

A critical and theoretical analysis of unbiased implicit variational inference

Kenny Chiu

April 7, 2022

1 Critical analysis

1.1 Introduction

Titsias and Ruiz [20] introduced *unbiased implicit variational inference* (UIVI) as a variational inference method with a flexible variational family that addresses the issues of the existing methods that it is built on. In this analysis, we summarize the work of Titsias and Ruiz [20] in the context of the literature and discuss the strengths and limitations of UIVI.

This analysis is organized as follows: Section 1.2 introduces the problem context and existing methods; Section 1.3 explains how UIVI works, how it addresses the limitations of previous methods, and its own limitations; and Section 1.4 highlights related work in the recent literature and discusses the general direction that the literature is moving towards.

1.2 Context and previous work

Variational inference (VI) [6] is a method for Bayesian inference that formulates the problem of finding the posterior distribution $p(\mathbf{z}|\mathbf{x})$ of latent variables \mathbf{z} given data \mathbf{x} as an optimization problem. VI posits a variational family $\mathcal{Q} = \{q_\theta\}$ of distributions indexed by variational parameters θ , and the goal is to identify the variational distribution $q_\theta(\mathbf{z}) \in \mathcal{Q}$ that best approximates the posterior distribution. In standard VI, the selected distribution q_θ is the one that minimizes the Kullback-Leibler (KL) divergence of q_θ and $p(\mathbf{z}|\mathbf{x})$, or equivalently, the one that maximizes the evidence lower bound (ELBO) defined as

$$\mathcal{L}(\theta) = \mathbb{E}_{q_\theta(\mathbf{z})} [\log p(\mathbf{x}, \mathbf{z}) - \log q_\theta(\mathbf{z})] .$$

Standard VI maximizes the ELBO using a coordinate ascent algorithm, which requires strong restrictions on the choice of the model and the variational family. These restrictions include (1) a mean-field assumption where the latent variables \mathbf{z} are marginally independent and the variational distribution factorizes as $q_\theta(\mathbf{z}) = \prod_{i=1}^d q_{\theta_i}(\mathbf{z}_i)$, and (2) a conjugate model where $p(\mathbf{z}_i)$ and $p(\mathbf{z}_i|\mathbf{x}, \mathbf{z}_{-i})$ are from the same distribution family. The consequence of these restrictions is a variational family that is often limited to be analytical (i.e., a subfamily of the exponential family) and in which marginal dependencies between the latent variables could not be modeled.

An important development after standard VI was black box VI (BBVI) [14], which relaxed the restrictive assumptions by optimizing the ELBO using a different approach. By rewriting the ELBO gradient in terms of an expectation, the gradient can be estimated unbiasedly and cheaply using Monte Carlo samples. The optimization approach of BBVI trades the restrictions of standard VI for the assumption that one can sample from the variational distribution $q_\theta(\mathbf{z})$. This weaker assumption expanded the possibilities for the choice of the variational family. One such proposed family was the hierarchical variational model (HVM) [15] containing distributions of the form $q_\theta(\mathbf{z}) = \int q_\theta(\mathbf{z}|\boldsymbol{\varepsilon})q_\theta(\boldsymbol{\varepsilon})d\boldsymbol{\varepsilon}$. An advantage of these hierarchical distributions over other variational distributions is the ease in being able to capture marginal dependencies between latent variables through the mixing distribution $q_\theta(\boldsymbol{\varepsilon})$.

Further pushing the assumption that one only needs to be able to sample from the variational distribution, one trend following the introduction of the HVM was the incorporation of deep neural networks to expand the modeling capacity of the hierarchical variational family. These models took various forms, such as through normalizing flows [17] or through implicit distributions [9] involving deep networks in which the density cannot be evaluated. Though the implicit models are flexible, the log density ratio in the ELBO is intractable in these models. Some works proposed using density ratio estimation to tackle this problem [e.g., 5, 9], but this approach is known to struggle in high-dimensional regimes [19].

The method that precedes UIVI and that was proposed to address the challenges of using implicit distributions in hierarchical variational models is *semi-implicit* VI (SIVI) [21]. SIVI makes use of a semi-implicit variational distribution in which (1) the variational conditional $q_\theta(\mathbf{z}|\boldsymbol{\varepsilon}) = q(\mathbf{z}|\boldsymbol{\varepsilon})$ is required to be reparameterizable and explicit, and (2) the mixing distribution $q_\theta(\boldsymbol{\varepsilon})$ is also required to be reparameterizable but possibly implicit.

SIVI then avoids the density ratio estimation problem by instead optimizing a lower bound for the ELBO calculated using samples that is only exact as the number of samples goes to infinity [10, 21].

1.3 Current work

Titsias and Ruiz [20] proposed UIVI as an alternative to SIVI that directly maximizes the ELBO as an objective rather than a surrogate lower bound. The motivation is that by directly optimizing the ELBO objective, UIVI should be more efficient than SIVI and therefore should result in faster convergence to the optimal variational approximation. UIVI allows for an unbiased ELBO objective by rewriting the ELBO gradient in terms of two expectations. One expectation is easily estimated using Monte Carlo samples, while the other expectation is over an inverse conditional for which UIVI estimates based on samples obtained through Markov chain Monte Carlo (MCMC) methods.

1.3.1 Unbiased implicit variational inference

Like in SIVI, UIVI starts with a hierarchical variational model setup where the variational distribution is

$$q_\theta(\mathbf{z}) = \int q_\theta(\mathbf{z}|\boldsymbol{\varepsilon})q(\boldsymbol{\varepsilon})d\boldsymbol{\varepsilon}.$$

UIVI requires the variational conditional $q_\theta(\mathbf{z}|\boldsymbol{\varepsilon})$ to be reparameterizable, i.e., that any sample $\mathbf{z} \sim q_\theta(\mathbf{z}|\boldsymbol{\varepsilon})$ can be rewritten as

$$\mathbf{z} = h_\theta(\mathbf{u}; \boldsymbol{\varepsilon}) := h_{\boldsymbol{\psi}=g_\theta(\boldsymbol{\varepsilon})}(\mathbf{u})$$

where h_ψ is some reparameterization function with parameters $\boldsymbol{\psi}$ that are the output of some arbitrarily complex function g_θ that depends on variational parameters θ and input $\boldsymbol{\varepsilon}$. To sample from $q_\theta(\mathbf{z})$, noise variables $\mathbf{u} \sim q(\mathbf{u})$ and $\boldsymbol{\varepsilon} \sim q(\boldsymbol{\varepsilon})$ are first sampled from fixed auxiliary distributions and then fed through h_θ . UIVI also requires that $q_\theta(\mathbf{z}|\boldsymbol{\varepsilon})$ and its log-gradient $\nabla_{\mathbf{z}} \log q_\theta(\mathbf{z}|\boldsymbol{\varepsilon})$ can be evaluated, which holds for common reparameterizable distributions such as Gaussian.

Under these assumptions, the ELBO can be rewritten as an expectation with respect to the noise distributions $q(\mathbf{u})$ and $q(\boldsymbol{\varepsilon})$ through a change of variables, and its gradient can be decomposed into two terms given by

$$\nabla_\theta \mathcal{L}(\theta) = \mathbb{E}_{q(\boldsymbol{\varepsilon})q(\mathbf{u})} \left[\nabla_{\mathbf{z}} \log p(\mathbf{x}, \mathbf{z}) \Big|_{\mathbf{z}=h_\theta(\mathbf{u}; \boldsymbol{\varepsilon})} \nabla_\theta h_\theta(\mathbf{u}; \boldsymbol{\varepsilon}) \right] - \mathbb{E}_{q(\boldsymbol{\varepsilon})q(\mathbf{u})} \left[\nabla_{\mathbf{z}} \log q_\theta(\mathbf{z}) \Big|_{\mathbf{z}=h_\theta(\mathbf{u}; \boldsymbol{\varepsilon})} \nabla_\theta h_\theta(\mathbf{u}; \boldsymbol{\varepsilon}) \right].$$

The first expectation can be estimated using samples from $q(\boldsymbol{\varepsilon})$ and $q(\mathbf{u})$. The second expectation is more difficult as $\nabla_{\mathbf{z}} \log q_\theta(\mathbf{z})$ may not be tractable if $q_\theta(\mathbf{z})$ is implicit. The first key trick in UIVI is to rewrite the gradient in the second term as an expectation of the form

$$\nabla_{\mathbf{z}} \log q(\mathbf{z}) = \mathbb{E}_{q_\theta(\boldsymbol{\varepsilon}|\mathbf{z})} [\nabla_{\mathbf{z}} \log q_\theta(\mathbf{z}|\boldsymbol{\varepsilon})]$$

which then allows for Monte Carlo estimation using samples from $q_\theta(\boldsymbol{\varepsilon}|\mathbf{z}) \propto q_\theta(\mathbf{z}|\boldsymbol{\varepsilon})q(\boldsymbol{\varepsilon})$. An MCMC sampler is used to sample from $q_\theta(\boldsymbol{\varepsilon}|\mathbf{z})$, and the second key trick in UIVI is to reuse the sample $(\mathbf{z}_i, \boldsymbol{\varepsilon}_i)$ from estimating the outer expectation as an initial point in the MCMC sampler. As the initial point is a sample from the same joint distribution $q_\theta(\mathbf{z}, \boldsymbol{\varepsilon})$, no burn-in is necessary and the only purpose of the MCMC is to break the dependency between samples that are used to estimate the inner expectation and that are used estimate the outer expectation. Thus, the gradient of the ELBO is estimated by

$$\hat{\nabla}_\theta \mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \left(\nabla_{\mathbf{z}} \log p(\mathbf{x}, \mathbf{z}) \Big|_{\mathbf{z}=h_\theta(\mathbf{u}_i; \boldsymbol{\varepsilon}_i)} - \frac{1}{m} \sum_{j=1}^m \nabla_{\mathbf{z}} \log q_\theta(\mathbf{z}|\boldsymbol{\varepsilon}'_j) \Big|_{\mathbf{z}=h_\theta(\mathbf{u}_i; \boldsymbol{\varepsilon}_i)} \right) \nabla_\theta h_\theta(\mathbf{u}_i; \boldsymbol{\varepsilon}_i)$$

where $\boldsymbol{\varepsilon}_i \sim q(\boldsymbol{\varepsilon})$, $\mathbf{u}_i \sim q(\mathbf{u})$, $\boldsymbol{\varepsilon}'_j \sim q_\theta(\boldsymbol{\varepsilon}|\mathbf{z})$ and with $n = 1$, $m = 5$ said to be used in practice. This gradient estimate can then be used in a standard stochastic gradient descent procedure to optimize the ELBO.

1.3.2 Other contributions

Aside from the UIVI algorithm, other contributions of the paper by Titsias and Ruiz [20] include the empirical evaluations of UIVI on synthetic and benchmark datasets. Using a Gaussian conditional with a neural network for the mean parameter, Hamiltonian Monte Carlo (HMC) [12] for the MCMC estimation of the ELBO gradient, and otherwise a fairly standard setup, UIVI is shown to be able to visually approximate various synthetic 2D distributions. Under a similar setup, UIVI is shown to be able to achieve better predictive performance than SIVI on the MNIST and HAPT [16] datasets while being comparable in terms of time per iteration. Finally, Titsias and Ruiz [20] show that for a variational autoencoder (VAE) [7] with a semi-implicit variational distribution, UIVI achieves a greater marginal log-likelihood on the test set compared to standard VAE and SIVI on the MNIST and Fashion-MNIST datasets.

1.3.3 Limitations

The paper by Titsias and Ruiz [20] has a few limitations. The main limitation is the lack of theoretical guarantees for the performance and convergence of UIVI. For example, it is unclear how flexible the UIVI variational family is without resorting to overparameterized deep networks, and no guidance is given on how to construct the model such that the true posterior is in or at least well-approximated by a member of the variational family. It is also unclear how good of an approximation UIVI is able to guarantee through its optimization procedure. However, we recognize that this is a common problem across the VI literature and generally stems from the challenge of analyzing general purpose methods that may include intractable and non-analytic components.

Another notable limitation of the paper is the missing discussion of the limitations of UIVI. In particular, the showcased experiments do not stress test UIVI, and there is no mention of possible directions for improving or extending UIVI. Related work published after the paper by Titsias and Ruiz [20] reported limited scalability with the number of latent parameters [8, 11]. This is likely a consequence of the stochastic optimization of the ELBO as well as the use of MCMC, both for which generally require an increasing number of samples in order to maintain the approximation quality as the dimensionality increases. The MCMC sampling in the UIVI optimization procedure may also result in greater variance of the ELBO gradient estimates [1] and further contribute to non-scalability by complicating potential parallelization of the algorithm [18]. In some instances, other issues common to MCMC approaches, such as poor mixing over different modes, appear to inhibit the performance of UIVI [18].

1.4 Other related work

While UIVI was proposed as an improved alternative to SIVI, there does not appear to be follow-up work in the literature that directly extends UIVI. As mentioned in the previous section, the inefficiency of MCMC in high-dimensional regimes is often cited as the main problem of UIVI [8, 11]. It appears that rather than trying to address this issue in UIVI, recent work in the literature return to SIVI and propose methods that either improve the quality of its approximation or allow it to scale more efficiently to high dimensions.

Several strategies for improving the SIVI approximation have been proposed in the literature around the time of or after the work by Titsias and Ruiz [20]. Molchanov et al. [10] proposed *doubly* SIVI (DSIVI) that expands the flexibility of standard SIVI by allowing both the posterior and prior to be semi-implicit. Sobolev and Vetrov [18] introduced *importance weighted hierarchical* VI (IWHVI), which optimizes a SIVI-like lower bound that incorporates elements from the bound used in importance weighted autoencoders [2]. SIVI, DSIVI and HVM can be seen as special cases of IWHVI and so the bound in IWHVI has the capacity to result in a tighter lower bound [18].

Recent work in the literature have focused more on improving the scalability of SIVI to high dimensions. Molchanova et al. [11] proposed *structured* SIVI where the high-dimensional semi-implicit distribution is assumed to factorize into low-dimensional semi-implicit distributions. Moens et al. [8] introduced *compositional implicit* VI, which adopts various mechanisms into SIVI including an adaptive solver for addressing the bias in the objective and sketch-based approximations that keep the method computationally practical for

high-dimensional regimes.

Though the developments in the related literature are mostly methodological, there have been some recent forays into the more theoretical side that attempt to provide statistical guarantees and insights for implicit VI. In particular, Plummer et al. [13] derive posterior contraction results for simple *non-linear latent variable models* by drawing connections to Gaussian convolutions. The NL-LVM has a structure that can be seen as a particular choice of the variational conditional and mixing distributions in UIVI, and so we suspect that the work by Plummer et al. [13] may provide a reasonable starting point for a theoretical analysis of UIVI.

2 Project report

An attempt at analyzing the theoretical properties of unbiased implicit variational inference

Abstract

Titsias and Ruiz [20] introduced unbiased implicit variational inference (UIVI) as an efficient alternative to semi-implicit variational inference. However, the theoretical properties and guarantees of UIVI are largely unknown and only conjectured in follow-up work based on empirical findings. We show that for a particular choice of the conditional and mixing distributions that make up the variational distribution in UIVI, the UIVI approximation is able to get arbitrarily close to the true posterior distribution under certain assumptions. We also discuss potential convergence issues of UIVI in the case of high dimensions and provide an initial attempt in analyzing the variance of the ELBO gradient estimator. We then suggest several directions for future work that may lead to a proper analysis of the theoretical properties of UIVI.

2.1 Introduction

Variational inference (VI) transforms the problem of computing the posterior distribution $p(\mathbf{z}|\mathbf{x})$ into a problem of minimizing the KL divergence (or equivalently, maximizing the evidence lower bound (ELBO)) between $p(\mathbf{z}|\mathbf{x})$ and a simpler variational distribution $q_{\theta}(\mathbf{z})$ belonging to some variational family \mathcal{Q}_{θ} indexed by variational parameters θ [6]. The performance of VI depends on the flexibility of the family \mathcal{Q}_{θ} as well as the ability to optimize θ over this family. Yin and Zhou [21] introduced *semi-implicit variational inference* (SIVI) in which arbitrarily complex functions (e.g., deep networks) could be used as components in the semi-implicit variational distribution to expand its modeling capacity. However, SIVI relies on optimizing a lower bound of the ELBO objective which may lead to slower convergence rates.

Titsias and Ruiz [20] introduced *unbiased implicit variational inference* (UIVI) as an alternative method to SIVI. Like SIVI, UIVI posits a flexible semi-implicit variational family to approximate the true posterior distribution. In contrast to SIVI, however, UIVI uses an unbiased estimator of the ELBO gradient which should result in faster convergence to the optimal variational approximation in theory. While the experiments conducted by Titsias and Ruiz [20] show promising results, theoretical analyses and guarantees are absent in their paper.

In this project, we discuss our attempts at analyzing the theoretical performance of UIVI in terms of its approximation quality and the variance of its ELBO gradient estimator. Although we are unable to make much progress due to time constraints, we suggest possible directions for continuing what is started in this work that may potentially lead to meaningful results. We also attempt to empirically evaluate UIVI, though we obtain unintuitive results that would require further investigation in order to explain.

This report is organized as follows: Section 2.2 provides a brief overview of UIVI; Section 2.3 introduces other notation used in this report; Section 2.4 presents our analysis of the approximation quality of UIVI for a particular choice of the variational conditional and mixing distributions; Section 2.5 describes our attempts at analyzing the variance of the ELBO gradient in UIVI; and Section 2.6 summarizes our findings and the possible directions of future work.

2.2 Background

UIVI approximates the true posterior distribution $p(\mathbf{z}|\mathbf{x})$ using a variational distribution of the form

$$q_{\theta}(\mathbf{z}) = \int q_{\theta}(\mathbf{z}|\varepsilon)q(\varepsilon)\lambda(d\varepsilon)$$

where the variational conditional $q_{\theta}(\mathbf{z}|\varepsilon)$ is required to be reparameterizable and explicit but where the dependency on θ is allowed to be arbitrarily complex. The auxiliary mixing distribution $q(\varepsilon)$ is taken to

be fixed, and let λ denote the Lebesgue measure on the specified support of ε . UIVI also requires that the log-gradient $\nabla_{\mathbf{z}} \log q_{\theta}(\mathbf{z}|\varepsilon)$ can be evaluated. Under these assumptions, the gradient of the ELBO in UIVI can be written as

$$\nabla_{\theta} \mathcal{L}(\theta) = \mathbb{E}_{q(\varepsilon)q(\mathbf{u})} \left[\nabla_{\mathbf{z}} \log p(\mathbf{x}, \mathbf{z}) \Big|_{\mathbf{z}=h_{\theta}(\mathbf{u}; \varepsilon)} \nabla_{\theta} h_{\theta}(\mathbf{u}; \varepsilon) \right] - \mathbb{E}_{q(\varepsilon)q(\mathbf{u})} \left[\nabla_{\mathbf{z}} \log q_{\theta}(\mathbf{z}) \Big|_{\mathbf{z}=h_{\theta}(\mathbf{u}; \varepsilon)} \nabla_{\theta} h_{\theta}(\mathbf{u}; \varepsilon) \right]$$

where \mathbf{z} is reparameterized through the form $\mathbf{z} = h_{\theta}(\mathbf{u}; \varepsilon)$ with $\mathbf{u} \sim q(\mathbf{u})$. The first term in the gradient can be estimated using Monte Carlo samples from $q(\varepsilon)$ and $q(\mathbf{u})$. For the second term, as $\nabla_{\mathbf{z}} \log q_{\theta}(\mathbf{z})$ may not necessarily be explicit, UIVI instead estimates its equivalent identity

$$\nabla_{\mathbf{z}} \log q_{\theta}(\mathbf{z}) = \mathbb{E}_{q_{\theta}(\varepsilon|\mathbf{z})} [\nabla_{\mathbf{z}} \log q_{\theta}(\mathbf{z}|\varepsilon)]$$

using Monte Carlo samples drawn from the reverse conditional $q_{\theta}(\varepsilon|\mathbf{z}) \propto q(\mathbf{z}|\varepsilon)q(\varepsilon)$ via Markov chain Monte Carlo (MCMC) methods. To avoid needing a burn-in phase, UIVI reuses samples from estimating the outer expectation as initial points in the MCMC such that the chain starts from stationarity.

2.3 Other notation

We briefly introduce other notation used in this report. Let \mathcal{B} denote the Borel σ -algebra of \mathbb{R}^d where d is the dimension of the latent variable \mathbf{z} . Let ϕ_{σ} denote the density of a $\mathcal{N}(0, \sigma^2 \mathbf{I}_d)$ distribution. We will often refer to both a distribution and its density by q and let the context distinguish between the one at hand. For densities f and q , let $f * q(z) = \int_{\mathcal{X}} f(z - x)q(x)dx$ denote the convolution of f and q .

2.4 Quality of approximation

Titsias and Ruiz [20] empirically show that UIVI is able to match the implicit variational distribution to various synthetic datasets and is able to better approximate several models compared to SIVI. However, they do not provide theoretical guarantees nor quantify the quality of the UIVI approximation. In this section, we aim to address this limitation of the paper by discussing particular cases of when UIVI is able or unable to theoretically recover the true distribution. We first show a simple example where the true posterior distribution is not in the variational family due to misspecification. We then show that under certain assumptions and choices of the reparameterization and mixing distribution, UIVI is able to approximate the true posterior arbitrarily closely.

2.4.1 Normal-normal example

We first illustrate through a simple example that a hierarchical variational family does not automatically guarantee a flexible variational approximation, and that the choice of the conditional and mixing distributions are still important in realizing the potential of UIVI.

Consider a univariate Gaussian posterior distribution $p(z|\mathbf{x})$ with mean $\mu_{Z|X}$ and known variance $\sigma_{Z|X}^2$. This is the case, for example, when we have a Gaussian likelihood with latent mean Z and known variance and a conjugate Gaussian prior for Z . Suppose that we choose to approximate $p(z|\mathbf{x})$ by a member of the UIVI variational family where the variational conditional is reparameterized as

$$z = h_{\theta}(u; \varepsilon) = \theta + \varepsilon + \sigma_{Z|X}u$$

with ε and u independent standard normal random variables and θ the variational parameter to optimize. This reparameterization corresponds to $\mu_{\theta}(\varepsilon) = \theta + \varepsilon$ being a simple additive function. Then it is easy to see that

$$\begin{aligned} q_{\theta}(z) &= \int_{\mathbb{R}} q_{\theta}(z|\varepsilon)q(\varepsilon)\lambda(d\varepsilon) \\ &= \int_{\mathbb{R}} \phi_{\sigma_{Z|X}}(\theta + \varepsilon)\phi_1(\varepsilon)\lambda(d\varepsilon) \\ &= \phi_{\sqrt{\sigma_{Z|X}^2+1}}(\theta) . \end{aligned}$$

In other words, our variational family is the set of univariate Gaussian distributions $\{\mathcal{N}(\theta, \sigma_{Z|X}^2 + 1) : \theta \in \mathbb{R}\}$. The true posterior distribution $\mathcal{N}(\mu_{Z|X}, \sigma_{Z|X}^2)$ is not in this variational family. The problem in this example is that our conditional distribution is too restrictive and misspecified. If we changed the reparameterization function to be

$$z = h'_{\theta}(u; \varepsilon) = \theta_1 + \varepsilon + \theta_2 u$$

where both θ_1 and θ_2 are learned parameters, then the new variational family corresponding to this reparameterization includes the true posterior distribution. While this example is very simple, it illustrates that relying on the hierarchical structure alone is insufficient for setting up a flexible variational family.

2.4.2 Flexible variational family

We now show that under certain assumptions and a particular choice of the conditional and mixing distributions, the induced UIVI variational family can be chosen to contain a distribution that approximates the true distribution arbitrarily closely. We do so using similar arguments that Plummer et al. [13] made for non-linear latent variable models (NL-LVM). Following their work, we assume that $\mathbf{z} \in \mathbb{R}^d$.

Consider a multivariate mixing distribution of the form $q(\varepsilon) = \prod_{i=1}^d q(\varepsilon_i)$ where $q(\varepsilon_i) = \text{Unif}(0, 1)$ for $i = 1, \dots, d$. Let the variational conditional $q_{\theta, \sigma}(\mathbf{z}|\varepsilon)$ be multivariate Gaussian with mean $\mu_{\theta}(\varepsilon)$ and covariance matrix $\Sigma_{\sigma} = \sigma^2 \mathbf{I}_d$ where $\mu_{\theta} : [0, 1]^d \rightarrow \mathbb{R}^d$ is some arbitrarily complex function. Note that we keep the variational parameters θ and the bandwidth σ separate for reasons that will be clear shortly. This distribution is reparameterizable through the form

$$\mathbf{z} = h_{\theta, \sigma}(\mathbf{u}; \varepsilon) = \mu_{\theta}(\varepsilon) + \Sigma_{\sigma}^{-\frac{1}{2}} \mathbf{u}$$

where $\mathbf{u} \sim \mathcal{N}(0, \mathbf{I}_d)$. Furthermore, the log-density of this distribution and its gradient are both evaluable and so this distribution satisfies the UIVI requirements. Note that this choice of $h_{\theta, \sigma}$ resembles the NL-LVM studied by Plummer et al. [13], which allows us to apply their results with slight modifications.

To study the approximation capabilities of the NL-LVM, the key insight of Plummer et al. [13] is that the marginal density of the reparameterized variable induced by the NL-LVM has the form of a convolution with a Gaussian kernel. This insight applies to our UIVI variational distribution as well where $q_{\theta, \sigma}(\mathbf{z})$ can be rewritten as

$$\begin{aligned} q_{\theta, \sigma}(\mathbf{z}) &= \int_{[0, 1]^m} q_{\theta, \sigma}(\mathbf{z}|\varepsilon) q(\varepsilon) \lambda(d\varepsilon) \\ &= \int_{[0, 1]^m} \phi_{\sigma}(\mathbf{z} - \mu_{\theta}(\varepsilon)) \lambda(d\varepsilon) \\ &= \int_{\mathbb{R}^d} \phi_{\sigma}(\mathbf{z} - \mathbf{t}) \nu_{\mu_{\theta}}(d\mathbf{t}) \end{aligned}$$

with $\nu_{\mu_{\theta}}(B) = \lambda(\mu_{\theta}^{-1}(B))$, $B \in \mathcal{B}$, being the image measure of λ under μ_{θ} . Using the approximation property of convolutions, we characterize the relationship between the true posterior distribution and this particular variational family through the following proposition.

Proposition 1. *For a fixed bandwidth σ , let \mathcal{Q}_{σ} denote the variational family of distributions $q_{\theta, \sigma}(\mathbf{z}|\varepsilon)$. Suppose that $\mu_{\theta}(\mathbf{t}) = F_{\mathbf{z}|\mathbf{x}}^{-1}(\mathbf{t})$ for all $\mathbf{t} \in [0, 1]^d$. Then $p(\mathbf{z}|\mathbf{x}) \in \mathcal{Q}_0$ (the limiting family as $\sigma \rightarrow 0$).*

Proof. If $\mu_{\theta}(\mathbf{t}) = F_{\mathbf{z}|\mathbf{x}}^{-1}(\mathbf{t})$ for all $\mathbf{t} \in [0, 1]$, then the variational distribution has the form $q_{\theta, \sigma}(\mathbf{z}) = \phi_{\sigma} * p(\mathbf{z}|\mathbf{x})$ by the change of variables argument above. The result then immediately follows from the consistency property of convolutions that $\int_{\mathcal{Z}} |\phi_{\sigma} * p(\mathbf{z}|\mathbf{x}) - p(\mathbf{z}|\mathbf{x})| d\mathbf{z} \rightarrow 0$ as $\sigma \rightarrow 0$ [4]. \square

Proposition 1 says that if the inverse CDF or quantile function $F_{\mathbf{z}|\mathbf{x}}^{-1}$ is in the set of functions $\{\mu_{\theta} : \theta \in \Theta\}$ that can be modeled by μ_{θ} , then the true posterior $p(\mathbf{z}|\mathbf{x})$ is a member of the limiting sequence of variational families as the bandwidth σ of the Gaussian kernel shrinks to zero. While this result also implies that the

true posterior is not in the variation family for any $\sigma > 0$, it suggests that for any measure of error, we can choose σ such that the best approximation will be close to the true posterior within a desired tolerance level.

What is not addressed by Proposition 1 is the question of whether we are able to achieve the best approximation for a given σ in practice. This depends on the functional form of μ_{θ} being flexible enough such that $F_{\mathbf{z}|\mathbf{x}}^{-1} \in \{\mu_{\theta} : \theta \in \Theta\}$ and whether our optimization procedure is able to identify the correct θ such that $\mu_{\theta} = F_{\mathbf{z}|\mathbf{x}}^{-1}$. At least with regards to the existence of a sufficiently flexible μ_{θ} , a universal approximation theorem (Theorem 3.1) in the work of Daniels and Velikova [3] states that for any continuous monotone non-decreasing function on a compact subset of \mathbb{R}^d , there exists a feedforward neural network with at most d hidden layers such that the pointwise error when approximating the function is within a set tolerance.

We note that the work by Plummer et al. [13] includes other results that may be of interest—namely, a posterior contraction result for NL-LVM. However, the results do not appear to generalize easily to UIVI. This is because the implicit model in NL-LVM is on the observable variable \mathbf{x} as opposed to the latent variable \mathbf{z} in UIVI. Ultimately, it is also likely more of interest to understand the convergence of the UIVI optimization procedure to the true posterior distribution $p(\mathbf{z}|\mathbf{x})$ for fixed \mathbf{x} rather than how the approximation changes with increasing amounts of observations.

2.4.3 Directions of future work for quantifying quality

As mentioned in the previous section, Proposition 1 does not address the question of how *good* the quality of the UIVI approximation is, and the other results of Plummer et al. [13] for NL-LVM do not generalize to UIVI. We briefly suggest some possible approaches that may allow quantifying the approximation quality of the convolution-based family using ideas from the previous section. We leave further exploration of these ideas for future work due to time constraints on this project.

One idea is to characterize the approximation quality in terms of the error from the Gaussian convolution for a given bandwidth $\sigma > 0$. Such a result would quantify the error in the best approximation possible for any fixed σ . The structure of such a result would look similar to those for kernel density estimation [e.g., 4, Section 9.5, Chapter 11] where the ratio of the expected empirical error and the theoretical error for some bandwidth is upper-bounded by some constant that depends on the true distribution. Being able to obtain such a result would likely require additional assumptions on the smoothness of $p(\mathbf{z}|\mathbf{x})$, e.g., similar to those made on the true density in the analysis of NL-LVM [13].

Another idea is to characterize the quality in terms of the error from approximating the quantile function by μ_{θ} for a fixed θ . Such a result would be necessary for obtaining a convergence rate of the UIVI optimization procedure for this variational family. However, this approach appears to be more challenging than the previous idea in that smoothness assumptions will need to be made on the inverse CDF and that the form of the function μ_{θ} will need to be specified. If such a result could be obtained, then combining this result with a result from the previous idea may allow for a characterization of the convergence rate to the best UIVI variational approximation.

2.5 Variance of gradient

It is known that the performance of standard SIVI suffers when the dimensionality d of the latent variable increases, with Molchanova et al. [11] citing the reason being that the variational distribution $q_{\theta}(\mathbf{z})$ is approximated using $K + 1$ Gaussian distributions (where K is the mixture size) when estimating the ELBO gradient. As the number of dimensions increases, the mixture size K needs to increase exponentially in order for the mixture to maintain a reasonable approximation of the posterior distribution. In this section, we attempt to show that UIVI with a Gaussian conditional inherently suffers from a similar problem where the ELBO gradient is estimated using a finite mixture of Gaussians. We also discuss our attempt at testing this empirically which returned questionable results.

Let $q_{\theta}(\mathbf{z}|\varepsilon)$ be multivariate Gaussian with mean $\mu_{\theta}(\varepsilon)$ and covariance matrix $\Sigma_{\theta}(\varepsilon)$. An unbiased Monte Carlo estimator of the ELBO gradient in UIVI is given by

$$\widehat{\nabla}_{\theta} \mathcal{L}(\theta; \mathbf{u}_{1:n}, \varepsilon_{1:n}, \varepsilon'_{1:m}) = \frac{1}{n} \sum_{i=1}^n \left(\nabla_{\mathbf{z}} \log p(\mathbf{x}, \mathbf{z})|_{\mathbf{z}=h_{\theta}(\mathbf{u}_i; \varepsilon_i)} - \frac{1}{m} \sum_{j=1}^m \nabla_{\mathbf{z}} \log q_{\theta}(\mathbf{z}|\varepsilon'_j)|_{\mathbf{z}=h_{\theta}(\mathbf{u}_i; \varepsilon_i)} \right) \nabla_{\theta} h_{\theta}(\mathbf{u}_i; \varepsilon_i)$$

where $\mathbf{u}_{1:n}$ are n i.i.d. samples from $\mathcal{N}(0, \mathbf{I}_d)$, $\varepsilon_{1:n}$ are n i.i.d. samples from $q(\varepsilon)$, and $\varepsilon'_{1:m}$ are m MCMC samples drawn from $q(\varepsilon|\mathbf{z})$. It is said that $n = 1$ is used in practice and m is kept small ($m = 5$ in the experiments by Titsias and Ruiz [20]). The interpretation of approximating the ELBO gradient with a Gaussian mixture appears in the second term where each sampled ε'_j corresponds to a single sampled Gaussian distribution. From this perspective, the number of MCMC samples m plays a similar role to the mixture size K in SIVI. However, it is not immediately obvious that this interpretation of sampling Gaussians necessarily impacts the performance of UIVI in high dimensions as the contribution of each sampled Gaussian to the gradient estimate is through their score function. Therefore, we aim to better understand the properties of this gradient estimator by studying its variance.

In particular, it is of interest to understand how the variance of this estimator behaves as the variational distribution approaches the true posterior distribution. If the true distribution lives in a high dimensional space, we may intuitively expect that as our variational distribution conforms to the true distribution, our sampled Gaussian distributions may also become more varied. Such an issue would pose a problem for convergence, and so we attempt to analyze the variance in this setting. In our analysis, we assume that θ is fixed and is exactly or close to the value such that $q_{\theta}(\mathbf{z}) = p(\mathbf{z}|\mathbf{x})$.

Deriving the variance of the gradient estimator is not straightforward due to the various sampled components and the arbitrarily complex mappings in $h_{\theta}(\mathbf{u}; \varepsilon)$. To simplify the analysis, we consider the case $n = 1$. We also apply the idea underlying Rao-Blackwellization [14] and define the functions

$$\begin{aligned} L_{\theta}(\mathbf{u}, \varepsilon, \varepsilon'_{1:m}) &= \widehat{\nabla}_{\theta} \mathcal{L}(\theta; \mathbf{u}, \varepsilon, \varepsilon'_{1:m}) , \\ L_{\theta}(\varepsilon'_{1:m}) &= \mathbb{E}_{\mathbf{u}, \varepsilon} [L_{\theta}(\mathbf{u}, \varepsilon, \varepsilon'_{1:m}) | \varepsilon'_{1:m}] . \end{aligned}$$

It then follows that

$$\text{Var}(L_{\theta}(\varepsilon'_{1:m})) = \text{Var}(L_{\theta}(\mathbf{u}, \varepsilon, \varepsilon'_{1:m})) + \mathbb{E}_{\mathbf{u}, \varepsilon, \varepsilon'} \left[(L_{\theta}(\mathbf{u}, \varepsilon, \varepsilon'_{1:m}) - L_{\theta}(\varepsilon'_{1:m}))^2 \right]$$

and so

$$\text{Var}(L_{\theta}(\varepsilon'_{1:m})) \leq \text{Var}(L_{\theta}(\mathbf{u}, \varepsilon, \varepsilon'_{1:m})) .$$

The function $L_{\theta}(\varepsilon'_{1:m})$ is easier to analyze and examining its variance is equivalent to examining a lower bound of the variance of the full gradient estimator. The function has the form

$$L_{\theta}(\varepsilon'_{1:m}) = \mathbb{E}_{\mathbf{u}, \varepsilon} \left[\nabla_{\mathbf{z}} \log p(\mathbf{x}, \mathbf{z})|_{\mathbf{z}=h_{\theta}(\mathbf{u}; \varepsilon)} \nabla_{\theta} h_{\theta}(\mathbf{u}; \varepsilon) \right] - \frac{1}{m} \sum_{j=1}^m \mathbb{E}_{\mathbf{u}, \varepsilon} \left[\nabla_{\mathbf{z}} \log q_{\theta}(\mathbf{z}|\varepsilon'_j)|_{\mathbf{z}=h_{\theta}(\mathbf{u}; \varepsilon)} \nabla_{\theta} h_{\theta}(\mathbf{u}; \varepsilon) \right] .$$

Note that the first term is constant. The second term retains the interpretation of sampling individual Gaussians through sampling ε'_j .

As the gradient of the ELBO is a vector, we consider the variance of the estimator for the gradient of the i -th parameter given by

$$\begin{aligned} \text{Var}(L_{\theta}(\varepsilon'_{1:m})_i) &= \frac{1}{m} \text{Var} \left(\mathbb{E}_{\mathbf{u}, \varepsilon} \left[\nabla_{\mathbf{z}} \log q_{\theta}(\mathbf{z}|\varepsilon')|_{\mathbf{z}=h_{\theta}(\mathbf{u}; \varepsilon)} \nabla_{\theta_i} h_{\theta}(\mathbf{u}; \varepsilon) \right] \right) \\ &= \frac{1}{m} \mathbb{E}_{\varepsilon'} \left[\mathbb{E}_{\mathbf{u}, \varepsilon} \left[\nabla_{\mathbf{z}} \log q_{\theta}(\mathbf{z}|\varepsilon')|_{\mathbf{z}=h_{\theta}(\mathbf{u}; \varepsilon)} \nabla_{\theta_i} h_{\theta}(\mathbf{u}; \varepsilon) \right]^2 \right] \end{aligned}$$

which follows from the fact that the expected score function is zero under the assumption that $q_{\theta}(\mathbf{z}) \approx p(\mathbf{z}|\mathbf{x})$, i.e.,

$$\mathbb{E}_{\mathbf{u}, \varepsilon, \varepsilon'} \left[\nabla_{\mathbf{z}} \log q_{\theta}(\mathbf{z}|\varepsilon') \Big|_{\mathbf{z}=h_{\theta}(\mathbf{u}; \varepsilon)} \nabla_{\theta_i} h_{\theta}(\mathbf{u}; \varepsilon) \right] = \mathbb{E}_{\mathbf{u}, \varepsilon, \varepsilon'} \left[\nabla_{\theta_i} \log q_{\theta}(\mathbf{z}|\varepsilon') \Big|_{\mathbf{z}=h_{\theta}(\mathbf{u}; \varepsilon)} \right] \approx 0 .$$

We can simplify the variance expression further by plugging in the Gaussian score function, which gives

$$\begin{aligned} \text{Var}(L_{\theta}(\varepsilon'_{1:m})_i) &= \frac{1}{m} \mathbb{E}_{\varepsilon'} \left[\mathbb{E}_{\mathbf{u}, \varepsilon} \left[\left(\Sigma_{\theta}(\varepsilon')^{-1} (h_{\theta}(\mathbf{u}; \varepsilon) - \mu_{\theta}(\varepsilon')) \right)^{\top} \nabla_{\theta_i} h_{\theta}(\mathbf{u}; \varepsilon) \right]^2 \right] \\ &= \frac{1}{m} \mathbb{E}_{\varepsilon'} \left[\mathbb{E}_{\mathbf{u}, \varepsilon} \left[\left(\Sigma_{\theta}(\varepsilon')^{-1} ((\mu_{\theta}(\varepsilon) + \Sigma_{\theta}(\varepsilon)\mathbf{u}) - \mu_{\theta}(\varepsilon')) \right)^{\top} \nabla_{\theta_i} h_{\theta}(\mathbf{u}; \varepsilon) \right]^2 \right] \\ &= \frac{1}{m} \mathbb{E}_{\varepsilon'} \left[\mathbb{E}_{\mathbf{u}, \varepsilon} \left[\left(\Sigma_{\theta}(\varepsilon')^{-1} (\mu_{\theta}(\varepsilon) - \mu_{\theta}(\varepsilon')) + \Sigma_{\theta}(\varepsilon')^{-1} \Sigma_{\theta}(\varepsilon) \mathbf{u} \right)^{\top} \nabla_{\theta_i} h_{\theta}(\mathbf{u}; \varepsilon) \right]^2 \right] . \end{aligned}$$

At this point, we find that it is difficult to proceed without further assumptions on Σ_{θ} and μ_{θ} . One idea is to impose a smoothness assumption on μ_{θ} such that the difference in evaluations at ε and ε' can be bounded. Another idea is to assume that μ_{θ_1} and Σ_{θ_2} do not share parameters and so $\nabla_{\theta_i} h_{\theta}(\mathbf{u}; \varepsilon)$ can be simplified in terms of only the function that θ_i is used in.

Though we are unable to make further progress, we can reason a few (perhaps obvious) conclusions from the above derivations. The lower bound to the variance of the original gradient estimator decreases as the number of samples m increases, which decreases the variation due to sampling the individual Gaussians. Also, the effect of the dimension of \mathbf{z} can only come into play in this Gaussian sampling through the (squared) score. If the score function depends on d at a rate greater than that of $O(\sqrt{d})$, then m needs to scale faster than linear with d in order to maintain the lower bound. Otherwise, m can compensate for increases in dimension by scaling linearly with d .

2.5.1 Experiments

For this project, we aimed to have a set of experiments that at least provide intuition of how the variance of the UIVI ELBO gradient estimator behaves in high dimensions. Due to time constraints, we were only able to carry out a subset of them with precarious results. In this section, we briefly discuss what we tried, the preliminary results that we obtained, and the other investigations that we had planned.

The code repository¹ owned by [Titsias and Ruiz](#) includes a MATLAB implementation of UIVI and SIVI. A number of the experiments presented in their paper are implemented as part of a UIVI demo. We decided to focus on the banana distribution as it seemed that it would be sufficiently difficult to approximate using a finite sample of Gaussians while still being simple enough for computational convenience. We extended their synthetic two-dimensional banana distribution experiment to one in d -dimensions. Samples from our d -dimensional banana distribution are generated through the following procedure:

1. Sample $\mathbf{w} \sim \mathcal{N}(0, \Sigma)$ where

$$\Sigma = \begin{bmatrix} 0.9 & \dots & 0.9 \\ \vdots & \ddots & \vdots \\ 0.9 & \dots & 0.9 \end{bmatrix} + 0.1 \mathbf{I}_d .$$

2. Apply transformation

$$\mathbf{z} = \begin{bmatrix} w_1 \\ w_2 + w_1^2 + 1 \\ \vdots \\ w_d + w_1^2 + 1 \end{bmatrix} .$$

The resulting distribution is a high-dimensional “banana” distribution. For example, Figure 1 shows the two-dimensional contours of the density (and samples) when $d = 3$.

¹<https://github.com/franrruiz/uivi>

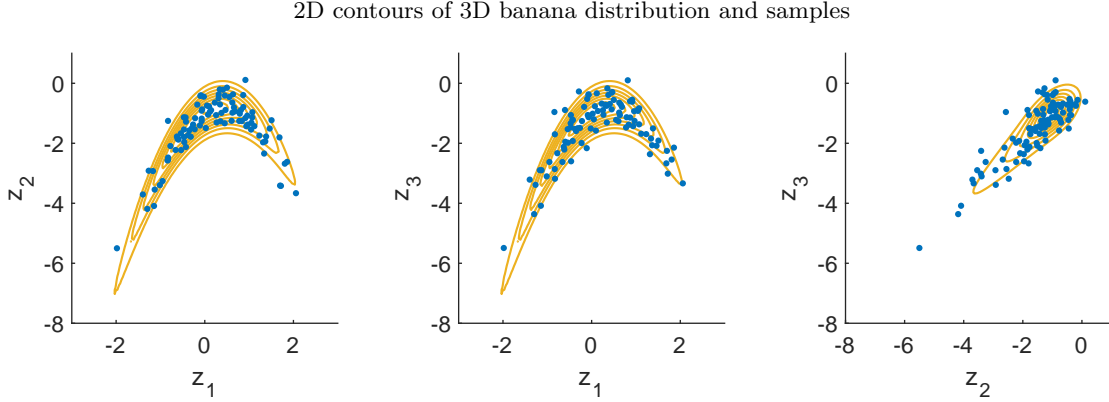


Figure 1: Two-dimensional contours of the density for a three-dimensional “banana” distribution and 100 samples drawn from the UIVI variational approximation ($m = 5$) after 50,000 SGD iterations.

Our experimental setup is as follows. The variational mixing distribution $q(\varepsilon)$ is taken to be univariate standard Gaussian, and the variational conditional distribution $q_{\theta}(\mathbf{z}|\varepsilon)$ is taken to be multivariate Gaussian with mean $\mu_{\mathbf{W},\mathbf{b}}(\varepsilon)$ and covariance matrix $\Sigma_{\sigma} = \text{diag}(\sigma)$. The mean function $\mu_{\mathbf{W},\mathbf{b}}$ is a neural network consisting of two hidden layers of 50 ReLU units each (with weights \mathbf{W} and biases \mathbf{b}). We run 50,000 iterations of stochastic gradient descent (SGD) to optimize the ELBO with respect to $\theta = (\mathbf{W}, \mathbf{b}, \sigma)$. Each iteration runs $5 + m$ Hamiltonian Monte Carlo (HMC) iterations (5 burn-in steps and m MCMC samples) with 5 leapfrog steps [12]. We use the same step size and learning rates as Titsias and Ruiz [20].

Our experiments with the d -dimensional banana distribution did not return intuitive results nor are they easy to interpret. Following the intuition described in Section 2.5, we expected the following results:

1. As the dimension d increases, the variance of the ELBO gradient estimator increases.
2. As the number of MCMC samples m increases, the variance of ELBO gradient estimator decreases.

However, we found that the variance of the ELBO gradient estimator does not change in an obvious way when changing the dimension d and the number of MCMC samples m . Figure 2 shows the estimates of ∇W_{11} , ∇b_{11} and $\nabla \sigma_1$ as d changes with fixed $m = 5$, and Figure 3 shows the same thing but with m changing and fixed $d = 10$. Based on Figure 4, the ELBO (estimated by the log-density of the banana distribution) appears to generally converge to a stable value by 1000 iterations. Sample variances computed using gradients computed in the final 49,000 SGD iterations for each configuration of d and m are shown in Table 1.

	d	$\text{Var}(\hat{\nabla} W_{11})$			$\text{Var}(\hat{\nabla} b_{11})$			$\text{Var}(\hat{\nabla} \sigma_1)$		
		3	10	30	3	10	30	3	10	30
m	1	7.44	0.11	0.59	46.92	33.36	30.78	84.81	147.14	125.92
	5	13.97	0.01	0.55	35.47	41.04	36.72	57.93	166.37	142.41
	10	4.72	0.01	2.90	31.87	39.68	41.52	69.79	143.35	172.34

Table 1: Sample variance of UIVI gradient estimators calculated from gradients computed in the last 49,000 SGD iterations for various configurations of dimension d and number of MCMC samples m .

Curiously, running the above experiments with the same variational distribution in SIVI also returned unexpected results. Figure 5 shows that increasing K from 50 to 100 does not decrease the variance of the gradient estimators. Furthermore, Figure 6 shows that the estimated ELBO in UIVI with $m = 1$ is notably larger than that of SIVI with $K = 50$ or $K = 100$.

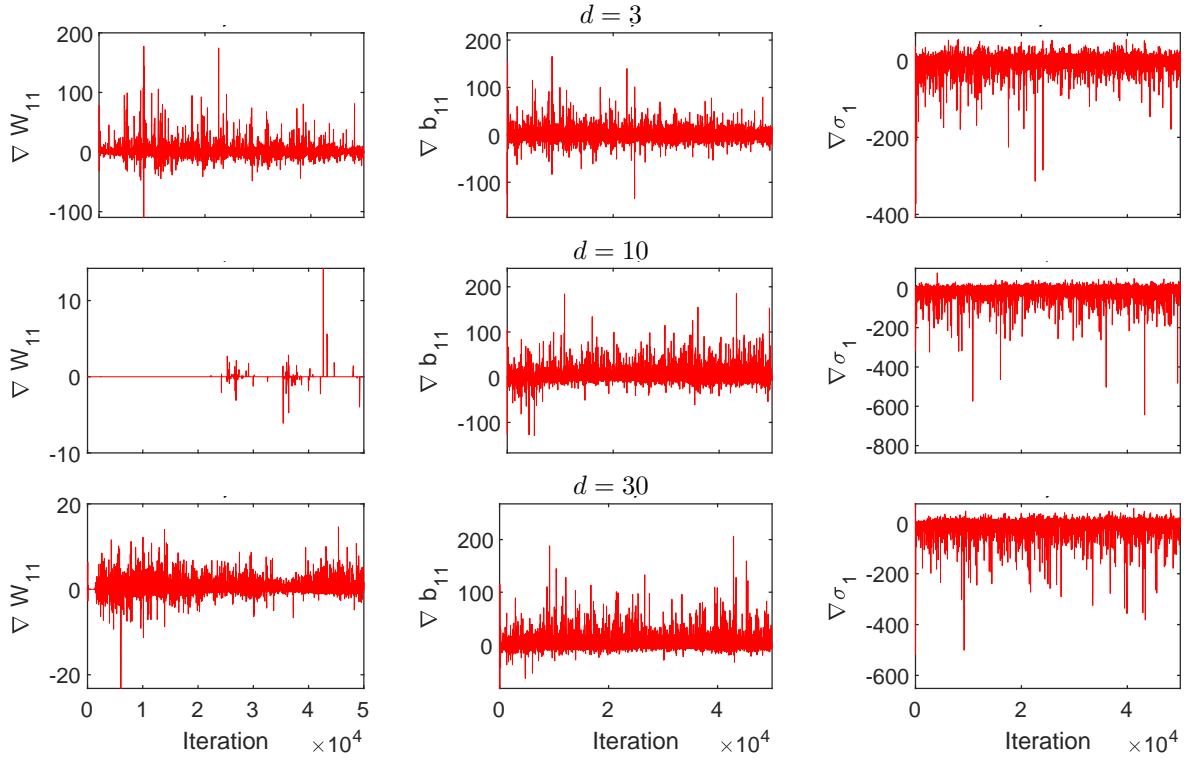


Figure 2: Estimated gradients of W_{11} , b_{11} , and σ_1 over 50,000 SGD iterations in UIVI for $d \in \{3, 10, 30\}$. In all cases, $m = 5$.

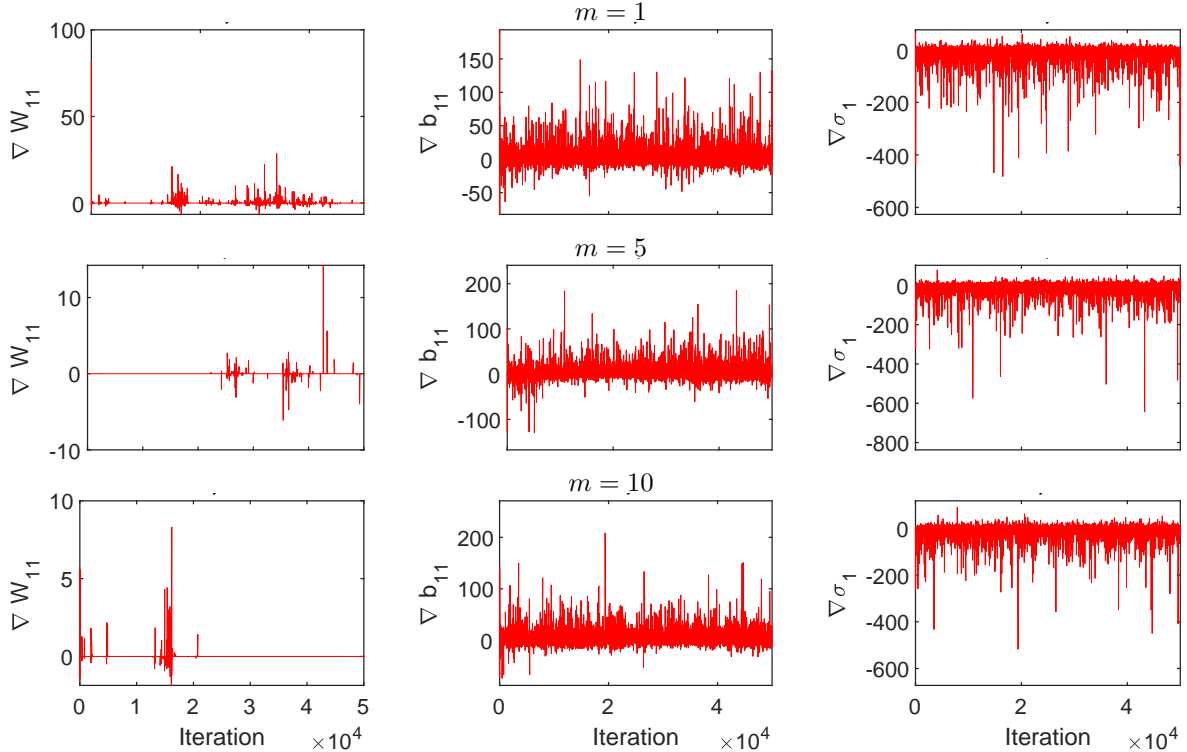


Figure 3: Estimated gradients of W_{11} , b_{11} , and σ_1 over 50,000 SGD iterations in UIVI for $m \in \{1, 5, 10\}$. In all cases, $d = 10$.

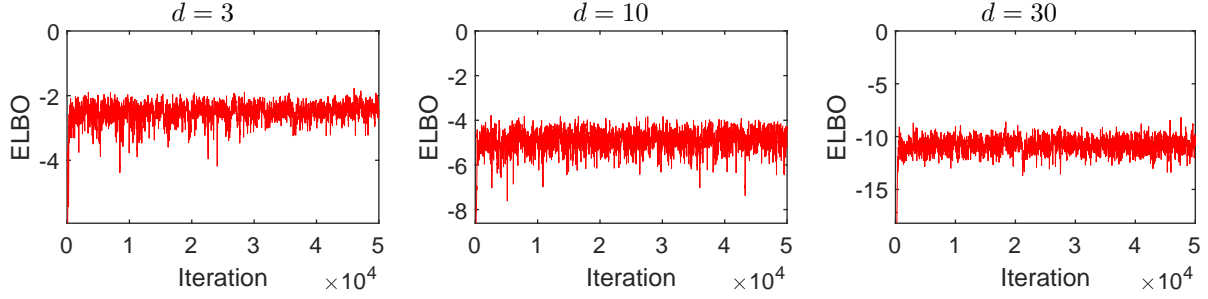


Figure 4: ELBO in UIVI estimated by the log-density of the banana distribution evaluated at a sample. Note that estimates are smoothed using a moving window of 50 iterations.

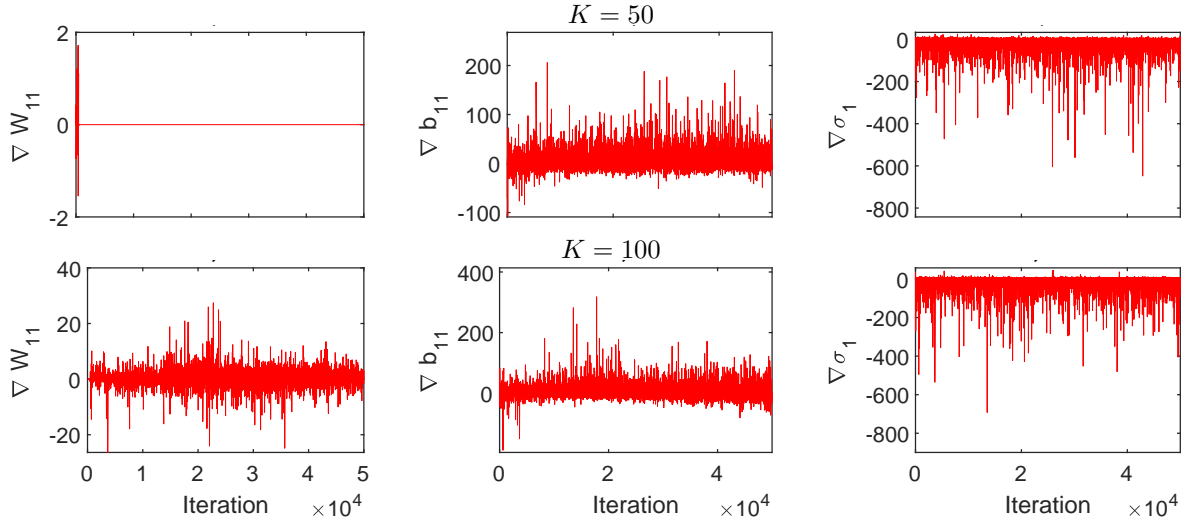


Figure 5: Estimated gradients of W_{11} , b_{11} , and σ_1 over 50,000 SGD iterations in SIVI for $K \in \{50, 100\}$. In both cases, $d = 10$.

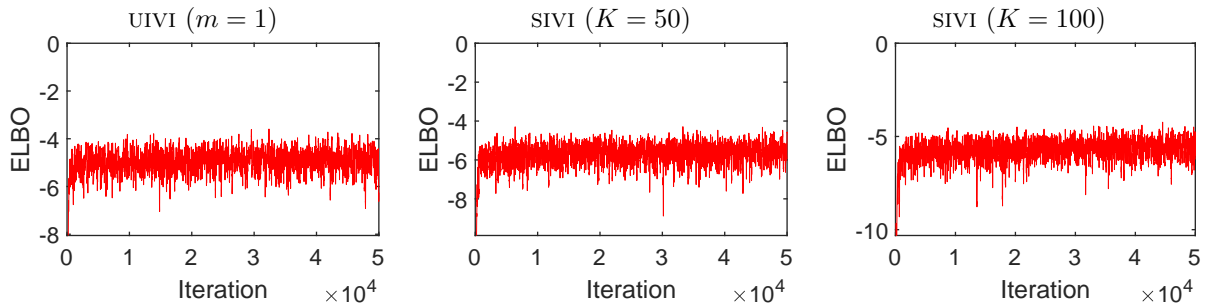


Figure 6: ELBO estimated by the log-density of the banana distribution evaluated at a sample. Note that estimates are smoothed using a moving window of 50 iterations.

Our results are unintuitive and the explanations for them are unclear. It may be that the high-dimensional banana distribution is not as difficult to approximate using Gaussians as we expected. It may also be that despite keeping the variational parameters the same across experiments, the changing experimental conditions lead to different interpretations and learned usages of the parameters. We also cannot rule out the possibility of bugs in the code. Given more time, we would investigate more thoroughly with the goal of explaining these results.

With more time, we would conduct other experiments in addition to a deeper investigation, such as:

1. repeating the above experiments with other high-dimensional synthetic distributions (e.g., Gaussian mixtures). If the results look relatively consistent across distributions, then it makes for a stronger case that our results are due to properties of UIVI rather than properties of the chosen banana distribution.
2. comparing the implicit VI methods (UIVI and SIVI) to non-implicit VI methods. Existing works have suggested that these issues related to dimensionality are not limited to only UIVI but also other implicit VI methods such as SIVI as well [8, 11]. If other methods that are specifically designed to address this issue (e.g., *structured* SIVI [11]) have better performance on the same experiments, then this would provide further evidence that the performance of UIVI suffers in high-dimensional problems.

2.6 Discussion

In this project, we attempted to analyze the theoretical performance of UIVI. Through a simple example, we showed that relying solely on the hierarchical form of the variational distribution is insufficient for obtaining a flexible variational family. We also showed that for a particular choice of the variational conditional and mixing distributions, the variational distribution in UIVI has an equivalent form as a Gaussian convolution. Then under the assumption that the map that transforms the noise is able to learn the inverse CDF of the posterior, we proved that the true posterior distribution is a member of the limiting variational family as the kernel bandwidth goes to zero.

We also attempted to analyze the variance of the ELBO gradient estimator. We discussed the intuition that UIVI can be seen as estimating the ELBO gradient using sampled Gaussian distributions and explained why this may potentially lead to convergence issues in high-dimensional problems. We then derived a lower bound of the estimator variance that may be useful if additional assumptions are imposed on the implicit functions.

The work in this project has been mainly exploratory, and a few ideas have been identified that may be of interest for follow-up work. The work by Plummer et al. [13] provides an interesting starting point for the theoretical analyses of implicit VI methods. While we have applied their approach to show the capacity of the variational family in UIVI as an initial result, the theory can likely be pushed further as discussed in Section 2.4.3. In terms of the variance of the ELBO gradient estimator, our theoretical analysis is incomplete and there is still much work to be done. Our experiments are also incomplete with the preliminary results being unintuitive and requiring deeper investigation.

References

- [1] Michael Betancourt. The fundamental incompatibility of scalable Hamiltonian Monte Carlo and naive data subsampling. In *International Conference on Machine Learning*, pages 533–540. PMLR, 2015.
- [2] Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*, 2015.
- [3] Hennie Daniels and Marina Velikova. Monotone and partially monotone neural networks. *IEEE Transactions on Neural Networks*, 21(6):906–917, 2010.
- [4] Luc Devroye and Gábor Lugosi. *Combinatorial methods in density estimation*. Springer Science & Business Media, 2001.
- [5] Ferenc Huszár. Variational inference using implicit distributions. *arXiv preprint arXiv:1702.08235*, 2017.
- [6] Michael I. Jordan, Zoubin Ghahramani, and et al. An introduction to variational methods for graphical models. In *Machine Learning*, pages 183–233. MIT Press, 1999.
- [7] Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [8] Vincent Moens, Hang Ren, Alexandre Maraval, Rasul Tutunov, Jun Wang, and Haitham Ammar. Efficient semi-implicit variational inference. *arXiv preprint arXiv:2101.06070*, 2021.
- [9] Shakir Mohamed and Balaji Lakshminarayanan. Learning in implicit generative models. *arXiv preprint arXiv:1610.03483*, 2016.
- [10] Dmitry Molchanov, Valery Kharitonov, Artem Sobolev, and Dmitry Vetrov. Doubly semi-implicit variational inference. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2593–2602. PMLR, 2019.
- [11] Iuliia Molchanova, Dmitry Molchanov, Novi Quadrianto, and Dmitry Vetrov. Structured semi-implicit variational inference. 2019.
- [12] Radford M Neal. MCMC using Hamiltonian dynamics. *Handbook of markov chain monte carlo*, 2(11): 2, 2011.
- [13] Sean Plummer, Shuang Zhou, Anirban Bhattacharya, David Dunson, and Debdeep Pati. Statistical guarantees for transformation based models with applications to implicit variational inference. In *International Conference on Artificial Intelligence and Statistics*, pages 2449–2457. PMLR, 2021.
- [14] Rajesh Ranganath, Sean Gerrish, and David Blei. Black box variational inference. In *Artificial Intelligence and Statistics*, pages 814–822. PMLR, 2014.
- [15] Rajesh Ranganath, Dustin Tran, and David Blei. Hierarchical variational models. In *International Conference on Machine Learning*, pages 324–333. PMLR, 2016.
- [16] Jorge-Luis Reyes-Ortiz, Luca Oneto, Alessandro Ghio, Albert Samà, Davide Anguita, and Xavier Parra. Human activity recognition on smartphones with awareness of basic activities and postural transitions. In *International conference on artificial neural networks*, pages 177–184. Springer, 2014.
- [17] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International Conference on Machine Learning*, pages 1530–1538. PMLR, 2015.
- [18] Artem Sobolev and Dmitry P Vetrov. Importance weighted hierarchical variational inference. *Advances in Neural Information Processing Systems*, 32, 2019.
- [19] Masashi Sugiyama, Taiji Suzuki, and Takafumi Kanamori. *Density ratio estimation in machine learning*. Cambridge University Press, 2012.

-
- [20] Michalis K Titsias and Francisco Ruiz. Unbiased implicit variational inference. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 167–176. PMLR, 2019.
 - [21] Mingzhang Yin and Mingyuan Zhou. Semi-implicit variational inference. In *International Conference on Machine Learning*, pages 5660–5669. PMLR, 2018.