
Variational (Gradient) Estimate of the Score Function in Energy-based Latent Variable Models

Fan Bao^{1,2}, Kun Xu¹, Chongxuan Li¹, Lanqing Hong², Jun Zhu^{*1}, Bo Zhang¹

¹Dept. of Comp. Sci. & Tech., Institute for AI, THBI Lab, BNRist Center,
State Key Lab for Intell. Tech. & Sys., Tsinghua University, Beijing, China

²Huawei Noah's Ark Lab

bf19@mails.tsinghua.edu.cn, {kunxu.thu, chongxuanli1991}@gmail.com,
honglanqing@huawei.com, {dcszj, dcszb}@tsinghua.edu.cn

Abstract

The learning and evaluation of energy-based latent variable models (EBLVMs) without any structural assumptions are highly challenging, because the true posteriors and the partition functions in such models are generally intractable. This paper presents variational estimates of the score function and its gradient with respect to the model parameters in a general EBLVM, referred to as *VaES* and *VaGES* respectively. The variational posterior is trained to minimize a certain divergence to the true model posterior and the bias in both estimates can be bounded by the divergence theoretically. With a minimal model assumption, *VaES* and *VaGES* can be applied to the *kernelized Stein discrepancy* (KSD) and *score matching* (SM)-based methods to learn EBLVMs. Besides, *VaES* can also be used to estimate the *exact Fisher divergence* between the data and general EBLVMs.

1 Introduction

An energy-based model (EBM) [17] firstly associates an energy $\mathcal{E}_\theta(\mathbf{v})$ with parameter θ to each configuration of visible variables \mathbf{v} , and then normalizes the energy to get a probability density as

$$p_\theta(\mathbf{v}) = \tilde{p}_\theta(\mathbf{v}) / \mathcal{Z}(\theta) = e^{-\mathcal{E}_\theta(\mathbf{v})} / \mathcal{Z}(\theta),$$

where $\tilde{p}_\theta(\mathbf{v})$ is the unnormalized density and $\mathcal{Z}(\theta) = \int e^{-\mathcal{E}_\theta(\mathbf{v})} d\mathbf{v}$ is called the partition function. As for the learning and evaluation, maximum likelihood estimate (MLE) [9, 31, 36, 24, 25, 2, 3, 7] and noise-contrastive estimation (NCE) [8, 26, 5] are commonly used strategies. However, MLE requires to estimate the partition function and NCE requires to design a proper noise distribution manually. In comparison, Fisher divergence based methods provide an efficient and effective alternative to train [13, 29, 34, 20] and evaluate [6] such models. The calculation of the Fisher divergence is based on the *score function* of an EBM, which is defined as $\nabla_{\mathbf{v}} \log p_\theta(\mathbf{v})$. Notably, the score function is independent of the partition function because $\nabla_{\mathbf{v}} \log p_\theta(\mathbf{v}) = \nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}) - \frac{\nabla_{\mathbf{v}} \log \mathcal{Z}(\theta)}{\mathcal{Z}(\theta)} = \nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v})$, and thereby tractable.

Energy-based latent variable models (EBLVMs) [9, 11, 27, 28, 16, 19, 1] incorporate hidden variables \mathbf{h} to EBMs, which increases the model capacity and enables feature extraction. Such a model defines a joint probability density over the visible variables \mathbf{v} and the latent variables \mathbf{h} as follows

$$p_\theta(\mathbf{v}, \mathbf{h}) = \tilde{p}_\theta(\mathbf{v}, \mathbf{h}) / \mathcal{Z}(\theta) = e^{-\mathcal{E}_\theta(\mathbf{v}, \mathbf{h})} / \mathcal{Z}(\theta),$$

where $\mathcal{E}_\theta(\mathbf{v}, \mathbf{h})$ is the energy function with θ as its parameters, $\tilde{p}_\theta(\mathbf{v}, \mathbf{h})$ is the unnormalized density and $\mathcal{Z}(\theta) = \int e^{-\mathcal{E}_\theta(\mathbf{v}, \mathbf{h})} d\mathbf{v} d\mathbf{h}$ is the partition function. At the cost of increasing expressiveness

*Corresponding author.

power, the posteriors of EBLVMs $p_\theta(\mathbf{h}|\mathbf{v})$ are generally intractable without a structural assumption [30, 33]. As a result, commonly used learning and evaluating algorithms in EBMs are not directly applicable to general nonstructural EBLVMs. To learn and evaluate general nonstructural EBLVMs, we propose variational estimates of the score function and its gradient w.r.t. the model parameters in such models, referred to as *VaES* and *VaGES*. Such estimates firstly introduce a tractable variational posterior $q_\phi(\mathbf{h}|\mathbf{v})$ to approximate the true one $p_\theta(\mathbf{h}|\mathbf{v})$, and then estimate the score function $\nabla_{\mathbf{v}} \log p_\theta(\mathbf{v})$ or its gradient $\frac{\partial \nabla_{\mathbf{v}} \log p_\theta(\mathbf{v})}{\partial \theta}$ by Monte Carlo methods based on samples from $q_\phi(\mathbf{h}|\mathbf{v})$. We show that under some assumptions, the bias introduced by the variational posterior can be bounded by the square root of the KL divergence or the Fisher divergence [15] between the variational posterior and the true posterior. With a minimal model assumption, VaES and VaGES can be applied to the *kernelized Stein discrepancy* (KSD) [22] and *score matching* (SM)-based methods [13, 34, 20] to learn EBLVMs. Besides, VaES can also be used to estimate the *exact Fisher divergence* between the data and general EBLVMs.

2 Method

In this paper, we want to estimate the score function and its gradient w.r.t. the model parameters in an EBLVM. Firstly recall [33] that the score function of an EBLVM can be expressed as

$$\nabla_{\mathbf{v}} \log p_\theta(\mathbf{v}) = \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})], \quad (1)$$

where $p_\theta(\mathbf{v}) = \int p_\theta(\mathbf{v}, \mathbf{h}) d\mathbf{h}$ is the marginal probability density and $p_\theta(\mathbf{h}|\mathbf{v})$ is the posterior probability density. We then show that the gradient of the score function w.r.t. the model parameters can be decomposed into a term of covariance and a term of expectation:

$$\begin{aligned} \frac{\partial \nabla_{\mathbf{v}} \log p_\theta(\mathbf{v})}{\partial \theta} &= \text{Cov}_{p_\theta(\mathbf{h}|\mathbf{v})}(\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}), \nabla_{\theta} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})) \\ &\quad + \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \theta} \right], \end{aligned} \quad (2)$$

where the covariance term is the covariance between two random vectors $\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})$, $\nabla_{\theta} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})$ with $\mathbf{h} \sim p_\theta(\mathbf{h}|\mathbf{v})$ and the expectation term is the expectation of a random matrix $\frac{\partial \nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \theta}$ with $\mathbf{h} \sim p_\theta(\mathbf{h}|\mathbf{v})$. A proof of Eqn. (2) is provided in Appendix A. For simplicity, we denote the covariance term and the expectation term as $c(\mathbf{v}; \theta)$ and $e(\mathbf{v}; \theta)$ respectively.

Eqn. (1,2) inspire us with a Monte Carlo estimation of $\nabla_{\mathbf{v}} \log p_\theta(\mathbf{v})$ and $\frac{\partial \nabla_{\mathbf{v}} \log p_\theta(\mathbf{v})}{\partial \theta}$ in an EBLVM. However, it needs samples from $p_\theta(\mathbf{h}|\mathbf{v})$, which requires the Markov chain Monte Carlo (MCMC) method [23] when only the energy function is accessible, leading to a high time cost. To efficiently get samples, we learn a variational posterior $q_\phi(\mathbf{h}|\mathbf{v})$ as an approximation of $p_\theta(\mathbf{h}|\mathbf{v})$ by

$$\min_{\phi} \mathcal{D}(q_\phi(\mathbf{h}|\mathbf{v}) || p_\theta(\mathbf{h}|\mathbf{v})), \quad (3)$$

where \mathcal{D} is the KL divergence when \mathbf{h} is discrete and the Fisher divergence when \mathbf{h} is continuous and the reasons will be specified in Sec. 2.2. In practice, we often need to estimate Eqn. (1) and Eqn. (2) on a minibatch of data $\mathbf{v}_{1:M}$, where M is the batch size. In this case, we don't initialize a new variational parameter ϕ for a new minibatch. We only maintain one parameter ϕ and update ϕ for K times on the minibatch with a prefixed learning rate scheme α by:

$$\phi \leftarrow \phi - \alpha \frac{1}{M} \sum_{i=1}^M \nabla_{\phi} \mathcal{D}(q_\phi(\mathbf{h}|\mathbf{v}_i) || p_\theta(\mathbf{h}|\mathbf{v}_i)). \quad (4)$$

2.1 Variational (Gradient) Estimate of the Score Function in EBLVMs

First, we consider estimating the score function. According to Eqn. (1), a naive variational estimate is

$$\frac{1}{L} \sum_{i=1}^L \nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}_i), \quad \mathbf{h}_i \stackrel{\text{i.i.d.}}{\sim} q_\phi(\mathbf{h}|\mathbf{v}), \quad (5)$$

where L is the number of samples from $q_\phi(\mathbf{h}|\mathbf{v})$. By noticing that

$$\mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} \nabla_{\mathbf{v}} \log q_\phi(\mathbf{h}|\mathbf{v}) = \int q_\phi(\mathbf{h}|\mathbf{v}) \frac{\nabla_{\mathbf{v}} q_\phi(\mathbf{h}|\mathbf{v})}{q_\phi(\mathbf{h}|\mathbf{v})} d\mathbf{h} = \nabla_{\mathbf{v}} \int q_\phi(\mathbf{h}|\mathbf{v}) d\mathbf{h} = \mathbf{0}, \quad (6)$$

we can subtract $r_v \log q(h|v)$ from the naive estimator without changing the value of the expectation, and the resulting variational estimate of the score function (VaES) is

$$\text{VaES}(v; \theta; \phi) = \frac{1}{L} \sum_{i=1}^L r_v \log \frac{p(v; h_i)}{q(h_i|v)}; \quad h_i \text{ i.i.d. } q(h|v); \quad (7)$$

When the variational posterior $q(h|v)$ is equal to the true posterior $p(h|v)$, $\text{VaES}(v; \theta; \phi) = r_v \log \frac{p(v; h)}{p(h|v)} = r_v \log p(v)$ is a deterministic variable and is exactly equal to the score function.

Then, we consider estimating the gradient of the score function in Eqn. (2). As for the covariance term $\alpha(v; \theta; \phi)$, we estimate it with the sample covariance matrix, and the resulting variational estimate is

$$\alpha(v; \theta; \phi) = \frac{1}{L-1} \sum_{i=1}^L r_v \log p(v; h_i) \frac{\partial \log p(v; h_i)}{\partial \theta} - \left(\frac{1}{L-1} \sum_{i=1}^L r_v \log p(v; h_i) \right) \frac{\partial \log p(v; h_i)}{\partial \theta}; \quad h_i \text{ i.i.d. } q(h|v); \quad (8)$$

As for the expectation term $\beta(v; \theta; \phi)$, we directly apply the Monte Carlo estimate and the resulting variational estimate is

$$\beta(v; \theta; \phi) = \frac{1}{L} \sum_{i=1}^L \frac{\partial^2 r_v \log p(v; h_i)}{\partial \theta^2}; \quad h_i \text{ i.i.d. } q(h|v); \quad (9)$$

Combining Eqn. (8) and Eqn. (9), the resulting variational gradient estimate of the score function (VaGES) is

$$\text{VaGES}(v; \theta; \phi) = \alpha(v; \theta; \phi) + \beta(v; \theta; \phi); \quad (10)$$

Remark: Although Eqn. (9) includes second derivatives, in practice we only need to estimate the product of the second derivatives with vectors and only two backpropagations are required for calculation (since $\frac{\partial^2 r_v \log p(v; h_i)}{\partial \theta^2} = \frac{\partial}{\partial \theta} \left(\frac{\partial r_v \log p(v; h_i)}{\partial \theta} \right)$).

2.2 Bounding the Bias

Notice that $\text{VaES}(v; \theta; \phi)$ and $\text{VaGES}(v; \theta; \phi)$ are actually estimating $E_{q(h|v)} [r_v \log p(v; h)]$ and

$$\text{Cov}_{q(h|v)}(r_v \log p(v; h); r_v \log p(v; h)) + E_{q(h|v)} \frac{\partial^2 r_v \log p(v; h)}{\partial \theta^2}$$

respectively and will introduce some bias. Firstly, we show that when h is discrete, the bias of $\text{VaES}(v; \theta; \phi)$ and $\text{VaGES}(v; \theta; \phi)$ can be bounded by the square root of the KL divergence between $q(h|v)$ and $p(h|v)$ under some assumptions on boundedness, as characterized in Theorem 1 and Theorem 2. So in Eqn. (3), we choose the KL divergence when h is discrete to learn the variational posterior $q(h|v)$.

Theorem 1. (h is discrete, VaES, proof in Appendix B) Suppose $\log p(v; h)$ is bounded w.r.t. $v; h$ and θ , then the bias of $\text{VaES}(v; \theta; \phi)$ can be bounded by the square root of the KL divergence between $q(h|v)$ and $p(h|v)$ up to multiplying a constant.

Theorem 2. (h is discrete, VaGES, proof in Appendix B) Suppose $\log p(v; h)$, $r_v \log p(v; h)$ and $\frac{\partial^2 r_v \log p(v; h)}{\partial \theta^2}$ are bounded w.r.t. $v; h$ and θ , then the bias of $\text{VaGES}(v; \theta; \phi)$ can be bounded by the square root of the KL divergence between $q(h|v)$ and $p(h|v)$ up to multiplying a constant.

Then, we show that under extra assumptions on the Stein regularity (see Def. 2) and boundedness of the Stein factors (see Def. 3), the bias of $\text{VaES}(v; \theta; \phi)$ and $\text{VaGES}(v; \theta; \phi)$ can be bounded by the square root of the Fisher divergence between $q(h|v)$ and $p(h|v)$, as characterized in Theorem 3 and Theorem 4. Although the boundedness of the Stein factors have only been verified under some simple cases [8], and haven't been extended to more complex cases, e.g., when θ is parameterized by a neural network, Theorem 3 and Theorem 4 still inspire us to choose the Fisher divergence in Eqn. (3) when h is continuous to learn the variational posterior $q(h|v)$.

Definition 1. [18] Suppose p is a probability density defined on \mathbb{R}^n and $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a function, we define g_f^p as a solution of the Stein equation $\mathbb{E}_p \nabla \cdot g = f - \mathbb{E}_p f$, where $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $S_p g(v) := \mathbb{E}_p \nabla \cdot \log p(v) \nabla g(v) + \text{Tr}(\mathbb{E}_p \nabla \nabla g(v))$.

Remark. The solution of the Stein equation exists. For example, let $f = \mathbb{E}_p f$, then

$$g_1(v) = \frac{1}{p(v)} \int_{-v_1}^{v_1} p(t; v_2; \dots; v_n) h(t; v_2; \dots; v_n) dt, \quad g_2(v) = \dots = g_n(v) = 0$$

is a solution.

Definition 2. Suppose p, q are probability densities defined on \mathbb{R}^n and $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a function, we say f satisfies the Stein regular condition w.r.t. q iff $\exists \delta \in \mathbb{Z} \setminus [1; m]; \lim_{\|v\| \rightarrow 1} q(v) g_{f_i}^p(v) = 0$.

Definition 3. [18] Suppose p, q are probability densities defined on \mathbb{R}^n and $f_S: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a function satisfying f_S satisfies the Stein regular condition w.r.t. q , we define $f_S^{p,q} := \mathbb{E}_{q(x)} \sum_{i=1}^m \|g_{f_i}^p(x)\|_2^2$,

referred to as the Stein factor of f w.r.t. p, q

Theorem 3. (h is continuous, VaES, proof in Appendix C) Suppose $(v; \dots; \dots)$, $\mathbb{E}_p \nabla \cdot \log p(v; h)$ as a function of h satisfies the Stein regular condition w.r.t. $p(h; v; \dots)$ and $q(h; v; \dots)$ and (2) the Stein factor of $\mathbb{E}_p \nabla \cdot \log p(v; h)$ as a function of h w.r.t. $p(h; v; \dots); q(h; v; \dots)$ is bounded w.r.t. v ; and \dots , then the bias of VaES($v; \dots; \dots$) can be bounded by the square root of the Fisher divergence between $q(h; v)$ and $p(h; v)$ up to multiplying a constant.

Theorem 4. (h is continuous, VaGES, proof in Appendix C) Suppose $(v; \dots; \dots)$, $\mathbb{E}_p \nabla \cdot \log p(v; h)$, $\mathbb{E}_q \nabla \cdot \log q(v; h)$, $\mathbb{E}_p \nabla \cdot \log p(v; h) \frac{\partial \log p(v; h)}{\partial v}$ and $\mathbb{E}_q \nabla \cdot \log q(v; h) \frac{\partial \log q(v; h)}{\partial v}$ as functions of h satisfy the Stein regular condition w.r.t. $p(h; v; \dots)$ and $q(h; v; \dots)$ and (2) the Stein factors of $\mathbb{E}_p \nabla \cdot \log p(v; h)$, $\mathbb{E}_q \nabla \cdot \log q(v; h)$, $\mathbb{E}_p \nabla \cdot \log p(v; h) \frac{\partial \log p(v; h)}{\partial v}$ and $\mathbb{E}_q \nabla \cdot \log q(v; h) \frac{\partial \log q(v; h)}{\partial v}$ as functions of h w.r.t. $p(h; v; \dots); q(h; v; \dots)$ are bounded w.r.t. v ; and \dots , (3) $\mathbb{E}_p \nabla \cdot \log p(v; h)$ and $\mathbb{E}_q \nabla \cdot \log q(v; h)$ are bounded w.r.t. $v; h$ and \dots , then the bias of VaGES($v; \dots; \dots$) can be bounded by the square root of the Fisher divergence between $q(h; v)$ and $p(h; v)$ up to multiplying a constant.

3 Learning EBLVMs

In this section, we show that VaES and VaGES can extend Kernelized Stein Discrepancy (KSD) [22] and methods [4, 20] based on score matching (SM) [3] to learn nonstructural EBLVMs. Besides, we empirically validate our methods and further details of experimental settings can be found in Appendix E.

3.1 Learning EBLVMs with KSD

Method. The KSD between a target probability density $p_t(v)$ and the model probability density $p(v)$ is defined as

$$\text{KSD}(p_t; p) = \mathbb{E}_{v; v^0 \sim p_t(v)} \mathbb{E}_p \nabla \cdot \log p(v) \nabla k(v; v^0) + \mathbb{E}_p \nabla \cdot \log p(v^0) \nabla k(v; v^0) + \text{Tr}(\mathbb{E}_p \nabla \nabla k(v; v^0)) + \mathbb{E}_p \nabla \cdot k(v; v^0) \nabla \log p(v^0) + \mathbb{E}_p \nabla \cdot k(v; v^0) \nabla \log p(v); \quad (11)$$

which properly measures the difference between p_t and p under some mild assumptions [22]. To learn an EBLVM from the target density, we can use gradient-based optimization to minimize $\text{KSD}(p_t; p)$, where the gradient w.r.t. is

$$\frac{\partial \text{KSD}(p_t; p)}{\partial \theta} = 2 \mathbb{E}_{v; v^0 \sim p_t(v)} \nabla \cdot (k(v; v^0) \nabla \log p(v) + \nabla k(v; v^0)) \frac{\partial \log p(v^0)}{\partial \theta}; \quad (12)$$

Estimating $\mathbb{E}_p \nabla \cdot \log p(v)$ with VaES and estimating $\mathbb{E}_p \nabla \cdot \log p(v^0) \frac{\partial \log p(v^0)}{\partial \theta}$ with VaGES, the variational stochastic gradient estimate of KSD (VaGES-KSD) is

$$\frac{1}{M} \sum_{i=1}^M (k(v_i; v_i^0) \text{VaES}(v_i; \dots) + \nabla k(v_i; v_i^0) \nabla \text{VaGES}(v_i^0; \dots)); \quad (13)$$

where the union of $v_{1:M}$ and $v_{1:M}^0$ is a minibatch from the target density, and $\text{VaES}(v_i; \dots)$ and $\text{VaGES}(v_i^0; \dots)$ are independent.

Model and Setting. To the best of our knowledge, KSD hasn't been shown feasible to scale up to natural images. So we illustrate the validity of VaGES-KSD by learning Gaussian restricted Boltzmann machines (GRBMs) [5, 10] on the 2-D checkerboard dataset, whose density is shown in Appendix F.1. The energy function of a GRBM is

$$E(v; h; \theta) = \frac{1}{2} \sum_j v_j b_j - \sum_j v_j \sum_i h_i w_{ij} - \frac{1}{2} \sum_i h_i^2 c_i^2 \quad (14)$$

with learnable parameters $\theta = (b; W; c)$. Since GRBM has a tractable posterior, we can directly learn it using KSD (see Eqn. (12)). We also compare another baseline where $p(v)$ is estimated by importance

sampling (i.e. $p(v) \approx \frac{1}{L} \sum_{i=1}^L \frac{p(v; h_i)}{\text{unif}(h_i)}; h_i \stackrel{\text{i.i.d.}}{\sim} \text{unif}(h)$,

where unif means the uniform distribution) and we call it IS-KSD. We generate 60,000 samples for training and 10,000 samples for testing. We use the RBF kernel $k(v; v^0) = \exp(-\frac{\|v - v^0\|_2^2}{2\sigma^2})$ and $\sigma = 0.1$.

Result. As shown in Fig. 1, VaGES-KSD outperforms IS-KSD on all L (i.e., the number of sampled v) and is comparable to the KSD baseline (the test log-likelihood curves of VaGES-KSD are very close to KSD when $L = 5; 10$).

3.2 Learning EBLVMs with Score Matching

Method. Recall that the Fisher divergence [15] between the target density p and the model density p_θ is

$$D_F(p_\theta(v) \| p(v)) = \frac{1}{2} E_{p_\theta(v)} \left\| \sum_j r_j \log p(v) - \sum_j r_j \log p_\theta(v) \right\|_2^2; \quad (15)$$

which measures the difference between the score functions of p and p_θ . The denoising score matching (DSM) [4] minimizes the Fisher divergence between a perturbed target density and the model density:

$$J_{\text{DSM}}(\theta), D_F(p_\theta(v) \| p_\sigma(v)) = \frac{1}{2} E_{p_\theta(w) p_\sigma(v|w)} \left\| \sum_j r_j \log p_\theta(v) - \sum_j r_j \log p_\sigma(v|w) \right\|_2^2; \quad (16)$$

where $p_\sigma(v) = \int p_\theta(w) p_\sigma(v|w) dw$ is the perturbed target density, $p_\sigma(v|w) = N(v|w; \frac{\sigma^2}{2} I)$ is the Gaussian perturbation and σ is a fixed noise level. The multiscale denoising score matching (MDSM) [20] is a variant of DSM which uses different levels of noise to learn an EBM in high-dimensional spaces:

$$J_{\text{MDSM}}(\theta), \frac{1}{2} E_{p_\theta(w) p(\cdot) p_\sigma(v|w)} \left\| \sum_j r_j \log p_\theta(v) - \sum_j r_j \log p_\sigma(v|w) \right\|_2^2; \quad (17)$$

where $p(\cdot)$ is a prior distribution over the fixable noise level and p_σ is a fixed noise level.

We extend the above two score matching methods (Eqn. (16) and Eqn. (17)) to learn EBLVMs. Firstly write Eqn. (16) and Eqn. (17) in a general form for simplicity

$$J(\theta) = \frac{1}{2} E_{p_\theta(w; v)} \left\| \sum_j r_j \log p_\theta(v) - \sum_j r_j \log p_\sigma(v|w) \right\|_2^2; \quad (18)$$

where $p_\theta(w; v)$ is the joint distribution of w and v (specifically, $p_\theta(w; v) = p_\theta(w) p_\sigma(v|w)$ for Eqn. (16) and $p_\theta(w; v) = p_\theta(w) p(\cdot) p_\sigma(v|w)$ for Eqn. (17)). We use gradient-based optimization to minimize $J(\theta)$ and its gradient w.r.t. θ is

$$\frac{\partial J(\theta)}{\partial \theta} = E_{p_\theta(w; v)} \left(\sum_j r_j \log p_\theta(v) - \sum_j r_j \log p_\sigma(v|w) \right) \frac{\partial \sum_j r_j \log p_\theta(v)}{\partial \theta}; \quad (19)$$

Estimating $r_v \log p(v)$ with VaES and estimating $\frac{\partial}{\partial \theta} \log p(v^0)$ with VaGES, the variational stochastic gradient estimate of the score matching methods (VaGES-SM) is

$$\frac{1}{M} \sum_{i=1}^M (\text{VaES}(v_i; \theta) - r_v \log p(v_i | w_i)) \text{VaGES}(v_i; \theta); \quad (20)$$

where $(w_{1:M}; v_{1:M})$ is a minibatch from $p_t(w; v)$, and $\text{VaES}(v_i; \theta)$ and $\text{VaGES}(v_i^0; \theta)$ are independent. We explicitly denote our methods as VaGES-DSM or VaGES-MDSM according to which score matching objective is used.

Result. To validate the effectiveness of VaGES-SM, we compare VaGES-SM with the state-of-art bi-level score matching (BiSM) [1] for learning EBLVMs, which reformulates Eqn. (16) and Eqn. (17) as bi-level optimization problems (BiDSM and BiMDSM) to learn EBLVMs and uses the gradient unrolling technique to solve the bi-level problems. We compare VaGES-DSM and BiDSM under different settings in GRBMs (see Eqn. (14)) on the Frey face dataset in Fig. 2. For BiSM, N means the number of samples for each v ; K means the times of updating parameters in the lower level problem given a minibatch; H means the number of steps for gradient unrolling in the higher level problem given a minibatch. Since the meaning of N and K in BiSM is similar to VaGES-SM, we compare them in the same dimension and leave another dimension when plotting Fig. 2 (a-d). VaGES-DSM outperforms BiDSM ($N=0$ or $N=2$) and has similar performance with BiDSM ($N=5$). Meanwhile it requires the least memory and the time consuming is similar to BiDSM ($N=2$).

(a) Log-likelihood (b) Fisher divergence (c) Time (s) (d) Memory (MB)

Figure 2: The comparison between VaGES-DSM and BiDSM in GRBMs on the Frey face dataset. The log-likelihood and the Fisher divergence (subtracted by the same unknown constant only relevant to the data) are evaluated on the testing dataset according to the best validation performance. The time is the training time of 2,000 iterations. VaGES-DSM have similar performance with BiDSM ($N=5$), and meanwhile requires less time and memory.

Besides, we compare VaGES-MDSM with BiMDSM in a deep EBLVM with energy function $E(v; h) = g_3(g_2(g_1(v; \theta_1); h); \theta_2)$, where g_1 and g_3 are two neural networks, g_2 is an additive coupling layer and θ is continuous. We evaluate them on the MNIST dataset. The BiMDSM has better visual quality than VaGES-MDSM. It indicates that there are still challenges to apply VaGES-MDSM to high-dimensional data, which is analyzed in Sec. 5.

(a) BiMDSM (b) VaGES-MDSM ($L=2$) (c) VaGES-MDSM ($L=200$)

Figure 3: Samples from deep EBLVMs trained by BiMDSM and VaGES-MDSM on the MNIST.

4 Evaluating EBLVMs

Method. In this setting, we are given an EBLVM $p(v; h) = e^{-E(v; h)} / Z(h)$ and a set of samples $\{v_i\}_{i=1}^n$ from the target density $p(v)$. We want to measure how well the model $p(v)$ approximates

the target density $p_t(v)$, which needs an absolute value representing the difference between p_t and p . This task is more difficult than model comparison, which only needs relative values to compare between models. We use the Fisher divergence of the maximum likelihood to measure how well p_t approximates p :

$$D_F^m(p_t(v) \| p(v)), \max_{f \in \mathcal{F}} E_{p_t(v)} E_{p(v)} [r_v \log p(v) - f(v) + r_v f(v)] - \frac{1}{2} \| \nabla f(v) \|_2^2; \quad (21)$$

where \mathcal{F} is the set of functions $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$, s.t. $\lim_{\|v\| \rightarrow 1} p_t(v) f(v) = 0$ and $p(\cdot)$ is a noise distribution (e.g., Gaussian distribution) introduced for computation efficiency. Under some mild assumptions (see Appendix D), we have $D_F(p_t \| p) = D_F^m(p_t \| p)$. In practice, \mathcal{F} is approximated by a neural network $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$, where θ is the parameter and Θ is the parameter space. We optimize the right hand side of Eqn. (21) and the gradient w.r.t.

$$E_{p_t(v)} E_{p(v)} [r_v \log p(v) - f(v) + r_v f(v)] - \frac{1}{2} \| \nabla f(v) \|_2^2;$$

Estimating $r_v \log p(v)$ with VaES, the variational stochastic gradient estimate is

$$\frac{1}{M} \sum_{i=1}^M \text{VaES}(v_i; \cdot) [r_v \log p(v) - f(v) + r_v f(v)] - \frac{1}{2} \| \nabla f(v) \|_2^2;$$

where $v_{1:M}$ is a minibatch from $p_t(v)$ and $\epsilon_{1:M}$ is a minibatch of noise from $p(\cdot)$. We refer to this method as VaGES-Fisher. In practice, Eqn. (21) is optimized on the training data, validated on the validation data (if there is one) and tested on the testing data.

Result. We validate the effectiveness of VaGES-Fisher in GRBMs. We initialize a GRBM $p(v)$, perturb its weight with increasing noise and provide the Fisher divergence between the initial GRBM and the perturbed one. The dimensions v and h are same and we experiment on dimensions of 200 and 500. Further details of experimental settings can be found in Appendix E. The result is shown in Fig. 4. We compare our estimated Fisher divergence with the accurate Fisher divergence (Fisher). Under both dimensions, our estimated one is close to the accurate one.

Figure 4: The Fisher divergence estimated by VaGES-Fisher (VaGES-Fisher) v.s. the accurate Fisher divergence (Fisher).

5 Challenges in High-Dimensional Spaces and Deep EBLVMs

Although our methods achieve promising results when the data is simple (e.g., the checkerboard dataset) or when the model is simple (e.g., GRBMs), there are still challenges in high-dimensional spaces or in deep EBLVMs. Consider the VaGES estimator, it's a matrix of size $\dim(v) \times \dim(v)$. So either more complex data or models will increase the size of the matrix, making it harder to estimate. In our estimate, the bias might matter more than the variance, since it doesn't increase the visual quality of samples, as shown in Fig. 3. It is an initial hypothesis and requires further experiments for validation. Besides, VaGES is based on the sample covariance matrix, which is relatively simple and there are might be advanced substitutes with better properties.

References

- [1] Fan Bao, Chongxuan Li, Kun Xu, Hang Su, Jun Zhu, and Bo Zhang. Bi-level score matching for learning energy-based latent variable models. arXiv preprint arXiv:2010.07856, 2020.
- [2] Chao Du, Kun Xu, Chongxuan Li, Jun Zhu, and Bo Zhang. Learning implicit generative models by teaching explicit ones. arXiv preprint arXiv:1807.03870, 2018.
- [3] Yilun Du and Igor Mordatch. Implicit generation and generalization in energy-based models. arXiv preprint arXiv:1903.08689, 2019.

- [4] Jianqing Fan, Yuan Liao, and Han Liu. An overview of the estimation of large covariance and precision matrices, 2016.
- [5] Ruiqi Gao, Erik Nijkamp, Diederik P Kingma, Zhen Xu, Andrew M Dai, and Ying Nian Wu. Flow contrastive estimation of energy-based models. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7518–7528, 2020.
- [6] Will Grathwohl, Kuan-Chieh Wang, Jörn-Henrik Jacobsen, David Duvenaud, and Richard Zemel. Learning the stein discrepancy for training and evaluating energy-based models without sampling.
- [7] Will Grathwohl, Kuan-Chieh Wang, Jörn-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. *arXiv preprint arXiv:1912.03263*, 2019.
- [8] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 297–304, 2010.
- [9] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002.
- [10] Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507, 2006.
- [11] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.
- [12] Tianyang Hu, Zixiang Chen, Hanxi Sun, Jincheng Bai, Mao Ye, and Guang Cheng. Stein neural sampler. *arXiv preprint arXiv:1810.03545*, 2018.
- [13] Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(Apr):695–709, 2005.
- [14] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- [15] Oliver Thomas Johnson. *Information theory and the central limit theorem*. World Scientific, 2004.
- [16] Volodymyr Kuleshov and Stefano Ermon. Neural variational inference and learning in undirected graphical models. In *Advances in Neural Information Processing Systems*, pages 6734–6743, 2017.
- [17] Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, and F Huang. A tutorial on energy-based learning. *Predicting structured data* (0), 2006.
- [18] Christophe Ley, Yvik Swan, et al. Stein's density approach and information inequalities. *Electronic Communications in Probability*, 8, 2013.
- [19] Chongxuan Li, Chao Du, Kun Xu, Max Welling, Jun Zhu, and Bo Zhang. To relieve your headache of training an mrf, take advantage. *arXiv preprint arXiv:1901.08400*, 2019.
- [20] Zengyi Li, Yubei Chen, and Friedrich T Sommer. Annealed denoising score matching: Learning energy-based models in high-dimensional spaces. *arXiv preprint arXiv:1910.07762*, 2019.
- [21] Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. In *Advances in neural information processing systems*, pages 2378–2386, 2016.
- [22] Qiang Liu, Jason Lee, and Michael Jordan. A kernelized stein discrepancy for goodness-of-fit tests. In *International conference on machine learning*, pages 276–284, 2016.
- [23] Radford M. Neal. *Probabilistic inference using markov chain monte carlo methods*, 1993.

- [24] Erik Nijkamp, Mitch Hill, Tian Han, Song-Chun Zhu, and Ying Nian Wu. On the anatomy of mcmc-based maximum likelihood learning of energy-based models. arXiv preprint arXiv:1903.12370, 2019.
- [25] Erik Nijkamp, Song-Chun Zhu, and Ying Nian Wu. On learning non-convergent short-run mcmc toward energy-based models. arXiv preprint arXiv:1904.09770, 2019.
- [26] Benjamin Rhodes and Michael U. Gutmann. Variational noise-contrastive estimation. In Kamalika Chaudhuri and Masashi Sugiyama, editors, Proceedings of Machine Learning Research volume 89 of Proceedings of Machine Learning Research, pages 2741–2750. PMLR, 16–18 Apr 2019. URL <http://proceedings.mlr.press/v89/rhodes19a.html>.
- [27] Ruslan Salakhutdinov and Geoffrey Hinton. Deep Boltzmann machines. Proceedings of the twelfth international conference on artificial intelligence and statistics, 2009.
- [28] Ruslan Salakhutdinov and Hugo Larochelle. Efficient learning of deep Boltzmann machines. In Proceedings of the thirteenth international conference on artificial intelligence and statistics, pages 693–700, 2010.
- [29] Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach to density and score estimation. arXiv preprint arXiv:1905.07088, 2019.
- [30] Kevin Swersky, Marc'Aurelio Ranzato, David Buchman, Nando D Freitas, and Benjamin M Marlin. On autoencoders and score matching for energy based models. Proceedings of the 28th international conference on machine learning (ICML), pages 1201–1208, 2011.
- [31] Tijmen Tieleman. Training restricted boltzmann machines using approximations to the likelihood gradient. In Proceedings of the 25th international conference on Machine learning, pages 1064–1071. ACM, 2008.
- [32] Alexandre B Tsybakov. Introduction to nonparametric estimation. Springer Science & Business Media, 2008.
- [33] Eszter Vértés, UCL Gatsby Unit, and Maneesh Sahani. Learning doubly intractable latent variable models via score matching. <http://approximateinference.org/2016/accepted/VertesSahani2016.pdf>.
- [34] Pascal Vincent. A connection between score matching and denoising autoencoders. Neural computation, 23(7):1661–1674, 2011.
- [35] Max Welling, Michal Rosen-Zvi, and Geoffrey E Hinton. Exponential family harmoniums with an application to information retrieval. Advances in neural information processing systems, pages 1481–1488, 2005.
- [36] Qiu Yixuan, Zhang Lingsong, and Wang Xiao. Unbiased contrastive divergence algorithm for training energy-based latent variable models. International Conference on Learning Representation, 2020.

A Proof of Equation (3)

Proof. Firstly we have

$$\begin{aligned} E_{p(h_{jv}; \cdot)} [r_v \log p(h_{jv}; \cdot)] &= E_{p(h_{jv}; \cdot)} \frac{r_v p(h_{jv}; \cdot)}{p(h_{jv}; \cdot)} = \int_Z p(h_{jv}; \cdot) \frac{r_v p(h_{jv}; \cdot)}{p(h_{jv}; \cdot)} dh \\ &= \int_Z r_v p(h_{jv}; \cdot) dh = r_v \int_Z p(h_{jv}; \cdot) dh = r_v 1 = 0; \end{aligned} \quad (22)$$

and similarly we have $E_{p(h_{jv}; \cdot)} [r_h \log p(h_{jv}; \cdot)] = 0$. Thereby, we have

$$\begin{aligned} r_v \log p(v; \cdot) &= r_v \log p(v; \cdot) + E_{p(h_{jv}; \cdot)} [r_v \log p(h_{jv}; \cdot)] = E_{p(h_{jv}; \cdot)} [r_v \log p(v; h; \cdot)]; \\ \text{and similarly we have } r_h \log p(v; \cdot) &= E_{p(h_{jv}; \cdot)} [r_h \log p(v; h; \cdot)]. \end{aligned}$$

Taking derivatives to Eqn. (22) w.r.t., we have

$$E_{p(h_{jv}; \cdot)} [r_v \log p(h_{jv}; \cdot)] \frac{\partial \log p(h_{jv}; \cdot)}{\partial v} + E_{p(h_{jv}; \cdot)} \frac{\partial r_v \log p(h_{jv}; \cdot)}{\partial v} = 0; \quad (23)$$

The leftmost term of Eqn. (23) can be written as

$$\begin{aligned} &E_{p(h_{jv}; \cdot)} [r_v \log p(h_{jv}; \cdot)] \frac{\partial \log p(h_{jv}; \cdot)}{\partial v} \\ &= E_{p(h_{jv}; \cdot)} [r_v \log p(v; h; \cdot)] \frac{\partial \log p(h_{jv}; \cdot)}{\partial v} + E_{p(h_{jv}; \cdot)} [r_v \log p(v; \cdot)] \frac{\partial \log p(h_{jv}; \cdot)}{\partial v} \\ &= E_{p(h_{jv}; \cdot)} [r_v \log p(v; h; \cdot)] \frac{\partial \log p(h_{jv}; \cdot)}{\partial v} \\ &= E_{p(h_{jv}; \cdot)} [r_v \log p(v; h; \cdot)] \frac{\partial \log p(v; h; \cdot)}{\partial v} + E_{p(h_{jv}; \cdot)} [r_v \log p(v; h; \cdot)] \frac{\partial \log p(v; \cdot)}{\partial v} \\ &= E_{p(h_{jv}; \cdot)} [r_v \log p(v; h; \cdot)] \frac{\partial \log p(v; h; \cdot)}{\partial v} + E_{p(h_{jv}; \cdot)} [r_v \log p(v; h; \cdot)] \frac{\partial \log p(v; \cdot)}{\partial v} \\ &= E_{p(h_{jv}; \cdot)} [r_v \log p(v; h; \cdot)] \frac{\partial \log p(v; h; \cdot)}{\partial v} + E_{p(h_{jv}; \cdot)} [r_v \log p(v; h; \cdot)] \frac{\partial \log p(v; \cdot)}{\partial v} \\ &= E_{p(h_{jv}; \cdot)} [r_v \log p(v; h; \cdot)] \frac{\partial \log p(v; h; \cdot)}{\partial v} \\ &= E_{p(h_{jv}; \cdot)} [r_v \log p(v; h; \cdot)] E_{p(h_{jv}; \cdot)} \frac{\partial \log p(v; h; \cdot)}{\partial v} \\ &= \text{Cov}_{p(h_{jv}; \cdot)} (r_v \log p(v; h; \cdot); r_h \log p(v; h; \cdot)); \end{aligned}$$

Thereby, we have

$$\begin{aligned} \frac{\partial r_v \log p(v; \cdot)}{\partial v} &= \frac{\partial r_v \log p(v; \cdot)}{\partial v} + E_{p(h_{jv}; \cdot)} [r_v \log p(h_{jv}; \cdot)] \frac{\partial \log p(h_{jv}; \cdot)}{\partial v} \\ &\quad + E_{p(h_{jv}; \cdot)} \frac{\partial r_v \log p(h_{jv}; \cdot)}{\partial v} \\ &= E_{p(h_{jv}; \cdot)} [r_v \log p(h_{jv}; \cdot)] \frac{\partial \log p(h_{jv}; \cdot)}{\partial v} \\ &\quad + E_{p(h_{jv}; \cdot)} \frac{\partial r_v \log p(v; h; \cdot)}{\partial v} \\ &= \text{Cov}_{p(h_{jv}; \cdot)} (r_v \log p(v; h; \cdot); r_h \log p(v; h; \cdot)) \\ &\quad + E_{p(h_{jv}; \cdot)} \frac{\partial r_v \log p(v; h; \cdot)}{\partial v}; \end{aligned}$$

□

B Proof of Theorem 1 and Theorem 2

Lemma 1. Suppose p, q are probability mass functions defined on Ω and $f : \Omega \rightarrow \mathbb{R}^m$, then the approximation error of estimating $\mathbb{E}_p f$ using $\mathbb{E}_q f$ can be bounded as

$$\| \mathbb{E}_p f - \mathbb{E}_q f \|_2 \leq \| f \|_1 \sqrt{2 D_{KL}(q \| p)};$$

where $\| f \|_1 := \sup_{\omega \in \Omega} \| f(\omega) \|_2$.

Proof. Firstly, we have

$$\| \mathbb{E}_p f - \mathbb{E}_q f \|_2 = \left\| \sum_{\omega \in \Omega} (p(\omega) - q(\omega)) f(\omega) \right\|_2 \leq \sum_{\omega \in \Omega} |p(\omega) - q(\omega)| \| f(\omega) \|_2 \leq \sum_{\omega \in \Omega} |p(\omega) - q(\omega)| \| f \|_1 = 2 D_{KL}(q \| p).$$

According to Pinsker's inequality [32], we have $\sum_{\omega \in \Omega} |p(\omega) - q(\omega)| \leq \sqrt{2 D_{KL}(p \| q)}$. Thereby, $\| \mathbb{E}_p f - \mathbb{E}_q f \|_2 \leq \sqrt{2 D_{KL}(p \| q)} \| f \|_1$. \square

Theorem 1. (h is discrete, VaES) Suppose $\log p(v; h)$ is bounded w.r.t. $v; h$ and σ , then the bias of $\text{VaES}(v; \sigma)$ can be bounded by the square root of the KL divergence between $q(h|v)$ and $p(h|v)$ up to multiplying a constant.

Proof. According to Lemma 1, we have

$$\| \mathbb{E}_{q(h|v)} [r - v \log p(v; h)] - \mathbb{E}_{p(h|v)} [r - v \log p(v; h)] \|_2 \leq \sqrt{2 D_{KL}(q(h|v) \| p(h|v))} \| r - v \log p(v; h) \|_2$$

By the boundedness of $r - v \log p(v; h)$, $\| r - v \log p(v; h) \|_2 \leq A$. Let $C = \sqrt{2}A$, then

$$\begin{aligned} & \| \mathbb{E}_{q(h|v)} [r - v \log p(v; h)] - \mathbb{E}_{p(h|v)} [r - v \log p(v; h)] \|_2 \\ &= \sqrt{2} \sqrt{D_{KL}(q(h|v) \| p(h|v))} \| r - v \log p(v; h) \|_2 \\ &\leq C \sqrt{D_{KL}(q(h|v) \| p(h|v))} \end{aligned}$$

Definition 1. Suppose A is a matrix, we define $\| A \|_2 := \sqrt{\lambda_{\max}(A^T A)}$.

Lemma 2. Suppose a, b are two vectors, then $\| ab^T \|_2 = \| a \|_2 \| b \|_2$.

Proof. $\| ab^T \|_2 = \sqrt{\lambda_{\max}(a a^T b b^T)} = \| a \|_2 \| b \|_2$. \square

Theorem 2. (h is discrete, VaGES) Suppose $\log p(v; h)$, $r - \log p(v; h)$ and $\frac{\partial}{\partial v} \log p(v; h)$ are bounded w.r.t. $v; h$ and σ , then the bias of $\text{VaGES}(v; \sigma)$ can be bounded by the square root of the KL divergence between $q(h|v)$ and $p(h|v)$ up to multiplying a constant.

Proof. According to Thm. 19 $C_1 < 1$, s.t.

$$\| \mathbb{E}_{q(h|v)} [r - v \log p(v; h)] - \mathbb{E}_{p(h|v)} [r - v \log p(v; h)] \|_2 \leq C_1 \sqrt{D_{KL}(q(h|v) \| p(h|v))};$$

Similarly, $C_2 < 1$, s.t.

$$\| \mathbb{E}_{q(h|v)} [r - \log p(v; h)] - \mathbb{E}_{p(h|v)} [r - \log p(v; h)] \|_2 \leq C_2 \sqrt{D_{KL}(q(h|v) \| p(h|v))};$$

and $C_3 < 1$, s.t.

$$\|E_q(r \vee \log p(v; h)) - E_p(r \vee \log p(v; h))\|_{jj_2} \leq C_3 \frac{q}{D_{KL}(q(hjv)jjp(hjv))}:$$

By the boundedness of $r \vee \log p(v; h)$ and $\log p(v; h)$, $r \vee \log p(v; h) \frac{q}{D_{KL}(q(hjv)jjp(hjv))}$ is also bounded. Therefore $C_4 < 1$, s.t.

$$\|E_q(r \vee \log p(v; h) \frac{q}{D_{KL}(q(hjv)jjp(hjv))}) - E_p(r \vee \log p(v; h) \frac{q}{D_{KL}(q(hjv)jjp(hjv))})\|_{jj_2} \leq C_4 \frac{q}{D_{KL}(q(hjv)jjp(hjv))}:$$

By the boundedness of $r \vee \log p(v; h)$ and $\log p(v; h)$, we can assume $C < 1$ is a constant that bounds $r \vee \log p(v; h)jj_2$ and $jj_2 \log p(v; h)jj_2$. Then, by the triangle inequality and Lemma. 2, we have

$$\begin{aligned} & \|E_q(r \vee \log p(v; h)) - E_p(r \vee \log p(v; h))\|_{jj_2} \\ & \leq \|E_q(r \vee \log p(v; h)) - E_p(r \vee \log p(v; h))\|_{jj_2} \\ & \quad + \|E_q(r \vee \log p(v; h)) - E_p(r \vee \log p(v; h))\|_{jj_2} \\ & = \|E_q(r \vee \log p(v; h)) - E_p(r \vee \log p(v; h))\|_{jj_2} \\ & \quad + \|E_q(r \vee \log p(v; h)) - E_p(r \vee \log p(v; h))\|_{jj_2} \\ & \leq \|E_q(r \vee \log p(v; h)) - E_p(r \vee \log p(v; h))\|_{jj_2} \\ & \quad + \|E_q(r \vee \log p(v; h)) - E_p(r \vee \log p(v; h))\|_{jj_2} \\ & \leq (C_2 + C_1) \frac{q}{D_{KL}(q(hjv)jjp(hjv))}: \end{aligned}$$

Therefore,

$$\begin{aligned} & \|Cov_q(r \vee \log p(v; h); r \log p(v; h)) - Cov_p(r \vee \log p(v; h); r \log p(v; h))\|_{jj_2} \\ & \leq \|E_q(r \vee \log p(v; h) \frac{q}{D_{KL}(q(hjv)jjp(hjv))}) - E_p(r \vee \log p(v; h) \frac{q}{D_{KL}(q(hjv)jjp(hjv))})\|_{jj_2} \\ & \quad + \|E_q(r \vee \log p(v; h)) - E_p(r \vee \log p(v; h))\|_{jj_2} \\ & \quad + \|E_q(r \vee \log p(v; h)) - E_p(r \vee \log p(v; h))\|_{jj_2} \\ & \leq (C_4 + C_2 + C_1) \frac{q}{D_{KL}(q(hjv)jjp(hjv))}: \end{aligned}$$

As a result,

$$\begin{aligned}
 & \int \int \text{Cov}_{q(hjv)}(r_{\nu} \log p(v; h); r_{\nu} \log p(v; h)) + E_{q(hjv)} \frac{\int r_{\nu} \log p(v; h)}{\int r_{\nu} \log p(v; h)} \\
 & \quad \frac{\int r_{\nu} \log p(v; h)}{\int r_{\nu} \log p(v; h)} \int \int_2 \\
 & = \int \int \text{Cov}_{q(hjv)}(r_{\nu} \log p(v; h); r_{\nu} \log p(v; h)) + E_{q(hjv)} \frac{\int r_{\nu} \log p(v; h)}{\int r_{\nu} \log p(v; h)} \\
 & \quad \text{Cov}_{p(hjv)}(r_{\nu} \log p(v; h); r_{\nu} \log p(v; h)) + E_{p(hjv)} \frac{\int r_{\nu} \log p(v; h)}{\int r_{\nu} \log p(v; h)} \int \int_2 \\
 & \int \int \text{Cov}_{q(hjv)}(r_{\nu} \log p(v; h); r_{\nu} \log p(v; h)) \\
 & \quad \text{Cov}_{p(hjv)}(r_{\nu} \log p(v; h); r_{\nu} \log p(v; h)) \int \int_2 \\
 & + \int \int E_{q(hjv)} \frac{\int r_{\nu} \log p(v; h)}{\int r_{\nu} \log p(v; h)} + E_{p(hjv)} \frac{\int r_{\nu} \log p(v; h)}{\int r_{\nu} \log p(v; h)} \int \int_2 \\
 & (C_4 + CC_2 + CC_1) \frac{D_{KL}(q(hjv) \parallel p(hjv))}{q} + C_3 \frac{D_{KL}(q(hjv) \parallel p(hjv))}{q} \\
 & = (C_4 + CC_2 + CC_1 + C_3) \frac{D_{KL}(q(hjv) \parallel p(hjv))}{q}
 \end{aligned}$$

□

C Proof of Theorem 3 and Theorem 4

Definition 2. Suppose p is a probability density on \mathbb{R}^n and $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$, we define $S_p f(v) := r_{\nu} \log p(v)^{\top} f(v) + \text{Tr}(r_{\nu} f(v))$.

Lemma 3. [21] Suppose p is a probability density on \mathbb{R}^n and $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a function satisfying $\lim_{\|v\| \rightarrow 1} p(v)f(v) = 0$, then $E_{p(v)}[S_p f(v)] = 0$.

Proof.

$$\begin{aligned}
 0 & = \int r_{\nu} (p(v)f(v)) dv = \int p(v) r_{\nu} f(v) + p(v) f(v) r_{\nu} \log p(v)^{\top} dv \\
 & = E_{p(v)} [r_{\nu} f(v) + f(v) r_{\nu} \log p(v)^{\top}]
 \end{aligned}$$

Thereby,

$$\begin{aligned}
 0 & = \text{Tr}(E_{p(v)} [r_{\nu} f(v) + f(v) r_{\nu} \log p(v)^{\top}]) = E_{p(v)} [\text{Tr}(r_{\nu} f(v)) + r_{\nu} \log p(v)^{\top} f(v)] \\
 & = E_{p(v)} [S_p f(v)]
 \end{aligned}$$

□

Lemma 4. Suppose q, p are probability densities on \mathbb{R}^n and $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ satisfies $\lim_{\|v\| \rightarrow 1} q(v)f(v) = 0$, we have

$$\int E_q S_p f = \frac{q}{E_{q(v)} \int \int f(v) \int \int_2} \frac{p}{D_F(q \parallel p)}$$

Proof. By Lemma 3, we have $E_q S_q f = 0$. Thereby,

$$\begin{aligned}
 \int E_q S_p f & = \int E_q S_p f - E_q S_q f = \int E_{q(v)} f(v)^{\top} (r_{\nu} \log p(v) - r_{\nu} \log q(v)) \\
 & \quad \frac{E_{q(v)} \int \int f(v) \int \int_2 r_{\nu} \log p(v) - r_{\nu} \log q(v) \int \int_2}{q} \\
 & = \frac{q}{E_{q(v)} \int \int f(v) \int \int_2} \frac{p}{D_F(q \parallel p)}
 \end{aligned}$$

□

Definition 3. [18] Suppose p is a probability density defined on \mathbb{R}^n and $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a function, we define g_f^p as a solution of the Stein equation $\mathbb{E}_p g = f - \mathbb{E}_p f$.

Remark. The solution of the Stein equation exists. For example, let h be $\mathbb{E}_p f$, then

$$g_1(v) = \frac{1}{p(v)} \int_1^{Z_{v_1}} p(t; v_2; \dots; v_n) h(t; v_2; \dots; v_n) dt; \quad g_2(v) = \dots = g_n(v) = 0$$

is a solution.

Definition 4. Suppose p, q are probability densities defined on \mathbb{R}^n and $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a function, we say f satisfies the Stein regular condition w.r.t. q iff $\exists i \in \{1, \dots, m\}; \lim_{j \rightarrow \infty} q(v) g_{f_i}^p(v) = 0$.

Definition 5. [18] Suppose p, q are probability densities defined on \mathbb{R}^n and $f_S: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a function satisfying f_S satisfies the Stein regular condition w.r.t. q , we define $f^{p;q} := \mathbb{E}_{q(x)} \sum_{i=1}^m \|g_{f_i}^p(x)\|_2^2$,

referred to as the Stein factor of f w.r.t. p, q .

Lemma 5. Suppose p, q are probability densities defined on \mathbb{R}^n and $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a function satisfying the Stein regular condition w.r.t. q , then the approximation error of estimating $\mathbb{E}_p f$ using $\mathbb{E}_q f$ can be bounded as

$$\| \mathbb{E}_q f - \mathbb{E}_p f \|_2 \leq \sqrt{f^{p;q} D_F(q||p)}$$

Proof. By Lemma 4, we have

$$\| \mathbb{E}_q f_i - \mathbb{E}_p f_i \|_2 = \| \mathbb{E}_q (f_i - \mathbb{E}_p f_i) \|_2 = \sqrt{\mathbb{E}_q S_p g_{f_i}^p} \leq \sqrt{\mathbb{E}_{q(v)} \|g_{f_i}^p(v)\|_2^2} \sqrt{D_F(q||p)}$$

Thereby, we have

$$\| \mathbb{E}_q f - \mathbb{E}_p f \|_2 = \sqrt{\sum_{i=1}^m \mathbb{E}_q \|f_i - \mathbb{E}_p f_i\|_2^2} \leq \sqrt{\sum_{i=1}^m \mathbb{E}_{q(v)} \|g_{f_i}^p(v)\|_2^2} \sqrt{D_F(q||p)} = \sqrt{f^{p;q} D_F(q||p)}$$

where $f^{p;q} := \mathbb{E}_{q(v)} \sum_{i=1}^m \|g_{f_i}^p(v)\|_2^2$. □

Theorem 3. (h is continuous, VaES) Suppose (1) $\delta(v; \cdot, \cdot), r_v \log p(v; h)$ as a function of h satisfies the Stein regular condition w.r.t. $p(h|v; \cdot)$ and $q(h|v; \cdot)$ and (2) the Stein factor of $r_v \log p(v; h)$ as a function of h w.r.t. $p(h|v; \cdot); q(h|v; \cdot)$ is bounded w.r.t. v ; and $\delta, r_v \log p(v; h)$ and $p(h|v)$ up to multiplying a constant.

Proof. It can be directly derived from Lemma 5. □

Theorem 4. (h is continuous, VaGES) Suppose (1) $\delta(v; \cdot, \cdot), r_v \log p(v; h), r \log p(v; h), r_v \log p(v; h) \frac{\partial \log p(v; h)}{\partial v}$ and $\frac{\partial r_v \log p(v; h)}{\partial v}$ as functions of h satisfy the Stein regular condition w.r.t. $p(h|v; \cdot)$ and $q(h|v; \cdot)$ and (2) the Stein factors of $r_v \log p(v; h), r \log p(v; h), r_v \log p(v; h) \frac{\partial \log p(v; h)}{\partial v}$ and $\frac{\partial r_v \log p(v; h)}{\partial v}$ as functions of h w.r.t. $p(h|v; \cdot); q(h|v; \cdot)$ are bounded w.r.t. v ; and $\delta, r_v \log p(v; h)$ and $r \log p(v; h)$ are bounded w.r.t. $v; h$ and $\delta, r_v \log p(v; h)$ and $r \log p(v; h)$ are bounded w.r.t. $v; h$ and $\delta, r_v \log p(v; h)$ and $r \log p(v; h)$ are bounded w.r.t. $v; h$, then the bias of VaGES($v; \cdot, \cdot$) can be bounded by the square root of the Fisher divergence between $q(h|v)$ and $p(h|v)$ up to multiplying a constant.

Proof. According to Lemma 5, $C_1 < 1$, s.t.

$$\| \mathbb{E}_q (r_v \log p(v; h)) - r_v \log p(v) \|_2 \leq C_1 \sqrt{D_F(q(h|v)||p(h|v))}$$

$9C_2 < 1$, s.t.

$$\| \mathbb{E}_q (r \log p(v; h)) - r \log p(v) \|_2 \leq C_2 \sqrt{D_F(q(h|v)||p(h|v))}$$

$\exists C_3 < \infty$, s.t.

$$\begin{aligned} & \left\| \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \right\|_2 \\ & \leq C_3 \sqrt{\mathcal{D}_F(q_\phi(\mathbf{h}|\mathbf{v})||p_\theta(\mathbf{h}|\mathbf{v}))}, \end{aligned}$$

and $\exists C_4 < \infty$, s.t.

$$\begin{aligned} & \left\| \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} \left[\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}) \frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}) \frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \right\|_2 \\ & \leq C_4 \sqrt{\mathcal{D}_F(q_\phi(\mathbf{h}|\mathbf{v})||p_\theta(\mathbf{h}|\mathbf{v}))}. \end{aligned}$$

By the boundedness of $\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})$ and $\nabla_{\boldsymbol{\theta}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})$, we can assume $C < \infty$ is a constant that bounds $\|\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})\|_2$ and $\|\nabla_{\boldsymbol{\theta}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})\|_2$. After establishing the above bounds w.r.t. the Fisher divergence, the rest proof is exactly the same as Theorem 2. For completeness, we restate the proof as follows. By the triangle inequality and Lemma. 2, we have

$$\begin{aligned} & \left\| \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})] \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \right. \\ & \quad \left. - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})] \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \right\|_2 \\ & \leq \left\| \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})] \left(\mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \right) \right\|_2 \\ & \quad + \left\| (\mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})] - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})]) \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \right\|_2 \\ & = \left\| \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})] \right\|_2 \left\| \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \right\|_2 \\ & \quad + \left\| \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})] - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})] \right\|_2 \left\| \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \right\|_2 \\ & \leq C \left\| \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \right\|_2 \\ & \quad + C \left\| \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})] - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})] \right\|_2 \\ & \leq CC_2 \sqrt{\mathcal{D}_F(q_\phi(\mathbf{h}|\mathbf{v})||p_\theta(\mathbf{h}|\mathbf{v}))} + CC_1 \sqrt{\mathcal{D}_F(q_\phi(\mathbf{h}|\mathbf{v})||p_\theta(\mathbf{h}|\mathbf{v}))} \\ & = (CC_2 + CC_1) \sqrt{\mathcal{D}_F(q_\phi(\mathbf{h}|\mathbf{v})||p_\theta(\mathbf{h}|\mathbf{v}))}. \end{aligned}$$

Thereby,

$$\begin{aligned} & \left\| \text{Cov}_{q_\phi(\mathbf{h}|\mathbf{v})} (\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}), \nabla_{\boldsymbol{\theta}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})) \right. \\ & \quad \left. - \text{Cov}_{p_\theta(\mathbf{h}|\mathbf{v})} (\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}), \nabla_{\boldsymbol{\theta}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})) \right\|_2 \\ & \leq \left\| \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} \left[\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}) \frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}) \frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \right\|_2 \\ & \quad + \left\| \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})] \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \right. \\ & \quad \left. - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} [\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})] \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \boldsymbol{\theta}} \right] \right\|_2 \\ & \leq C_4 \sqrt{\mathcal{D}_F(q_\phi(\mathbf{h}|\mathbf{v})||p_\theta(\mathbf{h}|\mathbf{v}))} + (CC_2 + CC_1) \sqrt{\mathcal{D}_F(q_\phi(\mathbf{h}|\mathbf{v})||p_\theta(\mathbf{h}|\mathbf{v}))} \\ & = (C_4 + CC_2 + CC_1) \sqrt{\mathcal{D}_F(q_\phi(\mathbf{h}|\mathbf{v})||p_\theta(\mathbf{h}|\mathbf{v}))}. \end{aligned}$$

As a result,

$$\begin{aligned}
& \|\text{Cov}_{q_\phi(\mathbf{h}|\mathbf{v})}(\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}), \nabla_{\theta} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})) + \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \theta} \right] \\
& \quad - \frac{\partial \nabla_{\mathbf{v}} \log \tilde{p}(\mathbf{v}; \theta)}{\partial \theta} \|_2 \\
& = \|\text{Cov}_{q_\phi(\mathbf{h}|\mathbf{v})}(\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}), \nabla_{\theta} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})) + \mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \theta} \right] \\
& \quad - \text{Cov}_{p_\theta(\mathbf{h}|\mathbf{v})}(\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}), \nabla_{\theta} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})) - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \theta} \right] \|_2 \\
& \leq \|\text{Cov}_{q_\phi(\mathbf{h}|\mathbf{v})}(\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}), \nabla_{\theta} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})) \\
& \quad - \text{Cov}_{p_\theta(\mathbf{h}|\mathbf{v})}(\nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h}), \nabla_{\theta} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})) \|_2 \\
& \quad + \|\mathbb{E}_{q_\phi(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \theta} \right] - \mathbb{E}_{p_\theta(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \nabla_{\mathbf{v}} \log \tilde{p}_\theta(\mathbf{v}, \mathbf{h})}{\partial \theta} \right] \|_2 \\
& \leq (C_4 + CC_2 + CC_1) \sqrt{\mathcal{D}_F(q_\phi(\mathbf{h}|\mathbf{v}) \| p_\theta(\mathbf{h}|\mathbf{v}))} + C_3 \sqrt{\mathcal{D}_F(q_\phi(\mathbf{h}|\mathbf{v}) \| p_\theta(\mathbf{h}|\mathbf{v}))} \\
& = (C_4 + CC_2 + CC_1 + C_3) \sqrt{\mathcal{D}_F(q_\phi(\mathbf{h}|\mathbf{v}) \| p_\theta(\mathbf{h}|\mathbf{v}))}.
\end{aligned}$$

□

D Consistency between \mathcal{D}_F^m and \mathcal{D}_F

Theorem 5. Suppose $\lim_{\|\mathbf{v}\| \rightarrow \infty} p_t(\mathbf{v})(\nabla_{\mathbf{v}} \log p_\theta(\mathbf{v}) - \nabla_{\mathbf{v}} \log p_t(\mathbf{v})) = \mathbf{0}$ and $\mathbb{E}_{p(\epsilon)} [\epsilon \epsilon^\top] = \mathbf{I}$, then $\mathcal{D}_F^m(p_t \| p_\theta) = \mathcal{D}_F(p_t \| p_\theta)$.

Proof. Notice that $\mathbb{E}_{p(\epsilon)} [\epsilon^\top \nabla_{\mathbf{v}} \mathbf{f}(\mathbf{v}) \epsilon] = \text{Tr}(\nabla_{\mathbf{v}} \mathbf{f}(\mathbf{v}))$ by the assumption $\mathbb{E}_{p(\epsilon)} [\epsilon \epsilon^\top] = \mathbf{I}$ and $\mathcal{D}_F^m(p_t \| p_\theta)$ can be simplified as

$$\mathcal{D}_F^m(p_t \| p_\theta) = \max_{\mathbf{f} \in \mathcal{F}} \mathbb{E}_{p_t(\mathbf{v})} \left[\nabla_{\mathbf{v}} \log p_\theta(\mathbf{v})^\top \mathbf{f}(\mathbf{v}) + \text{Tr}(\nabla_{\mathbf{v}} \mathbf{f}(\mathbf{v})) - \frac{1}{2} \|\mathbf{f}(\mathbf{v})\|_2^2 \right].$$

Suppose $\mathbf{f} \in \mathcal{F}$, i.e., \mathbf{f} is a function from \mathbb{R}^d to \mathbb{R}^d and $\lim_{\|\mathbf{v}\| \rightarrow \infty} p_t(\mathbf{v}) \mathbf{f}(\mathbf{v}) = \mathbf{0}$, by the Stein's identity, we have $\mathbb{E}_{p_t(\mathbf{v})} [\nabla_{\mathbf{v}} \log p_t(\mathbf{v})^\top \mathbf{f}(\mathbf{v}) + \text{Tr}(\nabla_{\mathbf{v}} \mathbf{f}(\mathbf{v}))] = \mathbf{0}$. Thereby, we have

$$\begin{aligned}
& \mathbb{E}_{p_t(\mathbf{v})} \left[\nabla_{\mathbf{v}} \log p_\theta(\mathbf{v})^\top \mathbf{f}(\mathbf{v}) + \text{Tr}(\nabla_{\mathbf{v}} \mathbf{f}(\mathbf{v})) - \frac{1}{2} \|\mathbf{f}(\mathbf{v})\|_2^2 \right] \\
& = \mathbb{E}_{p_t(\mathbf{v})} \left[\nabla_{\mathbf{v}} \log p_\theta(\mathbf{v})^\top \mathbf{f}(\mathbf{v}) + \text{Tr}(\nabla_{\mathbf{v}} \mathbf{f}(\mathbf{v})) - \frac{1}{2} \|\mathbf{f}(\mathbf{v})\|_2^2 \right] - \\
& \quad \mathbb{E}_{p_t(\mathbf{v})} [\nabla_{\mathbf{v}} \log p_t(\mathbf{v})^\top \mathbf{f}(\mathbf{v}) + \text{Tr}(\nabla_{\mathbf{v}} \mathbf{f}(\mathbf{v}))] \\
& = \mathbb{E}_{p_t(\mathbf{v})} \left[(\nabla_{\mathbf{v}} \log p_\theta(\mathbf{v}) - \nabla_{\mathbf{v}} \log p_t(\mathbf{v}))^\top \mathbf{f}(\mathbf{v}) - \frac{1}{2} \|\mathbf{f}(\mathbf{v})\|_2^2 \right] \\
& \leq \frac{1}{2} \mathbb{E}_{p_t(\mathbf{v})} [\|\nabla_{\mathbf{v}} \log p_\theta(\mathbf{v}) - \nabla_{\mathbf{v}} \log p_t(\mathbf{v})\|_2^2] = \mathcal{D}_F(p_t \| p_\theta).
\end{aligned}$$

The equality is achieved when $\mathbf{f}(\mathbf{v}) = \nabla_{\mathbf{v}} \log p_\theta(\mathbf{v}) - \nabla_{\mathbf{v}} \log p_t(\mathbf{v})$, which is a function in \mathcal{F} by assumption. As a result, $\mathcal{D}_F^m(p_t \| p_\theta) = \mathcal{D}_F(p_t \| p_\theta)$.

□

E Experimental Settings

For experiments on GRBMs, $q_\phi(\mathbf{h}|\mathbf{v})$ is always a Bernoulli distribution parameterized by a fully connected layer with the sigmoid activation and we use the Gumbel-Softmax trick [14] for reparameterization of $q_\phi(\mathbf{h}|\mathbf{v})$ with 0.1 as the temperature. By default, the variational parameter ϕ is updated for $K = 5$ times for each minibatch. The batch size is 100 for all experiments. We use the Adam optimizer for all experiments.

E.1 Learning EBLVMs with Kernelized Stein Discrepancy

The density of the checkerboard dataset is shown in Fig. 5 (a). The dimension of \mathbf{h} is 4. The learning rate is 0.001.

E.2 Learning EBLVMs with Score Matching

Experiments on GRBMs: The Frey face dataset² consists of gray-scaled face images of size 20×28 . Following BiSM, we split 1,400 images for training, 300 images for validation and 265 images for testing; the dimension of \mathbf{h} is 400; the learning rate is 0.0002.

Experiments on deep EBLVMs: The MNIST dataset consists of gray-scaled hand-written digits of size 28×28 . The dimension of \mathbf{h} is 20. Following BiSM, we split 60,000 samples for training and 10,000 samples for testing on MNIST; g_1 consists of a 12-layer ResNet and a linear layer which aligns the output dimension of g_1 with the dimension of \mathbf{h} ; g_3 is a fully connected layer; $q_\phi(\mathbf{h}|\mathbf{v})$ is a Gaussian distribution parameterized by a 3-layer convolutional neural network (CNN). For BiMDSM, the times of updating parameters in the lower level problem is $K = 5$ and the number of steps for gradient unrolling in the higher level problem is $N = 0$.

E.3 Evaluating EBLVMs

f_η is a multilayer perceptron (MLP) with 2 hidden layers and each layer has the same width. The GRBM is initialized as a standard Gaussian distribution by letting $\mathbf{b} = \mathbf{0}$, $\mathbf{c} = \mathbf{0}$, $\mathbf{W} = \mathbf{0}$, $\sigma = 1$, so we can get accurate samples from it. We get 20,000 samples from the initial GRBM, and split 16,000 samples for training, 2,000 samples for validation and 2,000 samples for testing. The learning rate is 0.0002. The number of \mathbf{h} sampled for each \mathbf{v} is $L = 1$.

²<http://www.cs.nyu.edu/~roweis/data.html>

