\lambda},\pi\right) \quad \text{is the KL divergence,} \qquad \qquad \mathbb{H} \quad \text{is the differential entropy,} \\ \pi \quad \quad \text{is the (target) posterior distribution, and} \quad q_{\lambda} \quad \text{is the variational distribution,} \\ \end{array}
```

While alternative approaches to VI (Dieng *et al.*, 2017; Hernandez-Lobato *et al.*, 2016; Kim *et al.*, 2022; Naesseth *et al.*, 2020) exist, so far, exclusive KL minimization has been the most successful. We thus use "exclusive KL minimization" as a synonym for VI, following convention.

Equivalently, one minimizes the negative evidence lower bound (ELBO, Jordan et al., 1999) F:

$$\underset{\lambda \in \Lambda}{\text{minimize}} F(\lambda) \triangleq \mathbb{E}_{\mathbf{z} \sim q_{\lambda}} - \log p(\mathbf{z}, \mathbf{x}) - \mathbb{H}(q_{\lambda}),$$

where  $\log p(\mathbf{z}, \mathbf{x})$  is the *joint likelihood*, which is proportional to the posterior as  $\pi(\mathbf{z}) \propto p(\mathbf{z}, \mathbf{x}) = p(\mathbf{x} \mid \mathbf{z}) p(\mathbf{z})$ , where  $p(\mathbf{x} \mid \mathbf{z})$  is the likelihood and  $p(\mathbf{z})$  is the prior.

#### 2.2 Variational Family

In this work, we focus on the following variational family. ( $\stackrel{d}{=}$  is equivalence in distribution.)

**Definition 1** (Reparameterized Family). Let  $\varphi$  be some d-variate distribution. Then,  $q_{\lambda}$  that can be equivalently represented as

$$\mathbf{z} \sim q_{\lambda} \quad \Leftrightarrow \quad \mathbf{z} \stackrel{\mathrm{d}}{=} \mathcal{T}_{\lambda}(\mathbf{u}); \quad \mathbf{u} \sim \varphi,$$
 is said to be part of a reparameterized family generated by the base distribution  $\varphi$  and the reparameterization function  $\mathcal{T}_{\lambda}$ .

**Definition 2** (Location-Scale Reparameterization Function).  $\mathcal{T}_{\lambda}: \mathbb{R}^d \to \mathbb{R}^d$  defined as

$$\mathcal{T}_{\lambda}(u) \triangleq Cu + m$$

with  $\lambda$  containing the parameters for forming the location  $m \in \mathbb{R}^d$  and scale  $C = C(\lambda) \in \mathbb{R}^{d \times d}$  is called the location-scale reparameterization function.

The location-scale family enables detailed theoretical analysis, as demonstrated by (Domke, 2019, 2020; Fujisawa & Sato, 2021; Kim *et al.*, 2023), and includes the most widely used variational families such as the Student-t, elliptical, and Gaussian families (Titsias & Lázaro-Gredilla, 2014).

**Handling Constrained Support** For common choices of the base distribution  $\varphi$ , the support of  $q_{\lambda}$  is the whole  $\mathbb{R}^d$ . Therefore, special treatment is needed when the support of  $\pi$  is constrained. Kucukelbir *et al.* (2017) proposed to handle this by applying diffeomorphic transformation denoted with  $\psi$ , often called *bjectors* (Dillon *et al.*, 2017; Fjelde *et al.*, 2020; Leger, 2023), to  $q_{\lambda}$  such that

$$\boldsymbol{\zeta} \sim q_{\psi,\lambda} \qquad \Leftrightarrow \qquad \boldsymbol{\zeta} \stackrel{d}{=} \psi^{-1}(\mathbf{z}); \quad \mathbf{z} \sim q_{\lambda},$$

such that the support of  $q_{\psi,\lambda}$  matches that of  $\pi$ . For example, when the support of  $\pi$  is  $\mathbb{R}_+$ , one can choose  $\psi^{-1} = \exp$ . This approach, known as automatic differentiation VI (ADVI), is now standard in most modern PPLs.

Why focus on posteriors with unconstrained supports? When bijectors are used, the entropy of  $q_{\lambda}$ ,  $\mathbb{H}(q_{\lambda})$ , needs to be adjusted by the Jacobian of  $\psi$  (Kucukelbir *et al.*, 2017),  $J_{\phi^{-1}}$ . However, applying the transformation to  $\pi$  instead of  $q_{\lambda}$  is mathematically equivalent and more convenient. In fact, bijectors can be automatically incorporated into our notation by implicitly setting

$$p(\mathbf{x} \mid \mathbf{z}) = \widetilde{p}(\mathbf{x} \mid \psi^{-1}(\mathbf{z}))$$
 and  $p(\mathbf{z}) = \widetilde{p}(\psi^{-1}(\mathbf{z})) |\mathbf{J}_{\psi^{-1}}(\mathbf{z})|$ ,

such that  $\widetilde{\pi}(\zeta) \propto \widetilde{p}(x \mid \zeta) \widetilde{p}(\zeta)$ , where  $\widetilde{\pi}$  is the constrained posterior that we are actually interested in. Therefore, our setup in Section 2.1, where the domain of z is taken to be the unconstrained  $\mathbb{R}^d$ , already encompasses constrained posteriors through ADVI.

Lastly, we impose light assumptions on the base distribution  $\varphi$ , which are already satisfied by most variational families used in practice. (*i.i.d.*: independently and identically distributed.)

**Assumption 1** (Base Distribution).  $\varphi$  is a d-variate distribution such that  $u \sim \varphi$  and  $u = (u_1, \dots, u_d)$  with i.i.d. components. Furthermore,  $\varphi$  is (i) symmetric and standardized such that  $\mathbb{E} u_i = 0$ ,  $\mathbb{E} u_i^2 = 1$ ,  $\mathbb{E} u_i^3 = 0$ , and (ii) has finite kurtosis  $\mathbb{E} u_i^4 = k_{\varphi} < \infty$ .

The assumptions on the variational family we will use throughout this work are collectively summarized in the following assumption:

**Assumption 2.** The variational family is the location-scale family formed by Definitions 1 and 2 with the base distribution  $\varphi$  satisfying Assumption 1.

#### 2.3 Scale Parameterizations

For the "scale" matrix  $C(\lambda)$  in the location-scale family, any parameterization that results in a positive-definite covariance  $CC^{\mathsf{T}} \in \mathbb{S}^{d}_{++}$  is valid. However, for the ELBO to ever be convex, the entropy  $\mathbb{H}(q_{\lambda})$  must be convex, which requires the mapping  $\lambda \mapsto CC^{\mathsf{T}}$  to be convex. To ensure this, we restrict C to (lower) triangular matrices with strictly positive eigenvalues, essentially, Cholesky factors. This leaves two of the most common parameterizations:

#### **Definition 3 (Mean-Field Family.).**

$$C = D_{\phi}(s)$$

where the d elements of s forms the diagonal and  $\lambda \in \Lambda$  such that

$$\Lambda = \{ (\boldsymbol{m}, \boldsymbol{s}) \mid \boldsymbol{m} \in \mathbb{R}^d, \boldsymbol{s} \in \mathcal{S} \}.$$

#### Definition 4 (Full-Rank Cholesky Family).

$$C=D_{\phi}\left( s\right) +L,$$

where the d elements of s forms the diagonal, L is d-by-d strictly lower triangular, and  $\lambda \in \Lambda$  such that

$$\Lambda = \{ (\boldsymbol{m}, \boldsymbol{s}, \boldsymbol{L}) \mid \boldsymbol{m} \in \mathbb{R}^d, \boldsymbol{s} \in \mathcal{S}, \text{vec}(\boldsymbol{L}) \in \mathbb{R}^{(d+1)d/2} \}.$$

Here, S is discussed in the next paragraph,  $\mathbf{D}_{\phi}(\mathbf{s}) \in \mathbb{R}^{d \times d}$  is a diagonal matrix such that  $\mathbf{D}_{\phi}(\mathbf{s}) \triangleq \text{diag}(\phi(\mathbf{s})) = \text{diag}(\phi(s_1), \dots, \phi(s_d))$ , and  $\phi$  is a function we call a *diagonal conditioner*.

**Linear v.s. Nonlinear Parameterizations** When the diagonal conditioner is a linear function  $\phi(x) = x$ , we say that the covariance parameterization is *linear*. In this case, to ensure that C is a Cholesky factor, the domain of s is set as  $S = \mathbb{R}^d_+$ . On the other hand, by choosing a nonlinear conditioner  $\phi : \mathbb{R} \to \mathbb{R}_+$ , we can make the domain of s to be the unconstrained  $S = \mathbb{R}^d$ . Because of this, nonlinear conditioners such as the softplus  $(x) \triangleq \log(1 + \exp(x))$  (Dugas *et al.*, 2000) are frequently used in practice, especially for mean-field. (See Table 1 by Kim *et al.*, 2023).

#### 2.4 Problem Structure of Black-Box Variational Inference

Exclusive KL minimization VI is fundamentally a composite (regularized) optimization problem

$$F(\lambda) = f(\lambda) + h(\lambda),$$
 (ELBO)

where  $f(\lambda) \triangleq \mathbb{E}_{\mathbf{z} \sim q_{\lambda}} \ell(\mathbf{z})$  is the *energy term*,  $\ell(\mathbf{z}) \triangleq -\log p(\mathbf{z}, \mathbf{x})$  is the negative joint log-likelihood, and  $h(\lambda) \triangleq -\mathbb{H}(q_{\psi,\lambda})$  is the *entropic regularizer*. From here, BBVI introduces more structure.

An illustration of the taxonomy is shown in Figure 1. In particular, BBVI has an *infinite sum* structure (IS). That is, it cannot be represented as a sum of finite subcomponents as in ERM. Furthermore.

$$F(\lambda) = \mathbb{E}_{\boldsymbol{u} \sim \varphi} f(\lambda; \boldsymbol{u}) + h(\lambda) \qquad (CP \cap IS)$$
$$= \mathbb{E}_{\boldsymbol{u} \sim \varphi} \ell(\mathcal{T}_{\lambda}(\boldsymbol{u})) + h(\lambda), \qquad (CP \cap IS \cap RP)$$

where  $f(\lambda; \mathbf{u}) \triangleq \ell(\mathcal{T}_{\lambda}(\mathbf{u}))$ .

**Theoretical Challenges** The structure of BBVI has multiple challenges that have hindered its theoretical analysis: (i) the stochasticity of the Jacobian of  $\mathcal{F}$  and (ii) The infinite sum structure.

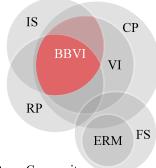
For Item (i), we can see that in

$$\nabla_{\lambda} \ell \left( \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right) \right) = \frac{\partial \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right)}{\partial \lambda} \nabla \ell \left( \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right) \right) = \frac{\partial \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right)}{\partial \lambda} g \left( \boldsymbol{\lambda} ; \boldsymbol{u} \right),$$

where  $g(\lambda; \boldsymbol{u}) \triangleq (\nabla \ell \circ \mathcal{T}_{\lambda})(\boldsymbol{u})$ , both the Jacobian of  $\mathcal{T}_{\lambda}$  and the gradient of the log-likelihood, g, depend on the randomness  $\boldsymbol{u}$ . Effectively decoupling the two is a major challenge to analyzing the properties of the ELBO and its gradient estimators (Domke, 2019, 2020).

For Item (ii), the problem is that recent analyses of SGD (Garrigos & Gower, 2023; Gower *et al.*, 2019; Nguyen *et al.*, 2018; Vaswani *et al.*, 2019) have increasingly been relying on the assumption that  $f(\lambda; \boldsymbol{u})$  is smooth for all  $\boldsymbol{u}$  such that

 $\|\nabla_{\lambda} f(\lambda; \boldsymbol{u}) - \nabla_{\lambda} f(\lambda'; \boldsymbol{u})\| \le L \|\lambda - \lambda'\|$  for some  $L < \infty$ . This is sensible if the support of  $\boldsymbol{u}$  is bounded, which is true for the ERM setting but not for the class of infinite sum (IS) problems. Previous works circumvented this issue by assuming (i) that the support of  $\boldsymbol{u}$  is bounded (Fujisawa & Sato, 2021) which implicitly changes the variational family, or (ii) that the gradient  $\nabla f$  is bounded by a constant (Buchholz *et al.*, 2018; Liu & Owen, 2021) which contradicts strong convexity (Nguyen *et al.*, 2018).



CP Composite

IS Infinite Sum

RP Reparameterized

FS Finite Sum

ERM Empirical Risk Minimization

Figure 1: **Taxonomy of variational inference**. Within BBVI, this work only considers the reparameterization gradient (BBVI  $\cap$  RP, shown in **dark red**). This leaves out BBVI with the score gradient (BBVI  $\setminus$  RP, shown in **light red**). The set VI  $\cap$  FS includes sparse variational Gaussian processes (Titsias, 2009), while the remaining set VI  $\setminus$  (FS  $\cup$  IS  $\cup$  RP) includes coordinate ascent VI (Blei *et al.*, 2017).

#### 3 The Evidence Lower Bound Under Nonlinear Scale Parameterizations

Under the linear parameterization ( $\phi(x) = x$ ), the properties of the ELBO, such as smoothness and convexity, have been previously analyzed by Challis & Barber (2013); Domke (2020); Titsias & Lázaro-Gredilla (2014). We generalize these results to nonlinear conditioners.

#### 3.1 Technical Assumptions

Let  $g_i(\lambda; \mathbf{u})$  be the *i*th coordinate of  $\mathbf{g}(\lambda; \mathbf{u})$  and recall that  $u_i$  denote the *i*th element of  $\mathbf{u}$ . Establishing convexity and smoothness of the ELBO under nonlinear parameterizations depends on a pair of necessary and sufficient assumptions. To establish smoothness:

**Assumption 3.** The gradient of  $\ell$  under reparameterization, g, satisfies

$$|\mathbb{E}g_i(\lambda; \boldsymbol{u}) u_i \phi''(s_i)| \leq L_s$$

for every coordinate i = 1, ... d, any  $\lambda \in \Lambda$ , and some  $0 < L_s < \infty$ .

Here,  $\phi''$  is the second derivative of  $\phi$ . The next one is required to establish convexity:

**Assumption 4.** The gradient of  $\ell$  under reparameterization, g, satisfies

$$\mathbb{E}g_i(\lambda; \mathbf{u}) u_i \geq 0$$

for every coordinate i = 1, ... d.

Intuitively, these assumption control how much  $\nabla \ell$  and  $\mathcal{T}_{\lambda}$  rotate the randomness  $\boldsymbol{u}$ . (Notice that the assumptions are closely related to the matrix  $\operatorname{Cov}(g(\lambda;\boldsymbol{u}),\boldsymbol{u})$ , the covariance between g and  $\boldsymbol{u}$ .) However, the peculiar aspect of these assumptions is that they are not implied by the convexity and smoothness of  $\ell$ . Especially, Assumption 3 strongly depends on the internals of  $\nabla \ell$ .

#### 3.2 Smoothness of the Entropy

Under the linear parameterization, Domke (2020) has previously shown that the entropic regularizer term h is not smooth. This fact immediately implies the ELBO is not smooth. However, certain nonlinear conditioners do result in a smooth regularizer.

**Lemma 1.** If the diagonal conditioner  $\phi$  is  $L_h$ -log-smooth, then the entropic regularizer  $h(\lambda)$  is  $L_h$ -smooth.

*Proof.* See the *full proof* in page 24.

**Example 1.** The following diagonal conditioners result in a smooth entropic regularizer:

- 1. Let  $\phi(x) = \text{softplus}(x)$ . Then, h is  $L_h$ -smooth with  $L_h \approx 0.167096$ .
- 2. Let  $\phi(x) = \exp(x)$ . Then, h is  $L_h$ -smooth for arbitrarily small  $L_h$ .

This might initially suggest that diagonal conditioners are a promising way of making the ELBO globally smooth. Unfortunately, the properties of the energy, f, change unfavorably.

#### 3.3 Smoothness of the Energy

**Inapplicability of Existing Proof Strategy** Previously, Domke (2020, Theorem 1) have proven that the energy is smooth when  $\phi$  is linear. The key step was to use Bessel's inequality based on the observation that the partial derivatives of the reparameterization function  $\mathcal{T}$  form unit bases in expectation. That is,

$$\mathbb{E}\left\langle \frac{\partial \mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)}{\partial \lambda_{i}}, \frac{\partial \mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)}{\partial \lambda_{j}} \right\rangle = \mathbb{1}_{i=j},$$

where  $\mathbb{1}_{i=j}$  is an indicator function that is 1 only when i=j and 0 otherwise.

Unfortunately, when  $\phi$  is nonlinear, the partial derivatives  $\partial \mathcal{F}_{\lambda}(u)/\partial \lambda_i$  for  $i=1,\ldots,p$  no longer form unit bases: while they are still orthogonal in expectation, the *lengths* change nonlinearly depending on  $\lambda$ . This leaves Bessel's inequality inapplicable. To circumvent this challenge, we establish a replacement for Bessel's inequality:

**Lemma 2.** Let **H** be a  $n \times n$  symmetric random matrix, where it is bounded as  $\|\mathbf{H}\|_2 \leq L < \infty$  almost surely. Also, let **J** be an  $m \times n$  random matrix such that  $\|\mathbb{E}\mathbf{J}^{\mathsf{T}}\mathbf{J}\|_2 < \infty$ . Then,

$$\|\mathbb{E} \mathbf{J}^{\top} \mathbf{H} \mathbf{J}\|_{2} \leq L \|\mathbb{E} \mathbf{J}^{\top} \mathbf{J}\|_{2}.$$

*Proof.* See the *full proof* in page 24.

**Remark 1.** By assuming that the joint log-likelihood  $\ell$  is smooth and twice-differentiable, we retrieve Theorem 1 of Domke (2020) by setting J to be the Jacobian of  $\mathcal{T}$ , and H to be the Hessian of  $\ell$  under reparameterization.

**Remark 2.** While our reparameterization function's partial derivatives still form orthogonal bases, they need not be; unlike Bessel's inequality, Lemma 2 does not require this. This implies that Lemma 2 is a strategy more general than Bessel's inequality.

Equipped with Lemma 2, we present our main result on smoothness:

**Theorem 1.** Let  $\ell$  be  $L_{\ell}$ -smooth and twice differentiable. Then, the following results hold:

- (i) If  $\phi$  is linear, the energy f is  $L_{\ell}$ -smooth.
- (ii) If  $\phi$  is 1-Lipschitz, the energy  $\ell$  is  $(L_{\ell} + L_s)$ -smooth if and only if Assumption 3 holds.

*Proof.* See the *full proof* in page 27.

Combined with Lemma 1, this directly implies that the overall ELBO is smooth.

**Corollary 1** (Smoothness of the ELBO). Let  $\ell$  be  $L_{\ell}$ -smooth and Assumption 3 hold. Furthermore, let the diagonal conditioner be 1-Lipschitz continuous, and  $L_{\phi}$ -log-smooth. Then, the ELBO is  $(L_{\ell} + L_{s} + L_{\phi})$ -smooth.

The increase of the smoothness constant implies that we need to use a smaller stepsize to guarantee convergence when using a nonlinear  $\phi$ . Furthermore, even on simple *L*-smooth examples Assumption 3 may not hold:

**Example 2.** Let  $\ell(z) = (1/2) z^{\mathsf{T}} A z$  and the diagonal conditioner be  $\phi(x) = \text{softplus}(x)$ . Then,

- (i) if A is dense and the variational family is the mean-field family or
- (ii) if A is diagonal and the variational family is the Cholesky family,

Assumption 3 holds with  $L_s \approx 0.26034 (\max_{i=1,...,d} A_{ii})$ .

(iii) If A is dense but the Cholesky family is used, Assumption 3 does not hold.

*Proof.* See the *full proof* in page 29.

Example 2 illustrates that establishing the smoothness of the energy becomes non-trivial under nonlinear parameterizations. Even when smoothness does hold, the increased smoothness constant implies that BBVI will be less robust to initialization and stepsizes. Furthermore, in the next section, we will show a much more grave problem: nonlinear parameterizations may affect the convergence *rate*.

#### 3.4 Convexity of the Energy

The convexity of the ELBO under linear parameterizations has first been established by Titsias & Lázaro-Gredilla (2014, Proposition 1) and Domke (2020, Theorem 9). In particular, Domke (2020) show that, when  $\phi$  is linear, if  $\ell$  is  $\mu$ -strongly convex, the energy is also  $\mu$ -strongly convex. However, when using a nonlinear  $\phi$  with a co-domain of  $\mathbb{R}_+$ , which is the whole point of using a nonlinear conditioner, strong convexity of  $\ell$  never transfers to f.

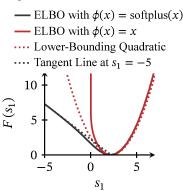


Figure 2: **Optimization landscape resulting from different**  $\phi$  **on a strongly-convex**  $\ell$ .  $\ell$  is the counter-example of Proposition 1 Item (ii).  $\phi(x) = x$  preserves strong convexity as shown by the lower-bounding quadratic (red dotted line ....).  $\phi$  = softplus violates the first-order condition of convexity (black dotted line ....).

**Theorem 2.** Let  $\ell$  be  $\mu$ -strongly convex. Then, we have the following:

- (i) If  $\phi$  is linear, the energy f is  $\mu$ -strongly convex.
- (ii) If  $\phi$  is convex, the energy f is convex if and only if Assumption 4 holds.
- (iii) If  $\phi$  is such that  $\phi \in C^1(\mathbb{R}, \mathbb{R}_+)$ , the energy f is not strongly convex.

*Proof.* See the *full proof* in page 33.

The following proposition provides some conditions for Assumption 4 to hold or not hold.

**Proposition 1.** We have the following:

- (i) If  $\ell$  is convex, then for the mean-field family, Assumption 4 holds.
- (ii) For the Cholesky family, there exists a convex  $\ell$  where Assumption 4 does not hold.

*Proof.* See the *full proof* in page 31.

For any continuous, differentiable nonlinear conditioner that maps only to non-negative reals, the strong convexity of  $\ell$  does lead to a strongly-convex ELBO. This phenomenon is visualized in Figure 2. The loss surface becomes flat near the optimal scale parameter. This problem becomes more noticeable as the optimal scale becomes smaller.

Nonlinear conditioners are suboptimal. As the dataset grows, Bayesian posteriors are known to "contract" as characterized by the Bernstein-von Mises theorem (van der Vaart, 1998). That is, the posterior variance becomes close to 0. This behavior also applies to misspecified variational posteriors as shown by Wang & Blei (2019). Thus, for large datasets, nonlinear conditioners mostly operate in the regime where they are suboptimal (locally less strongly convex). But linear conditioners result in a non-smooth entropy (Domke, 2020). This dilemma originally motivated Domke to consider proximal SGD, which we analyze in Section 4.2.

#### **Convergence Analysis of Black-Box Variational Inference**

#### **Black-Box Variational Inference**

BBVI with SGD repeats the steps:

$$\lambda_{t+1} = \lambda_t - \gamma_t \left( \widehat{\nabla f} \left( \lambda_t \right) + \nabla h \left( \lambda_t \right) \right), \quad \text{where} \quad \widehat{\nabla f} \left( \lambda_t \right) = \frac{1}{M} \sum_{m=1}^{M} \nabla_{\lambda} \ell \left( \mathcal{T}_{\lambda} \left( \mathbf{u}_m \right) \right)$$
 (1)

with  $u_m \sim \varphi$  is the M-sample reparameterization gradient estimator and  $\gamma_t$  is the stepsize. (See Kucukelbir et al., 2017 for algorithmic details.)

With our results in Section 3 and the results of Khaled & Richtárik (2023); Kim et al. (2023), we obtain a convergence guarantee. To apply the result of Kim et al. (2023), which bounds the gradient variance, we require an additional assumption.

**Assumption 5.** The negative log-likelihood  $\ell_{\text{like}}(z) \triangleq -\log p(x \mid z)$  is  $\mu$ -quadratically growing for all  $z \in \mathbb{R}^d$  such that

 $\frac{\mu}{2} \| \boldsymbol{z} - \bar{\boldsymbol{z}}_{\text{like}} \|_2^2 \le \ell_{\text{like}}(\boldsymbol{z}) - \ell_{\text{like}}^*,$  where  $\bar{\boldsymbol{z}}_{\text{like}}$  is the projection of  $\boldsymbol{z}$  to the set of minimizers of  $\ell_{\text{like}}$ , and  $\ell_{\text{like}}^* = \inf_{\boldsymbol{z} \in \mathbb{R}^d} \ell_{\text{like}}(\boldsymbol{z}).$ 

This assumption is weaker than assuming that the likelihood satisfies the Polyak-Łojasiewicz inequality (Karimi et al., 2016).

**Theorem 3.** Let Assumption 2 hold, the likelihood satisfy Assumption 5, and the assumptions of Corollary 1 hold such that the ELBO F is  $L_F$ -smooth with  $L_F = L_\ell + L_\phi + L_s$ . Then, the iterates generated by BBVI through Equation (1) and the M-sample reparameterization gradient include an  $\epsilon$ -stationary point such that  $\min_{0 \le t \le T-1} \mathbb{E} \| \nabla F(\lambda_t) \|_2 \le \epsilon$  for any  $\epsilon > 0$  if

$$T \ge \mathcal{O}\left(\frac{\left(F\left(\lambda_{0}\right) - F^{*}\right)^{2} L_{F} L_{\ell}^{2} C\left(d, k_{\varphi}\right)}{\mu M \epsilon^{4}}\right)$$

for some fixed stepsize  $\gamma$ , where  $C(d,\varphi) = d + k_{\varphi}$  for the Cholesky family and  $C(d,\varphi) = d + k_{\varphi}$  $2k_{\infty}\sqrt{d+1}$  for the mean-field family.

*Proof.* See the *full proof* in page 35.

**Remark 3.** Finding an  $\epsilon$ -stationary point of the ELBO has an iteration complexity of  $\mathcal{O}\left(dL_{\ell}^{2}\kappa M^{-1}\varepsilon^{-4}\right)$  for the Cholesky family and  $\mathcal{O}\left(\sqrt{d}L_{\ell}^{2}\kappa M^{-1}\varepsilon^{-4}\right)$  for the mean-field family.

#### 4.2 Black-Box Variational Inference with Proximal SGD

**Proximal SGD** For a composite objective F = f + h, proximal SGD repeats the steps:

$$\lambda_{t+1} = \operatorname{prox}_{\gamma_t,h} \left( \lambda_t - \gamma_t \widehat{\nabla f} \left( \lambda_t \right) \right) = \underset{\lambda \in \Lambda}{\operatorname{arg\,min}} \left[ \left\langle \widehat{\nabla f} \left( \lambda_t \right), \lambda \right\rangle + h \left( \lambda \right) + \frac{1}{2\gamma_t} \| \lambda - \lambda_t \|_2^2 \right], \quad (2)$$
 where prox is known as the *proximal* operator and  $\gamma_1, \dots, \gamma_T$  is a stepsize schedule.

In the context of VI, proximal SGD has previously been considered by Altosaar et al. (2018); Diao et al. (2023); Khan et al. (2016, 2015). Their overall focus has been on developing alternative algorithms by generalizing  $\|\lambda - \lambda^*\|$  to other metrics. In contrast, Domke (2020) considered proximal SGD with the regular Euclidean metric  $\|\lambda - \lambda^*\|_2$  for overcoming the non-smoothness of h under

linear parameterizations. Here, we prove the convergence of this scheme and show that it retrieves the fastest known convergence rates in stochastic first-order optimization.

**Proximal Operator for BBVI** In our context, h is the entropy of  $q_{\lambda}$  in the location-scale family. For this, Domke (2020) show that the proximal update for  $s_1, \dots, s_d$ , is

$$\operatorname{prox}_{\gamma_t,h}(s_i) = s_i + \frac{1}{2} \left( \sqrt{s_i^2 + 4\gamma_t} - s_i \right).$$

For other parameters, the proximal operator is the regular gradient descent update in Equation (1).

**Gradient Variance Bound** We first establish a bound on the gradient variance. In ERM, contemporary strategies do this by exploiting the finite sum structure of the objective (Section 2.4). Here, we establish a variance bound for RP estimator that does not rely on the finite sum assumption.

**Lemma 3** (Convex Expected Smoothness). Let  $\ell$  be  $L_{\ell}$ -smooth and  $\mu$ -strongly convex with the variational family satisfying Assumption 2 with the linear parameterization. Then,

$$\mathbb{E}\|\nabla_{\lambda}f(\lambda;\mathbf{u}) - \nabla_{\lambda'}f(\lambda';\mathbf{u})\|_{2}^{2} \leq 2L_{\ell}\kappa C(d,\varphi) B_{f}(\lambda,\lambda')$$

holds, where  $B_f(\lambda, \lambda') \triangleq f(\lambda) - f(\lambda') - \langle \nabla f(\lambda'), \lambda - \lambda' \rangle$  is the Bregman divergence,  $\kappa = L_\ell/\mu$  is the condition number,  $C(d, \varphi) = d + k_\varphi$  for the Cholesky family, and  $C(d, \varphi) = 2k_\varphi\sqrt{d} + 1$  for the mean-field family.

*Proof.* See the *full proof* in page 36.

Furthermore, the gradient variance at the optimum must be bounded:

**Lemma 4** (Domke, 2019; Kim et al., 2023). Let  $\ell$  be  $L_{\ell}$ -smooth with the variational family satisfying Assumption 2 and a 1-Lipschitz diagonal conditioner  $\phi$ . Then, the gradient variance at the optimum  $\lambda^* \in \arg\min_{\lambda \in \Lambda} F(\lambda)$  is bounded as

$$\sigma^{2} \leq \frac{1}{M} C\left(d, \varphi\right) L_{\ell}^{2} \left(\left\|\bar{\boldsymbol{z}} - \boldsymbol{m}^{*}\right\|_{2}^{2} + \left\|\boldsymbol{C}^{*}\right\|_{F}^{2}\right),$$

where  $\bar{z}$  is a stationary point of  $\ell$ ,  $m^*$  and  $C^*$  are the location and scale formed by  $\lambda^*$ , the constants are  $C(d, \varphi) = d + k_{\varphi}$  for the Cholesky family and  $C(d, \varphi) = 2k_{\varphi}\sqrt{d} + 1$  for the mean-field family,  $k_{\varphi}$  is the kurtosis of  $\varphi$  as defined in Assumption 1.

*Proof.* The full-rank case is proven by Domke (2019, Theorem 3), while the mean-field case is a basic corollary of the result by Kim *et al.* (2023, Lemma 2).

**Remark 4.** The dimensional dependence in the complexity of BBVI is transferred from the variance bound in Lemma 4. Unfortunately, for the Cholesky family, this dimensional dependence in the variance bound is tight (Domke, 2019).

**Main Result** With the gradient variance bounds, we now present our complexity result. The proof is identical to Theorem 3.2 by Gower *et al.* (2019), where they use a 2-stage decreasing stepsize schedule: the stepsize is initially held constant and then reduced in a 1/t rate.

**Theorem 4.** Let  $\ell$  be  $L_{\ell}$ -smooth and  $\mu$ -strongly convex. Then, for any  $\epsilon > 0$ , BBVI with proximal SGD in Equation (2), the M-sample reparameterization gradient estimator, a variational family satisfying Assumption 2 with the linear parameterization guarantees  $\mathbb{E}\|\lambda_T - \lambda^*\|_2^2 \le \epsilon$  if

$$\gamma_{t} = \begin{cases} \frac{M}{2L_{\ell}\kappa C(d,\varphi)} & \text{for } t \leq 4T_{\kappa} \\ \frac{2t+1}{(t+1)^{2}\mu} & \text{for } t > 4T_{\kappa}, \end{cases} \qquad T \geq \max\left(\frac{8\sigma^{2}}{\mu^{2}\epsilon} + \frac{4T_{\kappa}\|\lambda_{0} - \lambda^{*}\|_{2}}{e\sqrt{\epsilon}}, 4T_{\kappa}\right)$$

where  $\sigma^2$  is defined in Lemma 4,  $T_{\kappa} = [\kappa^2 C(d, \varphi) M^{-1}]$ ,  $\kappa = L_{\ell}/\mu$  is the condition number,  $\ell$  is Euler's constant,  $\ell$  =  $\ell$  arg  $\min_{\lambda \in \Lambda} F(\lambda)$ ,  $C(d, \varphi) = d + k_{\varphi}$  for the Cholesky family, and  $C(d, \varphi) = 2k_{\varphi}\sqrt{d} + 1$  for the mean-field family.

*Proof.* See the *full proof* in page 38.

**Remark 5.** BBVI with proximal SGD on  $\mu$ -strongly convex and  $L_{\ell}$ -smooth  $\ell$  has a complexity  $\mathcal{O}\left(\kappa^2 dM^{-1} \, \epsilon^{-1}\right)$  for the Cholesky family and  $\mathcal{O}\left(\kappa^2 \sqrt{d}M^{-1} \, \epsilon^{-1}\right)$  for the mean-field family.

**Remark 6.** We also provide a similar result with a fixed stepsize in Theorem 7 of Appendix F.3.2. In this case, the complexity is  $\mathcal{O}\left(\kappa^2 dM^{-1} \epsilon^{-1} \log \epsilon^{-1}\right)$  for the Cholesky family and  $\mathcal{O}\left(\kappa^2 \sqrt{d}M^{-1} \epsilon^{-1} \log \epsilon^{-1}\right)$  for the mean-field family.

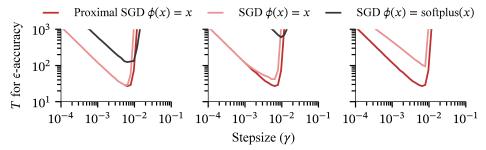


Figure 3: Stepsize versus the number of iterations for vanilla SGD and proximal SGD to achieve  $D_{KL}(q_{\lambda}, \pi) \leq \varepsilon = 1$  under different initializations for Gaussian posteriors. The initializations  $C(\lambda_0)$  are I,  $10^{-3}I$ ,  $10^{-5}I$  from left to right, respectively. The average suboptimality at iteration t was estimated from 10 independent runs. For each run, the target posterior was a 10-dimensional Gaussian with a covariance with a condition number  $\kappa = 10$  and a smoothness of L = 100.

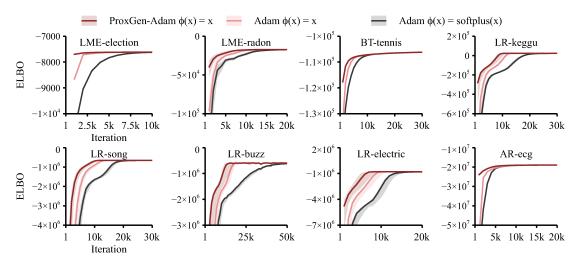


Figure 4: Comparison of BBVI convergence speed (ELBO v.s. Iteration) of different optimization algorithms. The error bands are the 80% quantiles estimated from 20 (10 for AR-eeg) independent replications. The results shown used a base stepsize of  $\gamma = 10^{-3}$ , while the initial point was  $m_0 = 0$ ,  $C_0 = I$ . Details on the setup can be found in the text of Section 5.2 and Appendix G.

#### 5 Experiments

#### 5.1 Synthetic Problem

**Setup** We first compare proximal SGD against vanilla SGD with linear and nonlinear parameterizations on a synthetic problem, which is log-smooth, strongly log-concave, and the exact solution is known. While a similar experiment was already conducted by Domke (2020), here we include nonlinear parameterizations, which were not originally considered. We run all algorithms with a fixed stepsize to infer a multivariate Gaussian with a full-rank covariance matrix. The variational approximation is a full-rank Gaussian formed by  $\varphi = \mathcal{N}(0,1)$  and the Cholesky parameterization.

**Results** The results are shown in Figure 3. Proximal SGD is clearly the most robust against initialization. Also, SGD with the nonlinear parameterization  $\phi(x) = \text{softplus}(x)$  is much slower to converge under all initializations. This confirms that linear parameterizations are indeed superior for both robustness against initializations and convergence speed.

#### 5.2 Realistic Problems

**Setup** We now evaluate proximal SGD on realistic problems. In practice, Adam (Kingma & Ba, 2015) is observed to be robust against stepsize choices (Zhang *et al.*, 2019). The reason why Adam performs well on non-smooth, non-convex problems is still under investigation (Kunstner *et al.*, 2023; Reddi *et al.*, 2023; Zhang *et al.*, 2022). Nonetheless, to compare fairly against Adam, we implement a recently proposed variant of proximal SGD called ProxGen (Yun *et al.*, 2021), which

includes an Adam-like update rule. The probabilitic models and datasets are fully described in Appendix G. We implement these models and BBVI on top of the Turing (Ge et al., 2018) probabilistic programming framework. Due to the size of these datasets, we implement doubly stochastic subsampling (Titsias & Lázaro-Gredilla, 2014) with a batch size of B = 100 (B = 500 for BT-tennis) with M = 10 Monte Carlo samples. For batch subsampling, we implement random-reshuffling, which is faster than independent subsampling both empirically (Bottou, 2009) and theoretically (Ahn et al., 2020; Haochen & Sra, 2019; Mishchenko et al., 2020; Nagaraj et al., 2019). We also observe that doubly stochastic BBVI benefits from reshuffling, but leave a detailed investigation to future works.

**Results** Representative results are shown in Figure 4, with additional results in Appendix H. Both ProxGen-Adam and Adam with linear parameterizations converge faster than Adam with nonlinear parameterization. Furthermore, for the case of election and buzz, Adam with the nonlinear parameterization converges much slower than the alternatives. When using linear parameterizations, ProxGen-Adam appears to be generally faster than Adam. We note, however, that due to the difference in the update rule between ProxGen-Adam and Adam, proximal operators alone might not fully explain the performance difference. Nevertheless, the results of our experiment do conclusively suggest that linear parameterizations are superior.

#### 6 Discussions

**Conclusions** In this work, we have proven the convergence of BBVI. Our assumptions encompass implementations that are actually used in practice, and our theoretical analysis revealed limitations in some of the popular design choices (mainly the use of nonlinear conditioners). To resolve this issue, we re-evaluated the utility of proximal SGD both theoretically and practically, where it achieved the strongest theoretical guarantees in stochastic first-order optimization.

**Related Works** To prove the convergence of BBVI, early works have *a-priori* "assumed" the regularity of the ELBO and the gradient estimator (Alquier & Ridgway, 2020; Buchholz *et al.*, 2018; Khan *et al.*, 2016, 2015; Liu & Owen, 2021; Regier *et al.*, 2017). Towards a more rigorous understanding, Domke (2019); Fan *et al.* (2015); Kim *et al.* (2023); Xu *et al.* (2019) studied the reparameterization gradient, Xu & Campbell (2022) studied the asymptotics of the ELBO, Challis & Barber (2013); Domke (2020); Titsias & Lázaro-Gredilla (2014) established convexity, and Domke (2020) established smoothness. On the other hand, Bhatia *et al.* (2022); Hoffman & Ma (2020) established rigorous convergence guarantees by considering simplified variant of BBVI where only the scale is optimized, and Fujisawa & Sato (2021) assumed that the support of  $\varphi$  is bounded almost surely. Meanwhile, under similar assumptions to ours, Diao *et al.* (2023); Lambert *et al.* (2022) recently established convergence guarantees for proximal SGD BBVI with a Bures-Wasserstein metric. Their computational properties differ from BBVI as they require Hessian evaluations. Also, understanding BBVI, which is VI with a Euclidean metric, is an important problem due to its practical relevance.

**Limitations** Our work has multiple limitations: (i) Our results are restricted to the location-scale family, (ii) the reparameterization gradient, and (iii) smooth joint log-likelihoods. However, the location-scale family with the reparameterization gradient is the most widely used combination in practice, and replacing the smoothness assumption is an active area of research in stochastic optimization. For our results on proximal SGD, we further assume that the joint log-likelihood is  $\mu$ -strongly convex (equivalently strongly log-concave posteriors). It is unclear how to extend the guarantees to only smooth but non-log-concave joint log-likelihoods.

**Open Problems** Although we have proven that the mean-field dimensional family has a dimension dependence of  $\mathcal{O}(\sqrt{d})$ , empirical results suggest room for improvement (Kim *et al.*, 2023). Therefore, we pose the following conjecture:

**Conjecture 1.** Under mild assumptions, BBVI for the mean-field variational family converges with only logarithmic dimensional dependence or no explicit dimensional dependence at all.

This would put mean-field BBVI in a regime clearly faster than approximate MCMC (Freund *et al.*, 2022). Also, it is unknown whether the  $\mathcal{O}(\kappa^2)$  condition number dependence dependence is tight. In fact, for proximal SGD BBVI in Bures-Wasserstien space, Diao *et al.* (2023) report a dependence of  $\mathcal{O}(\kappa)$ . Lastly, it would be interesting to see whether natural gradient VI (NGVI; Amari, 1998; Khan & Lin, 2017) can achieve similar convergence guarantees. While it is empirically known that NGVI often converges faster (Lin *et al.*, 2019), theoretical evidence has yet to follow.

#### Acknowledgments and Disclosure of Funding

The authors would like to thank Justin Domke for discussions on the concurrent results, Javier Burroni for pointing out a mistake in the earlier version of this work, and the anonymous reviewers for their constructive comments.

K. Kim and J. R. Gardner were funded by the National Science Foundation Award [IIS-2145644], while Y.-A. Ma was funded by the National Science Foundation Grants [NSF-SCALE MoDL-2134209] and [NSF-CCF-2112665 (TILOS)], the U.S. Department Of Energy, Office of Science, and the Facebook Research award.

#### References

- Agrawal, Abhinav, & Domke, Justin. 2021. Amortized Variational Inference for Simple Hierarchical Models. *Pages 21388–21399 of: Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc. (page 2)
- Agrawal, Abhinav, Sheldon, Daniel R, & Domke, Justin. 2020. Advances in Black-Box VI: Normalizing Flows, Importance Weighting, and Optimization. *Pages 17358–17369 of: Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc. (page 1)
- Ahn, Kwangjun, Yun, Chulhee, & Sra, Suvrit. 2020. SGD with Shuffling: Optimal Rates without Component Convexity and Large Epoch Requirements. *Pages 17526–17535 of: Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc. (page 10)
- Alquier, Pierre, & Ridgway, James. 2020. Concentration of Tempered Posteriors and of Their Variational Approximations. *The Annals of Statistics*, **48**(3), 1475–1497. (page 10)
- Altosaar, Jaan, Ranganath, Rajesh, & Blei, David. 2018. Proximity Variational Inference. *Pages* 1961–1969 of: Proceedings of the International Conference on Artificial Intelligence and Statistics. PMLR, vol. 84. JMLR. (page 7)
- Amari, Shun-ichi. 1998. Natural Gradient Works Efficiently in Learning. *Neural Computation*, **10**(2), 251–276. (page 10)
- Bertin-Mahieux, Thierry, Ellis, Daniel P.W., Whitman, Brian, & Lamere, Paul. 2011. The Million Song Dataset. *In: Proceedings of the International Conference on Music Information*. (page 40)
- Bhatia, Kush, Kuang, Nikki Lijing, Ma, Yi-An, & Wang, Yixin. 2022 (July). *Statistical and Computational Trade-Offs in Variational Inference: A Case Study in Inferential Model Selection*. arXiv Preprint arXiv:2207.11208. arXiv. (pages 1, 10)
- Bingham, Eli, Chen, Jonathan P., Jankowiak, Martin, Obermeyer, Fritz, Pradhan, Neeraj, Karaletsos, Theofanis, Singh, Rohit, Szerlip, Paul, Horsfall, Paul, & Goodman, Noah D. 2019. Pyro: Deep Universal Probabilistic Programming. *Journal of Machine Learning Research*, **20**(28), 1–6. (page 1)
- Blei, David M., Kucukelbir, Alp, & McAuliffe, Jon D. 2017. Variational Inference: A Review for Statisticians. *Journal of the American Statistical Association*, **112**(518), 859–877. (pages 2, 4)
- Bottou, Léon. 1999. On-Line Learning and Stochastic Approximations. *Pages 9–42 of: On-Line Learning in Neural Networks*, first edn. Cambridge University Press. (page 2)
- Bottou, Léon. 2009. Curiously Fast Convergence of Some Stochastic Gradient Descent Algorithms. (page 10)
- Buchholz, Alexander, Wenzel, Florian, & Mandt, Stephan. 2018. Quasi-Monte Carlo Variational Inference. *Pages 668–677 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 80. JMLR. (pages 4, 10)
- Carpenter, Bob, Gelman, Andrew, Hoffman, Matthew D., Lee, Daniel, Goodrich, Ben, Betancourt, Michael, Brubaker, Marcus, Guo, Jiqiang, Li, Peter, & Riddell, Allen. 2017. Stan: A Probabilistic Programming Language. *Journal of Statistical Software*, **76**(1). (page 1)
- Carvalho, Carlos M., Polson, Nicholas G., & Scott, James G. 2009. Handling Sparsity via the Horseshoe. *Pages 73–80 of: Proceedings of the International Conference on Artificial Intelligence and Statistics*. PMLR, vol. 5. JMLR. (page 41)
- Carvalho, Carlos M., Polson, Nicholas G., & Scott, James G. 2010. The Horseshoe Estimator for Sparse Signals. *Biometrika*, **97**(2), 465–480. (page 41)

- Challis, Edward, & Barber, David. 2013. Gaussian Kullback-Leibler Approximate Inference. *Journal of Machine Learning Research*, **14**(68), 2239–2286. (pages 1, 5, 10)
- Christmas, Jacqueline, & Everson, Richard. 2011. Robust Autoregression: Student-T Innovations Using Variational Bayes. *IEEE Transactions on Signal Processing*, **59**(1), 48–57. (page 41)
- Dhaka, Akash Kumar, Catalina, Alejandro, Andersen, Michael R, ns Magnusson, Må, Huggins, Jonathan, & Vehtari, Aki. 2020. Robust, Accurate Stochastic Optimization for Variational Inference. *Pages 10961–10973 of: Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc. (page 1)
- Dhaka, Akash Kumar, Catalina, Alejandro, Welandawe, Manushi, Andersen, Michael R., Huggins, Jonathan, & Vehtari, Aki. 2021. Challenges and Opportunities in High Dimensional Variational Inference. *Pages* 7787–7798 of: Advances in Neural Information Processing Systems, vol. 34. Curran Associates, Inc. (page 1)
- Diao, Michael Ziyang, Balasubramanian, Krishna, Chewi, Sinho, & Salim, Adil. 2023. Forward-Backward Gaussian Variational Inference via JKO in the Bures-Wasserstein Space. *Pages 7960–7991 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 202. JMLR. (pages 7, 10)
- Dieng, Adji Bousso, Tran, Dustin, Ranganath, Rajesh, Paisley, John, & Blei, David. 2017. Variational Inference via χ Upper Bound Minimization. *Pages 2729–2738 of: Advances in Neural Information Processing Systems*, vol. 30. Curran Associates, Inc. (page 2)
- Dillon, Joshua V., Langmore, Ian, Tran, Dustin, Brevdo, Eugene, Vasudevan, Srinivas, Moore, Dave, Patton, Brian, Alemi, Alex, Hoffman, Matt, & Saurous, Rif A. 2017 (Nov.). *TensorFlow Distributions*. arXiv Preprint arXiv:1711.10604. arXiv. (pages 1, 3)
- Domke, Justin. 2019. Provable Gradient Variance Guarantees for Black-Box Variational Inference. *Pages 329–338 of: Advances in Neural Information Processing Systems*, vol. 32. Curran Associates, Inc. (pages 1, 3, 4, 8, 10, 21, 22, 36)
- Domke, Justin. 2020. Provable Smoothness Guarantees for Black-Box Variational Inference. *Pages 2587–2596 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 119. JMLR. (pages 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 20, 25, 27, 33)
- Domke, Justin, Garrigos, Guillaume, & Gower, Robert. 2023. Provable Convergence Guarantees for Black-Box Variational Inference. *In: Advances in Neural Information Processing Systems (to Appear)*. New Orleans, LA, USA: arXiv. (pages 2, 17, 21)
- Dua, Dheeru, & Graff, Casey. 2017. UCI Machine Learning Repository. (page 40)
- Duchi, John, Hazan, Elad, & Singer, Yoram. 2011. Adaptive Subgradient Methods for Online Learning and Stochastic Optimization. *Journal of Machine Learning Research*, **12**(Jul), 2121–2159. (page 20)
- Dugas, Charles, Bengio, Yoshua, Bélisle, François, Nadeau, Claude, & Garcia, René. 2000. Incorporating Second-Order Functional Knowledge for Better Option Pricing. *In: Advances in Neural Information Processing Systems*, vol. 13. MIT Press. (page 4)
- Dwivedi, Raaz, Chen, Yuansi, Wainwright, Martin J., & Yu, Bin. 2019. Log-Concave Sampling: Metropolis-Hastings Algorithms Are Fast. *Journal of Machine Learning Research*, **20**(183), 1–42. (page 2)
- Fan, Kai, Wang, Ziteng, Beck, Jeff, Kwok, James, & Heller, Katherine A. 2015. Fast Second Order Stochastic Backpropagation for Variational Inference. *Pages 1387–1395 of: Advances in Neural Information Processing Systems*, vol. 28. Curran Associates, Inc. (page 10)
- Fjelde, Tor Erlend, Xu, Kai, Tarek, Mohamed, Yalburgi, Sharan, & Ge, Hong. 2020. Bijectors.jl: Flexible Transformations for Probability Distributions. *Pages 1–17 of: Proceedings of The Symposium on Advances in Approximate Bayesian Inference*. PMLR, vol. 118. JMLR. (page 3)
- Freund, Yoav, Ma, Yi-An, & Zhang, Tong. 2022. When Is the Convergence Time of Langevin Algorithms Dimension Independent? A Composite Optimization Viewpoint. *Journal of Machine Learning Research*, **23**(214), 1–32. (page 10)
- Fujisawa, Masahiro, & Sato, Issei. 2021. Multilevel Monte Carlo Variational Inference. *Journal of Machine Learning Research*, **22**(278), 1–44. (pages 3, 4, 10)

- Garrigos, Guillaume, & Gower, Robert M. 2023 (Feb.). Handbook of Convergence Theorems for (Stochastic) Gradient Methods. arXiv Preprint arXiv:2301.11235. arXiv. (pages 2, 4, 21, 37, 38)
- Ge, Hong, Xu, Kai, & Ghahramani, Zoubin. 2018. Turing: A Language for Flexible Probabilistic Inference. *Pages 1682–1690 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 84. JMLR. (pages 1, 10)
- Gelman, Andrew, & Hill, Jennifer. 2007. Data Analysis Using Regression and Multilevel/Hierarchical Models. Analytical Methods for Social Research. Cambridge; New York: Cambridge University Press. (page 40)
- Giordano, Ryan, Broderick, Tamara, & Jordan, Michael I. 2018. Covariances, Robustness, and Variational Bayes. *Journal of Machine Learning Research*, **19**(51), 1–49. (page 1)
- Giordano, Ryan, Ingram, Martin, & Broderick, Tamara. 2023 (Apr.). *Black Box Variational Inference with a Deterministic Objective: Faster, More Accurate, and Even More Black Box*. arXiv Preprint arXiv:2304.05527. arXiv. (page 41)
- Goldberger, Ary L., Amaral, Luis A. N., Glass, Leon, Hausdorff, Jeffrey M., Ivanov, Plamen Ch., Mark, Roger G., Mietus, Joseph E., Moody, George B., Peng, Chung-Kang, & Stanley, H. Eugene. 2000. PhysioBank, PhysioToolkit, and PhysioNet: Components of a New Research Resource for Complex Physiologic Signals. *Circulation*, **101**(23). (page 41)
- Gorbunov, Eduard, Hanzely, Filip, & Richtarik, Peter. 2020. A Unified Theory of SGD: Variance Reduction, Sampling, Quantization and Coordinate Descent. *Pages 680–690 of: Proceedings of the International Conference on Artificial Intelligence and Statistics*. PMLR, vol. 108. JMLR. (pages 21, 37)
- Gower, Robert Mansel, Loizou, Nicolas, Qian, Xun, Sailanbayev, Alibek, Shulgin, Egor, & Richtárik, Peter. 2019. SGD: General Analysis and Improved Rates. *Pages 5200–5209 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 97. JMLR. (pages 4, 8, 21, 38)
- Haochen, Jeff, & Sra, Suvrit. 2019. Random Shuffling Beats SGD after Finite Epochs. *Pages 2624–2633 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 97. JMLR. (page 10)
- Hernandez-Lobato, Jose, Li, Yingzhen, Rowland, Mark, Bui, Thang, Hernandez-Lobato, Daniel, & Turner, Richard. 2016. Black-Box Alpha Divergence Minimization. *Pages 1511–1520 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 48. JMLR. (page 2)
- Hoffman, Matthew, & Ma, Yian. 2020. Black-Box Variational Inference as a Parametric Approximation to Langevin Dynamics. *Pages 4324–4341 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 119. JMLR. (pages 1, 10)
- Jager, F., Taddei, A., Moody, G. B., Emdin, M., Antolič, G., Dorn, R., Smrdel, A., Marchesi, C., & Mark, R. G. 2003. Long-Term ST Database: A Reference for the Development and Evaluation of Automated Ischaemia Detectors and for the Study of the Dynamics of Myocardial Ischaemia. *Medical and Biological Engineering and Computing*, 41(2), 172–182. (pages 40, 41)
- Jordan, Michael I., Ghahramani, Zoubin, Jaakkola, Tommi S., & Saul, Lawrence K. 1999. An Introduction to Variational Methods for Graphical Models. *Machine Learning*, **37**(2), 183–233. (pages 2, 3)
- Karimi, Hamed, Nutini, Julie, & Schmidt, Mark. 2016. Linear Convergence of Gradient and Proximal-Gradient Methods under the Polyak-Łojasiewicz Condition. *Pages 795–811 of: Machine Learning and Knowledge Discovery in Databases*. Lecture Notes in Computer Science. Cham: Springer International Publishing. (page 7)
- Kawala, François, Douzal-Chouakria, Ahlame, Gaussier, Eric, & Diemert, Eustache. 2013. Prédictions d'activité Dans Les Réseaux Sociaux En Ligne. Page 16 of: Actes de La Conférence Sur Les Modèles et L'Analyse Des Réseaux : Approches Mathématiques et Informatique. (page 40)
- Khaled, Ahmed, & Richtárik, Peter. 2023. Better Theory for SGD in the Nonconvex World. *Transactions of Machine Learning Research*. (pages 7, 21, 34, 35)
- Khan, Mohammad, & Lin, Wu. 2017. Conjugate-Computation Variational Inference: Converting Variational Inference in Non-Conjugate Models to Inferences in Conjugate Models. *Pages 878–887 of: Proceedings of the International Conference on Artificial Intelligence and Statistics*. PMLR, vol. 54. JMLR. (page 10)

- Khan, Mohammad Emtiyaz, Babanezhad, Reza, Lin, Wu, Schmidt, Mark, & Sugiyama, Masashi. 2016. Faster Stochastic Variational Inference Using Proximal-Gradient Methods with General Divergence Functions. *Pages 319–328 of: Proceedings of the Conference on Uncertainty in Artificial Intelligence*. UAI'16. Arlington, Virginia, USA: AUAI Press. (pages 7, 10)
- Khan, Mohammad Emtiyaz E, Baque, Pierre, Fleuret, François, & Fua, Pascal. 2015. Kullback-Leibler Proximal Variational Inference. *In: Advances in Neural Information Processing Systems*, vol. 28. Curran Associates, Inc. (pages 7, 10)
- Kim, Kyurae, Oh, Jisu, Gardner, Jacob, Dieng, Adji Bousso, & Kim, Hongseok. 2022. Markov Chain Score Ascent: A Unifying Framework of Variational Inference with Markovian Gradients. *Pages 34802–34816 of: Advances in Neural Information Processing Systems*, vol. 35. Curran Associates, Inc. (page 2)
- Kim, Kyurae, Wu, Kaiwen, Oh, Jisu, & Gardner, Jacob R. 2023. Practical and Matching Gradient Variance Bounds for Black-Box Variational Bayesian Inference. *Pages 16853–16876 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 202. Honolulu, HI, USA: JMLR. (pages 1, 3, 4, 7, 8, 10, 21, 22, 34, 35, 36)
- Kingma, Diederik P., & Ba, Jimmy. 2015. Adam: A Method for Stochastic Optimization. *In: Proceedings of the International Conference on Learning Representations*. (pages 2, 9, 20)
- Kingma, Diederik P., & Welling, Max. 2014 (Apr.). Auto-Encoding Variational Bayes. *In: Proceedings of the International Conference on Learning Representations*. (page 2)
- Kucukelbir, Alp, Tran, Dustin, Ranganath, Rajesh, Gelman, Andrew, & Blei, David M. 2017. Automatic Differentiation Variational Inference. *Journal of Machine Learning Research*, **18**(14), 1–45. (pages 1, 3, 7, 21)
- Kunstner, Frederik, Chen, Jacques, Lavington, Jonathan Wilder, & Schmidt, Mark. 2023 (Feb.). Noise Is Not the Main Factor behind the Gap between Sgd and Adam on Transformers, but Sign Descent Might Be. *In: Proceedings of the International Conference on Learning Representations*. (page 9)
- Lambert, Marc, Chewi, Sinho, Bach, Francis, Bonnabel, Silvère, & Rigollet, Philippe. 2022. Variational Inference via Wasserstein Gradient Flows. *Pages 14434–14447 of: Advances in Neural Information Processing Systems*, vol. 35. Curran Associates, Inc. (page 10)
- Leger, Jean-Benoist. 2023 (Jan.). Parametrization Cookbook: A Set of Bijective Parametrizations for Using Machine Learning Methods in Statistical Inference. arXiv Preprint arXiv:2301.08297. arXiv. (page 3)
- Lin, Wu, Khan, Mohammad Emtiyaz, & Schmidt, Mark. 2019. Fast and Simple Natural-Gradient Variational Inference with Mixture of Exponential-Family Approximations. *Pages 3992–4002 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 97. JMLR. (page 10)
- Liu, Sifan, & Owen, Art B. 2021. Quasi-Monte Carlo Quasi-Newton in Variational Bayes. *Journal of Machine Learning Research*, **22**(243), 1–23. (pages 4, 10)
- Magnusson, Måns, Bürkner, Paul, & Vehtari, Aki. 2022 (Nov.). Posteriordb: A Set of Posteriors for Bayesian Inference and Probabilistic Programming. (page 40)
- Mishchenko, Konstantin, Khaled, Ahmed, & Richtarik, Peter. 2020. Random Reshuffling: Simple Analysis with Vast Improvements. *Pages 17309–17320 of: Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc. (page 10)
- Naesseth, Christian, Lindsten, Fredrik, & Blei, David. 2020. Markovian Score Climbing: Variational Inference with KL(p||q). *Pages 15499–15510 of: Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc. (page 2)
- Nagaraj, Dheeraj, Jain, Prateek, & Netrapalli, Praneeth. 2019. SGD without Replacement: Sharper Rates for General Smooth Convex Functions. *Pages 4703–4711 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 97. JMLR. (page 10)
- Nemirovski, A., Juditsky, A., Lan, G., & Shapiro, A. 2009. Robust Stochastic Approximation Approach to Stochastic Programming. *SIAM Journal on Optimization*, **19**(4), 1574–1609. (page 2)
- Nguyen, Lam, Nguyen, Phuong Ha, van Dijk, Marten, Richtarik, Peter, Scheinberg, Katya, & Takac, Martin. 2018. SGD and Hogwild! Convergence without the Bounded Gradients Assumption.

- Pages 3750–3758 of: Proceedings of the International Conference on Machine Learning. PMLR, vol. 80. JMLR. (page 4)
- Patil, Anand, Huard, David, & Fonnesbeck, Christopher. 2010. PyMC: Bayesian Stochastic Modelling in Python. *Journal of Statistical Software*, **35**(4). (page 1)
- Ranganath, Rajesh, Gerrish, Sean, & Blei, David. 2014. Black Box Variational Inference. *Pages 814–822 of: Proceedings of the International Conference on Artificial Intelligence and Statistics*. PMLR, vol. 33. JMLR. (page 1)
- Reddi, Sashank J., Kale, Satyen, & Kumar, Sanjiv. 2023 (May). On the Convergence of Adam and Beyond. *In: Proceedings of the International Conference on Learning Representations*. (page 9)
- Regier, Jeffrey, Jordan, Michael I, & McAuliffe, Jon. 2017. Fast Black-Box Variational Inference through Stochastic Trust-Region Optimization. *Pages 2399–2408 of: Advances in Neural Information Processing Systems*, vol. 30. Curran Associates, Inc. (page 10)
- Robbins, Herbert, & Monro, Sutton. 1951. A Stochastic Approximation Method. The Annals of Mathematical Statistics, 22(3), 400–407. (page 2)
- Roeder, Geoffrey, Wu, Yuhuai, & Duvenaud, David K. 2017. Sticking the Landing: Simple, Lower-Variance Gradient Estimators for Variational Inference. *Pages 6928–6937 of: Advances in Neural Information Processing Systems*, vol. 30. Curran Associates, Inc. (pages 2, 21)
- Shannon, Paul, Markiel, Andrew, Ozier, Owen, Baliga, Nitin S., Wang, Jonathan T., Ramage, Daniel, Amin, Nada, Schwikowski, Benno, & Ideker, Trey. 2003. Cytoscape: A Software Environment for Integrated Models of Biomolecular Interaction Networks. *Genome Research*, **13**(11), 2498–2504. (page 40)
- Titsias, Michalis. 2009. Variational Learning of Inducing Variables in Sparse Gaussian Processes. Pages 567–574 of: Proceedings of the International Conference on Artificial Intelligence and Statistics. PMLR, vol. 5. JMLR. (page 4)
- Titsias, Michalis, & Lázaro-Gredilla, Miguel. 2014. Doubly Stochastic Variational Bayes for Non-Conjugate Inference. *Pages 1971–1979 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 32. JMLR. (pages 1, 2, 3, 5, 6, 10, 21)
- van der Vaart, A. W. 1998. Asymptotic Statistics. First edn. Cambridge University Press. (page 7)
- Vaswani, Sharan, Bach, Francis, & Schmidt, Mark. 2019. Fast and Faster Convergence of SGD for Over-Parameterized Models and an Accelerated Perceptron. *Pages 1195–1204 of: Proceedings of the International Conference on Artificial Intelligence and Statistics*. PMLR, vol. 89. JMLR. (page 4)
- Wang, Yixin, & Blei, David. 2019. Variational Bayes under Model Misspecification. *In: Advances in Neural Information Processing Systems*, vol. 32. Curran Associates, Inc. (page 7)
- Welandawe, Manushi, Andersen, Michael Riis, Vehtari, Aki, & Huggins, Jonathan H. 2022 (Mar.). *Robust, Automated, and Accurate Black-Box Variational Inference*. arXiv Preprint arXiv:2203.15945. arXiv. (page 1)
- Wright, Stephen J., & Recht, Benjamin. 2021. *Optimization for Data Analysis*. New York: Cambridge University Press. (page 21)
- Xu, Ming, Quiroz, Matias, Kohn, Robert, & Sisson, Scott A. 2019. Variance Reduction Properties of the Reparameterization Trick. *Pages 2711–2720 of: Proceedings of the International Conference on Artificial Intelligence and Statistics*. PMLR, vol. 89. JMLR. (page 10)
- Xu, Zuheng, & Campbell, Trevor. 2022. The Computational Asymptotics of Gaussian Variational Inference and the Laplace Approximation. *Statistics and Computing*, **32**(4), 63. (page 10)
- Yao, Yuling, Vehtari, Aki, Simpson, Daniel, & Gelman, Andrew. 2018. Yes, but Did It Work?: Evaluating Variational Inference. *Pages 5581–5590 of: Proceedings of the International Conference on Machine Learning*. PMLR, vol. 80. JMLR. (page 1)
- Yun, Jihun, Lozano, Aurelie C, & Yang, Eunho. 2021. Adaptive Proximal Gradient Methods for Structured Neural Networks. *Pages 24365–24378 of: Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc. (pages 2, 9, 20)
- Zhang, Cheng, Butepage, Judith, Kjellstrom, Hedvig, & Mandt, Stephan. 2019. Advances in Variational Inference. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **41**(8), 2008–2026. (pages 2, 9)

Zhang, Yushun, Chen, Congliang, Shi, Naichen, Sun, Ruoyu, & Luo, Zhi-Quan. 2022. Adam Can Converge without Any Modification on Update Rules. *Pages 28386–28399 of: Advances in Neural Information Processing Systems*, vol. 35. Curran Associates, Inc. (page 9)

# On the Convergence of Black-Box Variational Inference *Appendix*

	Table of Contents				
1	Introduction				
2	Background2.1Black-Box Variational Inference2.2Variational Family				
3	The Evidence Lower Bound Under Nonlinear Scale Parameterizations 3.1 Technical Assumptions				
4	Convergence Analysis of Black-Box Variational Inference 4.1 Black-Box Variational Inference				
5	Experiments 5.1 Synthetic Problem				
6	Discussions				
A	Computational Resources				
В	Nomenclature				
C	Definitions				
D	ProxGen Adam for Black-Box Variational Inference				
E	Detailed Comparison Against Domke et al. (2023)				
F	Proofs F.1 Auxiliary Lemmas F.2 Properties of the Evidence Lower Bound F.2.1 Smoothness F.2.2 Convexity F.3 Convergence of Black-Box Variational Inference F.3.1 Vanilla Black-Box Variational Inference F.3.2 Proximal Black-Box Variational Inference				
G	Details of Experimental Setup				
H	Additional Experimental Results				
	•				

### **A** Computational Resources

Table 1: Computational Resources

Type	Model and Specifications		
System Topology	2 nodes with 2 sockets each with 24 logical threads (total 48 threads)		
Processor	1 Intel Xeon Silver 4310, 2.1 GHz (maximum 3.3 GHz) per socket		
Cache	1.1 MiB L1, 30 MiB L2, and 36 MiB L3		
Memory	250 GiB RAM		
Accelerator	1 NVIDIA RTX A5000 per node, 2 GHZ, 24GB RAM		

Running the experiments took approximately a week.

#### **B** Nomenclature

Symbol	Definition	Description	Section
λ		Variational parameters	2.1
z		Parameters of the target model $\pi$	2.1
$\mathcal{T}_{\pmb{\lambda}}$	$\triangleq C(\lambda) u + m$	location-scale reparameterization function	2.2
ü		Random vector before reparameterization	2.2
arphi		Base distribution of <b>u</b>	2.2
m		Location parameter (part of $\lambda$ )	2.2
$\boldsymbol{C}$		Scale parameter (part of $\lambda$ )	2.2 and 2.3
$\phi \ k_{arphi}$		Diagonal conditioner	2.3
$k_{oldsymbol{arphi}}$		Kurtosis (non-central 4th moment) of <b>u</b>	2.3
$D_{\phi}(s)$		Diagonal of $C$ using the diagonal conditioner $\phi$	2.3
$\hat{m{L}}$		Strictly lower triangular part of <i>C</i>	2.3
S		Elements forming the diagonal of <i>C</i>	2.3
$\ell\left( oldsymbol{z} ight)$	$\triangleq -\log p(\mathbf{z}, \mathbf{x})$	Negative joint likelihood	2.4
$f(\lambda)$	$\triangleq \mathbb{E}_{\mathbf{z} \sim q_{\lambda}} \ell \left( \mathbf{z} \right)$	Energy	2.4
$h(\lambda)$	$\triangleq -\mathbb{H}(q_{\lambda})$	Negative entropy	2.4
$F(\lambda)$	$\triangleq f(\lambda) + h(\lambda)$	Negative ELBO	2.1 and 2.4
$f(\boldsymbol{\lambda}; \boldsymbol{u})$	$\triangleq \ell\left(\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)\right)$	Negative Log-Likelihood under reparameterization	2.4
$g(\lambda; \boldsymbol{u})$	$\triangleq \nabla \ell \left( \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right) \right)$	Gradient of the Log-likelihood under reparameterization	2.4
M		Number of Monte Carlo samples	2.1
$\gamma_t$		Stepsize of (proximal) SGD at iteration t	4.1 and 4.2
$\widehat{ abla f}$		Reparameterization gradient estimator of the energy	4.1

#### **C** Definitions

For completeness, we provide formal definitions for some of the terms we used throughout the paper.

**Definition 5 (Smoothness).** A function  $f: \mathcal{Z} \to \mathbb{R}$  is said to be L-smooth if the inequality

$$\|\nabla f(z) - \nabla f(z')\| \le \|z - z'\|$$

holds for all  $z, z' \in \mathcal{Z}$ .

This assumption, also occasionally called Lipschitz smoothness, restricts the amount the gradient can change for a given distance. When f is twice differentiable, an equivalent condition is the Hessian to be bounded:

**Definition 6 (Smoothness).** A twice differentiable function  $f: \mathcal{Z} \to \mathbb{R}$  is said to be L-smooth if the inequality

$$\|\nabla^2 f(\mathbf{z})\| \le L$$

holds for all  $z \in \mathcal{Z}$ .

**Remark 7.** Assuming a function f is smooth is equivalent to assuming that f can be upper bounded by a quadratic function everywhere.

**Remark 8.** When the log-density  $\log \pi$  of a probability measure  $\Pi$  is L-smooth,  $\log \pi$  can be upper bounded everywhere by the log-density of a Gaussian.

**Definition 7 (Strong Convexity).** A twice differentiable function  $f: \mathbb{R}^d \to \mathbb{R}$  is said to be  $\mu$ -strongly convex if the inequality

$$\frac{\mu}{2}{\left\| {\boldsymbol{z} - \boldsymbol{z}'} \right\|^2} + \left\langle {\nabla f\left( {\boldsymbol{z}} \right),\boldsymbol{z} - \boldsymbol{z}'} \right\rangle + f\left( {\boldsymbol{z}} \right) \le f\left( {\boldsymbol{z}'} \right)$$

holds for all  $z, z' \in \mathbb{R}^d$  and some  $\mu > 0$ .

**Remark 9.** If Definition 7 holds only for  $\mu = 0$ , f is said to be (non-strongly) convex.

**Remark 10.** Assuming a function f is strongly convex is equivalent to assuming that f can be lower bounded by a quadratic.

**Definition 8** (**Strongly Log-Concave Measures**). For a probability measure Π in a Euclidean measurable space ( $\mathbb{R}^d$ ,  $\mathcal{B}(\mathbb{R}^d)$ ,  $\mathbb{P}$ ), where  $\mathcal{B}(\mathbb{R}^d)$  is the  $\sigma$ -algebra of Borel-measurable subsets of  $\mathbb{R}^d$ ,  $\mathbb{P}$  is the Lebesgue measure, we say Π is  $\mu$ -strongly log-concave if its log-density log  $\pi(z)$ :  $\mathbb{R}^d \to \mathbb{R}$  is  $\mu$ -strongly convex for some  $\mu > 0$ .

**Remark 11.** If Definition 8 holds only for  $\mu = 0$ ,  $\Pi$  is said to be (non-strongly) log-concave.

**Remark 12.** When  $\Pi$  is  $\mu$ -strongly log-concave,  $\log \pi$  can be lower bounded everywhere by the log-density of a Gaussian.

#### D ProxGen Adam for Black-Box Variational Inference

#### Algorithm 1: ProxGen-Adam for Black-Box Variational Inference

**Input:** Initial variational parameters  $\lambda_0$ , base stepsize  $\alpha$ , second moment stepsize  $\beta_2$ , momentum stepsize  $\{\beta_{1,t}\}_{t=1}^T$ , small positive constant  $\epsilon$ 

for  $t = 1, \dots, T$  do

estimate gradient of energy 
$$\widehat{\nabla f}$$

$$\begin{aligned}
\mathbf{g}_t &= \widehat{\nabla f}(\lambda) + \nabla h(\lambda) \\
\widehat{\lambda}_{t+1} &= \beta_{1,t} \widehat{\lambda}_t + (1 - \beta_{1,t}) \widehat{\lambda}_t \\
\mathbf{v}_{t+1} &= \beta_2 \mathbf{v}_t + (1 - \beta_2) \mathbf{g}_t^2 \\
\mathbf{\Gamma}_{t+1} &= \operatorname{diag} \left( \alpha / \left( \sqrt{\mathbf{v}_{t+1}} + \epsilon \right) \right) \\
\lambda_{t+1} &= \lambda_t - \mathbf{\Gamma}_{t+1} \widehat{\lambda}_{t+1} \\
\mathbf{s}_{t+1} &\leftarrow \operatorname{getscale} (\lambda_{t+1}) \\
\mathbf{s}_{t+1} &\leftarrow \mathbf{s}_{t+1} + \frac{1}{2} \left( \sqrt{\mathbf{s}_{t+1}^2 + 4 \gamma_{s,t+1}} - \mathbf{s}_{t+1} \right) \\
\lambda_{t+1} &\leftarrow \operatorname{setscale} (\lambda_{t+1}, \mathbf{s}_{t+1})
\end{aligned}$$

(By convention, all vector operations are elementwise.)

Adaptive and matrix-valued stepsize-variants of SGD such as Adam (Kingma & Ba, 2015), Ada-Grad (Duchi *et al.*, 2011) are widely used. The matrix stepsize of Adam at iteration *t* is given as

$$\Gamma_{t+1} = \operatorname{diag}\left(\alpha/\left(\sqrt{v_{t+1}} + \epsilon\right)\right),$$

where  $v_t$  is the exponential moving average of the second moment,  $\alpha$  is the "base stepsize." Furthermore, the matrix stepsize is applied to the moving average of the gradients, a scheme often called the (heavy-ball) momentum, denoted here as  $\overline{\lambda}_t$ .

Recently, Yun et al. (2021) have proven the convergence for these adaptive, momentum, and matrix-valued stepsize-based SGD methods with proximal steps. Then, the proximal operator is applied

$$\operatorname{prox}_{\Gamma_{t},h}\left(\lambda_{t}-\Gamma_{t}\overline{\lambda}_{t}\right)=\operatorname{arg\,min}_{\lambda}\left\{\left\langle\overline{\lambda}_{t},\lambda\right\rangle+h\left(\lambda\right)+\frac{1}{2}(\lambda-\lambda_{t})^{\mathsf{T}}\Gamma_{t}^{-1}\left(\lambda-\lambda_{t}\right)\right\}.$$

For Adam, the matrix-valued stepsize is a diagonal matrix. Thus, the proximal operator of Domke (2020) for each  $s_i$  forms independent 1-dimensional quadratic problems. Thus, the proximal step is given in the closed-form

$$\operatorname{prox}_{\Gamma_{t},h}(s_{i}) = s_{i} + \frac{1}{2} \left( \sqrt{s_{i}^{2} + 4\gamma_{s_{i}}} - s_{i} \right),$$

where, dropping the index t for clarity,  $\bar{s}_i$  is the element of  $\bar{\lambda}_t$  corresponding to  $s_i$ ,  $\gamma_{s_i}$  denotes the stepsize of  $s_i$  (a diagonal element of  $\Gamma_t$ ). Combined with the Adam-like stepsize rule, the algorithm is shown in Algorithm 1.

**Difference with Adam** In Algorithm 1, we can see the differences with vanilla Adam. Notably, ProxGenAdam does not perform bias correction of the estimated moments. Furthermore, while some implementations of Adam decay  $\beta_1$ , we keep it constant. It is possible that these differences could result in a different behavior from vanilla Adam. However, in this work, we follow the original implementation by Yun *et al.* (2021) as closely as possible and leave the comparison with vanilla Adam to future works.

#### E Detailed Comparison Against Domke et al. (2023)

In this section, we contrast our results against those of Domke *et al.* (2023). First, the main challenge to establishing a convergence guarantee for BBVI has been on bounding the gradient variance. In particular, Domke (2019) proved that the variance of the reparameterization gradient for the energy,  $\widehat{\nabla f}$ , is bounded as

$$\mathbb{E}\|\widehat{\nabla f}(\lambda)\|_{2}^{2} \le \alpha \|\lambda - \bar{\lambda}\|_{2}^{2} + \beta \tag{3}$$

for some finite positive constants  $\alpha$ ,  $\beta$  depending on the problem constants d, L,  $k_{\varphi}$ . Domke *et al.* (2023) call a gradient estimator satisfying this bound to be a "quadratic variance" estimator. Furthermore, they prove that the closed-form entropy (CFE; Kucukelbir *et al.*, 2017; Titsias & Lázaro-Gredilla, 2014) estimator:

$$\widehat{\nabla F}_{\text{CFE}}(\lambda) \triangleq \widehat{\nabla f}(\lambda) + \nabla h(\lambda)$$

and the STL estimator by Roeder et al. (2017):

$$\widehat{\nabla F}_{\mathrm{STL}}\left(\boldsymbol{\lambda}\right)\triangleq\frac{1}{M}\sum_{m=1}^{M}-\nabla_{\boldsymbol{\lambda}}\boldsymbol{\ell}\left(\boldsymbol{\mathcal{T}}_{\!\boldsymbol{\lambda}}\left(\boldsymbol{u}_{m}\right)\right)+\nabla_{\boldsymbol{\lambda}}\log q_{\boldsymbol{\lambda}}\left(\boldsymbol{\mathcal{T}}_{\!\boldsymbol{\nu}}\left(\boldsymbol{u}_{m}\right)\right)\Big|_{\boldsymbol{\nu}=\boldsymbol{\lambda}},$$

where  $\mathbf{u} \sim \varphi$ , also qualify as quadratic variance estimators.

Unfortunately, it has been unknown whether SGD is guaranteed to converge with a quadratic variance estimator except for strongly convex objectives (Wright & Recht, 2021, p. 85). Domke *et al.* (2023) expand the boundaries of SGD and prove that projected and proximal SGD with a quadratic variance estimator converges for both convex and strongly convex objectives. In particular, for the location-scale variational family, the linear parameterization, and log-concave objectives, they prove a complexity of  $\mathcal{O}(1/\epsilon^2)$ , and for strongly log-concave objectives, they prove a complexity of  $\mathcal{O}(1/\epsilon)$ .

On the other hand, Kim et al. (2023) and Lemma 12 developed the bound in Equation (3) to be of the form of

$$\mathbb{E}\|\widehat{\nabla F}(\lambda)\|_{2}^{2} \le A\left(F(\lambda) - F^{*}\right) + \|\nabla F(\lambda)\|_{2}^{2} + C\tag{4}$$

for some positive finite constants A, B, C, for which the convergence of SGD for convex, strongly convex (Garrigos & Gower, 2023; Gorbunov *et al.*, 2020), and non-convex objectives (Khaled & Richtárik, 2023) have already been proven. Applying these results to log-smooth and log-quadratically growing objectives, we prove a complexity of  $\mathcal{O}(1/\epsilon^4)$ , while for strong log-concave objectives, we also prove a complexity of  $\mathcal{O}(1/\epsilon)$ .

Overall, both approaches can be summarized as follows: we focused on establishing gradient variance bounds of known convergence proofs, while Domke *et al.* (2023) aimed to prove that the bound by Domke (2019) is sufficient to guarantee convergence. Note that, for strongly log-concave objectives, Equation (3) immediately implies Equation (4). Therefore, both approaches intersect in the case of strongly log-concave objectives. Indeed, Theorem 8 and the analogous result of Domke *et al.* (2023) are both based on the same proof strategy by Gower *et al.* (2019).

#### F Proofs

#### F.1 Auxiliary Lemmas

**Lemma 5.** Let  $\phi(x) = x$ . Then, the parameterization is linear in the sense that  $\mathcal{T}_{\lambda}$  is a bilinear function such that

$$\mathcal{T}_{\lambda-\lambda'}(\mathbf{u}) = \mathcal{T}_{\lambda}(\mathbf{u}) - \mathcal{T}_{\lambda'}(\mathbf{u}).$$

for any  $\lambda, \lambda' \in \Lambda$ .

Proof.

$$\mathcal{T}_{\lambda-\lambda'}(u) = (C(\lambda-\lambda'))u + (m-m')$$
  
=  $(D_{\phi}(s-s') + (L-L'))u + (m-m'),$ 

using the fact that  $\phi$  is the identity function,

$$= (D_{\phi}(s) - D_{\phi}(s') + (L - L'))u + (m - m')$$

$$= (C(\lambda) - C(\lambda'))u + (m + m')$$

$$= (C(\lambda)u + m) - (C(\lambda')u + m')$$

$$= \mathcal{T}_{\lambda}(u) - \mathcal{T}_{\lambda'}(u).$$

The linearity with respect to  $\boldsymbol{u}$  is obvious.

**Lemma 6.** Let the linear parameterization be used. Then, for any  $\lambda, \lambda' \in \Lambda$ , the inner product of the Jacobian of the reparameterization function satisfies the following equalities for any  $\mathbf{u} \in \mathbb{R}^d$ .

(i) For the Cholesky family (Domke, 2019, Lemma 8),

$$\left(\frac{\partial \mathcal{T}_{\lambda}(\boldsymbol{u})}{\partial \lambda}\right)^{\mathsf{T}} \frac{\partial \mathcal{T}_{\lambda}(\boldsymbol{u})}{\partial \lambda} = \left(1 + \|\boldsymbol{u}\|_{2}^{2}\right) \mathbf{I}$$

(ii) For the mean-field family (Kim et al., 2023, Lemma 1),

$$\left(\frac{\partial \mathcal{T}_{\lambda}(\boldsymbol{u})}{\partial \lambda}\right)^{\mathsf{T}} \frac{\partial \mathcal{T}_{\lambda}(\boldsymbol{u})}{\partial \lambda} = \left(1 + \left\|\boldsymbol{U}^{2}\right\|_{\mathsf{F}}\right) \mathbf{I},$$

where  $\mathbf{U} = \operatorname{diag}(u_1, \dots, u_d)$ 

**Lemma 7.** Let the linear parameterization be used. Then, for any  $\lambda \in \Lambda$  and any  $z \in \mathbb{R}^d$ , the following relationships hold.

(i) For the Cholesky family (Domke, 2019, Lemma 2),

$$\mathbb{E}\left(1+\left\|\boldsymbol{u}\right\|_{2}^{2}\right)\left\|\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)-\boldsymbol{z}\right\|_{2}^{2}=\left(d+1\right)\left\|\boldsymbol{m}-\boldsymbol{z}\right\|_{2}^{2}+\left(d+k_{\varphi}\right)\left\|\boldsymbol{C}\right\|_{\mathrm{F}}^{2}$$

(ii) For the mean-field family (Kim et al., 2023, Lemma 2),

$$\mathbb{E}\left(1+\left\|\boldsymbol{U}^{2}\right\|_{\mathrm{F}}\right)\left\|\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)-\boldsymbol{z}\right\|_{2}^{2}\leq\left(\sqrt{dk_{\varphi}}+k_{\varphi}\sqrt{d}+1\right)\left\|\boldsymbol{m}-\boldsymbol{z}\right\|_{2}^{2}+\left(2k_{\varphi}\sqrt{d}+1\right)\left\|\boldsymbol{C}\right\|_{\mathrm{F}}^{2}.$$

**Corollary 2.** Let the linear parameterization be used and  $\lambda, \lambda' \in \Lambda$  be any pair of variational parameters.

(i) For the Cholesky family,

$$\mathbb{E}\left(1+\left\|\mathbf{u}\right\|_{2}^{2}\right)\left\|\mathcal{T}_{\lambda'}\left(\mathbf{u}\right)-\mathcal{T}_{\lambda}\left(\mathbf{u}\right)\right\|_{2}^{2}\leq\left(k_{\varphi}+d\right)\left\|\lambda-\lambda'\right\|_{2}^{2}$$

(ii) For the mean-field family,

$$\mathbb{E}\left(1+\left\|\boldsymbol{U}^{2}\right\|_{\mathrm{F}}\right)\left\|\mathcal{T}_{\lambda'}\left(\boldsymbol{u}\right)-\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)\right\|_{2}^{2}\leq\left(2k_{\varphi}\sqrt{d}+1\right)\left\|\lambda-\lambda'\right\|_{2}^{2}$$

*Proof.* The results are a direct consequence of Lemma 7 and Lemma 5.

**Proof of (i)** We start from Lemma 7 as

$$\mathbb{E}\left(1+\left\|\boldsymbol{u}\right\|_{2}^{2}\right)\left\|\mathcal{T}_{\lambda-\lambda'}\left(\boldsymbol{u}\right)-\boldsymbol{z}\right\|_{2}^{2}=\left(d+1\right)\left\|\left(\boldsymbol{m}-\boldsymbol{m'}\right)-\boldsymbol{z}\right\|_{2}^{2}+\left(d+k_{\varphi}\right)\left\|\boldsymbol{C}\left(\lambda\right)-\boldsymbol{C}\left(\lambda'\right)\right\|_{\mathrm{F}}^{2},$$
 setting  $\boldsymbol{z}=\boldsymbol{0},$ 

$$=\left(d+1\right)\left\|\boldsymbol{m}-\boldsymbol{m}'\right\|_{2}^{2}+\left(d+k_{\varphi}\right)\left\|\boldsymbol{C}\left(\boldsymbol{\lambda}\right)-\boldsymbol{C}\left(\boldsymbol{\lambda}'\right)\right\|_{\mathrm{F}}^{2},$$

and since  $k_{\varphi} \ge 3$  by the property of the kurtosis,

$$\leq (d + k_{\varphi}) (\|\boldsymbol{m} - \boldsymbol{m}'\|_{2}^{2} + \|\boldsymbol{C}(\lambda) - \boldsymbol{C}(\lambda')\|_{F}^{2})$$
  
=  $(d + k_{\varphi}) \|\lambda - \lambda'\|_{2}^{2}$ .

**Proof of (ii)** Similarly, for the mean-field family, we can apply Lemma 7 as

$$\mathbb{E}\left(1+\left\|\boldsymbol{U}^{2}\right\|_{\mathrm{F}}\right)\left\|\mathcal{T}_{\boldsymbol{\lambda}-\boldsymbol{\lambda}'}\left(\boldsymbol{u}\right)-\bar{\boldsymbol{z}}\right\|_{2}^{2}\leq\left(\sqrt{dk_{\varphi}}+k_{\varphi}\sqrt{d}+1\right)\left\|\left(\boldsymbol{m}-\boldsymbol{m}'\right)-\bar{\boldsymbol{z}}\right\|_{2}^{2}+\left(2k_{\varphi}\sqrt{d}+1\right)\left\|\boldsymbol{C}-\boldsymbol{C}'\right\|_{\mathrm{F}}^{2},$$
 setting  $\boldsymbol{z}=\boldsymbol{0},$ 

$$= \left(\sqrt{dk_{\varphi}} + k_{\varphi}\sqrt{d} + 1\right) \left\|\boldsymbol{m} - \boldsymbol{m}'\right\|_{2}^{2} + \left(2k_{\varphi}\sqrt{d} + 1\right) \left\|\boldsymbol{C} - \boldsymbol{C}'\right\|_{F}^{2},$$

and since  $k_{\varphi} \geq 3$  by the property of the kurtosis,

$$\begin{split} & \leq \left(2k_{\varphi}\sqrt{d}+1\right)\left(\left\|\boldsymbol{m}-\boldsymbol{m}'\right\|_{2}^{2}+\left\|\boldsymbol{C}\left(\boldsymbol{\lambda}\right)-\boldsymbol{C}\left(\boldsymbol{\lambda}'\right)\right\|_{\mathrm{F}}^{2}\right) \\ & = \left(2k_{\varphi}\sqrt{d}+1\right)\left\|\boldsymbol{\lambda}-\boldsymbol{\lambda}'\right\|_{2}^{2}. \end{split}$$

Lemma 8. For the linear parameterization,

$$\mathbb{E}\|\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)-\mathcal{T}_{\lambda'}\left(\boldsymbol{u}\right)\|_{2}^{2}=\|\lambda-\lambda'\|_{2}^{2}$$

for any  $\lambda, \lambda' \in \Lambda$ .

*Proof.* First notice that, for linear parameterizations, we have

$$\mathbb{E}\|\mathcal{T}_{\lambda}(\boldsymbol{u})\|_{2}^{2} = \mathbb{E}\|\boldsymbol{C}\boldsymbol{u} + \boldsymbol{m}\|_{2}^{2}$$

$$= \mathbb{E}\boldsymbol{u}^{\mathsf{T}}\boldsymbol{C}^{\mathsf{T}}\boldsymbol{C}\boldsymbol{u} + \|\boldsymbol{m}\|_{2}^{2} + 2\boldsymbol{m}^{\mathsf{T}}\boldsymbol{C}\mathbb{E}\boldsymbol{u}$$

$$= \mathbb{E}\operatorname{tr}(\boldsymbol{u}^{\mathsf{T}}\boldsymbol{C}^{\mathsf{T}}\boldsymbol{C}\boldsymbol{u}) + \|\boldsymbol{m}\|_{2}^{2} + 2\boldsymbol{m}^{\mathsf{T}}\boldsymbol{C}\mathbb{E}\boldsymbol{u},$$

rotating the elements of the trace,

$$= \operatorname{tr}(\mathbf{C}^{\top} \mathbf{C} \mathbb{E} \mathbf{u} \mathbf{u}^{\top}) + \|\mathbf{m}\|_{2}^{2} + 2\mathbf{m}^{\top} \mathbf{C} \mathbb{E} \mathbf{u},$$

applying Assumption 1

$$= \operatorname{tr}(\boldsymbol{C}^{\mathsf{T}}\boldsymbol{C}) + \|\boldsymbol{m}\|_{2}^{2}$$
$$= \|\boldsymbol{C}\|_{F}^{2} + \|\boldsymbol{m}\|_{2}^{2}$$
$$= \|\boldsymbol{\lambda}\|_{2}^{2}.$$

Combined with Lemma 5, we have

$$\mathbb{E}\left\|\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)-\mathcal{T}_{\lambda'}\left(\boldsymbol{u}\right)\right\|_{2}^{2}=\mathbb{E}\left\|\mathcal{T}_{\lambda-\lambda'}\left(\boldsymbol{u}\right)\right\|_{2}^{2}=\left\|\lambda-\lambda'\right\|_{2}^{2}.$$

#### F.2 Properties of the Evidence Lower Bound

#### F.2.1 Smoothness

**Lemma 1.** If the diagonal conditioner  $\phi$  is  $L_h$ -log-smooth, then the entropic regularizer  $h(\lambda)$  is  $L_h$ -smooth.

*Proof.* The entropic regularizer is

$$h(\lambda) = -H(\varphi) - \sum_{i=1}^{d} \log \phi(s_i),$$

and depends only on the diagonal elements  $s_1, \ldots, s_d$  of C. The Hessian of h is then a diagonal matrix, where only the entries that correspond to  $s_1, \ldots, s_d$  are non-zero. The Lipschitz smoothness constant is then the constant  $L_h < \infty$  that satisfies

$$\frac{\partial^2 h(\lambda)}{\partial s_i^2} = -\frac{\mathrm{d}^2 \log \phi}{\mathrm{d} s_i^2} < L_h$$

for all i = 1, ..., d, which is the smoothness constant of  $s_i \mapsto \log \phi(s_i)$ .

**Lemma 2.** Let **H** be a  $n \times n$  symmetric random matrix, where it is bounded as  $\|\mathbf{H}\|_2 \leq L < \infty$  almost surely. Also, let **J** be an  $m \times n$  random matrix such that  $\|\mathbb{E}\mathbf{J}^{\mathsf{T}}\mathbf{J}\|_2 < \infty$ . Then,

$$\|\mathbb{E}\mathbf{J}^{\mathsf{T}}\mathbf{H}\mathbf{J}\|_{2} \leq L\|\mathbb{E}\mathbf{J}^{\mathsf{T}}\mathbf{J}\|_{2}.$$

*Proof.* By the property of the Rayleigh quotients, for a symmetric matrix A, its maximum eigenvalue is given in the variational form

$$\sup_{\|\boldsymbol{x}\| \leq 1} \boldsymbol{x}^{\mathsf{T}} \boldsymbol{H} \boldsymbol{x} = \sigma_{\max} \left( \boldsymbol{H} \right) \leq \sqrt{\sigma_{\max} \left( \boldsymbol{H} \right)^2} = \left\| \boldsymbol{H} \right\|_2,$$

where  $\sigma_{\max}(A)$  is the maximal eigenvalue of A. Notice the relationship with the  $\ell_2$ -operator norm. The inequality is strict only if all eigenvalues are negative.

From the property above,

$$\|\mathbb{E} \mathbf{J}^{\mathsf{T}} \mathbf{H} \mathbf{J}\|_{2} = \sup_{\|\mathbf{x}\|_{2} \le 1} \mathbf{x}^{\mathsf{T}} (\mathbb{E} \mathbf{J}^{\mathsf{T}} \mathbf{H} \mathbf{J}) \mathbf{x}.$$

By reparameterizing as y = Jx,

$$= \sup_{\|x\|_2 \le 1} \mathbb{E} \mathbf{y}^\mathsf{T} \mathbf{H} \mathbf{y},$$

and the property of the  $\ell_2$ -operator norm,

$$\leq \sup_{\|\mathbf{x}\|_{2} \leq 1} \mathbb{E}\|\mathbf{H}\|_{2}\|\mathbf{y}\|_{2}^{2} = \sup_{\|\mathbf{x}\|_{2} \leq 1} \mathbb{E}\|\mathbf{H}\|_{2}\|\mathbf{J}\mathbf{x}\|_{2}^{2}.$$

From our assumption about the maximal eigenvalue of H,

$$\leq L \sup_{\|\boldsymbol{x}\|_{2} \leq 1} \mathbb{E} \|\boldsymbol{J}\boldsymbol{x}\|_{2}^{2},$$

denoting the  $\ell_2$  vector norm as a quadratic form as,

$$= L \sup_{\|\boldsymbol{x}\|_2 \le 1} \boldsymbol{x}^\top (\mathbb{E} \boldsymbol{J}^\top \boldsymbol{J}) \boldsymbol{x},$$

again, by the property of the  $\ell_2$ -operator norm,

$$\leq L \|\mathbb{E} \mathbf{J}^{\mathsf{T}} \mathbf{J}\|_{2} \sup_{\|\mathbf{x}\|_{2} \leq 1} \|\mathbf{x}\|_{2}^{2}$$
$$= L \|\mathbb{E} \mathbf{J}^{\mathsf{T}} \mathbf{J}\|_{2}.$$

**Lemma 9.** For a 1-Lipschitz diagonal conditioner  $\phi$ , the Jacobian of the location-scale reparameterization function  $\mathcal{T}_{\lambda}$  satisfies

$$\left\| \mathbb{E} \left( \frac{\partial \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right)}{\partial \lambda} \right)^{\mathsf{T}} \frac{\partial \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right)}{\partial \lambda} \right\|_{2} \leq 1.$$

*Proof.* For notational clarity, we will occasionally represent  $\mathcal{T}_{\lambda}$  as

$$\mathcal{F}_{\lambda}(\mathbf{u}) = \mathcal{F}(\lambda; \mathbf{u}),$$

such that  $\mathcal{T}_i(\lambda; \mathbf{u})$  denotes the *i*th component of  $\mathcal{T}_{\lambda}$ .

From the definition of  $\mathcal{T}_{\lambda}$ , it is straightforward to notice that its Jacobian is the concatenation of 3 block matrices

$$J_{m} = \frac{\partial \mathcal{T}_{\lambda}(u)}{\partial m}, \qquad J_{s} = \frac{\partial \mathcal{T}_{\lambda}(u)}{\partial s}, \text{ and } J_{L} = \frac{\partial \mathcal{T}_{\lambda}(u)}{\partial \text{vec}(L)}.$$

The *m* block form a deterministic identity matrix

$$J_{m} = \frac{\partial \mathcal{T}_{\lambda}(u)}{\partial m} = \mathbf{I},$$

which is shown by (Domke, 2020, Lemma 4).

The proof strategy is as follows: we will directly compute the squared Jacobian through block matrix multiplication. The key is that, after expectation, the resulting matrix becomes diagonal. Then, the  $\ell_2$  operator norm, or maximal eigenvalue, follows trivially as the maximal diagonal element. First,

$$\mathbb{E}\left(\frac{\partial \mathcal{T}(\lambda; u)}{\partial \lambda}\right)^{\mathsf{T}} \frac{\partial \mathcal{T}(\lambda; u)}{\partial \lambda} = \mathbb{E}\begin{bmatrix}J_{J_{L}}^{\mathsf{T}}\\J_{L}^{\mathsf{T}}\\J_{L}^{\mathsf{T}}\end{bmatrix} \begin{bmatrix}J_{m} & J_{s} & J_{L}\end{bmatrix}$$

$$= \mathbb{E}\begin{bmatrix}J_{m}^{\mathsf{T}}J_{m} & J_{m}^{\mathsf{T}}J_{s} & J_{m}^{\mathsf{T}}J_{L}\\J_{L}^{\mathsf{T}}J_{m} & J_{L}^{\mathsf{T}}J_{s} & J_{L}^{\mathsf{T}}J_{L}\end{bmatrix}$$

$$= \begin{bmatrix}\mathbf{I} & \mathbb{E}J_{s} & \mathbb{E}J_{L}\\\mathbb{E}J_{L}^{\mathsf{T}} & \mathbb{E}J_{L}^{\mathsf{T}}J_{s} & \mathbb{E}J_{L}^{\mathsf{T}}J_{L}\end{bmatrix}.$$

$$= \begin{bmatrix}\mathbb{E}J_{L}^{\mathsf{T}} & \mathbb{E}J_{L}^{\mathsf{T}}J_{s} & \mathbb{E}J_{L}^{\mathsf{T}}J_{L}\\\mathbb{E}J_{L}^{\mathsf{T}} & \mathbb{E}J_{L}^{\mathsf{T}}J_{s} & \mathbb{E}J_{L}^{\mathsf{T}}J_{L}\end{bmatrix}.$$

For  $J_s$ , the entries are

$$\frac{\partial \mathcal{T}_{i}(\mathbf{u})}{\partial s_{i}} = \phi'(s_{i}) u_{j} \mathbb{1}_{i=j},$$

which is a diagonal matrix. Thus, by Assumption 1,

$$\mathbb{E} J_s = \mathbf{O}, \qquad \mathbb{E} J_s^{\mathsf{T}} J_s = \mathrm{diag} \left( \phi'(s) \right)^2.$$

For  $J_s$ , the entries are

$$\frac{\partial \mathcal{T}_i(\lambda; \boldsymbol{u})}{\partial L_{ik}} = u_k \mathbb{1}_{i=j}.$$

To gather some intuition, the case of d = 4 looks like the following:

$$\boldsymbol{J_L} = \begin{bmatrix} u_2 & u_3 & u_4 & & & & & \\ & & u_1 & u_3 & u_4 & & & & \\ & & & & u_1 & u_2 & u_4 & & & \\ & & & & & u_1 & u_2 & u_4 & & \\ & & & & & u_1 & u_2 & u_3 \end{bmatrix}.$$

It is crucial to notice that the *i*th row does *not* include  $u_i$ . This means that, the matrix  $J_s^{\mathsf{T}} J_L$  has entries that are either 0, or  $\phi'(s_i) u_i u_j$  for  $i \neq j$ , which is  $\mathbb{E} \phi'(s_i) u_i u_j = 0$  by Assumption 1. Therefore,

$$\mathbb{E} J_{\mathbf{s}}^{\mathsf{T}} J_{L} = \mathbf{O}.$$

Finally, the elements of  $J_L^T J_L$  are

$$\mathbb{E}\sum_{i=0}^{d}\frac{\partial \mathcal{T}_{i}(\lambda;\boldsymbol{u})}{\partial L_{jk}}\frac{\partial \mathcal{T}_{i}(\lambda;\boldsymbol{u})}{\partial L_{lm}}=\mathbb{E}\sum_{i=0}^{d}u_{k}u_{m}\mathbb{1}_{i=j}\mathbb{1}_{i=l}=\mathbb{1}_{j=l}(\mathbb{E}u_{k}u_{m})=\mathbb{1}_{j=l}\mathbb{1}_{k=m},$$

where the last equality follows from Assumption 1, which forms an identity matrix as

$$\mathbb{E} J_L^{\mathsf{T}} J_L = \mathbf{I}.$$

Therefore, the expected-squared Jacobian is now

$$\begin{split} \mathbb{E} \left( \frac{\partial \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right)}{\partial \boldsymbol{\lambda}} \right)^{\mathsf{T}} \frac{\partial \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right)}{\partial \boldsymbol{\lambda}} &= \begin{bmatrix} \mathbf{I} & \mathbb{E} \boldsymbol{J}_{s} & \mathbb{E} \boldsymbol{J}_{L} \\ \mathbb{E} \boldsymbol{J}_{s}^{\mathsf{T}} & \mathbb{E} \boldsymbol{J}_{s}^{\mathsf{T}} \boldsymbol{J}_{s} & \mathbb{E} \boldsymbol{J}_{s}^{\mathsf{T}} \boldsymbol{J}_{L} \\ \mathbb{E} \boldsymbol{J}_{L}^{\mathsf{T}} & \mathbb{E} \boldsymbol{J}_{L}^{\mathsf{T}} \boldsymbol{J}_{s} & \mathbb{E} \boldsymbol{J}_{L}^{\mathsf{T}} \boldsymbol{J}_{L} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I} & \\ \operatorname{diag} \left( \boldsymbol{\phi} \left( \boldsymbol{s} \right) \right)^{2} \\ & \mathbf{I} \end{bmatrix}, \end{split}$$

which, conveniently, is a diagonal matrix. The maximal singular value of a block-diagonal matrix is the maximal singular value of each block. And since each block is diagonal with only positive entries, the largest element forms the maximal singular value. As we assume that  $\phi$  is 1-Lipchitz, the element of all blocks is lower-bounded by 0 and upper-bounded by 1. Therefore, the maximal singular value of the expected-squared Jacobian is bounded by 1.

**Theorem 1.** Let  $\ell$  be  $L_{\ell}$ -smooth and twice differentiable. Then, the following results hold:

- (i) If  $\phi$  is linear, the energy f is  $L_{\ell}$ -smooth.
- (ii) If  $\phi$  is 1-Lipschitz, the energy  $\ell$  is  $(L_{\ell} + L_s)$ -smooth if and only if Assumption 3 holds.

*Proof.* For notational clarity, we will occasionally represent  $\mathcal{T}_{\lambda}$  as

$$\mathcal{F}_{\lambda}(\mathbf{u}) = \mathcal{F}(\lambda; \mathbf{u}),$$

such that  $\mathcal{T}_i(\lambda; \boldsymbol{u})$  denotes the *i*th component of  $\mathcal{T}_{\lambda}$ .

By the Leibniz and chain rule, the Hessian of the energy f follows as

$$\nabla^{2} f(\lambda) = \mathbb{E} \nabla_{\lambda}^{2} \ell\left(\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)\right)$$

$$= \underbrace{\mathbb{E} \left(\frac{\partial \mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)}{\partial \lambda}\right)^{\mathsf{T}} \nabla^{2} \ell\left(\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)\right) \frac{\partial \mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)}{\partial \lambda}}_{\triangleq T_{\text{tim}}} + \underbrace{\mathbb{E} \sum_{i=1}^{d} D_{i} \ell\left(\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)\right) \frac{\partial^{2} \mathcal{T}_{i}\left(\boldsymbol{\lambda};\boldsymbol{u}\right)}{\partial \lambda^{2}}}_{\triangleq T_{\text{top}}}.$$

When  $\mathcal{F}$  is linear with respect to  $\lambda$ , it is clear that we have

$$\frac{\partial^2 \mathcal{T}_i(\lambda; \mathbf{u})}{\partial \lambda^2} = \mathbf{0}.$$
 (5)

Then,  $T_{\text{non}}$  is zero. In contrast,  $T_{\text{lin}}$  appears for both the linear and nonlinear cases. Therefore,  $T_{\text{non}}$  fully characterizes the effect of nonlinearity in the reparameterization function.

Now, the triangle inequality yields

$$\|\nabla^2 f(\lambda)\|_2 = \|T_{\text{lin}} + T_{\text{non}}\|_2 \le \|T_{\text{lin}}\|_2 + \|T_{\text{non}}\|_2,$$

where equality is achieved when either term is 0. On the contrary, the reverse triangle inequality states that

$$\left|\left\|T_{\mathrm{lin}}\right\|_{2}-\left\|T_{\mathrm{non}}\right\|_{2}\right|\leq\left\|\nabla^{2}f\left(\boldsymbol{\lambda}\right)\right\|_{2}.$$

This implies that, if either  $T_{\rm lin}$  or  $T_{\rm non}$  is unbounded, the Hessian is not bounded. Thus, ensuring that  $T_{\rm lin}$  and  $T_{\rm non}$  are bounded is sufficient and necessary to establish that f is smooth.

**Proof of (i)** The bound on the linear part,  $T_{lin}$ , follows from Lemma 2 as

$$\begin{split} \left\| T_{\text{lin}} \right\|_2 &= \left\| \mathbb{E} \left( \frac{\partial \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right)}{\partial \lambda} \right)^{\mathsf{T}} \nabla^2 \ell \left( \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right) \right) \frac{\partial \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right)}{\partial \lambda} \right\|_2 \\ &\leq L_{\ell} \left\| \mathbb{E} \left( \frac{\partial \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right)}{\partial \lambda} \right)^{\mathsf{T}} \frac{\partial \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right)}{\partial \lambda} \right\|_2, \end{split}$$

and from the 1-Lipschitzness of  $\phi$ , Lemma 9 yields

$$\leq L_{\ell}$$
.

When  $\phi$  is linear, it immediately follows from Equation (5) that

$$\|\nabla^2 f(\lambda)\|_2 = \|T_{\text{lin}}\|_2 \le L_\ell,$$

which is tight as shown by Domke (2020, Theorem 6).

**Proof of (ii)** For the nonlinear part  $T_{\text{non}}$ , we use the fact that  $\mathcal{T}_i(\lambda; \boldsymbol{u})$  is given as

$$\mathcal{T}_{i}(\boldsymbol{\lambda};\boldsymbol{u})=m_{i}+\phi\left(s_{i}\right)u_{i}+\sum_{j\neq i}L_{ij}u_{j}.$$

The second derivative of  $\mathcal{T}_i$  is clearly non-zero only for the nonlinear part involving  $s_1, \dots, s_d$ . Thus,  $T_{\text{non}}$  follows as

$$T_{\text{non}} = \mathbb{E} \sum_{i=1}^{d} D_{i} \ell \left( \mathcal{T}_{\lambda} \left( \boldsymbol{u} \right) \right) \frac{\partial^{2} \mathcal{T}_{i} \left( \lambda; \boldsymbol{u} \right)}{\partial \lambda^{2}}$$

$$= \mathbb{E} \sum_{i=1}^{d} g_{i} \left( \lambda; \boldsymbol{u} \right) \frac{\partial^{2} \mathcal{T}_{i} \left( \lambda; \boldsymbol{u} \right)}{\partial \lambda^{2}}$$

$$= \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \mathbb{E} \sum_{i=1}^{d} g_{i} \left( \lambda; \boldsymbol{u} \right) \frac{\partial^{2} \mathcal{T}_{i} \left( \lambda; \boldsymbol{u} \right)}{\partial s^{2}} & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}.$$

Furthermore, the second-order derivatives with respect to  $s_1, \dots, s_d$  are given as

$$\frac{\partial^2 \mathcal{T}_i(\lambda; \mathbf{u})}{\partial s_i^2} = \mathbb{1}_{i=j} \phi''(s_j).$$

Considering this, the only non-zero block of  $T_{non}$  forms a diagonal matrix as

$$\mathbb{E} \sum_{i=1}^{d} g_{i}(\lambda; \boldsymbol{u}) \frac{\partial^{2} \mathcal{T}_{i}(\lambda; \boldsymbol{u})}{\partial s} = \begin{bmatrix} \mathbb{E} g_{1}(\lambda; \boldsymbol{u}) \frac{\partial^{2} \mathcal{T}_{1}(\lambda; \boldsymbol{u})}{\partial s_{1}^{2}} & & \\ & \ddots & \\ & & \mathbb{E} g_{d}(\lambda; \boldsymbol{u}) \frac{\partial^{2} \mathcal{T}_{d}(\lambda; \boldsymbol{u})}{\partial s_{d}^{2}} \end{bmatrix}$$
$$= \begin{bmatrix} \mathbb{E} g_{1}(\lambda; \boldsymbol{u}) \phi''(s_{1}) u_{1} & & \\ & \ddots & \\ & & \mathbb{E} g_{d}(\lambda; \boldsymbol{u}) \phi''(s_{d}) u_{d} \end{bmatrix}$$

This implies that the only non-zero entries of  $T_{\text{non}}$  lie on its diagonal. Since the  $\ell_2$  norm of a diagonal matrix is the value of the maximal diagonal element,

$$\|T_{\text{non}}\|_{2} \leq \max_{i=1,\dots,d} \mathbb{E}g_{i}(\lambda; \boldsymbol{u}) \phi''(s_{i}) u_{i} \leq L_{s},$$

where  $L_s$  is finite constant if Assumption 3 holds. On the contrary, if a finite  $L_s$  does not exist,  $\|T_{\text{non}}\|_2$  cannot be bounded. Therefore, the energy is smooth if and only if Assumption 3 holds. When it does, the energy f is  $L_f + L_s$  smooth.

**Example 2.** Let  $\ell(z) = (1/2) z^{\mathsf{T}} A z$  and the diagonal conditioner be  $\phi(x) = \text{softplus}(x)$ . Then,

- (i) if A is dense and the variational family is the mean-field family or
- (ii) if A is diagonal and the variational family is the Cholesky family,

Assumption 3 holds with  $L_s \approx 0.26034 (\max_{i=1,...,d} A_{ii})$ .

(iii) If A is dense but the Cholesky family is used, Assumption 3 does not hold.

*Proof.* Since the gradient is

$$\nabla \ell \left( \boldsymbol{z} \right) = \boldsymbol{A} \boldsymbol{z},$$

combined with reparameterization, we have

$$g(\lambda; \mathbf{u}) = A(C\mathbf{u} + \mathbf{m})$$

Then, for each coordinate i = 1, ..., d, we have

$$\mathbb{E}g_{i}(\lambda; \boldsymbol{u}) u_{i} \phi''(s_{i}) = \mathbb{E}\left(\sum_{j} \sum_{k \leq j} A_{ij} C_{jk} u_{k} + \sum_{j} A_{ij} m_{j}\right) u_{i} \phi''(s_{i})$$

$$= \sum_{i} \sum_{k \leq j} A_{ij} C_{jk} \mathbb{E}u_{k} u_{i} \phi''(s_{i}) + \sum_{i} A_{ij} m_{j} \mathbb{E}u_{i} \phi''(s_{i}),$$

and from Assumption 1,

$$\begin{split} &=\phi''(s_i)\sum_j\sum_{k\leq j}A_{ij}C_{jk}\mathbb{I}_{k=i}\\ &=\phi''(s_i)\sum_jA_{ij}C_{ji}. \end{split}$$

Furthermore, the diagonal of C involves  $\phi$  such that

$$\mathbb{E}g_{i}\left(\boldsymbol{\lambda};\boldsymbol{u}\right)u_{i}\phi''(s_{i}) = \underbrace{A_{ii}\phi\left(s_{i}\right)\phi''\left(s_{i}\right)}_{T_{\text{diag}}} + \underbrace{\sum_{j< i}A_{ij}C_{ji}\phi''(s_{i})}_{T_{\text{off}}}.$$

For the softplus function, we have

$$0 < \phi''(s) < 1$$

for any finite s, and we have

$$\sup_{s} \phi(s) \phi''(s) \approx 0.26034,$$

where the supremum was numerically approximated. Then, it is clear that  $T_{\rm diag}$  is finite as long as the diagonals of  $\boldsymbol{A}$  are finite. Furthermore, we have the following:

- (i) If A is diagonal, then  $T_{\text{off}}$  is 0.
- (ii) If A is dense but C is diagonal due to the use of the mean-field family,  $T_{\rm off}$  is again 0.
- (iii) However, when both  $\bf A$  and  $\bf C$  are not diagonal,  $T_{\rm off}$  can be made arbitrarily large.

#### F.2.2 Convexity

**Lemma 10.** Let  $\ell$  be convex. Then, for a convex nonlinear  $\phi$ , the inequality

$$\langle \nabla_{\lambda} f(\lambda), \lambda - \lambda' \rangle \leq \mathbb{E} \langle \nabla g(\lambda; \mathbf{u}), \mathcal{T}_{\lambda}(\mathbf{u}) - \mathcal{T}_{\lambda'}(\mathbf{u}) \rangle$$

holds for all  $\lambda \in \Lambda$  if and only if Assumption 4 holds. For the linear parameterization, the inequality becomes equality.

Proof. First, notice that the left-hand side is

$$\left\langle \nabla_{\lambda} f\left(\lambda\right), \lambda - \lambda' \right\rangle = \sum_{i=1}^{p} \left\langle \nabla_{\lambda} \mathbb{E} \ell\left(\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)\right), \lambda - \lambda' \right\rangle = \mathbb{E} \sum_{i=1}^{p} \left\langle \left(\frac{\partial \mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)}{\partial \lambda_{i}}\right) g\left(\lambda; \boldsymbol{u}\right), \lambda - \lambda' \right\rangle.$$

By restricting us to the location-scale family, we then get

$$= \mathbb{E}\left(\underbrace{\sum_{i} \left(\frac{\partial \mathcal{T}_{\lambda}\left(\mathbf{u}\right)}{\partial m_{i}}\right) g\left(\boldsymbol{\lambda}; \mathbf{u}\right) \left(m_{i} - m_{i}'\right)}_{\text{convexity with respect to } \mathbf{m}} + \underbrace{\sum_{ij} \left(\frac{\partial \mathcal{T}_{\lambda}\left(\mathbf{u}\right)}{\partial L_{ij}}\right) g\left(\boldsymbol{\lambda}; \mathbf{u}\right) \left(L_{ij} - L_{ij}'\right)}_{\text{convexity with respect to } \mathbf{L}}\right)$$

$$+\sum_{i}\mathbb{E}\left(\frac{\partial\mathcal{T}_{\lambda}(\boldsymbol{u})}{\partial s_{i}}\right)g(\lambda;\boldsymbol{u})(s_{i}-s'_{i}),$$

convexity with respect to  $\boldsymbol{s}$ 

and plugging the derivatives of the reparameterization function,

$$= \mathbb{E}\left(\sum_{i} g_{i}(\boldsymbol{\lambda}; \boldsymbol{u}) \left(m_{i} - m'_{i}\right) + \sum_{i} \sum_{j < i} u_{j} g_{i}(\boldsymbol{\lambda}; \boldsymbol{u}) \left(L_{ij} - L'_{ij}\right) + \sum_{i} \phi'\left(s_{i}\right) u_{i} g_{i}\left(\boldsymbol{\lambda}; \boldsymbol{u}\right) \left(s_{i} - s'_{i}\right)\right).$$

On the other hand, the right-hand side follows as

$$\mathbb{E}\left\langle \nabla g\left(\boldsymbol{\lambda};\boldsymbol{u}\right),\ \mathcal{T}_{\boldsymbol{\lambda}}\left(\boldsymbol{u}\right)-\mathcal{T}_{\boldsymbol{\lambda}'}\left(\boldsymbol{u}\right)\right\rangle$$

$$\begin{split} &= \mathbb{E}\left(\left\langle g\left(\boldsymbol{\lambda};\boldsymbol{u}\right),\,\boldsymbol{m}-\boldsymbol{m}'\right\rangle + \left\langle g\left(\boldsymbol{\lambda};\boldsymbol{u}\right),\left(\boldsymbol{L}-\boldsymbol{L}'\right)\boldsymbol{u}\right\rangle + \left\langle g\left(\boldsymbol{\lambda};\boldsymbol{u}\right),\,\left(\boldsymbol{\Phi}\left(\boldsymbol{s}\right)-\boldsymbol{\Phi}\left(\boldsymbol{s}'\right)\right)\boldsymbol{u}\right\rangle\right) \\ &= \mathbb{E}\left(\sum_{i}g_{i}\left(\boldsymbol{\lambda};\boldsymbol{u}\right)\left(m_{i}-m_{i}'\right) + \sum_{i}\sum_{j< i}u_{j}g_{i}\left(\boldsymbol{\lambda};\boldsymbol{u}\right)\left(L_{ij}-L_{ij}'\right) + \sum_{i}g_{i}\left(\boldsymbol{\lambda};\boldsymbol{u}\right)u_{i}\left(\phi\left(\boldsymbol{s}_{i}\right)-\phi\left(\boldsymbol{s}_{i}'\right)\right)\right). \end{split}$$

The convexity with respect to the m and L is clear from the first two terms; they are equal. The statement is now up to the last term. That is, the statement holds if

$$\mathbb{E}\sum_{i}g_{i}(\boldsymbol{\lambda};\boldsymbol{u})\ u_{i}\ \phi'(s_{i})\big(s_{i}-s_{i}'\big)\ \leq\ \mathbb{E}\sum_{i}g_{i}(\boldsymbol{\lambda};\boldsymbol{u})\ u_{i}\big(\phi(s_{i})-\phi(s_{i}')\big). \tag{6}$$

For this, we will show that Assumption 4 is both necessary and sufficient.

**Proof of sufficiency** Equation (6) holds if

$$\mathbb{E}g_i(\lambda; \mathbf{u}) \ u_i \geq 0$$

for all i = 1, ..., d, which is non other than Assumption 4.

**Proof of necessity** Suppose that the inequality

$$\langle \nabla_{\lambda} f(\lambda), \lambda - \lambda' \rangle \leq \mathbb{E} \langle g(\lambda; \mathbf{u}), \mathcal{T}_{\lambda} (\mathbf{u}) - \mathcal{T}_{\lambda'} (\mathbf{u}) \rangle$$

holds for all  $\lambda \in \Lambda$ , implying

$$\sum_{i} \mathbb{E} u_{i} g_{i}(\lambda; \boldsymbol{u}) \phi'(s_{i})(s_{i} - s'_{i}) \leq \sum_{i} \mathbb{E} u_{i} g_{i}(\lambda; \boldsymbol{u}) ((\phi(s_{i}) - \phi(s'_{i})).$$

For any  $\lambda$ , we are free to set any  $\lambda' \in \Lambda$  and check whether we can retrieve Assumption 4 for this specific  $\lambda$ . Now, for each axis i, set  $s'_i = s_i$  for all  $j \neq i$ , then

$$\mathbb{E} u_i g_i(\lambda; \mathbf{u}) \phi'(s_i)(s_i - s_i') \leq \mathbb{E} u_i g_i(\lambda; \mathbf{u}) (\phi(s_i) - \phi(s_i')).$$

Since  $\phi$  is assumed to be convex such that

$$\phi'(s_i)(s_i - s_i') \le \phi(s_i) - \phi(s_i').$$

it follows that

$$\mathbb{E}\,u_i\,g_i\left(\boldsymbol{\lambda};\boldsymbol{u}\right)\geq 0.\tag{7}$$

Therefore, for any  $\lambda \in \Lambda$  it must be that Assumption 4 holds.

#### **Proposition 1.** We have the following:

- (i) If  $\ell$  is convex, then for the mean-field family, Assumption 4 holds.
- (ii) For the Cholesky family, there exists a convex  $\ell$  where Assumption 4 does not hold.

*Proof.* For (i), the key property is the monotonicity of the gradient.

**Proof of (i)** For the mean-field family, recall that

$$C_{ii} = \phi(s_i)$$
.

Also, observe that

$$Cu + m = (C_{11}u_1 + m_1, ..., C_{dd}u_d + m_d).$$

By the property of convex functions,  $\nabla \ell$  is monotone such that

$$\langle \nabla \ell(\mathbf{z}) - \nabla \ell(\mathbf{z}'), \mathbf{z} - \mathbf{z}' \rangle \geq 0.$$

Now, by setting  $\mathbf{z} = \mathbf{C}\mathbf{u} + \mathbf{m}$  and  $\mathbf{z}' = \mathbf{C}\mathbf{u} + \mathbf{m} - C_{ii}u_i\mathbf{e}_i$ , we obtain

$$\langle \nabla \ell \left( \mathbf{C} \mathbf{u} + \mathbf{m} \right) - \nabla \ell \left( \mathbf{C} \mathbf{u} + \mathbf{m} - C_{ii} u_i \mathbf{e}_i \right), C_{ii} u_i \mathbf{e}_i \rangle \geq 0$$

for every i = 1, ..., d.

For the mean-field family,  $Cu + m - C_{ii}u_i\mathbf{e}_i$  is now independent of  $u_i$ . Thus,

$$\mathbb{E} C_{ii}u_i D_i \ell (\mathbf{C}\mathbf{u} + \mathbf{m}) \ge \mathbb{E} C_{ii}u_i D_i \ell (\mathbf{C}\mathbf{u} + \mathbf{m} - \mathbf{e}_i C_{ii}u_i)$$

$$= C_{ii} (\mathbb{E} u_i) (\mathbb{E} D_i f (\mathbf{C}\mathbf{u} + \mathbf{m} - \mathbf{e}_i C_{ii}u_i))$$

$$= 0,$$

where  $D_i f$  denotes the *i*th axis of  $\nabla f$ . Since  $C_{ii} > 0$  by design,

$$\mathbb{E} C_{ii} u_i D_i f(\mathbf{C}\mathbf{u} + \mathbf{m}) > 0 \quad \Leftrightarrow \quad \mathbb{E} u_i D_i f(\mathbf{C}\mathbf{u} + \mathbf{m}) > 0,$$

which is Assumption 4.

**Proof of (ii)** We provide an example that proves the statement. Let  $\ell(z) = \frac{1}{2} z^{T} A z$ . Then,

$$g(\lambda; u) = \ell(\mathcal{T}_{\lambda}(u)) = A(Cu + m) = ACu + Am.$$

Suppose that we choose  $\lambda$  such that

$$\boldsymbol{C} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$$

and m = 0. Also, setting

$$A = \begin{bmatrix} 1 & -2 \\ -2 & 5 \end{bmatrix},$$

we get a strongly convex function  $\ell$ . Then,

$$g(\lambda; \mathbf{u}) = \mathbf{ACu} = \begin{bmatrix} 1 & -2 \\ -2 & 5 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} -1 & -2 \\ 3 & 5 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} -u_1 - 2u_2 \\ 3u_1 + 5u_2 \end{bmatrix}$$

Finally, we have

$$\mathbb{E}g_1(\lambda; \mathbf{u}) u_1 = \mathbb{E}(-u_1 - 2u_2) u_1 = -1 < 0,$$

which violates Assumption 4.

**Lemma 11.** For any function  $f \in C^1(\mathbb{R}, \mathbb{R}_+)$ , there is no constant  $0 < L < \infty$  such that

$$|f(x) - f(y)| \ge L|x - y|.$$

*Proof.* Suppose for the sake of contradiction that such L > 0 exists. Letting  $y \to x$  gives  $|f'(x)| \ge L$  for all  $x \in \mathbb{R}$ . For each x, either  $f'(x) \le -L$  or  $f'(x) \ge L$  holds. We discuss two cases based on the value of f'(0).

If  $f'(0) \ge L$ , we claim that  $f'(x) \ge L$  for all  $x \in \mathbb{R}$ . Otherwise, f'(x) < L for some x implies  $f'(x) \le -L$ . By the intermediate value theorem (f' is continuous), there exists a point y between 0 and x that attains the value f'(y) = 0, which is a contradiction.

Now that  $f'(x) \ge L > 0$  for all x, f is an increasing function. For any x < 0, we have

$$f(x) = f(x) - f(0) + f(0)$$
  
= -|f(x) - f(0)| + f(0)  
\leq -L|x| + f(0).

Here, we can plug  $x' = -\frac{f(0)}{L}$  as

$$f(x') = -L\left|-\frac{f(0)}{L}\right| + f(0) = -|f(0)| + f(0) \le 0,$$

which implies that  $f(x') \notin \mathbb{R}_+$ , which is a contradiction.

Now we discuss the second case  $f'(0) \le -L$ . By a similar argument,  $f'(x) \le -L$  for all  $x \in \mathbb{R}$ . Thus, f is a decreasing function. For any x > 0, we have

$$f(x) = f(x) - f(0) + f(0)$$
  
= -|f(x) - f(0)| + f(0)  
\(\leq -Lx + f(0).

Picking  $x' = \frac{f(0)}{L}$  results in  $f(x') \notin \mathbb{R}_+$ , which is a contradiction.

**Theorem 2.** Let  $\ell$  be  $\mu$ -strongly convex. Then, we have the following:

- (i) If  $\phi$  is linear, the energy f is  $\mu$ -strongly convex.
- (ii) If  $\phi$  is convex, the energy f is convex if and only if Assumption 4 holds.
- (iii) If  $\phi$  is such that  $\phi \in C^1(\mathbb{R}, \mathbb{R}_+)$ , the energy f is not strongly convex.

*Proof.* The special case (i) is proven by Domke (2020, Theorem 9). We focus on the general statement (ii).

If  $\ell$  is  $\mu$ -strongly convex, the inequality

$$\ell(z) - \ell(z') \ge \langle \nabla \ell(z'), z - z' \rangle + \frac{\mu}{2} ||z - z'||_2^2$$
 (8)

holds, where the general convex case is obtained as a special case with  $\mu = 0$ . The goal is to relate this to the ( $\mu$ -strong-)convexity of the energy with respect to the variational parameters given by

$$f(\lambda) - f(\lambda') \ge \langle \nabla_{\lambda} f(\lambda'), \lambda - \lambda' \rangle + \frac{\mu}{2} \|\lambda - \lambda'\|_{2}^{2}.$$

**Proof of (ii)** Plugging the reparameterized latent variables to Equation (8) and taking the expectation, we have

$$\mathbb{E}\ell\left(\mathcal{T}_{\lambda}\left(\mathbf{u}\right)\right) - \mathbb{E}\ell\left(\mathcal{T}_{\lambda'}\left(\mathbf{u}\right)\right) \geq \mathbb{E}\left\langle\nabla\ell\left(\mathcal{T}_{\lambda'}\left(\mathbf{u}\right)\right), \mathcal{T}_{\lambda}\left(\mathbf{u}\right) - \mathcal{T}_{\lambda'}\left(\mathbf{u}\right)\right\rangle + \frac{\mu}{2}\mathbb{E}\|\mathcal{T}_{\lambda}\left(\mathbf{u}\right) - \mathcal{T}_{\lambda'}\left(\mathbf{u}\right)\|_{2}^{2}$$

$$\Leftrightarrow \qquad f\left(\lambda\right) - f\left(\lambda'\right) \geq \mathbb{E}\left\langle\nabla\ell\left(\mathcal{T}_{\lambda'}\left(\mathbf{u}\right)\right), \mathcal{T}_{\lambda}\left(\mathbf{u}\right) - \mathcal{T}_{\lambda'}\left(\mathbf{u}\right)\right\rangle + \frac{\mu}{2}\mathbb{E}\|\mathcal{T}_{\lambda}\left(\mathbf{u}\right) - \mathcal{T}_{\lambda'}\left(\mathbf{u}\right)\|_{2}^{2}$$

$$\Leftrightarrow \qquad f\left(\lambda\right) - f\left(\lambda'\right) \geq \mathbb{E}\left\langle\nabla g\left(\lambda';\mathbf{u}\right), \mathcal{T}_{\lambda}\left(\mathbf{u}\right) - \mathcal{T}_{\lambda'}\left(\mathbf{u}\right)\right\rangle + \frac{\mu}{2}\mathbb{E}\|\mathcal{T}_{\lambda}\left(\mathbf{u}\right) - \mathcal{T}_{\lambda'}\left(\mathbf{u}\right)\|_{2}^{2}$$

Thus, the energy is convex if and only if

$$\mathbb{E}\left\langle g\left(\boldsymbol{\lambda};\boldsymbol{u}\right),\mathcal{T}_{\boldsymbol{\lambda}}\left(\boldsymbol{u}\right)-\mathcal{T}_{\boldsymbol{\lambda}'}\left(\boldsymbol{u}\right)\right\rangle \geq\left\langle \nabla f\left(\boldsymbol{\lambda}\right),\boldsymbol{\lambda}-\boldsymbol{\lambda}'\right\rangle$$

holds. This is established by Lemma 10.

**Proof of (iii)** We now prove that, under the nonlinear parameterization, the energy cannot be strongly convex. When the energy is convex, it is also strongly convex if and only if

$$\frac{\mu}{2}\mathbb{E}\left\|\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)-\mathcal{T}_{\lambda'}\left(\boldsymbol{u}\right)\right\|_{2}^{2}\geq\frac{\mu}{2}\left\|\lambda-\lambda'\right\|_{2}^{2}$$

From the proof of Domke (2020, Lemma 5), it follows that

$$\mathbb{E}\left\|\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right)-\mathcal{T}_{\lambda'}\left(\boldsymbol{u}\right)\right\|_{2}^{2}=\left\|\boldsymbol{C}-\boldsymbol{C}'\right\|_{F}^{2}+\left\|\boldsymbol{m}-\boldsymbol{m}'\right\|_{2}^{2}$$

Furthermore, under nonlinear parameterizations,

$$\|C - C'\|_{F}^{2} + \|m - m'\|_{2}^{2}$$

$$= \|(D_{\phi}(s) - D_{\phi}(s')) - (L - L')\|_{F}^{2} + \|m - m'\|_{2}^{2},$$

expanding the quadratic,

$$= \left\| \boldsymbol{D_{\phi}}\left(\boldsymbol{s}\right) - \boldsymbol{D_{\phi}}\left(\boldsymbol{s'}\right) \right\|_{\mathrm{F}}^{2} + \left\| \boldsymbol{L} - \boldsymbol{L'} \right\|_{\mathrm{F}}^{2} - 2\left\langle \boldsymbol{D_{\phi}}\left(\boldsymbol{s}\right) - \boldsymbol{D_{\phi}}\left(\boldsymbol{s'}\right), \boldsymbol{L} - \boldsymbol{L'}\right\rangle_{\mathrm{F}} + \left\| \boldsymbol{m} - \boldsymbol{m'} \right\|_{2}^{2},$$

and since  $D_{\phi}(s)$  and L reside in different sub-spaces, they are orthogonal. Thus,

$$= \| \mathbf{D}_{\phi}(s) - \mathbf{D}_{\phi}(s') \|_{F}^{2} + \| \mathbf{L} - \mathbf{L}' \|_{F}^{2} + \| \mathbf{m} - \mathbf{m}' \|_{2}^{2}$$

$$= \| \phi(s) - \phi(s') \|_{2}^{2} + \| \mathbf{L} - \mathbf{L}' \|_{F}^{2} + \| \mathbf{m} - \mathbf{m}' \|_{2}^{2}.$$
(9)

For the energy term to be strongly convex, Equation (9) must be bounded *below* by  $\|\lambda - \lambda'\|_2^2$ . Evidently, this implies that a necessary and sufficient condition is that

$$|\phi(s_{ii}) - \phi(s'_{ii})| \ge L |s_{ii} - s'_{ii}|$$

by some constant  $0 < L < \infty$ . Notice that the direction of the inequality is reversed from the Lipschitz condition. Unfortunately, there is no such continuous and differentiable function  $\phi : \mathbb{R} \to \mathbb{R}_+$ , as established by Lemma 11. Thus, for any diagonal conditioner  $\phi \in C^1(\mathbb{R}, \mathbb{R}_+)$ , the energy cannot be strongly convex.

#### F.3 Convergence of Black-Box Variational Inference

#### F.3.1 Vanilla Black-Box Variational Inference

**Theorem 5.** Let the variational family satisfy Assumption 2, the likelihood satisfy Assumption 5, and the assumptions of Corollary 1 hold such that the ELBO, F, is  $L_F$ -smooth with  $L_F = L_\ell + L_\phi + L_s$ . Then, if the stepsize satisfy  $\gamma < 1/L_F$ , the iterates of BBVI with SGD and the M-sample reparameterization gradient estimator satisfy

$$\begin{split} \min_{0 \leq t \leq T-1} \mathbb{E} \left\| \nabla F\left(\boldsymbol{\lambda}_{t}\right) \right\|_{2}^{2} &\leq \gamma \frac{2L_{F}L_{\ell}\kappa\,C\left(d,\varphi\right)}{M} \, \left( \left\| \bar{\boldsymbol{z}}_{\text{joint}} - \bar{\boldsymbol{z}}_{\text{like}} \right\|_{2}^{2} + 2\left(F^{*} - f_{\text{L}}^{*}\right) \right) \\ &+ \frac{2}{\gamma T} \bigg( 1 + \gamma^{2} \frac{4L_{F}L_{\ell}\,\kappa}{M} \, C\left(d,\varphi\right) \bigg)^{T} \left( F\left(\boldsymbol{\lambda}_{0}\right) - F^{*} \right). \end{split}$$

where

$$\begin{split} \bar{\boldsymbol{z}}_{\text{joint}} &= \operatorname{proj}_{\ell}(\boldsymbol{z}) & \text{is the projection of } \boldsymbol{z} \text{ onto set of minimizers of } \ell \\ \bar{\boldsymbol{z}}_{\text{like}} &= \operatorname{proj}_{\ell_{\text{like}}}(\boldsymbol{z}) & \text{is the projection of } \boldsymbol{z} \text{ onto set of minimizers of } \ell_{\text{like}}, \\ \kappa &= L_{\ell}/\mu & \text{is the condition number}, \\ F^* &= \inf_{\boldsymbol{\lambda} \in \Lambda} F(\boldsymbol{\lambda}), \\ \ell^*_{\text{like}} &= \inf_{\boldsymbol{\lambda} \in \mathbb{R}^d} \ell_{\text{like}}(\boldsymbol{z}), \\ C(d, \varphi) &= d + k_{\varphi} & \text{for the Cholesky nonlinear}, \end{split}$$

$$C(d,\varphi) = 2k_{\varphi}\sqrt{d} + 1$$
 for the mean-field nonlinear,

M is the number of Monte Carlo samples.

*Proof.* Khaled & Richtárik (2023, Theorem 2) show that, if the objective function F is  $L_F$ -smooth and the stochastic gradients satisfy the ABC given as

$$\mathbb{E}\|\widehat{\nabla F}\left(\boldsymbol{\lambda}\right)\|_{2}^{2} \leq A\left(F\left(\boldsymbol{\lambda}\right) - F^{*}\right) + B\|\nabla F\|_{2}^{2} + C$$

for some  $0 < A, B, C < \infty$ , SGD guarantees

$$\min_{0 \leq t \leq T-1} \mathbb{E} \left\| \nabla F\left(\boldsymbol{\lambda}_{t}\right) \right\|_{2}^{2} \leq L_{F} C \gamma + \frac{2 \left(1 + L_{F} \gamma^{2} A\right)^{T}}{\gamma T} \left(F\left(\boldsymbol{\lambda}_{0}\right) - F^{*}\right).$$

Under the conditions of Corollary 1, F is  $L_F$ -smooth with  $L_F = L_\ell + L_s + L_\phi$ . Furthermore, under Assumption 5, Kim *et al.* (2023) show that the Monte Carlo gradient estimates satisfy

$$\begin{split} \mathbb{E} \|\widehat{\nabla F}\left(\pmb{\lambda}\right)\|_{2}^{2} &\leq \frac{4L_{\ell}^{2}C\left(d,\varphi\right)}{\mu M}\left(F\left(\pmb{\lambda}\right) - F^{*}\right) + B\|\nabla F\|_{2}^{2} \\ &+ \frac{2L_{\ell}^{2}C\left(d,\varphi\right)}{\mu M}\|\bar{\pmb{z}}_{\mathrm{joint}} - \bar{\pmb{z}}_{\mathrm{like}}\|_{2}^{2} + \frac{4L_{\ell}^{2}C\left(d,\varphi\right)}{\mu M}\left(F^{*} - \ell_{\mathrm{like}}^{*}\right), \end{split}$$

This means that the ABC condition is satisfied with constants

$$A = \frac{4L_{\ell}^{2}}{\mu M} C(d, \varphi), \qquad B = 1, \qquad C = \frac{2L_{\ell}^{2}}{\mu M} C(d, \varphi) \left\| \bar{\mathbf{z}}_{\text{joint}} - \bar{\mathbf{z}}_{\text{like}} \right\|_{2}^{2} + \frac{4L_{\ell}^{2}}{\mu M} C(d, \varphi) \left( F^{*} - \ell_{\text{like}}^{*} \right).$$

Plugging these constants in, we obtain

$$\begin{split} \min_{0 \leq t \leq T-1} \mathbb{E} \left\| \nabla f \left( \lambda_t \right) \right\|_2^2 & \leq \gamma \frac{2 L_F L_\ell^2 C \left( d, \varphi \right)}{\mu M} \left( \left\| \bar{\mathbf{z}}_{\text{joint}} - \bar{\mathbf{z}}_{\text{like}} \right\|_2^2 + 2 \left( F^* - \ell_{\text{like}}^* \right) \right) \\ & + \frac{2}{\gamma T} \left( 1 + \gamma^2 L_F \frac{4 L_\ell^2}{\mu M} C \left( d, \varphi \right) \right)^T \left( F \left( \lambda_0 \right) - F^* \right). \end{split}$$

Substituting the condition number yields the stated result.

**Theorem 3.** Let Assumption 2 hold, the likelihood satisfy Assumption 5, and the assumptions of Corollary 1 hold such that the ELBO F is  $L_F$ -smooth with  $L_F = L_\ell + L_\phi + L_s$ . Then, the iterates generated by BBVI through Equation (1) and the M-sample reparameterization gradient include an  $\epsilon$ -stationary point such that  $\min_{0 \le t \le T-1} \mathbb{E} \|\nabla F(\lambda_t)\|_2 \le \epsilon$  for any  $\epsilon > 0$  if

$$T \ge \mathcal{O}\left(\frac{\left(F\left(\lambda_{0}\right) - F^{*}\right)^{2} L_{F} L_{\ell}^{2} C\left(d, k_{\varphi}\right)}{\mu M \epsilon^{4}}\right)$$

for some fixed stepsize  $\gamma$ , where  $C(d, \varphi) = d + k_{\varphi}$  for the Cholesky family and  $C(d, \varphi) = 2k_{\varphi}\sqrt{d} + 1$  for the mean-field family.

*Proof.* As a corollary to Theorem 5, Khaled & Richtárik (2023, Corollary 1) show that, for an  $L_F$ -smooth objective function F, a gradient estimator satisfying the ABC condition, an  $\epsilon$ -stationary point can be encountered if

$$\gamma = \min\left(\frac{1}{\sqrt{L_FAT}}, \frac{1}{L_FB}, \frac{\epsilon}{2L_FC}\right), \qquad T \geq \frac{12\left(F\left(\lambda_0\right) - F^*\right)L_F}{\epsilon^2} \max\left(B, \frac{12\left(F\left(\lambda_0\right) - F^*\right)A}{\epsilon^2}, \frac{2C}{\epsilon^2}\right).$$

Under Assumption 5, Kim et al. (2023) show that the Monte Carlo gradient estimates satisfy

$$\begin{split} \mathbb{E}\|\widehat{\nabla F}\left(\boldsymbol{\lambda}\right)\|_{2}^{2} &\leq \frac{4L_{\ell}^{2}C\left(d,\varphi\right)}{\mu M}\left(F\left(\boldsymbol{\lambda}\right) - F^{*}\right) + B\|\nabla F\|_{2}^{2} \\ &+ \frac{2L_{\ell}^{2}C\left(d,\varphi\right)}{\mu M}\|\bar{\boldsymbol{z}}_{\text{joint}} - \bar{\boldsymbol{z}}_{\text{like}}\|_{2}^{2} + \frac{4L_{\ell}^{2}C\left(d,\varphi\right)}{\mu M}\left(F^{*} - \ell_{\text{like}}^{*}\right), \end{split}$$

This means that the ABC condition is satisfied with constants

$$A = \frac{4L_f^2}{\mu M} C\left(d,\varphi\right), \qquad B = 1, \qquad C = \frac{2L_f^2}{\mu M} C\left(d,\varphi\right) \left(\left\|\bar{\mathbf{z}}_{\text{joint}} - \bar{\mathbf{z}}_{\text{like}}\right\|_2^2 + 2\left(F^* - f_{\text{L}}^*\right)\right).$$

where

$$ar{z}_{
m joint} = \operatorname{proj}_{\ell}(oldsymbol{z})$$
 is the projection of  $oldsymbol{z}$  onto set of minimizers of  $\ell$  is the projection of  $oldsymbol{z}$  onto set of minimizers of  $\ell$  like, 
$$F^* = \inf_{oldsymbol{\lambda} \in \Lambda} F(oldsymbol{\lambda}),$$
 
$$\ell^*_{
m like} = \inf_{oldsymbol{\lambda} \in \mathbb{R}^d} \ell_{
m like}(oldsymbol{z}),$$
 
$$C(d, \varphi) = d + k_{\varphi} \qquad \text{for the Cholesky family,}$$
 
$$C(d, \varphi) = 2k_{\varphi}\sqrt{d} + 1 \qquad \text{for the mean-field family,}$$
 is the number of Monte Carlo samples.

Plugging these constants in, we obtain

$$T \geq \frac{12\left(F\left(\lambda_{0}\right) - F^{*}\right)L_{F}}{\epsilon^{2}} \max\left(1, \frac{48\left(F\left(\lambda_{0}\right) - F^{*}\right)L_{\ell}^{2}C\left(d, \varphi\right)}{\mu M \epsilon^{2}}, \frac{8L_{\ell}^{2}C\left(d, \varphi\right)\left(\left\|\bar{\mathbf{z}}_{\text{joint}} - \bar{\mathbf{z}}_{\text{like}}\right\|_{2}^{2} + \left(F^{*} - \ell_{\text{like}}^{*}\right)\right)}{\mu M \epsilon^{2}}\right)$$

$$= \mathcal{O}\left(\frac{\left(F\left(\lambda_{0}\right) - F^{*}\right)^{2}L_{F}L_{\ell}^{2}C\left(d\right)}{\mu M \epsilon^{4}}\right),$$

where we omitted the dependence on  $k_{\varphi}$  and the minimizers of  $\ell$  and  $\ell_{\rm like}$ .

#### F.3.2 Proximal Black-Box Variational Inference

**Lemma 3 (Convex Expected Smoothness).** Let  $\ell$  be  $L_{\ell}$ -smooth and  $\mu$ -strongly convex with the variational family satisfying Assumption 2 with the linear parameterization. Then,

$$\mathbb{E}\|\nabla_{\lambda}f(\lambda;\mathbf{u}) - \nabla_{\lambda'}f(\lambda';\mathbf{u})\|_{2}^{2} \leq 2L_{\ell}\kappa C(d,\varphi) B_{f}(\lambda,\lambda')$$

holds, where  $B_f(\lambda, \lambda') \triangleq f(\lambda) - f(\lambda') - \langle \nabla f(\lambda'), \lambda - \lambda' \rangle$  is the Bregman divergence,  $\kappa = L_\ell/\mu$  is the condition number,  $C(d, \varphi) = d + k_\varphi$  for the Cholesky family, and  $C(d, \varphi) = 2k_\varphi\sqrt{d} + 1$  for the mean-field family.

Proof. First, we have

$$\mathbb{E}\|\nabla_{\lambda}f(\lambda;\boldsymbol{u}) - \nabla_{\lambda'}f(\lambda';\boldsymbol{u})\|_{2}^{2} = \mathbb{E}\|\nabla_{\lambda}\ell\left(\mathcal{T}_{\lambda}(\boldsymbol{u})\right) - \nabla_{\lambda'}\ell\left(\mathcal{T}_{\lambda'}(\boldsymbol{u})\right)\|_{2}^{2}$$

$$= \mathbb{E}\left\|\frac{\partial\mathcal{T}_{\lambda}(\boldsymbol{u})}{\partial\lambda}g(\lambda,\boldsymbol{u}) - \frac{\partial\mathcal{T}_{\lambda'}(\boldsymbol{u})}{\partial\lambda'}g(\lambda',\boldsymbol{u})\right\|_{2}^{2}.$$

For the linear parameterization, the Jacobian of  $\mathcal{T}_{\lambda}$  does not depend on  $\lambda$ . Therefore,

$$= \mathbb{E} \left\| \frac{\partial \mathcal{T}_{\lambda}(\mathbf{u})}{\partial \lambda} \left( g(\lambda, \mathbf{u}) - g(\lambda', \mathbf{u}) \right) \right\|_{2}^{2}$$

and Lemma 6 yields

$$=J_{\mathcal{T}}(\boldsymbol{u})\mathbb{E}\|g(\boldsymbol{\lambda},\boldsymbol{u})-g(\boldsymbol{\lambda}',\boldsymbol{u})\|_{2}^{2},$$

where

$$J_{\mathcal{T}}(\boldsymbol{u}) = 1 + \|\boldsymbol{u}\|_{2}^{2}$$
 for the Cholesky family and  $J_{\mathcal{T}}(\boldsymbol{u}) = 1 + \|\boldsymbol{U}^{2}\|_{F}$  for the mean-field family.

From now on, we apply the strategy of Domke (2019, Theorem 3) for resolving the randomness u. That is,

$$\mathbb{E}J_{\mathcal{T}}(\boldsymbol{u}) \| \mathbf{g}(\boldsymbol{\lambda}, \boldsymbol{u}) - \mathbf{g}(\boldsymbol{\lambda}', \boldsymbol{u}) \|_{2}^{2} = J_{\mathcal{T}}(\boldsymbol{u}) \| \nabla \ell \left( \mathcal{T}_{\boldsymbol{\lambda}}(\boldsymbol{u}) \right) - \nabla \ell \left( \mathcal{T}_{\boldsymbol{\lambda}'}(\boldsymbol{u}) \right) \|_{2}^{2}$$

from the  $L_{\ell}$ -smoothness of f,

$$\leq L_{\ell}^{2} \mathbb{E} J_{\mathcal{T}}(\boldsymbol{u}) \left\| \mathcal{T}_{\lambda}(\boldsymbol{u}) - \mathcal{T}_{\lambda'}(\boldsymbol{u}) \right\|_{2}^{2},$$

and applying Corollary 2,

$$\leq L_{\ell}^{2} C(d, \varphi) \|\lambda - \lambda'\|_{2}^{2}$$

The last step follows the approach of Kim *et al.* (2023), where we convert the quadratic bound into a bound involving the energy. Recall that the  $\mu$ -strongly convexity of  $\ell$  implies

$$\frac{\mu}{2} \| \mathbf{z}' - \mathbf{z} \|_{2}^{2} \le \ell(\mathbf{z}) - \ell(\mathbf{z}') - \langle \nabla \ell(\mathbf{z}'), \mathbf{z} - \mathbf{z}' \rangle. \tag{10}$$

From Lemma 8, we have

$$L_{\ell}^{2} C\left(d,\varphi\right) \left\|\lambda - \lambda'\right\|_{2}^{2} = L_{f}^{2} C\left(d,\varphi\right) \mathbb{E}\left\|\mathcal{T}_{\lambda}\left(\boldsymbol{u}\right) - \mathcal{T}_{\lambda'}\left(\boldsymbol{u}\right)\right\|_{2}^{2},$$

and by  $\mu$ -strongly convexity,

$$\leq \frac{2L_{\ell}^{2}}{\mu} C(d, \varphi) \mathbb{E}\left(\ell\left(\mathcal{T}_{\lambda}(\boldsymbol{u})\right) - \ell\left(\mathcal{T}_{\lambda'}(\boldsymbol{u})\right) - \langle \nabla \ell\left(\mathcal{T}_{\lambda'}(\boldsymbol{u})\right), \mathcal{T}_{\lambda}(\boldsymbol{u}) - \mathcal{T}_{\lambda'}(\boldsymbol{u})\rangle\right) \\
= \frac{2L_{\ell}^{2}}{\mu} C(d, \varphi) \mathbb{E}\left(f\left(\lambda; \boldsymbol{u}\right) - f\left(\lambda'; \boldsymbol{u}\right) - \langle g\left(\lambda'; \boldsymbol{u}\right), \mathcal{T}_{\lambda}(\boldsymbol{u}) - \mathcal{T}_{\lambda'}(\boldsymbol{u})\rangle\right) \\
= \frac{2L_{\ell}^{2}}{\mu} C(d, \varphi) \left(f\left(\lambda\right) - f\left(\lambda'\right) - \mathbb{E}\left\langle g\left(\lambda'; \boldsymbol{u}\right), \mathcal{T}_{\lambda}(\boldsymbol{u}) - \mathcal{T}_{\lambda'}(\boldsymbol{u})\rangle\right).$$

Finally, by applying the equality in Lemma 10,

$$=\frac{2L_{\ell}^{2}}{u}C(d,\varphi)\big(f(\lambda)-f(\lambda')-\langle\nabla f(\lambda'),\lambda-\lambda'\rangle\big).$$

**Lemma 12** (Variance Transfer). Let  $\ell$  be  $L_{\ell}$ -smooth and  $\mu$ -strongly convex with the variational family satisfying Assumption 2 with the linear parameterization. Also, let  $\widehat{\nabla f}$  be an M-sample gradient estimator of the energy. Then,

$$\operatorname{tr} \mathbb{V} \, \widehat{\nabla f} \, (\lambda) \leq \frac{4 L_{\ell} \kappa \, C \, (d, \varphi)}{M} \, \mathrm{B}_{f} \, (\lambda, \lambda') + 2 \operatorname{tr} \mathbb{V} \, \widehat{\nabla f} \, (\lambda') \,,$$

 $\kappa = L_{\ell}/\mu$  is the condition number,  $B_f$  is the Bregman divergence defined in Lemma 3,  $C(d, \varphi) = d + k_{\varphi}$  for the Cholesky family, and  $C(d, \varphi) = 2k_{\varphi}\sqrt{d} + 1$  for the mean-field family.

*Proof.* First, the M-sample gradient estimator is defined as

$$\widehat{\nabla f}(\lambda) = \frac{1}{M} \sum_{m=1}^{M} \nabla_{\lambda} f(\lambda; \mathbf{u}_{m}),$$

where  $u_m \sim \varphi$ . Since  $u_1, \dots, u_m$  are independent and identically distributed, we have

$$\operatorname{tr} \mathbb{V} \widehat{\nabla f}(\lambda) = \frac{1}{M} \operatorname{tr} \mathbb{V} \nabla_{\lambda} f(\lambda; \boldsymbol{u}).$$

From here, given Lemma 3, the proof is identical with that of Garrigos & Gower (2023, Lemma 8.20), except for the constants.

**Theorem 6.** Let  $\ell$  be  $L_\ell$ -smooth and  $\mu$ -strongly convex. Then, BBVI with proximal SGD in Equation (2), M-Monte Carlo samples, a variational family satisfying Assumption 2, the linear parameterization, and a fixed stepsize  $0 < \gamma \le \frac{M}{2L_\ell \, \kappa \, C(d, \varphi)}$ , the iterates satisfy

$$\mathbb{E}\|\boldsymbol{\lambda}_{T}-\boldsymbol{\lambda}^{*}\|_{2}^{2} \leq (1-\gamma\mu)^{T}\|\boldsymbol{\lambda}_{0}-\boldsymbol{\lambda}^{2}\|_{2}^{2}+\frac{2\gamma\sigma^{2}}{\mu},$$

where  $\kappa = L_{\ell}/\mu$  is the condition number,  $\sigma^2$  is defined in Lemma 4,  $\lambda^* = \arg\min_{\lambda \in \Lambda} F(\lambda)$ ,  $C(d, \varphi) = d + k_{\varphi}$  for the Cholesky family, and  $C(d, \varphi) = 2k_{\varphi}\sqrt{d} + 1$  for the mean-field family.

Proof. Provided that

**(A.6.1)** the energy f is  $\mu$ -strongly convex,

**(A.6.2)** the energy f is  $L_{\ell}$ -smooth,

(A.6.3) the regularizer h is convex,

(A.6.4) the regularizer h is lower semi-continuous,

(A.6.5) the convex expected smoothness condition holds,

(A.6.6) the variance transfer condition holds, and

(A.6.7) the gradient variance  $\sigma^2$  at the optimum is finite such that  $\sigma^2 < \infty$ ,

the proof is identical to that of Garrigos & Gower (2023, Theorem 11.9), which is based on the results of Gorbunov *et al.* (2020, Corollary A.2).

In our setting,

(A.6.1) is established by Theorem 2,

(A.6.2) is established by Theorem 1,

(A.6.3) is trivially satisfied since h is the negative entropy,

(A.6.4) is trivially satisfied since h is continuous,

(A.6.5) is established in Lemma 3,

(A.6.6) is established in Lemma 12,

(A.6.7) is established in Lemma 4.

The only difference is that, we replace the constant  $L_{\max}$  in the proof of Garrigos & Gower to  $L_{\ell} \kappa C(d, \varphi)/M$ . This stems from the different constants in the variance transfer condition.

**Theorem 7.** Let  $\ell$  be  $L_{\ell}$ -smooth and  $\mu$ -strongly convex. Then, for any  $\epsilon > 0$ , BBVI with proximal SGD in Equation (2), M-Monte Carlo samples, a variational family satisfying Assumption 2, and the linear parameterization guarantees  $\mathbb{E}\|\lambda_T - \lambda^*\|_2^2 \leq \epsilon$  if

$$\gamma = \min \left( \frac{\epsilon}{2} \frac{\mu}{2\sigma^2}, \frac{M}{2L_{\ell} \kappa C(d, \varphi)} \right), \qquad T \ge \max \left( \frac{1}{\epsilon} \frac{4\sigma^2}{\mu^2}, \frac{2\kappa^2 C(d, \varphi)}{M} \right) \log \left( \frac{2\|\lambda_0 - \lambda^*\|}{\epsilon} \right),$$

where  $\kappa = L_e/\mu$ ,  $\sigma^2$  is defined in Lemma 4,  $\lambda^* = \arg\min_{\lambda \in \Lambda} F(\lambda)$ ,  $C(d, \varphi) = d + k_{\varphi}$  for the Cholesky family, and  $C(d, \varphi) = 2k_{\varphi}\sqrt{d} + 1$  for the mean-field family.

*Proof.* This is a corollary of the fixed stepsize convergence guarantee in Theorem 6 as shown by Garrigos & Gower (2023, Corollary 11.10). They guarantee an  $\epsilon$ -accurate solution as long as

$$\gamma = \min\left(\frac{\epsilon}{2} \frac{2}{2\sigma_{\mathrm{F}}^*}, \frac{1}{2L_{\mathrm{max}}}\right), \quad T \geq \max\left(\frac{1}{\epsilon} \frac{4\sigma_{\mathrm{F}}^*}{\mu^2}, \frac{2L_{\mathrm{max}}}{\mu}\right) \log\left(\frac{2\|\lambda_0 - \lambda^*\|}{\epsilon}\right).$$

In our notation,  $\sigma_F^* = \sigma^2$  and  $L_{\text{max}} = L_{\ell} \kappa C(d, \varphi)/M$ .

**Theorem 8.** Let  $\ell$  be  $L_{\ell}$ -smooth and  $\mu$ -strongly convex. Then, BBVI with proximal SGD in Equation (2), the M-sample reparameterization gradient estimator, a variational family satisfying Assumption 2, the linear parameterization,  $T \geq 4T_{\kappa}$ , and a stepsize schedule of

$$\gamma_t = \begin{cases} \frac{M}{2L_{\ell}\kappa C(d,\varphi)} & for \quad t \leq 4T_{\kappa} \\ \frac{2t+1}{(t+1)^2\mu} & for \quad t > 4T_{\kappa}, \end{cases}$$

where  $T_{\kappa} = [\kappa^2 C(d, \varphi) M^{-1}]$ ,  $\kappa = L_{\ell}/\mu$  is the condition number,  $C(d, \varphi) = d + k_{\varphi}$  for the Cholesky family, and  $C(d, \varphi) = 2k_{\varphi}\sqrt{d} + 1$  for the mean-field family, then the iterates satisfy

$$\mathbb{E}\|\boldsymbol{\lambda}_{T}-\boldsymbol{\lambda}^{*}\|_{2}^{2} \leq \frac{16 T_{\kappa}^{2} \|\boldsymbol{\lambda}_{0}-\boldsymbol{\lambda}^{*}\|_{2}^{2}}{e^{2}T^{2}} + \frac{8\sigma^{2}}{\mu^{2}T}$$

where  $\sigma^2$  is defined in Lemma 4, e is Euler's constant, and  $\lambda^* = \arg\min_{\lambda \in \Lambda} F(\lambda)$ .

*Proof.* Under our assumptions, Theorem 6 holds, of which the proof is essentially obtaining the recursion

$$\mathbb{E}\|\boldsymbol{\lambda}_{t+1} - \boldsymbol{\lambda}^*\|_2^2 = (1 - \gamma_t \mu) \, \mathbb{E}\|\boldsymbol{\lambda}_t - \boldsymbol{\lambda}^*\|_2^2 + 2\gamma_t^2 \sigma^2.$$

Instead of a fixed stepsize, we can apply the decreasing stepsize rule in the proof statement, then which the proof becomes identical to that of Gower *et al.* (2019, Theorem 3.2). We only need to replace  $\mathcal{L}$  with  $L_{\max}$  in the proof of Garrigos & Gower (2023, Theorem 11.9). This, in our notation, is  $L_{\max} = L_{\ell} \kappa C(d, \varphi)/M$ .

**Theorem 4.** Let  $\ell$  be  $L_{\ell}$ -smooth and  $\mu$ -strongly convex. Then, for any  $\epsilon > 0$ , BBVI with proximal SGD in Equation (2), the M-sample reparameterization gradient estimator, a variational family satisfying Assumption 2 with the linear parameterization guarantees  $\mathbb{E}\|\lambda_T - \lambda^*\|_2^2 \le \epsilon$  if

satisfying Assumption 2 with the linear parameterization guarantees 
$$\mathbb{E}\|\lambda_{T} - \lambda^{*}\|_{2}^{2} \leq \epsilon$$
 if 
$$\gamma_{t} = \begin{cases} \frac{M}{2L_{\rho\kappa}C(d,\varphi)} & \text{for } t \leq 4T_{\kappa} \\ \frac{2i+1}{(t+1)^{2}\mu} & \text{for } t > 4T_{\kappa}, \end{cases} \qquad T \geq \max\left(\frac{8\sigma^{2}}{\mu^{2}\epsilon} + \frac{4T_{\kappa}\|\lambda_{0} - \lambda^{*}\|_{2}}{e\sqrt{\epsilon}}, 4T_{\kappa}\right)$$

where  $\sigma^2$  is defined in Lemma 4,  $T_{\kappa} = [\kappa^2 C(d, \varphi) M^{-1}]$ ,  $\kappa = L_{\ell}/\mu$  is the condition number, e is Euler's constant,  $\lambda^* = \arg\min_{\lambda \in \Lambda} F(\lambda)$ ,  $C(d, \varphi) = d + k_{\varphi}$  for the Cholesky family, and  $C(d, \varphi) = 2k_{\varphi}\sqrt{d} + 1$  for the mean-field family.

*Proof.* The computational complexity follows from the smallest number of iterations T such that

$$\mathbb{E}\|\boldsymbol{\lambda}_T - \boldsymbol{\lambda}^*\|_2^2 \le \frac{16T_\kappa^2\|\boldsymbol{\lambda}_0 - \boldsymbol{\lambda}^*\|_2^2}{\mathrm{e}^2T^2} + \frac{8\sigma^2}{\mu^2T} \le \epsilon$$

By multiplying both sides with  $T^2$  as

$$T^{2}\epsilon - \frac{8\sigma^{2}}{\mu^{2}}T - \frac{16T_{\kappa}^{2}\|\lambda_{0} - \lambda^{*}\|_{2}^{2}}{e^{2}} \ge 0,$$
(11)

we can see that we are looking for the smallest positive integer that is larger than the solution of a quadratic equation with respect to T. This is given as

$$T \geq \frac{\frac{8\sigma^2}{\mu^2} + \sqrt{\left(\frac{8\sigma^2}{\mu^2}\right)^2 + 64\epsilon \frac{T_K^2 \|\lambda_0 - \lambda^*\|_2^2}{e^2}}}{2\epsilon}.$$

Applying the inequality  $\sqrt{a+b} \le \sqrt{a} + \sqrt{b}$ ,

$$\frac{\frac{8\sigma^{2}}{\mu^{2}} + \sqrt{\left(\frac{8\sigma^{2}}{\mu^{2}}\right)^{2} + 64\varepsilon \frac{T_{\kappa}^{2} \|\lambda_{0} - \lambda^{*}\|_{2}^{2}}}{2\varepsilon}}{2\varepsilon} \leq \frac{\frac{8\sigma^{2}}{\mu^{2}} + \left(\frac{8\sigma^{2}}{\mu^{2}}\right) + \sqrt{64\varepsilon \frac{T_{\kappa}^{2} \|\lambda_{0} - \lambda^{*}\|_{2}^{2}}{e^{2}}}}{2\varepsilon}$$
$$= \frac{\frac{16\sigma^{2}}{\mu^{2}} + \sqrt{\varepsilon} \frac{8T_{\kappa} \|\lambda_{0} - \lambda^{*}\|_{2}}{e}}{2\varepsilon}$$
$$= \frac{8\sigma^{2}}{\mu^{2}\varepsilon} + \frac{4T_{\kappa} \|\lambda_{0} - \lambda^{*}\|_{2}}{e\sqrt{\varepsilon}}.$$

Thus,  $\mathbb{E}\|\pmb{\lambda}_T - \pmb{\lambda}^*\|_2^2 \leq \epsilon$  can be satisfied with a number of iterations at least

$$T \ge \max\left(\frac{8\sigma^2}{\mu^2 \epsilon} + \frac{4T_{\kappa} \|\lambda_0 - \lambda^*\|_2}{e\sqrt{\epsilon}}, \ 4T_{\kappa}\right).$$

#### **G** Details of Experimental Setup

Table 2: Summary of Datasets and Problems

Abbrev.	Model	Dataset	d	N
LME-election LME-radon	Linear Mixed Effects	1988 U.S. presidential election (Gelman & Hill, 2007) U.S. household radon levels (Gelman & Hill, 2007)	90 391	11,566 12,573
BT-tennis	Bradley-Terry	ATP World Tour tennis	6030	172,199
LR-keggu LR-song LR-buzz LR-electric	Linear Regression	KEGG-undirected (Shannon <i>et al.</i> , 2003) million songs (Bertin-Mahieux <i>et al.</i> , 2011) buzz in social media (Kawala <i>et al.</i> , 2013) household electric	31 94 81 15	63,608 515,345 583,250 2,049,280
AR-ecg	Sparse Autoregression	Long-term ST ECG (Jager et al., 2003)	63	20,642,000

Linear Regression (LR-\*) We consider a basic Bayesian hierarchical linear regression model

$$\begin{split} & \sigma_{\alpha} \sim \mathcal{N}_{+}\left(0, 10^{2}\right), \quad \sigma_{\beta} \sim \mathcal{N}_{+}\left(0, 10^{2}\right), \quad \sigma \sim \mathcal{N}_{+}\left(0, 0.3^{2}\right) \\ & \boldsymbol{\beta} \sim \mathcal{N}\left(\boldsymbol{0}, \sigma_{\beta}^{2} \boldsymbol{I}\right), \quad \boldsymbol{\alpha} \sim \mathcal{N}\left(0, \sigma_{\alpha}^{2}\right), \\ & y_{i} \sim \mathcal{N}\left(\boldsymbol{\beta}^{\mathsf{T}} \boldsymbol{x}_{i} + \boldsymbol{\alpha}, \, \sigma^{2}\right), \end{split}$$

where a weakly informative half-normal hyperprior  $\mathcal{N}_+$ , a normal distribution with the support restricted to  $\mathbb{R}_+$ , is assigned on the hyperparameters. For the datasets, we consider large-scale regression problems obtained from the UCI repository (Dua & Graff, 2017), shown in Table 2. For all datasets, we standardize the regressors  $x_i$  and the outcomes  $y_i$ .

**Radon Levels (MLE-radon)** MLE-radon is a radon level regression problem by Gelman & Hill (2007). It fits a hierarchical mixed-effects model for estimating household radon levels across different counties while considering the floor elevation of each site. The model is described as

$$\begin{split} \sigma &\sim \mathcal{N}_{+}\left(0,1^{2}\right), \quad \sigma_{\alpha} \sim \mathcal{N}_{+}\left(0,1^{2}\right), \quad \mu_{\alpha} \sim \mathcal{N}\left(0,10^{2}\right), \quad \boldsymbol{\varepsilon} \sim \mathcal{N}\left(\mathbf{0},10^{2}\mathbf{I}\right) \\ \beta_{1} &\sim \mathcal{N}\left(0,10^{2}\right), \quad \beta_{2} \sim \mathcal{N}\left(0,10^{2}\right) \\ \boldsymbol{\alpha} &= \mu_{\alpha} + \sigma_{\alpha}\boldsymbol{\varepsilon} \\ \mu_{i} &= \boldsymbol{\alpha}[\mathrm{county}_{i}] + \beta_{1}\log\left(\mathrm{uppm}_{i}\right) + \mathrm{floor}_{i}\,\beta_{2} \\ \log \mathrm{radon}_{i} &\sim \mathcal{N}\left(\mu_{i},\sigma^{2}\right), \end{split}$$

which uses variable slopes and intercepts with non-centered parameterization. The dataset was obtained from PosteriorDB (Magnusson *et al.*, 2022). Also, for the radon regression problem, the Minnesota subset is often used due to computational reasons. Here, we use the full national dataset.

**Presidential Election (MLE-election)** MLE-election is a model for studying the effects of sociological factors on the 1988 United States presidential election (Gelman & Hill, 2007). The model is described as

$$\begin{split} &\sigma_{\text{age}} \sim \mathcal{N}\left(0, 100^2\right), \quad \sigma_{\text{edu}} \sim \mathcal{N}\left(0, 100^2\right), \quad \sigma_{\text{age} \times \text{edu}} \sim \mathcal{N}\left(0, 100^2\right) \\ &\sigma_{\text{state}} \sim \mathcal{N}\left(0, 100^2\right), \quad \sigma_{\text{region}} \sim \mathcal{N}\left(0, 100^2\right), \\ &\boldsymbol{b}_{\text{age}} \sim \mathcal{N}\left(\boldsymbol{0}, \sigma_{\text{age}}^2 \mathbf{I}\right), \quad \boldsymbol{b}_{\text{edu}} \sim \mathcal{N}\left(\boldsymbol{0}, \sigma_{\text{edu}}^2 \mathbf{I}\right), \quad \boldsymbol{b}_{\text{age} \times \text{edu}} \sim \mathcal{N}\left(\boldsymbol{0}, \sigma_{\text{age} \times \text{edu}}^2 \mathbf{I}\right), \\ &\boldsymbol{b}_{\text{state}} \sim \mathcal{N}\left(\boldsymbol{0}, \sigma_{\text{state}}^2 \mathbf{I}\right), \quad \boldsymbol{b}_{\text{region}} \sim \mathcal{N}\left(\boldsymbol{0}, \sigma_{\text{region}}^2 \mathbf{I}\right) \\ &\boldsymbol{\beta} \sim \mathcal{N}\left(\boldsymbol{0}, 100^2 \mathbf{I}\right) \\ &\boldsymbol{\beta}_{i} = \boldsymbol{\beta}_1 + \boldsymbol{\beta}_2 \, \text{black}_i + \boldsymbol{\beta}_3 \, \text{female}_i + \boldsymbol{\beta}_4 \, \boldsymbol{v}_{\text{prev},i} + \boldsymbol{\beta}_5 \, \text{female}_i \, \text{black}_i \\ &+ \boldsymbol{b}_{\text{age}}[\text{age}_i] + \boldsymbol{b}_{\text{edu}}[\text{edu}_i] + \boldsymbol{b}_{\text{age} \times \text{edu}}[\text{age}_i \, \text{edu}_i] + \boldsymbol{b}_{\text{state}}[\text{state}_i] + \boldsymbol{b}_{\text{region}}[\text{region}_i] \\ &\boldsymbol{y}_i \sim \text{bernoulli}\left(\boldsymbol{p}_i\right). \end{split}$$

The dataset was obtained from PosteriorDB (Magnusson et al., 2022).

**Bradley-Terry (BT-Tennis)** BT-Tennis is a Bradley-Terry model for estimating the skill of professional tennis players used by Giordano *et al.* (2023). The model is described as

$$\sigma \sim \mathcal{N}_{+}(0, 1)$$

$$\theta \sim \mathcal{N}(\mathbf{0}, \sigma^{2}\mathbf{I})$$

$$p_{i} \sim \theta[\min_{i}] - \theta[\log_{i}]$$

$$v_{i} \sim \text{bernoulli}(p),$$

where  $win_i$ ,  $los_i$  are the indices of the winning and losing players for the *i*th game, respectively. While we subsample over the games i = 1, ..., N, each player's involvement is sparse in that each player plays only a handful of games. Consequently, the subsampling noise is substantial. Therefore, we use a larger batch size of 500. Similarly to Giordano *et al.* (2023), we use the ATP World Tour data publically available online <sup>1</sup>.

**Autoregression (AR-ecg)** AR-ecg is a linear autoregressive model. Here, we use a Student-t likelihood as originally proposed by Christmas & Everson (2011). While they originally imposed an automatic relevance detection prior on the autoregressive coefficients, we instead set a horseshoe shrinkage prior (Carvalho *et al.*, 2009, 2010). Since the horseshoe is known to result in complex posterior geometry, this should make the problem more challenging. The model is described as

```
\begin{split} &\alpha_d = 10^{-2}, \quad \beta_d = 10^{-2}, \quad \alpha_d = 10^{-2}, \quad \beta_d = 10^{-2}, \\ &d \sim \operatorname{gamma}\left(\alpha_d, \beta_d\right), \\ &\sigma^{-1} \sim \operatorname{inverse-gamma}\left(\alpha_\sigma, \beta_\sigma\right), \\ &\tau \sim \operatorname{cauchy}_+\left(0, 1\right), \\ &\lambda \sim \operatorname{cauchy}_+\left(\mathbf{0}, \mathbf{1}\right), \\ &\theta \sim \mathcal{N}\left(0, \tau \operatorname{diag}\left(\lambda\right)\right) \\ &y[n] \sim \operatorname{stduent-t}\left(d, \, \theta_1 y[n-1] + \theta_2 y[n-2] + \dots + \theta_P y[n-P], \sigma\right), \end{split}
```

where d is the degrees-of-freedom for the Student-t likelihood, cauchy  $_{+}$  is a half-Cauchy prior.

For the dataset, we use the long-term electrocardiogram measurements of Jager *et al.* (2003) obtained from Physionet (Goldberger *et al.*, 2000). The data instance we used has a duration of 23 hours sampled at 250 Hz with 12-bit resolution over a range of  $\pm 10$  millivolts. During the experiments, we observed that the hyperparameters suggested by Christmas & Everson are sensitive to the signal amplitude. Therefore, we scaled the signal amplitude to be  $\pm 10$ .

Ihttps://datahub.io/sports-data/atp-world-tour-tennis-data

#### H Additional Experimental Results

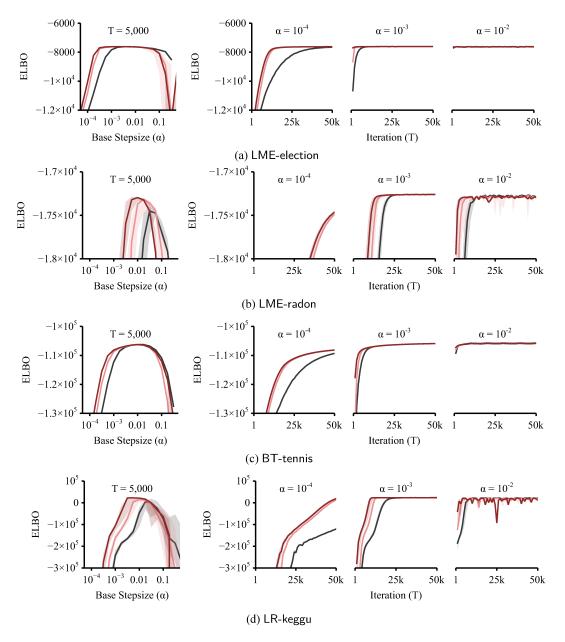


Figure 5: BBVI convergence speed (ELBO v.s. Iteration) and robustness against stepsize (ELBO at T=50,000 v.s. Base stepsize). The error bands are the 80% quantiles estimated from 20 (10 for AR-eeg) independent replications. The initial point was  $m_0 = 0$ ,  $C_0 = I$ .

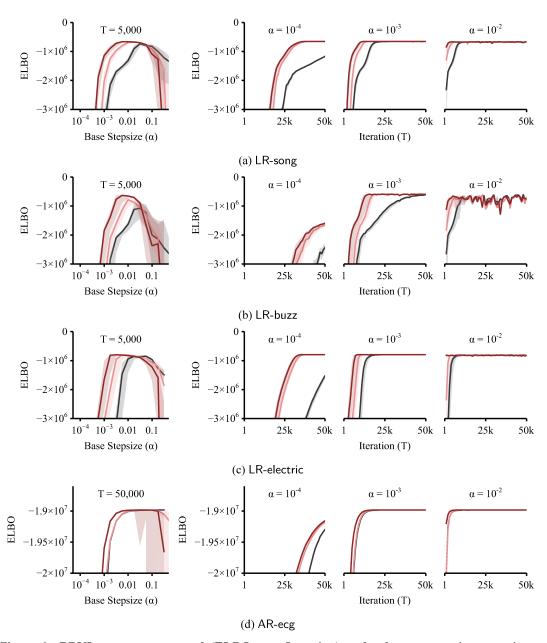


Figure 6: BBVI convergence speed (ELBO v.s. Iteration) and robustness against stepsize (ELBO at T=50,000 v.s. Base stepsize). The error bands are the 80% quantiles estimated from 20 (10 for AR-eeg) independent replications. The initial point was  $m_0 = 0$ ,  $C_0 = I$ .