
Bi-level Score Matching for Learning Energy-based Latent Variable Models: Appendix

Fan Bao*, **Chongxuan Li***, **Kun Xu**, **Hang Su[†]**, **Jun Zhu[†]**, **Bo Zhang**
 Dept. of Comp. Sci. & Tech., Institute for AI, THBI Lab, BNRist Center,
 State Key Lab for Intell. Tech. & Sys., Tsinghua University, Beijing, China
 bf19@mails.tsinghua.edu.cn, {chongxuanli1991, kunxu.thu}@gmail.com,
 {suhangss, dcszj, dcszb}@tsinghua.edu.cn

A Proofs and Derivations

A.1 Proof of Theorem 1 and Further Analysis of the Potential Bias

Theorem 1. *Assuming that $\forall \theta \in \Theta, \exists \phi \in \Phi$ such that $\mathcal{D}(q(\mathbf{h}|\mathbf{v}; \phi) || p(\mathbf{h}|\mathbf{v}; \theta)) = 0, \forall \mathbf{v} \in \text{supp}(q)$, we have $\nabla_{\theta} \mathcal{J}(\theta) = \nabla_{\theta} \mathcal{J}_{Bi}(\theta, \phi^*(\theta))$.*

Proof. Suppose $\theta \in \Theta, \phi_0 \in \Phi$ satisfies that $q(\mathbf{h}|\mathbf{v}; \phi_0) = p(\mathbf{h}|\mathbf{v}; \theta)$ for all $\mathbf{v} \in \text{supp}(q)$, then $\mathcal{G}(\theta, \phi_0) = \mathbb{E}_{q(\mathbf{v}, \epsilon)} \mathcal{D}(q(\mathbf{h}|\mathbf{v}; \phi_0) || p(\mathbf{h}|\mathbf{v}; \theta)) = 0$. By the definition of $\phi^*(\theta)$, we have $0 \leq \mathcal{G}(\theta, \phi^*(\theta)) \leq \mathcal{G}(\theta, \phi_0) = 0$, and thereby $\mathcal{G}(\theta, \phi^*(\theta)) = 0$. It means that $\phi^*(\theta)$ also satisfies that $q(\mathbf{h}|\mathbf{v}; \phi^*(\theta)) = p(\mathbf{h}|\mathbf{v}; \theta)$ for all $\mathbf{v} \in \text{supp}(q)$. Finally, we have

$$\begin{aligned}
 \mathcal{J}_{Bi}(\theta, \phi^*(\theta)) &= \mathbb{E}_{q(\mathbf{v}, \epsilon)} \mathbb{E}_{q(\mathbf{h}|\mathbf{v}; \phi)} \mathcal{F} \left(\nabla_{\mathbf{v}} \log \frac{\tilde{p}(\mathbf{v}, \mathbf{h}; \theta)}{q(\mathbf{h}|\mathbf{v}; \phi)}, \epsilon, \mathbf{v} \right) \Big|_{\phi=\phi^*(\theta)} \\
 &= \mathbb{E}_{q(\mathbf{v}, \epsilon)} \mathbb{E}_{p(\mathbf{h}|\mathbf{v}; \theta)} \mathcal{F} \left(\nabla_{\mathbf{v}} \log \frac{\tilde{p}(\mathbf{v}, \mathbf{h}; \theta)}{p(\mathbf{h}|\mathbf{v}; \theta)}, \epsilon, \mathbf{v} \right) \\
 &= \mathbb{E}_{q(\mathbf{v}, \epsilon)} \mathbb{E}_{q(\mathbf{h}|\mathbf{v}; \phi)} \mathcal{F} (\nabla_{\mathbf{v}} \log \tilde{p}(\mathbf{v}; \theta), \epsilon, \mathbf{v}) \\
 &= \mathbb{E}_{q(\mathbf{v}, \epsilon)} \mathcal{F} (\nabla_{\mathbf{v}} \log \tilde{p}(\mathbf{v}; \theta), \epsilon, \mathbf{v}) = \mathcal{J}(\theta),
 \end{aligned}$$

and thereby $\nabla_{\theta} \mathcal{J}(\theta) = \nabla_{\theta} \mathcal{J}_{Bi}(\theta, \phi^*(\theta))$. □

When the assumptions in Theorem 1 don't hold, we can still bound the bias between $\mathcal{J}(\theta)$ and $\mathcal{J}_{Bi}(\theta, \phi^*(\theta))$ by the minimum of $\mathcal{G}(\theta, \phi)$ up to a constant under the following surrogate assumptions:

1. There exists a set of conditional densities $R = \{r(\mathbf{h}|\mathbf{v}; \eta) : \eta \in H\}$ parameterized by η including both $\{p(\mathbf{h}|\mathbf{v}; \theta) | \theta \in \Theta\}$ and $\{q(\mathbf{h}|\mathbf{v}; \phi) | \phi \in \Phi\}$, and the divergence between two conditional densities in R can be bounded by the distance of their parameterizations from below, i.e., $\exists C_1 > 0, \forall \eta_1, \eta_2 \in H, C_1 \|\eta_1 - \eta_2\| \leq \mathbb{E}_{q(\mathbf{v}, \epsilon)} \mathcal{D}(r(\mathbf{h}|\mathbf{v}; \eta_1) || r(\mathbf{h}|\mathbf{v}; \eta_2))$.
2. $\mathcal{J}'_{Bi}(\theta, \eta) := \mathbb{E}_{q(\mathbf{v}, \epsilon)} \mathbb{E}_{r(\mathbf{h}|\mathbf{v}; \eta)} \mathcal{F} \left(\nabla_{\mathbf{v}} \log \frac{\tilde{p}(\mathbf{v}, \mathbf{h}; \theta)}{r(\mathbf{h}|\mathbf{v}; \eta)}, \epsilon, \mathbf{v} \right)$ is Lipschitz continuous w.r.t. η on H , with C_2 as its Lipschitz constant, and C_2 is independent of θ .

*Equal contribution. [†] Corresponding author.

Based on assumption 1, there exists a mapping T_p from Θ to H and a mapping T_q from Φ to H , s.t. $p(\mathbf{h}|\mathbf{v}; \boldsymbol{\theta}) = r(\mathbf{h}|\mathbf{v}; T_p(\boldsymbol{\theta}))$ and $q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi}) = r(\mathbf{h}|\mathbf{v}; T_q(\boldsymbol{\phi}))$. The bias can be bounded as

$$\begin{aligned}
|\mathcal{J}_{Bi}(\boldsymbol{\theta}, \boldsymbol{\phi}^*(\boldsymbol{\theta})) - \mathcal{J}(\boldsymbol{\theta})| &= |\mathcal{J}'_{Bi}(\boldsymbol{\theta}, T_q(\boldsymbol{\phi}^*(\boldsymbol{\theta}))) - \mathcal{J}'_{Bi}(\boldsymbol{\theta}, T_p(\boldsymbol{\theta}))| \\
&\leq C_2 \|T_q(\boldsymbol{\phi}^*(\boldsymbol{\theta})) - T_p(\boldsymbol{\theta})\| \\
&\leq \frac{C_2}{C_1} \mathbb{E}_{q(\mathbf{v}, \boldsymbol{\epsilon})} \mathcal{D}(r(\mathbf{h}|\mathbf{v}; T_q(\boldsymbol{\phi}^*(\boldsymbol{\theta}))) \| r(\mathbf{h}|\mathbf{v}; T_p(\boldsymbol{\theta}))) \\
&= \frac{C_2}{C_1} \mathbb{E}_{q(\mathbf{v}, \boldsymbol{\epsilon})} \mathcal{D}(q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi}^*(\boldsymbol{\theta})) \| p(\mathbf{h}|\mathbf{v}; \boldsymbol{\theta})) \\
&= \frac{C_2}{C_1} \mathcal{G}(\boldsymbol{\theta}, \boldsymbol{\phi}^*(\boldsymbol{\theta})) = \frac{C_2}{C_1} \min_{\boldsymbol{\phi} \in \Phi} \mathcal{G}(\boldsymbol{\theta}, \boldsymbol{\phi}).
\end{aligned}$$

Thereby, to ensure $|\mathcal{J}_{Bi}(\boldsymbol{\theta}, \boldsymbol{\phi}^*(\boldsymbol{\theta})) - \mathcal{J}(\boldsymbol{\theta})| < \delta$, it's enough to ensure $\min_{\boldsymbol{\phi} \in \Phi} \mathcal{G}(\boldsymbol{\theta}, \boldsymbol{\phi}) < \frac{C_1}{C_2} \delta$. We notice that the assumption does not necessarily hold, especially in the context of deep learning and we leave a deeper analysis for the future work.

A.2 Derivation of Divergences used in the Lower Level Problem

We now derive the equivalent forms of divergences used in the lower level optimization. If the KL divergence is used, we have:

$$\begin{aligned}
\mathcal{D}_{KL}(q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi}) \| p(\mathbf{h}|\mathbf{v}; \boldsymbol{\theta})) &= \mathbb{E}_{q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})} \log \frac{q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})}{p(\mathbf{h}|\mathbf{v}; \boldsymbol{\theta})} \\
&= \mathbb{E}_{q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})} \log \frac{q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi}) p(\mathbf{v}; \boldsymbol{\theta}) \mathcal{Z}(\boldsymbol{\theta})}{\tilde{p}(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})} \\
&= \mathbb{E}_{q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})} \left[\log \frac{q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})}{\tilde{p}(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})} \right] + \log p(\mathbf{v}; \boldsymbol{\theta}) + \log \mathcal{Z}(\boldsymbol{\theta}) \\
&\equiv \mathbb{E}_{q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})} \log \frac{q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})}{\tilde{p}(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})},
\end{aligned}$$

where the last equivalence holds because we optimize the divergence only with respect to $\boldsymbol{\phi}$.

If the Fisher divergence is used, we have:

$$\begin{aligned}
\mathcal{D}_F(q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi}) \| p(\mathbf{h}|\mathbf{v}; \boldsymbol{\theta})) &= \frac{1}{2} \mathbb{E}_{q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})} [\|\nabla_{\mathbf{h}} \log q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi}) - \nabla_{\mathbf{h}} \log p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})\|_2^2] \\
&= \frac{1}{2} \mathbb{E}_{q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})} [\|\nabla_{\mathbf{h}} \log q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi}) - \nabla_{\mathbf{h}} \log \tilde{p}(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) - \nabla_{\mathbf{h}} \log \mathcal{Z}(\boldsymbol{\theta})\|_2^2] \\
&= \frac{1}{2} \mathbb{E}_{q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})} [\|\nabla_{\mathbf{h}} \log q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi}) - \nabla_{\mathbf{h}} \log \tilde{p}(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})\|_2^2].
\end{aligned}$$

A.3 Some Mathematical Pre-knowledge for Proof of Theorem 2

Let \mathbf{x} be a vector in \mathbb{C}^n and $\|\mathbf{x}\|$ be the 2-norm of \mathbf{x} . Let $A \in \mathbb{C}^{n \times m}$ be a matrix and $\|A\| := \sup_{\mathbf{x} \in \mathbb{C}^m \setminus \{0\}} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|}$ be the natural norm of A induced by the 2-norm. Let $\|f\|_{Lip} := \sup_{\mathbf{x}_1 \neq \mathbf{x}_2 \in X} \frac{\|f(\mathbf{x}_2) - f(\mathbf{x}_1)\|}{\|\mathbf{x}_2 - \mathbf{x}_1\|}$ be the Lipschitz constant of a function f mapping from a normed vector space (or a subset of it) to another normed vector space, and $\|f\|_{\infty} := \sup_{\mathbf{x} \in X} \|f(\mathbf{x})\|$ be the norm superior of a function taking values in a normed vector space.

Lemma 1. *Suppose $A \in \mathbb{R}^{n \times n}$ is a symmetric positive semi-definite matrix, then $\|A\| = \sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} \langle A\mathbf{x}, \mathbf{x} \rangle$. Furthermore, if A is invertible, then $\|A^{-1}\| = \left(\inf_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} \langle A\mathbf{x}, \mathbf{x} \rangle \right)^{-1}$.*

Proof. By the property of Hermitian matrix, we have $\|A\| = \sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} |\langle A\mathbf{x}, \mathbf{x} \rangle|$. Since A is positive semi-definite, we have $\langle A\mathbf{x}, \mathbf{x} \rangle \geq 0$ and $\|A\| = \sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} \langle A\mathbf{x}, \mathbf{x} \rangle$.

If A is invertible, then

$$\begin{aligned} \|A^{-1}\| &= \sup_{\mathbf{x} \in \mathbb{C}^n, \mathbf{x} \neq 0} \frac{\|A^{-1}\mathbf{x}\|}{\|\mathbf{x}\|} = \sup_{\mathbf{x} \in \mathbb{C}^n, \mathbf{x} \neq 0} \frac{\|\mathbf{x}\|}{\|A\mathbf{x}\|} = \left(\inf_{\mathbf{x} \in \mathbb{C}^n, \mathbf{x} \neq 0} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} \right)^{-1} \\ &= \left(\inf_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} \|A\mathbf{x}\| \right)^{-1} = \left(\inf_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} (\langle A^2\mathbf{x}, \mathbf{x} \rangle)^{\frac{1}{2}} \right)^{-1} \\ &= (\lambda_{\min}(A^2))^{-\frac{1}{2}} = (\lambda_{\min}(A))^{-1} = \left(\inf_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} \langle A\mathbf{x}, \mathbf{x} \rangle \right)^{-1}. \end{aligned}$$

□

Lemma 2. Suppose $A \subset \mathbb{R}^{n \times n}$ is a symmetric positive semi-definite matrix and $\alpha > 0$, s.t. $\alpha\|A\| \leq 1$, then $\|I - \alpha A\| \leq 1$. Furthermore, if A is invertible, then $\|I - \alpha A\| = 1 - \alpha\|A^{-1}\|^{-1}$.

Proof. By the property of Hermitian matrix, we have

$$\|I - \alpha A\| = \sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} |\langle (I - \alpha A)\mathbf{x}, \mathbf{x} \rangle| = \sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} |1 - \langle \alpha A\mathbf{x}, \mathbf{x} \rangle|.$$

Since $\alpha\|A\| = \sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} |\langle \alpha A\mathbf{x}, \mathbf{x} \rangle| \leq 1$, we have

$$\sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} |1 - \langle \alpha A\mathbf{x}, \mathbf{x} \rangle| = \sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} 1 - \langle \alpha A\mathbf{x}, \mathbf{x} \rangle \leq 1.$$

As a result, $\|I - \alpha A\| \leq 1$. If A is invertible, by Lemma 1, we have

$$\sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} 1 - \langle \alpha A\mathbf{x}, \mathbf{x} \rangle = 1 - \alpha \inf_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|=1} \langle A\mathbf{x}, \mathbf{x} \rangle = 1 - \alpha\|A^{-1}\|^{-1}$$

□

A.4 Proof of Theorem 2

For clarity, we explicitly write $\hat{\phi}^n(\boldsymbol{\theta})$ as $\hat{\phi}^n(\boldsymbol{\theta}, \phi^0)$ to emphasize the dependence on ϕ^0 , and $\hat{\phi}^n(\boldsymbol{\theta}, \phi^0)$ is recursively defined as

$$\hat{\phi}^n(\boldsymbol{\theta}, \phi^0) = \hat{\phi}^{n-1}(\boldsymbol{\theta}, \phi^0) - \alpha \frac{\partial \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi} \Big|_{\phi=\hat{\phi}^{n-1}(\boldsymbol{\theta}, \phi^0)}, \quad (1)$$

where we slightly abuse the notation for simplicity and ϕ^0 is also denoted as $\hat{\phi}^0(\boldsymbol{\theta}, \phi^0)$.

Let $\hat{\mathcal{J}}_{B_i}^n(\boldsymbol{\theta}, \phi^0) := \hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \hat{\phi}^n(\boldsymbol{\theta}, \phi^0))$ be the surrogate loss, we firstly build the relationship between the surrogate loss and the accurate loss $\hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta}))$ by the following lemma.

Lemma 3. $\hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta})) = \hat{\mathcal{J}}_{B_i}^n(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta}))$ for all $n \geq 0$.

Proof. Since $\frac{\partial \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi} \Big|_{\phi=\hat{\phi}^*(\boldsymbol{\theta})} = 0$, we have $\hat{\phi}^1(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta})) = \hat{\phi}^0(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta})) = \hat{\phi}^*(\boldsymbol{\theta})$. Similarly, we have $\hat{\phi}^n(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta})) = \hat{\phi}^*(\boldsymbol{\theta})$ for all $n \geq 1$ by the mathematical induction. As a result, $\hat{\mathcal{J}}_{B_i}^n(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta})) = \hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \hat{\phi}^n(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta}))) = \hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta}))$. □

We can further bound the difference between the gradient of the surrogate loss $\frac{\partial \hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \hat{\phi}^n(\boldsymbol{\theta}, \phi^0))}{\partial \boldsymbol{\theta}}$ and the gradient of the true loss $\frac{\partial \hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}}$ as

$$\begin{aligned}
& \left\| \frac{\partial \hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \hat{\phi}^n(\boldsymbol{\theta}, \phi^0))}{\partial \boldsymbol{\theta}} - \frac{\partial \hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\| = \left\| \frac{\partial \hat{\mathcal{J}}_{B_i}^n(\boldsymbol{\theta}, \phi^0)}{\partial \boldsymbol{\theta}} - \frac{\partial \hat{\mathcal{J}}_{B_i}^n(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\| \\
&= \left\| \frac{\partial \hat{\mathcal{J}}_{B_i}^n(\boldsymbol{\theta}, \phi^0)}{\partial \boldsymbol{\theta}} - \frac{\partial \hat{\mathcal{J}}_{B_i}^n(\boldsymbol{\theta}, \phi^0)}{\partial \boldsymbol{\theta}} \Big|_{\phi^0 = \hat{\phi}^*(\boldsymbol{\theta})} - \frac{\partial \hat{\mathcal{J}}_{B_i}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \Big|_{\phi^0 = \hat{\phi}^*(\boldsymbol{\theta})} \frac{\partial \hat{\phi}^*(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right\| \\
&\leq \left\| \frac{\partial \hat{\mathcal{J}}_{B_i}^n(\boldsymbol{\theta}, \phi^0)}{\partial \boldsymbol{\theta}} - \frac{\partial \hat{\mathcal{J}}_{B_i}^n(\boldsymbol{\theta}, \phi^0)}{\partial \boldsymbol{\theta}} \Big|_{\phi^0 = \hat{\phi}^*(\boldsymbol{\theta})} \right\| + \left\| \frac{\partial \hat{\mathcal{J}}_{B_i}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \Big|_{\phi^0 = \hat{\phi}^*(\boldsymbol{\theta})} \frac{\partial \hat{\phi}^*(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right\| \quad (2)
\end{aligned}$$

The first term in Eqn. (2) has a bound in the form of $(A + Bn)\kappa^n \|\phi^0 - \hat{\phi}^*(\boldsymbol{\theta})\|$ and the second term in Eqn. (2) has a bound in the form of $C\kappa^n$, as shown in Theorem 2.

Theorem 2. Suppose the following assumptions hold:

1. Both Θ and Φ are compact and convex,
2. $\hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \phi) \in C^2(\Omega)$, $\hat{\mathcal{G}}(\boldsymbol{\theta}, \phi) \in C^3(\Omega)$, where Ω is an open set including $\Theta \times \Phi$ (i.e. $\hat{\mathcal{J}}_{B_i}$ and $\hat{\mathcal{G}}$ are second and third order continuously differentiable on Ω respectively),
3. $\hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)$ is strongly convex on Φ for all $\boldsymbol{\theta} \in \Theta$,
4. $\forall n \geq 0, \forall \boldsymbol{\theta} \in \Theta, \forall \phi^0 \in \Phi, \hat{\phi}^n(\boldsymbol{\theta}, \phi^0) \in \Phi$ and $\hat{\phi}^*(\boldsymbol{\theta}) \in \Phi$,

then when α is small enough, there exists $A, B, C > 0$ and $\kappa \in (0, 1)$ independent of $\boldsymbol{\theta}$ and ϕ^0 , s.t.,

$$\left\| \frac{\partial \hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \hat{\phi}^n(\boldsymbol{\theta}, \phi^0))}{\partial \boldsymbol{\theta}} - \frac{\partial \hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \hat{\phi}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\| \leq (A + Bn)\kappa^n \|\phi^0 - \hat{\phi}^*(\boldsymbol{\theta})\| + C\kappa^n,$$

for all $\boldsymbol{\theta} \in \Theta, \phi^0 \in \Phi$ and $n \geq 0$.

Proof. By assumptions 1 and 2, when $\boldsymbol{\theta} \in \Theta$ and $\phi \in \Phi$, the norms of k order ($0 \leq k \leq 2$) partial derivatives of $\hat{\mathcal{J}}_{B_i}(\boldsymbol{\theta}, \phi)$ can be bounded by a positive constant A_1 and the norms of k order ($0 \leq k \leq 3$) partial derivatives of $\hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)$ can be bounded by a positive constant A_2 . By assumption 2 and 3, $\frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi^2}$ is positive definite and thereby invertible for all $\boldsymbol{\theta} \in \Theta$ and $\phi \in \Phi$. By assumptions 1, 2, 3 and the smoothness of matrix inverse operator, we have $A_3 := \sup_{\boldsymbol{\theta} \in \Theta} \sup_{\phi \in \Phi} \|(\frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi^2})^{-1}\| < \infty$.

We choose the learning rate α s.t. $\alpha \leq \frac{1}{A_2}$. By Lemma 2 we have

$$\|I - \alpha \frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi^2}\| = 1 - \alpha \|(\frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi^2})^{-1}\|^{-1} \leq 1 - \alpha A_3^{-1}, \quad \forall \boldsymbol{\theta} \in \Theta, \forall \phi \in \Phi.$$

Taking partial derivative of Eqn. (1) w.r.t. ϕ^0 , we have

$$\begin{aligned}
\frac{\partial \hat{\phi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} &= \frac{\partial \hat{\phi}^{n-1}(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} - \alpha \frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi^2} \Big|_{\phi = \hat{\phi}^{n-1}(\boldsymbol{\theta}, \phi^0)} \frac{\partial \hat{\phi}^{n-1}(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \\
&= (I - \alpha \frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi^2} \Big|_{\phi = \hat{\phi}^{n-1}(\boldsymbol{\theta}, \phi^0)}) \frac{\partial \hat{\phi}^{n-1}(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0}, \\
\left\| \frac{\partial \hat{\phi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \right\| &\leq \left\| I - \alpha \frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi^2} \Big|_{\phi = \hat{\phi}^{n-1}(\boldsymbol{\theta}, \phi^0)} \right\| \left\| \frac{\partial \hat{\phi}^{n-1}(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \right\| \\
&\leq (1 - \alpha A_3^{-1}) \left\| \frac{\partial \hat{\phi}^{n-1}(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \right\|, \quad \forall \boldsymbol{\theta} \in \Theta, \forall \phi^0 \in \Phi, \forall n \geq 1.
\end{aligned}$$

Thereby, we have

$$\left\| \frac{\partial \hat{\phi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \right\| \leq (1 - \alpha A_3^{-1})^n \left\| \frac{\partial \hat{\phi}^0(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \right\| = (1 - \alpha A_3^{-1})^n, \quad \forall \boldsymbol{\theta} \in \Theta, \forall \phi^0 \in \Phi, \forall n \geq 0.$$

Taking partial derivative of Eqn. (1) w.r.t. θ , we have

$$\begin{aligned}
\frac{\partial \hat{\phi}^n(\theta, \phi^0)}{\partial \theta} &= \frac{\partial \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \theta} - \alpha \frac{\partial^2 \hat{\mathcal{G}}(\theta, \phi)}{\partial \theta \partial \phi} \Big|_{\phi=\hat{\phi}^{n-1}(\theta, \phi^0)} \\
&\quad - \alpha \frac{\partial^2 \hat{\mathcal{G}}(\theta, \phi)}{\partial \phi^2} \Big|_{\phi=\hat{\phi}^{n-1}(\theta, \phi^0)} \frac{\partial \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \theta} \\
&= (I - \alpha \frac{\partial^2 \hat{\mathcal{G}}(\theta, \phi)}{\partial \phi^2} \Big|_{\phi=\hat{\phi}^{n-1}(\theta, \phi^0)}) \frac{\partial \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \theta} \\
&\quad - \alpha \frac{\partial^2 \hat{\mathcal{G}}(\theta, \phi)}{\partial \theta \partial \phi} \Big|_{\phi=\hat{\phi}^{n-1}(\theta, \phi^0)}, \tag{3}
\end{aligned}$$

$$\begin{aligned}
\| \frac{\partial \hat{\phi}^n(\theta, \phi^0)}{\partial \theta} \| &\leq \| I - \alpha \frac{\partial^2 \hat{\mathcal{G}}(\theta, \phi)}{\partial \phi^2} \Big|_{\phi=\hat{\phi}^{n-1}(\theta, \phi^0)} \| \| \frac{\partial \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \theta} \| \\
&\quad + \alpha \| \frac{\partial^2 \hat{\mathcal{G}}(\theta, \phi)}{\partial \theta \partial \phi} \Big|_{\phi=\hat{\phi}^{n-1}(\theta, \phi^0)} \| \\
&\leq (1 - \alpha A_3^{-1}) \| \frac{\partial \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \theta} \| + \alpha A_2, \quad \forall \theta \in \Theta, \forall \phi^0 \in \Phi, \forall n \geq 1.
\end{aligned}$$

Thereby, we have

$$\begin{aligned}
\| \frac{\partial \hat{\phi}^n(\theta, \phi^0)}{\partial \theta} \| &\leq (1 - \alpha A_3^{-1})^n (\| \frac{\partial \phi^0(\theta, \phi)}{\partial \theta} \| - A_3 A_2) + A_3 A_2 \\
&= (1 - (1 - \alpha A_3^{-1})^n) A_3 A_2 \leq A_3 A_2, \quad \forall \theta \in \Theta, \forall \phi^0 \in \Phi, \forall n \geq 0.
\end{aligned}$$

Taking partial derivative of Eqn. (3) w.r.t. ϕ^0 , we have

$$\begin{aligned}
\frac{\partial^2 \hat{\phi}^n(\theta, \phi^0)}{\partial \phi^0 \partial \theta} &= (-\alpha \frac{\partial^3 \hat{\mathcal{G}}(\theta, \phi)}{\partial \phi^3} \Big|_{\phi=\hat{\phi}^{n-1}(\theta, \phi^0)} \frac{\partial \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \phi^0}) \frac{\partial \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \theta} \\
&\quad + (I - \alpha \frac{\partial^2 \hat{\mathcal{G}}(\theta, \phi)}{\partial \phi^2} \Big|_{\phi=\hat{\phi}^{n-1}(\theta, \phi^0)}) \frac{\partial^2 \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \phi^0 \partial \theta} \\
&\quad - \alpha \frac{\partial^3 \hat{\mathcal{G}}(\theta, \phi)}{\partial \phi \partial \theta \partial \phi} \Big|_{\phi=\hat{\phi}^{n-1}(\theta, \phi^0)} \frac{\partial \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \phi^0}, \\
\| \frac{\partial^2 \hat{\phi}^n(\theta, \phi^0)}{\partial \phi^0 \partial \theta} \| &\leq \alpha \| \frac{\partial^3 \hat{\mathcal{G}}(\theta, \phi)}{\partial \phi^3} \Big|_{\phi=\hat{\phi}^{n-1}(\theta, \phi^0)} \| \| \frac{\partial \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \phi^0} \| \| \frac{\partial \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \theta} \| \\
&\quad + \| I - \alpha \frac{\partial^2 \hat{\mathcal{G}}(\theta, \phi)}{\partial \phi^2} \Big|_{\phi=\hat{\phi}^{n-1}(\theta, \phi^0)} \| \| \frac{\partial^2 \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \phi^0 \partial \theta} \| \\
&\quad + \alpha \| \frac{\partial^3 \hat{\mathcal{G}}(\theta, \phi)}{\partial \phi \partial \theta \partial \phi} \Big|_{\phi=\hat{\phi}^{n-1}(\theta, \phi^0)} \| \| \frac{\partial \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \phi^0} \| \\
&\leq \alpha A_2 (1 - \alpha A_3^{-1})^{n-1} A_2 A_3 + (1 - \alpha A_3^{-1}) \| \frac{\partial^2 \hat{\phi}^{n-1}(\theta, \phi^0)}{\partial \phi^0 \partial \theta} \| \\
&\quad + \alpha A_2 (1 - \alpha A_3^{-1})^{n-1}, \quad \forall \theta \in \Theta, \forall \phi^0 \in \Phi, \forall n \geq 1.
\end{aligned}$$

Thereby, we have

$$\begin{aligned}
\| \frac{\partial^2 \hat{\phi}^n(\theta, \phi^0)}{\partial \phi^0 \partial \theta} \| &\leq n (1 - \alpha A_3^{-1})^{n-1} \alpha A_2 (A_2 A_3 + 1) + \| \frac{\partial^2 \phi^0(\theta, \phi)}{\partial \phi^0 \partial \theta} \| (1 - \alpha A_3^{-1}) \\
&= n (1 - \alpha A_3^{-1})^{n-1} \alpha A_2 (A_2 A_3 + 1), \quad \forall \theta \in \Theta, \forall \phi^0 \in \Phi, \forall n \geq 0.
\end{aligned}$$

The derivative of $\hat{\mathcal{J}}_{Bi}^n(\boldsymbol{\theta}, \phi^0)$ w.r.t. $\boldsymbol{\theta}$ is

$$\frac{\partial \hat{\mathcal{J}}_{Bi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \boldsymbol{\theta}} = \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \phi)}{\partial \boldsymbol{\theta}} \Big|_{\phi=\hat{\phi}^n(\boldsymbol{\theta}, \phi^0)} + \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \phi)}{\partial \phi} \Big|_{\phi=\hat{\phi}^n(\boldsymbol{\theta}, \phi^0)} \frac{\partial \hat{\phi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \boldsymbol{\theta}}.$$

Taking Lipschitz constant to both sides w.r.t. ϕ^0 on Φ and by the convexity of Φ , we have

$$\begin{aligned} \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}^n(\boldsymbol{\theta}, \cdot)}{\partial \boldsymbol{\theta}} \right\|_{Lip} &\leq \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \cdot)}{\partial \boldsymbol{\theta}} \right\|_{Lip} \|\hat{\phi}^n(\boldsymbol{\theta}, \cdot)\|_{Lip} + \sup_{\phi \in \Phi} \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \phi)}{\partial \phi} \right\| \left\| \frac{\partial \hat{\phi}^n(\boldsymbol{\theta}, \cdot)}{\partial \boldsymbol{\theta}} \right\|_{Lip} \\ &\quad + \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \cdot)}{\partial \phi} \right\|_{Lip} \|\hat{\phi}^n(\boldsymbol{\theta}, \cdot)\|_{Lip} \sup_{\phi^0 \in \Phi} \left\| \frac{\partial \hat{\phi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \boldsymbol{\theta}} \right\| \\ &\leq \sup_{\phi \in \Phi} \left\| \frac{\partial^2 \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \phi)}{\partial \phi \partial \boldsymbol{\theta}} \right\| \sup_{\phi^0 \in \Phi} \left\| \frac{\partial \hat{\phi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \right\| \\ &\quad + \sup_{\phi \in \Phi} \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \phi)}{\partial \phi} \right\| \sup_{\phi^0 \in \Phi} \left\| \frac{\partial^2 \hat{\phi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0 \partial \boldsymbol{\theta}} \right\| \\ &\quad + \sup_{\phi \in \Phi} \left\| \frac{\partial^2 \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \phi)}{\partial \phi^2} \right\| \sup_{\phi^0 \in \Phi} \left\| \frac{\partial \hat{\phi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \right\| \sup_{\phi^0 \in \Phi} \left\| \frac{\partial \hat{\phi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \boldsymbol{\theta}} \right\| \\ &\leq A_1(1 - \alpha A_3^{-1})^n + A_1 n(1 - \alpha A_3^{-1})^{n-1} \alpha A_2(A_2 A_3 + 1) \\ &\quad + A_1(1 - \alpha A_3^{-1})^n A_3 A_2, \quad \forall \boldsymbol{\theta} \in \Theta, \forall n \geq 0. \end{aligned} \tag{4}$$

As a result, we can bound the first term of Eqn. (2) as

$$\begin{aligned} &\left\| \frac{\partial \hat{\mathcal{J}}_{Bi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \boldsymbol{\theta}} - \frac{\partial \hat{\mathcal{J}}_{Bi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \boldsymbol{\theta}} \Big|_{\phi^0=\hat{\phi}^*(\boldsymbol{\theta})} \right\| \\ &\leq A_1(1 + A_2 A_3) \left(1 + \frac{\alpha A_2}{1 - \alpha A_3^{-1}} n\right) (1 - \alpha A_3^{-1})^n \|\phi^0 - \hat{\phi}^*(\boldsymbol{\theta})\|, \quad \forall \boldsymbol{\theta} \in \Theta, \forall \phi^0 \in \Phi, \forall n \geq 0. \end{aligned} \tag{5}$$

For the second term of Eqn. (2), the partial derivative $\frac{\partial \hat{\mathcal{J}}_{Bi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \Big|_{\phi^0=\hat{\phi}^*(\boldsymbol{\theta})}$ can be expanded as

$$\frac{\partial \hat{\mathcal{J}}_{Bi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \Big|_{\phi^0=\hat{\phi}^*(\boldsymbol{\theta})} = \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \phi)}{\partial \phi} \Big|_{\phi=\hat{\phi}^*(\boldsymbol{\theta})} \frac{\partial \hat{\phi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \Big|_{\phi^0=\hat{\phi}^*(\boldsymbol{\theta})},$$

and thereby

$$\begin{aligned} \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \Big|_{\phi^0=\hat{\phi}^*(\boldsymbol{\theta})} \right\| &\leq \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \phi)}{\partial \phi} \Big|_{\phi=\hat{\phi}^*(\boldsymbol{\theta})} \right\| \left\| \frac{\partial \hat{\phi}^n(\boldsymbol{\theta}, \phi^0)}{\partial \phi^0} \Big|_{\phi^0=\hat{\phi}^*(\boldsymbol{\theta})} \right\| \\ &\leq A_1(1 - \alpha A_3^{-1})^n, \quad \forall \boldsymbol{\theta} \in \Theta, \forall n \geq 0. \end{aligned}$$

For calculating $\frac{\partial \hat{\phi}^*(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$, we take partial derivative to $\frac{\partial \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi} \Big|_{\phi=\hat{\phi}^*(\boldsymbol{\theta})} = 0$ w.r.t. $\boldsymbol{\theta}$ and get

$$\frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \boldsymbol{\theta} \partial \phi} \Big|_{\phi=\hat{\phi}^*(\boldsymbol{\theta})} + \frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi^2} \Big|_{\phi=\hat{\phi}^*(\boldsymbol{\theta})} \frac{\partial \hat{\phi}^*(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = 0,$$

and thereby

$$\frac{\partial \hat{\phi}^*(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = - \left(\frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi^2} \Big|_{\phi=\hat{\phi}^*(\boldsymbol{\theta})} \right)^{-1} \frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \boldsymbol{\theta} \partial \phi} \Big|_{\phi=\hat{\phi}^*(\boldsymbol{\theta})},$$

$$\left\| \frac{\partial \hat{\phi}^*(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right\| \leq \left\| \left(\frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \phi^2} \Big|_{\phi=\hat{\phi}^*(\boldsymbol{\theta})} \right)^{-1} \right\| \left\| \frac{\partial^2 \hat{\mathcal{G}}(\boldsymbol{\theta}, \phi)}{\partial \boldsymbol{\theta} \partial \phi} \Big|_{\phi=\hat{\phi}^*(\boldsymbol{\theta})} \right\| \leq A_3 A_2, \quad \forall \boldsymbol{\theta} \in \Theta.$$

Thus, the second term of Eqn. (2) can be bounded as

$$\begin{aligned} \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}^n(\boldsymbol{\theta}, \boldsymbol{\phi}^0)}{\partial \boldsymbol{\phi}^0} \Big|_{\boldsymbol{\phi}^0 = \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta})} \frac{\partial \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right\| &\leq \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}^n(\boldsymbol{\theta}, \boldsymbol{\phi}^0)}{\partial \boldsymbol{\phi}^0} \Big|_{\boldsymbol{\phi}^0 = \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta})} \right\| \left\| \frac{\partial \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right\| \\ &\leq A_1(1 - \alpha A_3^{-1})^n A_3 A_2, \quad \forall \boldsymbol{\theta} \in \Theta, \forall n \geq 0. \end{aligned} \quad (6)$$

By Eqn. (2,5,6), we get

$$\begin{aligned} &\left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^n(\boldsymbol{\theta}, \boldsymbol{\phi}^0))}{\partial \boldsymbol{\theta}} - \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\| \\ &\leq A_1(1 + A_2 A_3) \left(1 + \frac{\alpha A_2}{1 - \alpha A_3^{-1}} n\right) (1 - \alpha A_3^{-1})^n \|\boldsymbol{\phi}^0 - \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta})\| \\ &\quad + A_1(1 - \alpha A_3^{-1})^n A_3 A_2, \quad \forall \boldsymbol{\theta} \in \Theta, \forall \boldsymbol{\phi}^0 \in \Phi, \forall n \geq 0. \end{aligned}$$

Let $A := A_1(1 + A_2 A_3)$, $B := A_1(1 + A_2 A_3) \frac{\alpha A_2}{1 - \alpha A_3^{-1}}$, $C := A_1 A_2 A_3$, $\kappa := 1 - \alpha A_3^{-1}$, then $A, B, C > 0$ and $\kappa \in (0, 1)$ are constants independent of $\boldsymbol{\theta}$ and $\boldsymbol{\phi}^0$ and

$$\begin{aligned} &\left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^n(\boldsymbol{\theta}, \boldsymbol{\phi}^0))}{\partial \boldsymbol{\theta}} - \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\| \\ &\leq (A + Bn)\kappa^n \|\boldsymbol{\phi}^0 - \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta})\| + C\kappa^n, \quad \forall \boldsymbol{\theta} \in \Theta, \forall \boldsymbol{\phi}^0 \in \Phi, \forall n \geq 0. \end{aligned}$$

□

A.5 Proof of Corollary 3

Corollary 3. (BiSM finds δ -stationary points) For any accuracy level $\delta > 0$, assuming Theorem 2 holds, using a sufficiently large N , i.e. asymptotically $\mathcal{O}(\log \frac{1}{\delta})$, and a proper learning rate scheme β [1], Algorithm 1 in the main text converges to a δ -stationary point of BiSM, namely,

$$\left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\| \leq \delta,$$

and further a δ -stationary point of SM if Theorem 1 also holds.

Proof. For any $\delta > 0$ and $\boldsymbol{\theta} \in \Theta$, assuming that

$$\left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\| > \delta,$$

we have

$$\begin{aligned} &2 \left\langle \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}}, \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^N(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\rangle \\ &= \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\|^2 + \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^N(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\|^2 - \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} - \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^N(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\|^2 \\ &\geq \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\|^2 - \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} - \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^N(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\|^2 \\ &> \delta^2 - \left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} - \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^N(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\|^2. \end{aligned}$$

If Theorem 2 holds, using a sufficiently large N , i.e. asymptotically $\mathcal{O}(\log \frac{1}{\delta})$, we have

$$\left\| \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} - \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^N(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\|^2 \leq \delta^2,$$

which implies that

$$\left\langle \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}}, \frac{\partial \hat{\mathcal{J}}_{Bi}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^N(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right\rangle > 0.$$

Therefore, using a proper learning rate scheme β such that $\sum_{k=1}^{\infty} \beta_k = \infty$, $\sum_{k=1}^{\infty} \beta_k^2 < \infty$ [1], Algorithm 1, i.e. stochastic gradient descent based on $\frac{\partial \hat{\mathcal{J}}_{\text{Bi}}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^N(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}}$, will decrease $\|\frac{\partial \hat{\mathcal{J}}_{\text{Bi}}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}}\|$ in expectation until it converges to a δ -stationary point of BiSM such that $\|\frac{\partial \hat{\mathcal{J}}_{\text{Bi}}(\boldsymbol{\theta}, \hat{\boldsymbol{\phi}}^*(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}}\| \leq \delta$, according to Corollary 4.12 in Bottou et al. [1], whose regularity conditions are covered by the assumptions in Theorem 2. Further, if a δ -stationary point of BiSM is also a δ -stationary point of SM if Theorem 1 also holds. \square

B Experimental Settings

B.1 GRBM

The batch size is 100 on both the checkerboard dataset and the Frey face dataset². We train 100,000 iterations on the checkerboard dataset and 20,000 iterations on the Frey face dataset. The noise level [14] of DSM and BiDSM is 0.05 on the checkerboard dataset and 0.3 for on the Frey face dataset. The type of random directions [12] of SSM and BiSSM is the multivariate Rademacher distribution on both datasets. We choose $q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})$ as a Bernoulli distribution for all BiSM methods and use the Gumbel-Softmax trick [7] for reparameterization of $q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})$ with 0.1 as the temperature.

On both datasets, we tune the learning rate in $\{10^{-4}, 3 \times 10^{-4}, 10^{-3}, 3 \times 10^{-3}, 10^{-2}\}$ according to the visual quality of density plots and samples respectively. On the checkerboard dataset, all methods achieve similar results with the learning rates 10^{-3} , 3×10^{-3} and 10^{-2} and we choose 10^{-3} as the default value. On the Frey face dataset, we find that both DSM and BiDSM can work on the learning rate 10^{-4} and 3×10^{-4} and we choose 2×10^{-4} as the final learning rate. We also split a validation dataset from the Frey face dataset to choose the best model according to their corresponding loss on the validation dataset. We run 10 evaluations of the validation dataset during training.

On the Frey face dataset, we tune the noise level in $\{0.01, 0.03, 0.1, 0.3, 1\}$ for DSM and BiDSM and both methods only work on the noise level 0.3, so we choose 0.3 as the final noise level.

We run 1,000 steps Gibbs sampling to sample from GRBM on both datasets and all methods.

B.2 Deep EBLVM

The batch size is 100 on both the MNIST, CIFAR10 and CelebA datasets. We scale the CelebA datasets to 32×32 and 64×64 and explicitly denote them as CelebA32 or CelebA64 when necessary. Following [9], we train 100,000 iterations on the MNIST dataset and 300,000 iterations on the CIFAR10 and the CelebA datasets; the noise level is geometrically distributed in the range $[0.1, 3.0]$ on the MNIST dataset and uniformly distributed in the range $[0.05, 1.2]$ on the CIFAR10 and the CelebA dataset; σ_0 (see [9]) is 0.1 on both datasets. We choose $q(\mathbf{h}|\mathbf{v}; \boldsymbol{\phi})$ as a Gaussian distribution parameterized by a 3-layer convolutional neural network for BiMDSM.

The energy function is $\mathcal{E}(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = g_3(g_2(g_1(\mathbf{v}; \boldsymbol{\theta}_1), \mathbf{h}); \boldsymbol{\theta}_2)$ for the deep EBLVM trained by BiMDSM and is $\mathcal{E}(\mathbf{v}; \boldsymbol{\theta}) = g_3(g_1(\mathbf{v}; \boldsymbol{\theta}_1); \boldsymbol{\theta}_2)$ for the fully visible deep EBM trained by the baseline MDSM. $g_1(\cdot)$ is a 12-layer ResNet for MNIST, a 18-layer ResNet for CIFAR10 and CelebA32 following [9], or a 24-layer ResNet for CelebA64. For the EBLVM, an extra fully connected layer is introduced in $g_1(\cdot)$ to match the dimension of \mathbf{h} . $g_2(\cdot, \cdot)$ is an additive coupling layer [3] to make the features output by $g_1(\cdot)$ and the latent variables strongly coupled. $g_3(\cdot)$ consists of a fully connected layer with an ELU activation function and use the square of 2-norm to output a scalar.

To sample from the deep EBLVM, we firstly resample data from the training dataset and inference their approximate posterior mean. We then sample from $p(\mathbf{v}|\mathbf{h})$ given \mathbf{h} equal to the approximate posterior mean. Although it introduces bias due to the difference between $p(\mathbf{h}|\mathbf{v})$ and $q(\mathbf{h}|\mathbf{v})$, we find this sampling procedure can increase sample quality and diversity compared to directly sampling from $p(\mathbf{v}, \mathbf{h})$. Besides, this sampling procedure is not to reconstruct the training data since $p(\mathbf{v}|\mathbf{h})$ is multimodal, as shown in Fig. 5. To sample from deep EBM, we directly sample from $p(\mathbf{v})$. We use the annealed Langevin dynamics [9] as our sampling technique. Following [9], we choose $[1, 100]$ as the range of temperature and 0.02 as the step length for annealed Langevin dynamics on both EBLVM and EBM.

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

For MNIST, MDSM spends about 4 hours training a deep EBM and BiMDSM spends about 8 hours training a deep EBLVM. For CIFAR10, MDSM spends about 32 hours training a deep EBM and BiMDSM spends about 48 hours training a deep EBLVM. For CelebA32, BiMDSM spends about 48 hours training a deep EBLVM. The above experiments on deep EBLVMs and deep EBMs are conducted on 1 GeForce RTX 2080 Ti GPU. For CelebA64, BiMDSM spends about a week training a deep EBLVM on 4 GeForce RTX 1080 Ti GPUs.

C Additional Results

C.1 GRBM

C.1.1 Sensitivity analysis on N

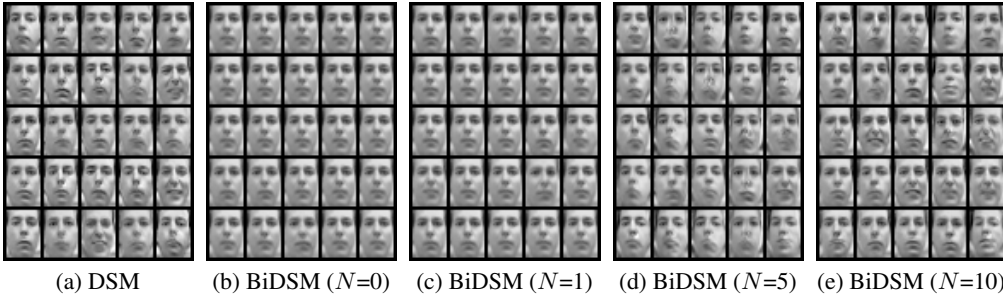
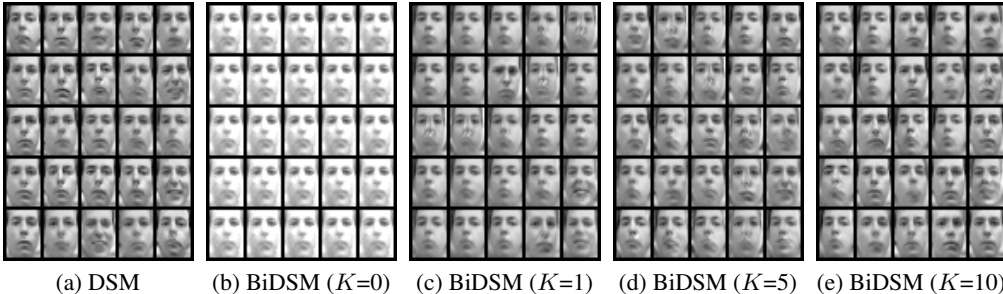


Figure 1: Samples from GRBMs trained by DSM and BiDSM on different N (0, 1, 5 and 10) according to the best validation performance on the Frey face dataset.

Fig. 1 shows samples from GRBMs trained by DSM and BiDSM on different N . The sample quality of BiDSM increases as N increases, and is comparable to DSM when $N=10$. The result is consistent with the test Fisher divergence quantitative results in Tab. 1 in the full paper.

C.1.2 Sensitivity analysis on K



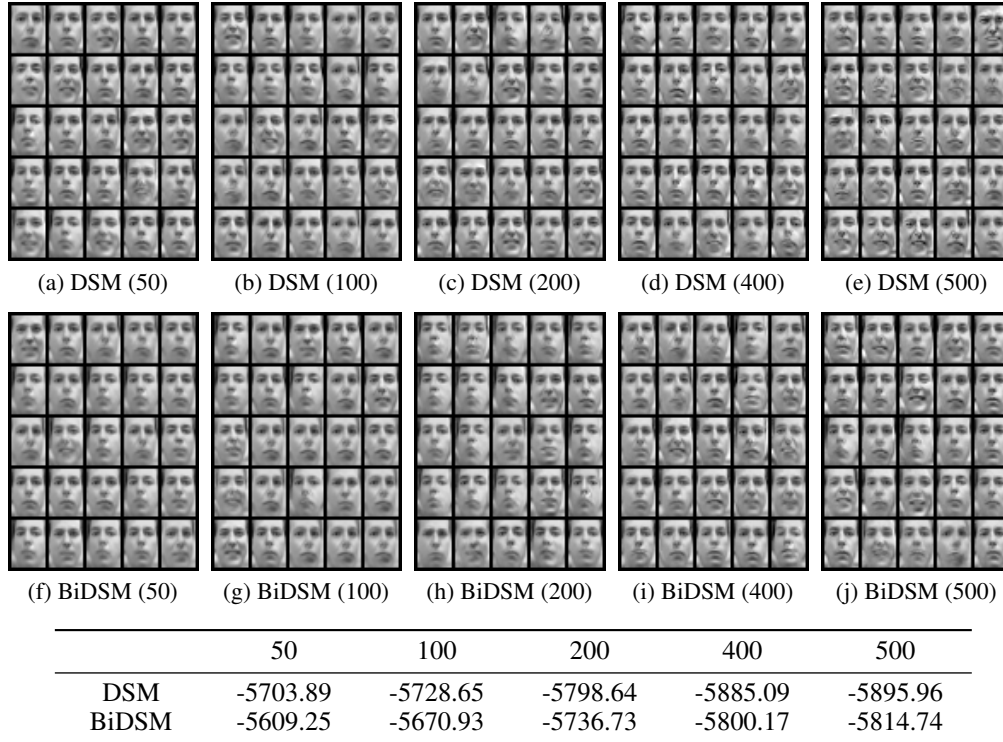
DSM	BiDSM ($K=0$)	BiDSM ($K=1$)	BiDSM ($K=5$)	BiDSM ($K=10$)
-5885.09	-3775.31	-5684.97	-5780.18	-5795.52

(f) Test Fisher divergence \downarrow (subtracted by the same unknown constant only relevant to the data)

Figure 2: Samples from GRBMs trained by DSM and BiDSM on different K (0, 1, 5 and 10) according to the best validation performance on the Frey face dataset.

Fig. 2 shows samples and test Fisher divergence from GRBMs trained by DSM and BiDSM on different K . The sample quality of BiDSM increases as K increases and the test Fisher divergence decreases as K increases. When $K = 0$, the variational posterior $q(\mathbf{h}|\mathbf{v}; \phi)$ doesn't change during training, leading to a much worse result than others.

C.1.3 Sensitivity analysis on dimensions of h



(k) Test Fisher divergence \downarrow (subtracted by the same unknown constant only relevant to the data)

Figure 3: Samples (a-j) and test Fisher divergence (k) of GRBMs trained by DSM and BiDSM on different dimensions of h (50, 100, 200, 400 and 500) according to the best validation performance on the Frey face dataset. N is 10 for BiDSM.

Fig. 3 shows samples and test Fisher divergence from GRBM trained by DSM and BiDSM on different dimensions of h . Both the sample quality and test Fisher divergence of BiDSM are comparable to DSM on different dimensions of h .

C.1.4 Time Complexity Comparison

Table 1: Time complexity comparison in GRBMs on the Frey face dataset. The time is the averaged training time of 100 iterations. All experiments are conducted on one GeForce GTX 1080 Ti GPU.

(a) Comparison on N				(b) Comparison on K			
Methods	Time (s)			Methods	Time (s)		
BiDSM ($N=0, K=5$)	4.35			BiDSM ($K=1, N=5$)	7.30		
BiDSM ($N=1, K=5$)	5.07			BiDSM ($K=2, N=5$)	7.75		
BiDSM ($N=5, K=5$)	8.61			BiDSM ($K=5, N=5$)	8.61		
BiDSM ($N=10, K=5$)	13.78			BiDSM ($K=10, N=5$)	10.82		

(c) Comparison between different methods						
Methods	CD-5	SSM	DSM	VNCE	BiDSM ($N=0, K=5$)	BiDSM ($N=K=5$)
Time (s)	1.59	1.51	1.33	4.36	4.35	8.61

According to Algorithm 1, the time complexity and space complexity in a training iteration is $\mathcal{O}(K + N)$ and $\mathcal{O}(N)$ respectively. In Tab. 1 (a-b), we show the time complexity comparison of BiDSM on different N (0, 1, 5 and 10) and K (0, 1, 5 and 10). The training time is approximately

linearly correlated to both N and K . In Tab. 1 (c), we show time complexity comparison between different methods. VNCE and our BiDSM are two methods of learning nonstructural EBLVMs, which require extra time to learn in a black-box manner compared to CD-5, SSM and DSM. While VNCE and BiDSM ($N=0, K=5$) have the similar time complexity, to the best of our knowledge VNCE hasn't been shown feasible to scale up to natural images, including the Frey face dataset (it doesn't hurt the time complexity comparison on this dataset). Besides, as stated in Appendix B.2, the training time of 100,000 iterations is 8h for BiMDSM in a deep EBLVM and 4h for MDSM in a deep EBM on MNIST; the training time of 300,000 iterations is 48h for BiMDSM in a deep EBLVM and 32h for MDSM in a deep EBM on CIFAR10. Thus, BiSM can learn general EBLVMs without a prohibitive cost.

C.2 Deep EBLVM

C.2.1 Sensitivity analysis on dimensions of h

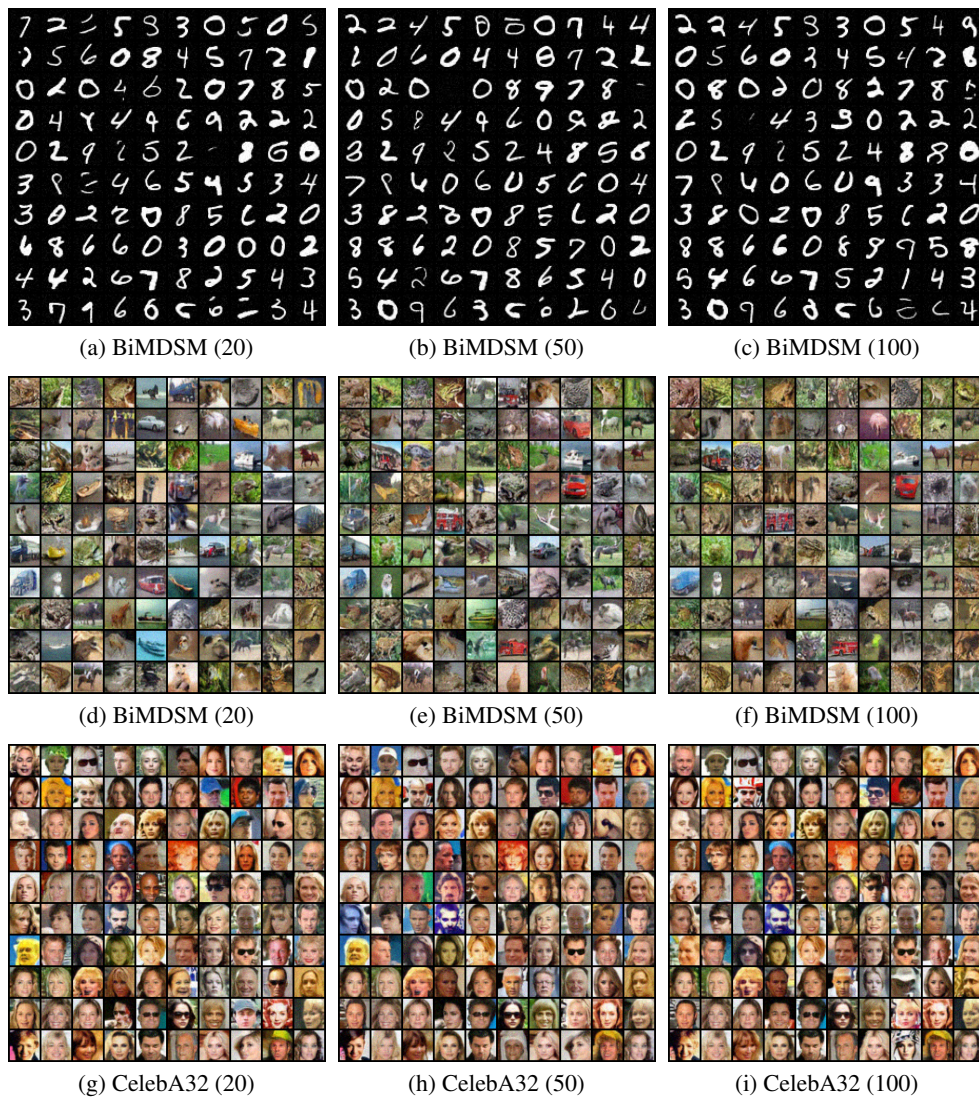


Figure 4: Samples from EBLVMs trained by BiMDSM on the MNIST, CIFAR10 and CelebA32 datasets with different dimensions of h (20, 50 and 100).

Fig. 4 shows samples from EBLVMs trained by BiMDSM on the MNIST, CIFAR10 and CelebA32 datasets with different dimensions of h . The EBLVMs can produce meaningful samples in all

settings. Notice that across different dimensions of \mathbf{h} on one dataset, samples at the same position are sometimes similar because we initialize the same random seeds for different dimensions of \mathbf{h} .

C.2.2 Conditionally Sampling

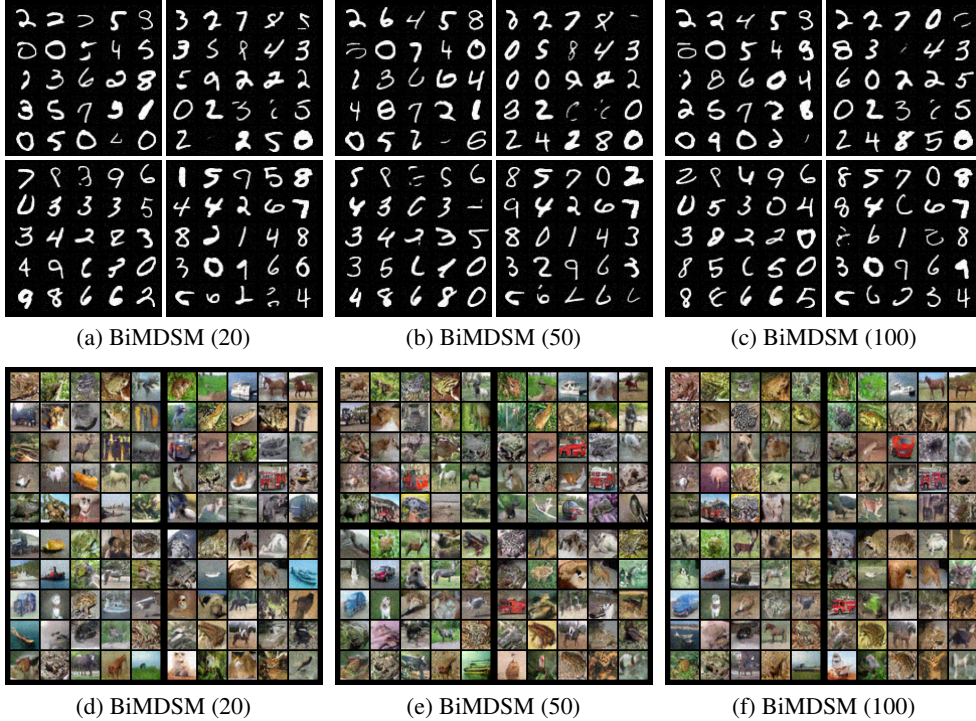


Figure 5: Samples from conditional distribution $p(\mathbf{v}|\mathbf{h})$ of EBLVMs trained by BiMDSM on the MNIST and CIFAR10 datasets with different dimensions of \mathbf{h} (20, 50 and 100).

Fig. 5 shows samples from conditional distribution $p(\mathbf{v}|\mathbf{h})$ of EBLVMs trained by BiMDSM on the MNIST and CIFAR10 datasets with different dimensions of \mathbf{h} . Each subfigure is split to four parts, and samples in the same part correspond to the same \mathbf{h} , which is inferred from a training data via the approximate posterior mean. On each dataset, we use the same four training data in all settings to infer \mathbf{h} . The samples from $p(\mathbf{v}|\mathbf{h})$ are highly diverse, suggesting that $p(\mathbf{v}|\mathbf{h})$ of an EBLVM is multimodal. Intrinsicly, this is because the deep EBLVMs used here defines the conditional distribution $p(\mathbf{v}|\mathbf{h})$ in a highly nonstructural way, in contrast to the hierarchical manner used in previous methods [8, 5, 6, 10]. Notice that it doesn't contradict the inference results in Sec. C.2.3, since $p(\mathbf{v}|\mathbf{h}) \propto p(\mathbf{v})p(\mathbf{h}|\mathbf{v})$ and $p(\mathbf{v}|\mathbf{h})$ can be dominated by $p(\mathbf{v})$.

C.2.3 Inference Results

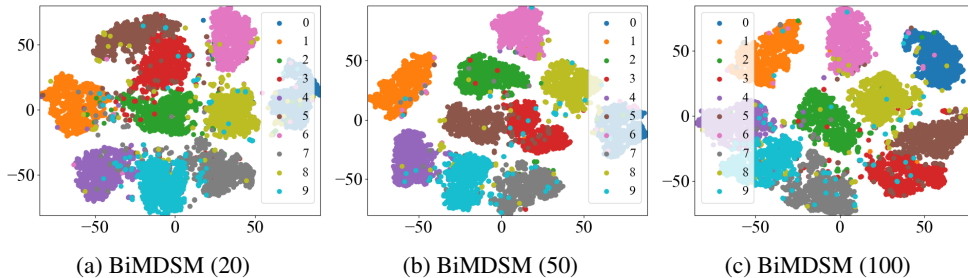


Figure 6: t-SNE [13] embedding of the approximate posterior mean for the test MNIST data on different dimensions of \mathbf{h} .

Table 2: Test classification accuracy (%) \uparrow of the approximate posterior mean of EBLVMs trained by BiMDSM. We show results on the MNIST and CIFAR10 datasets with different dimensions of h (20, 50 and 100). We use default linear SVM [4] implemented by sklearn as the classifier.

	BiMDSM (20)	BiMDSM (50)	BiMDSM (100)	Linear SVM on raw data
MNIST	93.85	97.39	97.75	91.58
CIFAR10	34.83	39.58	46.46	28.19

Fig. 6 shows the t-SNE [13] embedding of the approximate posterior mean for the test MNIST data on different dimensions of h . For MNIST, the embedding can be well separated on different dimensions of h . For CIFAR10, the intra-class distance can sometime be larger than the inner-class distance, and thereby the embedding can hardly be separated according to the class [2].

We train a linear SVM classifier³ using the posterior mean learned by BiMDSM as features. Tab. 2 shows the test classification accuracy. On both datasets, the accuracy increases as the dimension of h increases. The results of BiMDSM are better than a linear SVM classifier trained on raw data, suggesting that the features capture the underlying semantics of the images. We mention that previous EBLVMs [6, 10, 11] apply a supervised fine-tuning procedure to obtain better classification results. In contrast, we focus on a purely unsupervised learning setting here because our main goal is not to achieve the state-of-the-art classification results but to validate that the deep EBLVMs learned by BiSM can extract semantic features from natural images.

C.2.4 Results on CelebA64



Figure 7: Samples of 64×64 resolution from an EBLVM trained by BiMDSM on CelebA64. The dimension of h is 20.

Fig. 7 shows promising results on scaling to images of higher resolutions. We show samples of 64×64 resolution from an EBLVM trained by BiMDSM on CelebA64, which are of high diversity.

References

- [1] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *Siam Review*, 60(2):223–311, 2018.

³We use the default implementation provided by the sklearn package. Please see details in the online document: <https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html>

- [2] LI Chongxuan, Max Welling, Jun Zhu, and Bo Zhang. Graphical generative adversarial networks. In *Advances in neural information processing systems*, pages 6069–6080, 2018.
- [3] Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.
- [4] Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. Liblinear: A library for large linear classification. *Journal of machine learning research*, 9(Aug):1871–1874, 2008.
- [5] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [6] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.
- [7] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- [8] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [9] 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.
- [10] Ruslan Salakhutdinov and Geoffrey Hinton. Deep Boltzmann machines. In *Proceedings of the twelfth international conference on artificial intelligence and statistics*, 2009.
- [11] 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.
- [12] 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.
- [13] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008. URL <http://www.jmlr.org/papers/v9/vandermaaten08a.html>.
- [14] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.