

A. Proof of Theorem 1

Assumption 1. The function $q_\theta(\mathbf{x})f^*(T_\omega(\mathbf{x}))$ and its partial derivative $\nabla_\theta q_\theta(\mathbf{x})f^*(T_\omega(\mathbf{x}))$ are continuous w.r.t. θ and \mathbf{x} .

Theorem 1. For a θ -parametrized energy-based model $q_\theta(\mathbf{x}) = \frac{\exp(-E_\theta(\mathbf{x}))}{Z_\theta}$, under Assumption 1, the gradient of the variational representation of *f*-divergence (Equation (3)) w.r.t. θ can be written as:

$$-\nabla_\theta \mathbb{E}_{q_\theta(\mathbf{x})}[f^*(T_\omega(\mathbf{x}))] = \mathbb{E}_{q_\theta(\mathbf{x})}[\nabla_\theta E_\theta(\mathbf{x})f^*(T_\omega(\mathbf{x}))] - \mathbb{E}_{q_\theta(\mathbf{x})}[\nabla_\theta E_\theta(\mathbf{x})] \cdot \mathbb{E}_{q_\theta(\mathbf{x})}[f^*(T_\omega(\mathbf{x}))]$$

When we can obtain i.i.d. samples from q_θ , we can get an unbiased estimation of the gradient.

Proof. Under Assumption 1, with Leibniz integral rule, we have:

$$\begin{aligned} -\nabla_\theta \mathbb{E}_{q_\theta(\mathbf{x})}[f^*(T_\omega(\mathbf{x}))] &= -\int \nabla_\theta q_\theta(\mathbf{x})f^*(T_\omega(\mathbf{x}))d\mathbf{x} \\ &= -\int q_\theta(\mathbf{x})\frac{\nabla_\theta q_\theta(\mathbf{x})}{q_\theta(\mathbf{x})}f^*(T_\omega(\mathbf{x}))d\mathbf{x} \\ &= -\int q_\theta(\mathbf{x})\nabla_\theta \log q_\theta(\mathbf{x})f^*(T_\omega(\mathbf{x}))d\mathbf{x} \\ &= -\int q_\theta(\mathbf{x})[-\nabla_\theta E_\theta(\mathbf{x}) - \nabla_\theta \log Z_\theta]f^*(T_\omega(\mathbf{x}))d\mathbf{x} \\ &= \int q_\theta(\mathbf{x})\nabla_\theta E_\theta(\mathbf{x})f^*(T_\omega(\mathbf{x}))d\mathbf{x} + \nabla_\theta \log Z_\theta \cdot \int q_\theta(\mathbf{x})f^*(T_\omega(\mathbf{x}))d\mathbf{x} \\ &= \int q_\theta(\mathbf{x})\nabla_\theta E_\theta(\mathbf{x})f^*(T_\omega(\mathbf{x}))d\mathbf{x} + \frac{\int \exp(-E_\theta(\mathbf{x}))(-\nabla_\theta E_\theta(\mathbf{x}))d\mathbf{x}}{Z_\theta} \cdot \int q_\theta(\mathbf{x})f^*(T_\omega(\mathbf{x}))d\mathbf{x} \\ &= \mathbb{E}_{q_\theta(\mathbf{x})}[\nabla_\theta E_\theta(\mathbf{x})f^*(T_\omega(\mathbf{x}))] - \mathbb{E}_{q_\theta(\mathbf{x})}[\nabla_\theta E_\theta(\mathbf{x})] \cdot \mathbb{E}_{q_\theta(\mathbf{x})}[f^*(T_\omega(\mathbf{x}))] \end{aligned}$$

We can simply use i.i.d. samples from q_θ to get an unbiased estimation of the first term, and we introduce the following lemma for estimating the second term $\mathbb{E}_{q_\theta(\mathbf{x})}[\nabla_\theta E_\theta(\mathbf{x})] \cdot \mathbb{E}_{q_\theta(\mathbf{x})}[f^*(T_\omega(\mathbf{x}))]$.

Lemma 3. Let P denote a probability distribution over sample space \mathcal{X} and $g, h : \mathcal{X} \rightarrow \mathbb{R}$ are functions. The following empirical estimator is unbiased with respect to $\mathbb{E}_P[g(X)] \cdot \mathbb{E}_P[h(X)]$:

$$\frac{1}{nm} \sum_{i=1}^n g(x_i) \cdot \sum_{j=1}^m h(y_j) \tag{21}$$

where $\mathbf{x}_{1:n} \sim P$, $\mathbf{y}_{1:m} \sim P$ are two independent sets of i.i.d. samples from P .

Proof. Let $X_{1:n}$ and $Y_{1:m}$ denote the random variables corresponding to the sampling process for $\mathbf{x}_{1:n}$, $\mathbf{y}_{1:m}$ respectively. Let P_n and P_m denote the probability distribution for $X_{1:n}$ and $Y_{1:m}$. Then we have:

$$\mathbb{E}_{P_n} \mathbb{E}_{P_m} \left[\frac{1}{nm} \sum_{i=1}^n g(X_i) \cdot \sum_{j=1}^m h(Y_j) \right] \tag{22}$$

$$= \mathbb{E}_{P_n} \left[\frac{1}{n} \sum_{i=1}^n g(X_i) \cdot \mathbb{E}_{P_m} \left[\frac{1}{m} \sum_{j=1}^m h(Y_j) \right] \right] \tag{23}$$

$$= \mathbb{E}_{P_n} \left[\frac{1}{n} \sum_{i=1}^n g(X_i) \cdot \mathbb{E}_P[h(X)] \right] \tag{24}$$

$$= \mathbb{E}_P[g(X)] \cdot \mathbb{E}_P[h(X)] \tag{25}$$

where the first equality comes from the law of total expectation and the fact that $X_{1:n}$ and $Y_{1:m}$ are independent. Therefore, the estimator is unbiased. \square

Lemma 3 shows that we can use i.i.d. samples from q_θ to get an unbiased estimation of the second term. \square

B. Proof of Theorem 2

Theorem 2. Let P and Q be two probability measures over the Borel σ -algebra on domain \mathcal{X} with densities p, q and $P \ll Q$. Additionally, let q be an energy-based distribution, $q(\mathbf{x}) = \exp(-E(\mathbf{x}))/Z_q$. For any class of functions \mathcal{H} mapping from \mathcal{X} to \mathbb{R} such that the expectations in the following equation are finite and $\forall \mathbf{x} \in \mathcal{X}$, \mathcal{H} contains an element $H(\mathbf{x}) = \log(p(\mathbf{x})Z_q)$, then we have

$$D_f(P||Q) = \sup_{H \in \mathcal{H}} \mathbb{E}_{p(\mathbf{x})}[f'(\exp(H(\mathbf{x}) + E(\mathbf{x})))] - \mathbb{E}_{q(\mathbf{x})}[f^*(f'(\exp(H(\mathbf{x}) + E(\mathbf{x})))))]$$

where the supreme is attained at $H^*(\mathbf{x}) = \log(p(\mathbf{x})Z_q)$.

Proof. First, we define a function class $\hat{\mathcal{T}} = \{f'(r)|r : \mathcal{X} \rightarrow \mathbb{R}_+\}$. Since $\hat{\mathcal{T}}$ is a subset of \mathcal{T} (in Lemma 1) and the optimal $T^* = f'(p/q)$ is an element of $\hat{\mathcal{T}}$, from Lemma 1, we have:

$$D_f(P||Q) = \sup_{T \in \hat{\mathcal{T}}} \mathbb{E}_{p(\mathbf{x})}[T(\mathbf{x})] - \mathbb{E}_{q(\mathbf{x})}[f^*(T(\mathbf{x}))]. \quad (26)$$

Without loss of generality, we can reparametrize all functions $T \in \hat{\mathcal{T}}$ with H such that $T(\mathbf{x}) = f'(\exp(H(\mathbf{x}) + E(\mathbf{x})))$ for all $\mathbf{x} \in \mathcal{X}$, since the transformations $f'(\cdot)$, $\exp(\cdot)$ and addition by $E(\mathbf{x})$ are all bijections (note that f is strictly convex and differentiable). The optimal $H^*(\mathbf{x})$ then satisfies:

$$\forall \mathbf{x} \in \mathcal{X}, f'(\exp(H^*(\mathbf{x}) + E(\mathbf{x}))) = T^*(\mathbf{x}) = f'(p(\mathbf{x})/q(\mathbf{x}))$$

From the bijectivity of $f'(\cdot)$, we have:

$$\begin{aligned} H^*(\mathbf{x}) &= \log p(\mathbf{x}) - \log q(\mathbf{x}) - E(\mathbf{x}) \\ &= \log p(\mathbf{x}) + E(\mathbf{x}) + \log Z_q - E(\mathbf{x}) = \log(p(\mathbf{x})Z_q), \end{aligned}$$

□

C. Proof of Theorem 3

Assumption 2. The function $p(\mathbf{x})f'(\exp(H_\omega(\mathbf{x}) + E_\theta(\mathbf{x})))$ and its partial derivative $\nabla_\theta f'(\exp(H_\omega(\mathbf{x}) + E_\theta(\mathbf{x})))$ are continuous w.r.t. θ and \mathbf{x} .

Assumption 3. The function $q_\theta(\mathbf{x})f^*(f'(\exp(H_\omega(\mathbf{x}) + E_\theta(\mathbf{x}))))$ and its partial derivative $\nabla_\theta q_\theta(\mathbf{x})f^*(f'(\exp(H_\omega(\mathbf{x}) + E_\theta(\mathbf{x}))))$ are continuous w.r.t. θ and \mathbf{x} .

Theorem 3. For a θ -parametrized energy-based model $q_\theta(\mathbf{x}) = \frac{\exp(-E_\theta(\mathbf{x}))}{Z_\theta}$ and a fixed ω , under Assumptions 2 and 3, the gradient of $\mathcal{L}_{f\text{-EBM}}(\theta, \omega)$ (the objective in Equation (14)) with respect to θ can be written as:

$$\begin{aligned} \nabla_\theta \mathcal{L}_{f\text{-EBM}}(\theta, \omega) &= \mathbb{E}_{p(\mathbf{x})}[\nabla_\theta f'(\exp(f_\omega(\mathbf{x}) + E_\theta(\mathbf{x})))] + \mathbb{E}_{q_\theta(\mathbf{x})}[F_{\theta, \omega}(\mathbf{x})\nabla_\theta E_\theta(\mathbf{x})] - \\ &\quad \mathbb{E}_{q_\theta(\mathbf{x})}[\nabla_\theta F_{\theta, \omega}(\mathbf{x})] - \mathbb{E}_{q_\theta(\mathbf{x})}[\nabla_\theta E_\theta(\mathbf{x})] \cdot \mathbb{E}_{q_\theta(\mathbf{x})}[F_{\theta, \omega}(\mathbf{x})] \end{aligned}$$

where $F_{\theta, \omega}(\mathbf{x}) = f^*(f'(\exp(H_\omega(\mathbf{x}) + E_\theta(\mathbf{x}))))$. When we can obtain i.i.d. samples from p and q_θ , we can get an unbiased estimation of the gradient.

Proof. Under Assumption 2, with Leibniz integral rule, for the first term in Equation (14), we have:

$$\nabla_\theta \mathbb{E}_{p(\mathbf{x})}[f'(\exp(H_\omega(\mathbf{x}) + E_\theta(\mathbf{x})))] = \int p(\mathbf{x})\nabla_\theta f'(\exp(H_\omega(\mathbf{x}) + E_\theta(\mathbf{x})))d\mathbf{x} = \mathbb{E}_{p(\mathbf{x})}[\nabla_\theta f'(\exp(f_\omega(\mathbf{x}) + E_\theta(\mathbf{x})))]$$

For notational simplicity, let us use $F_{\theta, \omega}(\mathbf{x})$ to denote $f^*(f'(\exp(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x}))))$. Under Assumption 3, with Leibniz integral rule, for the second term, we have:

$$\begin{aligned}
 & -\nabla_{\theta} \mathbb{E}_{q_{\theta}(\mathbf{x})}[f^*(f'(\exp(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x})))))] \\
 = & -\int \nabla_{\theta}(q_{\theta}(\mathbf{x})F_{\theta, \omega}(\mathbf{x}))d\mathbf{x} \\
 = & -\int F_{\theta, \omega}(\mathbf{x})\nabla_{\theta}q_{\theta}(\mathbf{x})d\mathbf{x} - \int q_{\theta}(\mathbf{x})\nabla_{\theta}F_{\theta, \omega}(\mathbf{x})d\mathbf{x} \\
 = & -\int q_{\theta}(\mathbf{x})F_{\theta, \omega}(\mathbf{x})\frac{\nabla_{\theta}q_{\theta}(\mathbf{x})}{q_{\theta}(\mathbf{x})}d\mathbf{x} - \int q_{\theta}(\mathbf{x})\nabla_{\theta}F_{\theta, \omega}(\mathbf{x})d\mathbf{x} \\
 = & -\int q_{\theta}(\mathbf{x})F_{\theta, \omega}(\mathbf{x})\nabla_{\theta} \log q_{\theta}(\mathbf{x})d\mathbf{x} - \int q_{\theta}(\mathbf{x})\nabla_{\theta}F_{\theta, \omega}(\mathbf{x})d\mathbf{x} \\
 = & -\int q_{\theta}(\mathbf{x})F_{\theta, \omega}(\mathbf{x})[-\nabla_{\theta}E_{\theta}(\mathbf{x}) - \nabla_{\theta} \log Z_{\theta}]d\mathbf{x} - \int q_{\theta}(\mathbf{x})\nabla_{\theta}F_{\theta, \omega}(\mathbf{x})d\mathbf{x} \\
 = & \int q_{\theta}(\mathbf{x})F_{\theta, \omega}(\mathbf{x})\nabla_{\theta}E_{\theta}(\mathbf{x})d\mathbf{x} - \int q_{\theta}(\mathbf{x})\nabla_{\theta}F_{\theta, \omega}(\mathbf{x})d\mathbf{x} + \nabla_{\theta} \log Z_{\theta} \cdot \int q_{\theta}(\mathbf{x})F_{\theta, \omega}(\mathbf{x})d\mathbf{x} \\
 = & \int q_{\theta}(\mathbf{x})F_{\theta, \omega}(\mathbf{x})\nabla_{\theta}E_{\theta}(\mathbf{x})d\mathbf{x} - \int q_{\theta}(\mathbf{x})\nabla_{\theta}F_{\theta, \omega}(\mathbf{x})d\mathbf{x} + \frac{\int \exp(-E_{\theta}(\mathbf{x}))(-\nabla_{\theta}E_{\theta}(\mathbf{x}))d\mathbf{x}}{Z_{\theta}} \cdot \int q_{\theta}(\mathbf{x})F_{\theta, \omega}(\mathbf{x})d\mathbf{x} \\
 = & \mathbb{E}_{q_{\theta}(\mathbf{x})}[F_{\theta, \omega}(\mathbf{x})\nabla_{\theta}E_{\theta}(\mathbf{x})] - \mathbb{E}_{q_{\theta}(\mathbf{x})}[\nabla_{\theta}F_{\theta, \omega}(\mathbf{x})] - \mathbb{E}_{q_{\theta}(\mathbf{x})}[\nabla_{\theta}E_{\theta}(\mathbf{x})] \cdot \mathbb{E}_{q_{\theta}(\mathbf{x})}[F_{\theta, \omega}(\mathbf{x})]
 \end{aligned}$$

Similar to the proof in Appendix A, we can use *i.i.d.* samples from p and q_{θ} to get an unbiased estimation of the gradient. \square

D. Proof for the Local Convergence of *f*-EBM

D.1. Non-linear Dynamical Systems

In this section, we present a brief introduction of non-linear dynamical system theory (Hassan, 1996). For a comprehensive description, please refer to (Hassan, 1996; Nagarajan & Kolter, 2017). We further generalize some of the theories which will be useful for establishing the local convergence property of *f*-EBM later.

Consider a system consisting of variables $\phi \in \Phi \subseteq \mathbb{R}^n$ whose time derivative is defined by the vector field $v(\phi)$:

$$\dot{\phi} = v(\phi) \tag{27}$$

where $v : \Phi \rightarrow \mathbb{R}^n$ is a locally Lipschitz mapping from a domain Φ into \mathbb{R}^n .

Suppose ϕ^* is an equilibrium point of the system in Equation (27), *i.e.*, $v(\phi^*) = 0$. Let ϕ_t denote the state of the system at time t . To begin with, we introduce the following definition to characterize the stability of ϕ^* :

Definition 3 (Definition 4.1 in (Hassan, 1996)). *The equilibrium point ϕ^* for the system defined in Equation (27) is*

- *stable if for each $\epsilon > 0$, there is $\delta = \delta(\epsilon) > 0$ such that*

$$\|\phi_0 - \phi^*\| < \delta \implies \forall t \geq 0, \|\phi_t - \phi^*\| < \epsilon.$$

- *unstable if not stable.*
- *asymptotically stable if it is stable and $\delta > 0$ can be chosen such that*

$$\|\phi_0 - \phi^*\| < \delta \implies \lim_{t \rightarrow \infty} \phi_t = \phi^*.$$

- *exponentially stable if it is asymptotically stable and $\delta, k, \lambda > 0$ can be chosen such that:*

$$\|\phi_0 - \phi^*\| < \delta \implies \|\phi_t\| \leq k\|\phi_0\| \exp(-\lambda t).$$

The system is stable if for any value of ϵ , we can find a value of δ (possibly dependent on ϵ), such that a trajectory starting in a δ neighborhood of the equilibrium point will never leave the ϵ neighborhood of the equilibrium point. However, such a system may either converge to the equilibrium point or orbit within the ϵ ball. By contrast, asymptotic stability is a stronger notion of stability in the sense that trajectories starting in a δ neighborhood of the equilibrium will converge to the equilibrium point in the limit $t \rightarrow \infty$. Moreover, if ϕ^* is asymptotically stable, we call the algorithm obtained by iteratively applying the updates in Equation (27) *locally convergent* to ϕ^* . If ϕ^* is exponentially stable, we call the corresponding algorithm *linearly convergent* to ϕ^* .

Now we introduce the following theorem, which is central for studying the asymptotic stability of a system:

Theorem 4 (Theorem 4.15 in (Hassan, 1996)). *Let ϕ^* be an equilibrium point for the non-linear system*

$$\dot{\phi} = v(\phi) \tag{28}$$

where $v : \Phi \rightarrow \mathbb{R}^n$ is continuously differentiable and Φ is a neighborhood of ϕ^* . Let \mathbf{J} be the Jacobian of the system in Equation (28) at the equilibrium point:

$$\mathbf{J} = \left. \frac{\partial v(\phi)}{\partial \phi} \right|_{\phi=\phi^*} \tag{29}$$

Then, we have:

- The equilibrium point ϕ^* is asymptotically stable and exponentially stable if \mathbf{J} is a Hurwitz matrix, i.e., $\text{Re}(\lambda) < 0$ for all eigenvalues λ of \mathbf{J} .
- The equilibrium point ϕ^* is unstable if $\text{Re}(\lambda) > 0$ for one or more of the eigenvalues of \mathbf{J} .

Proof. See proof for Theorem 4.7, Theorem 4.15 and Corollary 4.3 in (Hassan, 1996). □

With Theorem 4, the stability of an equilibrium point can be analyzed by examining if all the eigenvalues of the Jacobian $v'(\phi)|_{\phi=\phi^*}$ have strictly negative real part.

In the following, we will use $\lambda_{\max}(\cdot)$ and $\lambda_{\min}(\cdot)$ to denote the largest and smallest eigenvalues of a non-zero positive semi-definite matrix. Now we introduce the following theorem to upper bound the real part of the eigenvalues of a matrix with a specific form:

Theorem 5. *Suppose $\mathbf{J} \in \mathbb{R}^{(m+n) \times (m+n)}$ is of the following form:*

$$\mathbf{J} = \begin{bmatrix} \mathbf{0} & \mathbf{P} \\ -\mathbf{P}^\top & -\mathbf{Q} \end{bmatrix}$$

where $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is a symmetric real positive definite matrix and $\mathbf{P}^\top \in \mathbb{R}^{n \times m}$ is a full column rank matrix. Then, for every eigenvalue λ of \mathbf{J} , $\text{Re}(\lambda) < 0$. More precisely, we have:

- When $\text{Im}(\lambda) = 0$,

$$\text{Re}(\lambda) \leq -\frac{\lambda_{\min}(\mathbf{Q})\lambda_{\min}(\mathbf{P}\mathbf{P}^\top)}{\lambda_{\min}(\mathbf{Q})\lambda_{\max}(\mathbf{Q}) + \lambda_{\min}(\mathbf{P}\mathbf{P}^\top)}$$

- When $\text{Im}(\lambda) \neq 0$,

$$\text{Re}(\lambda) \leq -\frac{\lambda_{\min}(\mathbf{Q})}{2}$$

Proof. See Lemma G.2 in (Nagarajan & Kolter, 2017). □

This theorem is useful for proving the local convergence of GANs. To prove the stability of *f*-EBM in the following sections, we need the following generalized theorem:

Theorem 6. Suppose $\mathbf{J} \in \mathbb{R}^{(m+n) \times (m+n)}$ is of the following form:

$$\mathbf{J} = \begin{bmatrix} -\mathbf{S} & \mathbf{P} \\ -\mathbf{P}^\top & -\mathbf{Q} \end{bmatrix}$$

where $\mathbf{S} \in \mathbb{R}^{m \times m}$ is a symmetric real positive semi-definite matrix, $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is a symmetric real positive definite matrix, $\mathbf{P}^\top \in \mathbb{R}^{n \times m}$ is a full column rank matrix. Then, for every eigenvalue λ of \mathbf{J} , $\text{Re}(\lambda) < 0$. More precisely, we have:

- When $\text{Im}(\lambda) = 0$,

$$\lambda_1 < -\frac{\lambda_{\min}(\mathbf{Q})\lambda_{\min}(\mathbf{P}\mathbf{P}^\top)}{\lambda_{\min}(\mathbf{Q})\lambda_{\max}(\mathbf{S}, \mathbf{Q}) + \lambda_{\min}(\mathbf{P}\mathbf{P}^\top)} < 0$$

where $\lambda_{\max}(\mathbf{S}, \mathbf{Q}) = \max(\lambda_{\max}(\mathbf{S}), \lambda_{\max}(\mathbf{Q}))$.

- When $\text{Im}(\lambda) \neq 0$,

$$\lambda_1 \leq -\frac{\lambda_{\min}(\mathbf{S}) + \lambda_{\min}(\mathbf{Q})}{2} < 0$$

Proof. We prove this theorem in a similar way to the proof for Theorem 5. Consider the following generic eigenvector equation:

$$\begin{bmatrix} -\mathbf{S} & \mathbf{P} \\ -\mathbf{P}^\top & -\mathbf{Q} \end{bmatrix} \begin{bmatrix} \mathbf{a}_1 + i\mathbf{a}_2 \\ \mathbf{b}_1 + i\mathbf{b}_2 \end{bmatrix} = (\lambda_1 + i\lambda_2) \begin{bmatrix} \mathbf{a}_1 + i\mathbf{a}_2 \\ \mathbf{b}_1 + i\mathbf{b}_2 \end{bmatrix}$$

where $\mathbf{a}_i, \mathbf{b}_i, \lambda_i$ are all real valued and the vector is normalized, i.e., $\|\mathbf{a}_1\|^2 + \|\mathbf{a}_2\|^2 + \|\mathbf{b}_1\|^2 + \|\mathbf{b}_2\|^2 = 1$. The above equation can be rewritten as:

$$\begin{bmatrix} -\mathbf{S}\mathbf{a}_1 + \mathbf{P}\mathbf{b}_1 + i(-\mathbf{S}\mathbf{a}_2 + \mathbf{P}\mathbf{b}_2) \\ -\mathbf{P}^\top\mathbf{a}_1 - \mathbf{Q}\mathbf{b}_1 + i(-\mathbf{P}^\top\mathbf{a}_2 - \mathbf{Q}\mathbf{b}_2) \end{bmatrix} = \begin{bmatrix} \lambda_1\mathbf{a}_1 - \lambda_2\mathbf{a}_2 + i(\lambda_1\mathbf{a}_2 + \lambda_2\mathbf{a}_1) \\ \lambda_1\mathbf{b}_1 - \lambda_2\mathbf{b}_2 + i(\lambda_1\mathbf{b}_2 + \lambda_2\mathbf{b}_1) \end{bmatrix}$$

By equating the real and complex parts, we get:

$$-\mathbf{S}\mathbf{a}_1 + \mathbf{P}\mathbf{b}_1 = \lambda_1\mathbf{a}_1 - \lambda_2\mathbf{a}_2 \tag{30}$$

$$-\mathbf{P}^\top\mathbf{a}_1 - \mathbf{Q}\mathbf{b}_1 = \lambda_1\mathbf{b}_1 - \lambda_2\mathbf{b}_2 \tag{31}$$

$$-\mathbf{S}\mathbf{a}_2 + \mathbf{P}\mathbf{b}_2 = \lambda_1\mathbf{a}_2 + \lambda_2\mathbf{a}_1 \tag{32}$$

$$-\mathbf{P}^\top\mathbf{a}_2 - \mathbf{Q}\mathbf{b}_2 = \lambda_1\mathbf{b}_2 + \lambda_2\mathbf{b}_1 \tag{33}$$

By multiplying the Equations (30), (31), (32), (33) by $\mathbf{a}_1^\top, \mathbf{b}_1^\top, \mathbf{a}_2^\top, \mathbf{b}_2^\top$ and adding them together, we get:

$$\begin{aligned} & -\mathbf{a}_1^\top\mathbf{S}\mathbf{a}_1 + \mathbf{a}_1^\top\mathbf{P}\mathbf{b}_1 - \mathbf{b}_1^\top\mathbf{P}^\top\mathbf{a}_1 - \mathbf{b}_1^\top\mathbf{Q}\mathbf{b}_1 - \mathbf{a}_2^\top\mathbf{S}\mathbf{a}_2 + \mathbf{a}_2^\top\mathbf{P}\mathbf{b}_2 - \mathbf{b}_2^\top\mathbf{P}^\top\mathbf{a}_2 - \mathbf{b}_2^\top\mathbf{Q}\mathbf{b}_2 \\ & = \lambda_1\mathbf{a}_1^\top\mathbf{a}_1 - \lambda_2\mathbf{a}_1^\top\mathbf{a}_2 + \lambda_1\mathbf{b}_1^\top\mathbf{b}_1 - \lambda_2\mathbf{b}_1^\top\mathbf{b}_2 + \lambda_1\mathbf{a}_2^\top\mathbf{a}_2 + \lambda_2\mathbf{a}_2^\top\mathbf{a}_1 + \lambda_1\mathbf{b}_2^\top\mathbf{b}_2 + \lambda_2\mathbf{b}_2^\top\mathbf{b}_1 \end{aligned}$$

which simplifies to

$$-\mathbf{a}_1^\top\mathbf{S}\mathbf{a}_1 - \mathbf{a}_2^\top\mathbf{S}\mathbf{a}_2 - \mathbf{b}_1^\top\mathbf{Q}\mathbf{b}_1 - \mathbf{b}_2^\top\mathbf{Q}\mathbf{b}_2 = \lambda_1(\mathbf{a}_1^\top\mathbf{a}_1 + \mathbf{b}_1^\top\mathbf{b}_1 + \mathbf{a}_2^\top\mathbf{a}_2 + \mathbf{b}_2^\top\mathbf{b}_2) = \lambda_1$$

Because $\mathbf{S} \succeq \mathbf{0}, \mathbf{Q} \succ \mathbf{0}$, $-\mathbf{a}_1^\top\mathbf{S}\mathbf{a}_1 - \mathbf{a}_2^\top\mathbf{S}\mathbf{a}_2 - \mathbf{b}_1^\top\mathbf{Q}\mathbf{b}_1 - \mathbf{b}_2^\top\mathbf{Q}\mathbf{b}_2 \leq 0$, where the equality holds only if $\mathbf{b}_1 = \mathbf{0}, \mathbf{b}_2 = \mathbf{0}$ (as well as $-\mathbf{a}_1^\top\mathbf{S}\mathbf{a}_1 - \mathbf{a}_2^\top\mathbf{S}\mathbf{a}_2 = 0$). Next we show that $\mathbf{b}_1 = \mathbf{0}, \mathbf{b}_2 = \mathbf{0}$ is contradictory with the condition that \mathbf{P}^\top is a full column rank matrix. First, when $-\mathbf{a}_1^\top\mathbf{S}\mathbf{a}_1 - \mathbf{a}_2^\top\mathbf{S}\mathbf{a}_2 - \mathbf{b}_1^\top\mathbf{Q}\mathbf{b}_1 - \mathbf{b}_2^\top\mathbf{Q}\mathbf{b}_2 = 0$, we have $\lambda_1 = 0$. Applying these to Equations (31) and (33), we get $\mathbf{P}^\top\mathbf{a}_1 = \mathbf{0}$ and $\mathbf{P}^\top\mathbf{a}_2 = \mathbf{0}$. Since one of \mathbf{a}_1 and \mathbf{a}_2 is non-zero (otherwise the eigenvector is zero), this implies \mathbf{P}^\top is not a full column rank matrix, which is contradictory with the condition that \mathbf{P}^\top has full column rank. Consequently, $-\mathbf{a}_1^\top\mathbf{S}\mathbf{a}_1 - \mathbf{a}_2^\top\mathbf{S}\mathbf{a}_2 - \mathbf{b}_1^\top\mathbf{Q}\mathbf{b}_1 - \mathbf{b}_2^\top\mathbf{Q}\mathbf{b}_2 = \lambda_1 < 0$.

Now we proceed to get a tighter upper bound on the real part of the eigenvalue λ_1 . By multiplying Equations (30) and (32) by $-\mathbf{a}_2^\top$ and \mathbf{a}_1^\top and adding them together, we get:

$$\mathbf{a}_2^\top\mathbf{S}\mathbf{a}_1 - \mathbf{a}_2^\top\mathbf{P}\mathbf{b}_1 - \mathbf{a}_1^\top\mathbf{S}\mathbf{a}_2 + \mathbf{a}_1^\top\mathbf{P}\mathbf{b}_2 = -\lambda_1\mathbf{a}_2^\top\mathbf{a}_1 + \lambda_2\mathbf{a}_2^\top\mathbf{a}_2 + \lambda_1\mathbf{a}_1^\top\mathbf{a}_2 + \lambda_2\mathbf{a}_1^\top\mathbf{a}_1$$

which simplifies to:

$$-a_2^\top P b_1 + a_1^\top P b_2 = \lambda_2 a_2^\top a_2 + \lambda_2 a_1^\top a_1 \quad (34)$$

Similarly, by multiplying Equations (31) and (33) by $-b_2^\top$ and b_1^\top and adding them together, we get:

$$b_2^\top P^\top a_1 + b_2^\top Q b_1 - b_1^\top P^\top a_2 - b_1^\top Q b_2 = -\lambda_1 b_2^\top b_1 + \lambda_2 b_2^\top b_2 + \lambda_1 b_1^\top b_2 + \lambda_2 b_1^\top b_1 \quad (35)$$

which simplifies to:

$$b_2^\top P^\top a_1 - b_1^\top P^\top a_2 = \lambda_2 b_2^\top b_2 + \lambda_2 b_1^\top b_1 \quad (36)$$

With Equation (34) and Equation (36), we have:

$$\lambda_2 (\|a_1\|^2 + \|a_2\|^2) = \lambda_2 (\|b_1\|^2 + \|b_2\|^2)$$

which implies that either $\|a_1\|^2 + \|a_2\|^2 = \|b_1\|^2 + \|b_2\|^2 = 1/2$ or $\lambda_2 = 0$.

In the first case ($\lambda_2 \neq 0$ and $\|a_1\|^2 + \|a_2\|^2 = \|b_1\|^2 + \|b_2\|^2 = 1/2$), from $-a_1^\top S a_1 - a_2^\top S a_2 - b_1^\top Q b_1 - b_2^\top Q b_2 = \lambda_1$, we get an upper bound:

$$\lambda_1 \leq -\frac{\lambda_{\min}(S) + \lambda_{\min}(Q)}{2}$$

This upper bound is strictly negative since $\lambda_{\min}(S) \geq 0$ and $\lambda_{\min}(Q) > 0$.

Now we introduce the following lemma which is useful for deriving the upper bound of λ_1 in the second case:

Lemma 4. *Let $S \succeq 0$ and $Q \succeq 0$ be two real symmetric matrices. If $a^\top S a + b^\top Q b = c$, then $a^\top S^\top S a + b^\top Q^\top Q b \in [c \cdot \min(\lambda_{\min}(S), \lambda_{\min}(Q)), c \cdot \max(\lambda_{\max}(S), \lambda_{\max}(Q))]$.*

Proof. Let $S = U_S \Lambda_S U_S^\top$ and $Q = U_Q \Lambda_Q U_Q^\top$ be the eigenvalue decompositions of S and Q . Then, we have

$$c = a^\top S a + b^\top Q b = a^\top U_S \Lambda_S U_S^\top a + b^\top U_Q \Lambda_Q U_Q^\top b = x^\top \Lambda_S x + y^\top \Lambda_Q y$$

where $x = U_S^\top a$ and $y = U_Q^\top b$. Therefore, we have:

$$c = \sum_i x_i^2 \lambda_S^i + \sum_j y_j^2 \lambda_Q^j$$

Similarly, we have:

$$\begin{aligned} a^\top S^\top S a + b^\top Q^\top Q b &= a^\top U_S \Lambda_S U_S^\top U_S \Lambda_S U_S^\top a + b^\top U_Q \Lambda_Q U_Q^\top U_Q \Lambda_Q U_Q^\top b \\ &= \sum_i x_i^2 (\lambda_S^i)^2 + \sum_j y_j^2 (\lambda_Q^j)^2 \end{aligned}$$

which differs from c by a multiplicative factor within $[\min(\lambda_{\min}(S), \lambda_{\min}(Q)), \max(\lambda_{\max}(S), \lambda_{\max}(Q))]$. □

In the second case, the imaginary part of the eigenvalue is zero ($\lambda_2 = 0$), which implies the imaginary part of the eigenvector must also be zero, $a_2 = b_2 = \mathbf{0}$. Applying this to Equations (30), (31), (32), (33), we get:

$$\begin{aligned} -S a_1 + P b_1 &= \lambda_1 a_1 \\ -P^\top a_1 - Q b_1 &= \lambda_1 b_1 \end{aligned}$$

Rearranging the above equations, we get:

$$\begin{aligned} P b_1 &= (\lambda_1 I + S) a_1 \\ -P^\top a_1 &= (\lambda_1 I + Q) b_1 \end{aligned}$$

Squaring both sides of the equations, we get:

$$\begin{aligned} \mathbf{b}_1^\top \mathbf{P}^\top \mathbf{P} \mathbf{b}_1 &= \mathbf{a}_1^\top (\lambda_1^2 \mathbf{I} + 2\lambda_1 \mathbf{S} + \mathbf{S}^\top \mathbf{S}) \mathbf{a}_1 = \lambda_1^2 \|\mathbf{a}_1\|^2 + 2\lambda_1 \mathbf{a}_1^\top \mathbf{S} \mathbf{a}_1 + \mathbf{a}_1^\top \mathbf{S}^\top \mathbf{S} \mathbf{a}_1 \\ \mathbf{a}_1^\top \mathbf{P} \mathbf{P}^\top \mathbf{a}_1 &= \mathbf{b}_1^\top (\lambda_1^2 \mathbf{I} + 2\lambda_1 \mathbf{Q} + \mathbf{Q}^\top \mathbf{Q}) \mathbf{b}_1 = \lambda_1^2 \|\mathbf{b}_1\|^2 + 2\lambda_1 \mathbf{b}_1^\top \mathbf{Q} \mathbf{b}_1 + \mathbf{b}_1^\top \mathbf{Q}^\top \mathbf{Q} \mathbf{b}_1 \end{aligned}$$

Summing these two equations together, using the fact that $-\mathbf{a}_1^\top \mathbf{S} \mathbf{a}_1 - \mathbf{b}_1^\top \mathbf{Q} \mathbf{b}_1 = \lambda_1$ and $\|\mathbf{a}_1\|^2 + \|\mathbf{b}_1\|^2 = 1$, we get:

$$\begin{aligned} \mathbf{b}_1^\top \mathbf{P}^\top \mathbf{P} \mathbf{b}_1 + \mathbf{a}_1^\top \mathbf{P} \mathbf{P}^\top \mathbf{a}_1 &= \lambda_1^2 (\|\mathbf{a}_1\|^2 + \|\mathbf{b}_1\|^2) - 2\lambda_1^2 + \mathbf{a}_1^\top \mathbf{S}^\top \mathbf{S} \mathbf{a}_1 + \mathbf{b}_1^\top \mathbf{Q}^\top \mathbf{Q} \mathbf{b}_1 \\ &= -\lambda_1^2 + \mathbf{a}_1^\top \mathbf{S}^\top \mathbf{S} \mathbf{a}_1 + \mathbf{b}_1^\top \mathbf{Q}^\top \mathbf{Q} \mathbf{b}_1 \end{aligned}$$

Let $\lambda_{\max}(\mathbf{S}, \mathbf{Q})$ denote $\max(\lambda_{\max}(\mathbf{S}), \lambda_{\max}(\mathbf{Q}))$. With Lemma 4, we have:

$$\mathbf{b}_1^\top \mathbf{P}^\top \mathbf{P} \mathbf{b}_1 + \mathbf{a}_1^\top \mathbf{P} \mathbf{P}^\top \mathbf{a}_1 \leq -\lambda_1^2 - \lambda_1 \lambda_{\max}(\mathbf{S}, \mathbf{Q}) \quad (37)$$

Furthermore, we have:

$$-\lambda_1 = \mathbf{a}_1^\top \mathbf{S} \mathbf{a}_1 + \mathbf{b}_1^\top \mathbf{Q} \mathbf{b}_1 \geq \lambda_{\min}(\mathbf{Q}) \|\mathbf{b}_1\|^2 = \lambda_{\min}(\mathbf{Q}) (1 - \|\mathbf{a}_1\|^2) \implies \|\mathbf{a}_1\|^2 \geq 1 + \frac{\lambda_1}{\lambda_{\min}(\mathbf{Q})}$$

Note that λ_1 can either satisfy $\lambda_1 \leq -\lambda_{\min}(\mathbf{Q})$ or $-\lambda_{\min}(\mathbf{Q}) < \lambda_1 < 0$. In the first scenario, we already obtain an upper bound: $\lambda_1 \leq -\lambda_{\min}(\mathbf{Q})$. So we focus on deriving an upper bound for the second scenario. From Equation (37), we have:

$$\begin{aligned} -\lambda_1^2 - \lambda_1 \lambda_{\max}(\mathbf{S}, \mathbf{Q}) &\geq \mathbf{a}_1^\top \mathbf{P} \mathbf{P}^\top \mathbf{a}_1 \geq \lambda_{\min}(\mathbf{P} \mathbf{P}^\top) \|\mathbf{a}_1\|^2 \geq \lambda_{\min}(\mathbf{P} \mathbf{P}^\top) \left(1 + \frac{\lambda_1}{\lambda_{\min}(\mathbf{Q})}\right) \\ \implies -\lambda_1 \left(\lambda_1 + \lambda_{\max}(\mathbf{S}, \mathbf{Q}) + \frac{\lambda_{\min}(\mathbf{P} \mathbf{P}^\top)}{\lambda_{\min}(\mathbf{Q})}\right) &\geq \lambda_{\min}(\mathbf{P} \mathbf{P}^\top) \\ \implies -\lambda_1 \left(\lambda_{\max}(\mathbf{S}, \mathbf{Q}) + \frac{\lambda_{\min}(\mathbf{P} \mathbf{P}^\top)}{\lambda_{\min}(\mathbf{Q})}\right) &> \lambda_{\min}(\mathbf{P} \mathbf{P}^\top) \end{aligned} \quad (38)$$

where the last implication uses the fact that $\lambda_1 < 0$. From Equation (38), we get an upper bound for λ_1 :

$$\lambda_1 < -\lambda_{\min}(\mathbf{Q}) \frac{\lambda_{\min}(\mathbf{P} \mathbf{P}^\top)}{\lambda_{\min}(\mathbf{Q}) \lambda_{\max}(\mathbf{S}, \mathbf{Q}) + \lambda_{\min}(\mathbf{P} \mathbf{P}^\top)}$$

Since the fraction in above equation lies in $(0, 1)$, we will use this as the upper bound for the second case ($\lambda_2 = 0$). Also note that this upper bound is strictly negative, since $\lambda_{\max}(\mathbf{S}, \mathbf{Q}) \geq \lambda_{\min}(\mathbf{Q}) > 0$ and $\lambda_{\min}(\mathbf{P} \mathbf{P}^\top) > 0$. □

D.2. Notations and Setup

First, we can reformulate the minimax game defined in Equation (14) as:

$$\min_{\boldsymbol{\theta}} \max_{\boldsymbol{\omega}} V(\boldsymbol{\theta}, \boldsymbol{\omega}) = \min_{\boldsymbol{\theta}} \max_{\boldsymbol{\omega}} \mathbb{E}_{p(\mathbf{x})} [A(H_{\boldsymbol{\omega}}(\mathbf{x}) + E_{\boldsymbol{\theta}}(\mathbf{x}))] - \mathbb{E}_{q_{\boldsymbol{\theta}}(\mathbf{x})} [B(H_{\boldsymbol{\omega}}(\mathbf{x}) + E_{\boldsymbol{\theta}}(\mathbf{x}))] \quad (39)$$

where the functions $A(u)$ and $B(u)$ are defined as:

$$A(u) = f'(\exp(u)), \quad B(u) = f^*(f'(\exp(u))) - f^*(f'(1)) \quad (40)$$

with $B(0) = 0$. Since $f^*(f'(1))$ is a constant, and $\mathbb{E}_{q_{\boldsymbol{\theta}}(\mathbf{x})} [B(H_{\boldsymbol{\omega}}(\mathbf{x}) + E_{\boldsymbol{\theta}}(\mathbf{x}))] = \mathbb{E}_{q_{\boldsymbol{\theta}}(\mathbf{x})} [f^*(f'(\exp(H_{\boldsymbol{\omega}}(\mathbf{x}) + E_{\boldsymbol{\theta}}(\mathbf{x}))))] - f^*(f'(1))$, the above formulation is equivalent to the original *f*-EBM minimax game in Equation (14) (up to a constant that does not depend on $\boldsymbol{\theta}$ and $\boldsymbol{\omega}$).

We have the following theorem to characterize the properties of functions $A(u)$ and $B(u)$:

Theorem 7. *For any *f*-divergence with closed, strictly convex and third-order differentiable generator function *f*, functions $A(u)$, $B(u)$ defined as:*

$$A(u) = f'(\exp(u)), \quad B(u) = f^*(f'(\exp(u))) - f^*(f'(1))$$

satisfy the following properties:

- $A'(0) = B'(0) = f''(1) > 0$
- $A''(0) - B''(0) = -f'''(1) < 0$
- $A''(0) - B''(0) + 2B'(0) = f''(1) = B'(0) > 0$

where A', B' are first-order derivative of A, B ; A'', B'' are second-order derivative of A, B .

Proof. First, we introduce the following lemma for convex conjugates and subgradients:

Lemma 5. *If f is closed and convex, then we have:*

$$y \in \partial f(x) \iff x \in \partial f^*(y)$$

Proof. If $y \in \partial f(x)$, then $f^*(y) = \sup_u (yu - f(u)) = yx - f(x)$. Therefore, we have:

$$\begin{aligned} f^*(v) &= \sup_u vu - f(u) \\ &\geq vx - f(x) \\ &= x(v - y) - f(x) + xy \\ &= f^*(y) + x(v - y) \end{aligned}$$

Because this holds for all v , we have $x \in \partial f^*(y)$. The reverse implication $x \in \partial f^*(y) \implies y \in \partial f(x)$ follows from $f^{**} = f$. \square

Let us use $g(y)$ to denote the conjugate function, $g(y) = f^*(y)$. Next we have the following lemma for the gradient of conjugate:

Lemma 6. *If f is strictly convex and differentiable, then we have: $g'(y) = \arg \max_x (yx - f(x))$, $g'(f'(u)) = u$.*

Proof. Since f is strictly convex, x maximizes $yx - f(x)$ if and only if $y \in \partial f(x)$. With Lemma 5, we know that

$$y \in \partial f(x) \iff x \in \partial f^*(y) = \{g'(y)\}$$

Therefore $g'(y) = \arg \max_x (yx - f(x))$. Then we have:

$$g'(f'(u)) = \arg \max_x f'(u)x - f(x) := \arg \max_x h(x)$$

Since $h'(x) = f'(u) - f'(x)$, $h''(x) = -f''(x) < 0$, we have $\arg \max_x f'(u)x - f(x) = u$. \square

Now we are ready to proof the properties of functions $A(u)$ and $B(u)$. Recall that $A(u) = f'(\exp(u))$, $B(u) = f^*(f'(\exp(u))) - f^*(f'(1))$, with Lemma 6, we have:

$$\begin{aligned} A'(u) &= f''(\exp(u)) \exp(u) \\ A''(u) &= f'''(\exp(u)) \exp(2u) + f''(\exp(u)) \exp(u) \\ B'(u) &= g'(f'(\exp(u))) f''(\exp(u)) \exp(u) \\ &= f''(\exp(u)) \exp(2u) \\ B''(u) &= f'''(\exp(u)) \exp(3u) + 2f''(\exp(u)) \exp(2u) \end{aligned}$$

Then we have:

$$\begin{aligned} A'(0) &= f''(1) \\ A''(0) &= f'''(1) + f''(1) \\ B'(0) &= f''(1) \\ B''(0) &= f'''(1) + 2f''(1) \end{aligned}$$

Since f is strictly convex with $\forall x \in \text{dom}(f), f''(x) > 0$, we have:

$$\begin{aligned} A'(0) &= B'(0) = f''(1) > 0 \\ A''(0) - B''(0) &= -f''(1) < 0 \\ A''(0) - B''(0) + 2B'(0) &= f''(1) = B'(0) > 0 \end{aligned}$$

□

Some examples of Theorem 7 can be found in Table 3.

Name	$A'(0)$	$B'(0)$	$A''(0)$	$B''(0)$
Kullback-Leibler	1	1	0	1
Reverse Kullback-Leibler	1	1	-1	0
Pearson χ^2	2	2	2	4
Neyman χ^2	2	2	-4	-2
Squared Hellinger	$\frac{1}{2}$	$\frac{1}{2}$	$-\frac{1}{4}$	$\frac{1}{4}$
Jensen-Shannon	$\frac{1}{2}$	$\frac{1}{2}$	$-\frac{1}{4}$	$\frac{1}{4}$
α -divergence ($\alpha \notin \{0, 1\}$)	1	1	$\alpha - 1$	α

Table 3. Some examples of f -divergences and the corresponding first- and second-order gradient values of functions $A(u)$ and $B(u)$ (defined in Equation (40)) at the equilibrium point.

Let us use (θ^*, ω^*) to denote the equilibrium point and we have the following realizability assumption:

Assumption 4 (Realizability). $\exists \theta^*, \omega^*$, such that $\forall \mathbf{x} \in \text{supp}(p), q_{\theta^*}(\mathbf{x}) = p(\mathbf{x})$ and $H_{\omega^*}(\mathbf{x}) = \log(p(\mathbf{x})Z_{\theta^*})$.

That is we assume the energy function and the variational functions are powerful enough such that at the equilibrium point, the model distribution q_{θ^*} matches the true data distribution p and the variational function achieves the optimal form in Theorem 2. Note that this also implies:

$$\begin{aligned} \forall \mathbf{x} \in \text{supp}(p), H_{\omega^*}(\mathbf{x}) + E_{\theta^*}(\mathbf{x}) &= \log(p(\mathbf{x})Z_{\theta^*}) - \log(\exp(-E_{\theta^*}(\mathbf{x}))) \\ &= \log(p(\mathbf{x})/q_{\theta^*}(\mathbf{x})) = 0 \end{aligned}$$

Furthermore, we have the following assumption on the energy function:

Assumption 5. $(\mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})] - 2\mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x})] \mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})])|_{\theta^*}$ is positive semi-definite.

In one-dimension case, this assumption is saying: $\mathbb{E}_{p(x)}[(\nabla_{\theta} E_{\theta}(x))^2] - 2(\mathbb{E}_{p(x)}[\nabla_{\theta} E_{\theta}(x)])^2 \geq 0$. Intuitively, this implies the sum of energies over the entire sample space changes smoothly such that the average change of energies is small compared to the sum of the squared change of energies. For a better understanding of this assumption and its rationality, in the following, we will introduce two simple examples that satisfy this assumption.

Example 1. Consider the following quadratic energy function:

$$E_{\theta}(x) = \frac{(x - \mu)^2}{2\sigma^2}, \quad \theta = [\mu, \sigma]^{\top}, \quad p = q_{\theta} = \frac{\exp(-\frac{(x-\mu)^2}{2\sigma^2})}{\sqrt{2\pi}\sigma}$$

which corresponds to a Gaussian distribution.

The first-order gradient is:

$$\frac{\partial E_{\theta}(x)}{\partial \mu} = \frac{x - \mu}{\sigma^2}, \quad \frac{\partial E_{\theta}(x)}{\partial \sigma} = -\frac{(x - \mu)^2}{\sigma^3}$$

Then, we have:

$$\begin{aligned}
 \int_{-\infty}^{+\infty} p(x) \frac{\partial E_{\theta}(x)}{\partial \mu} dx &= \int_{-\infty}^{+\infty} \frac{\exp(-\frac{(x-\mu)^2}{2\sigma^2})}{\sqrt{2\pi}\sigma} \frac{x-\mu}{\sigma^2} dx = 0 \\
 \int_{-\infty}^{+\infty} p(x) \frac{\partial E_{\theta}(x)}{\partial \sigma} dx &= - \int_{-\infty}^{+\infty} \frac{\exp(-\frac{(x-\mu)^2}{2\sigma^2})}{\sqrt{2\pi}\sigma} \frac{(x-\mu)^2}{\sigma^3} dx \\
 &= \frac{1}{\sqrt{2\pi}\sigma^4} \int_{-\infty}^{+\infty} \exp(-\frac{x^2}{2\sigma^2}) x^2 dx = \frac{1}{\sigma} \\
 \int_{-\infty}^{+\infty} p(x) \nabla_{\theta} E_{\theta}(x) dx &= \int_{-\infty}^{+\infty} p(x) \begin{bmatrix} \frac{\partial E_{\theta}(x)}{\partial \mu} \\ \frac{\partial E_{\theta}(x)}{\partial \sigma} \end{bmatrix} dx = \begin{bmatrix} 0 \\ \frac{1}{\sigma} \end{bmatrix} \\
 \int_{-\infty}^{+\infty} p(x) \left(\frac{\partial E_{\theta}(x)}{\partial \mu} \right)^2 dx &= \int_{-\infty}^{+\infty} \frac{\exp(-\frac{(x-\mu)^2}{2\sigma^2})}{\sqrt{2\pi}\sigma} \frac{(x-\mu)^2}{\sigma^4} dx \\
 &= \frac{1}{\sqrt{2\pi}\sigma^5} \int_{-\infty}^{+\infty} \exp(-\frac{x^2}{2\sigma^2}) x^2 dx = \frac{1}{\sigma^2} \\
 \int_{-\infty}^{+\infty} p(x) \frac{\partial E_{\theta}(x)}{\partial \mu} \frac{\partial E_{\theta}(x)}{\partial \sigma} dx &= - \int_{-\infty}^{+\infty} \frac{\exp(-\frac{(x-\mu)^2}{2\sigma^2})}{\sqrt{2\pi}\sigma} \frac{(x-\mu)^3}{\sigma^5} dx = 0 \\
 \int_{-\infty}^{+\infty} p(x) \left(\frac{\partial E_{\theta}(x)}{\partial \sigma} \right)^2 dx &= \int_{-\infty}^{+\infty} \frac{\exp(-\frac{(x-\mu)^2}{2\sigma^2})}{\sqrt{2\pi}\sigma} \frac{(x-\mu)^4}{\sigma^6} dx = 0 \\
 &= \frac{1}{\sqrt{2\pi}\sigma^7} \int_{-\infty}^{+\infty} \exp(-\frac{x^2}{2\sigma^2}) x^4 dx = \frac{3}{\sigma^2} \\
 \int_{-\infty}^{+\infty} p(x) \nabla_{\theta} E_{\theta}(x) \nabla_{\theta}^{\top} E_{\theta}(x) dx &= \int_{-\infty}^{+\infty} p(x) \begin{bmatrix} \left(\frac{\partial E_{\theta}(x)}{\partial \mu} \right)^2 & \frac{\partial E_{\theta}(x)}{\partial \mu} \frac{\partial E_{\theta}(x)}{\partial \sigma} \\ \frac{\partial E_{\theta}(x)}{\partial \mu} \frac{\partial E_{\theta}(x)}{\partial \sigma} & \left(\frac{\partial E_{\theta}(x)}{\partial \sigma} \right)^2 \end{bmatrix} dx = \begin{bmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{3}{\sigma^2} \end{bmatrix}
 \end{aligned}$$

Therefore, Assumption 5 holds:

$$\mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})] - 2\mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x})] \mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})] = \begin{bmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{\sigma^2} \end{bmatrix} \succ \mathbf{0}$$

Example 2. We consider a more powerful energy function:

$$E_{\theta}(\mathbf{x}) = -\log \left(\sum_{k=1}^K \pi_k \frac{\exp(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\top} \Sigma_k^{-1}(\mathbf{x}-\boldsymbol{\mu}_k))}{\sqrt{(2\pi)^n |\Sigma_k|}} \right), \quad \sum_{k=1}^K \pi_k = 1, \quad \mathbf{x} \in \mathbb{R}^n, \quad \boldsymbol{\mu}_k \in \mathbb{R}^n, \quad \Sigma_k \in S_{++}^n$$

where $\pi_1, \dots, \pi_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \Sigma_1, \dots, \Sigma_K$ are the learnable parameters θ ; S_{++}^n denotes the space of symmetric positive definite $n \times n$ matrices.

The partition function is:

$$Z_{\theta} = \int \exp(-E_{\theta}(\mathbf{x})) d\mathbf{x} = \int \sum_{k=1}^K \pi_k \frac{\exp(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\top} \Sigma_k^{-1}(\mathbf{x}-\boldsymbol{\mu}_k))}{\sqrt{(2\pi)^n |\Sigma_k|}} d\mathbf{x} = \sum_{k=1}^K \pi_k = 1$$

Therefore, Assumption 5 holds:

$$\begin{aligned}
 &\mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})] - 2\mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x})] \mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})] \\
 &= \mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})] - 2\nabla_{\theta} \log Z_{\theta} \nabla_{\theta}^{\top} \log Z_{\theta} \\
 &= \mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})]
 \end{aligned}$$

which is a positive semi-definite moment matrix.

Note that the model distribution induced by the energy function correspond to a Gaussian mixture model (GMM). With the universal approximation theorem of GMM (i.e., any smooth density can be approximated with any specific nonzero amount of error by a GMM with enough components (Goodfellow et al., 2016)), the above EBM can be used to fit any probability distribution when K is large enough.

Finally, the last assumption we need is:

Assumption 6. $(\mathbb{E}_{p(\mathbf{x})}[\nabla_{\omega} H_{\omega}(\mathbf{x}) \nabla_{\omega}^{\top} H_{\omega}(\mathbf{x})])|_{\omega^*}$ and $(\mathbb{E}_{p(\mathbf{x})}[\nabla_{\omega} H_{\omega}(\mathbf{x})] \mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})])|_{(\theta^*, \omega^*)}$ are full column rank.

This is similar to the assumption used in the local convergence analysis of GANs (Nagarajan & Kolter, 2017; Mescheder et al., 2018). Alternatively, we can replace it with another assumption that the rank deficiencies of the above matrices, if any, correspond to equivalent equilibria. In the following, we will use Assumption 6.

D.3. Local Convergence of Single-Step *f*-EBM

In the following, we will establish a theoretical proof of the local convergence property of Single-Step *f*-EBM. More specifically, instead of assuming the variational function H_{ω} is optimal at every update of the energy function, we consider a more realistic setting, where both the variational function H_{ω} and energy function E_{θ} take simultaneous gradient steps, with time derivatives defined as:

$$\begin{pmatrix} \dot{\theta} \\ \dot{\omega} \end{pmatrix} = \begin{pmatrix} -\nabla_{\theta} V(\theta, \omega) \\ \nabla_{\omega} V(\theta, \omega) \end{pmatrix} \quad (41)$$

That is we use gradient descent to update θ and gradient ascent to update ω with respect to $V(\theta, \omega)$ at the same frequency. In practice we can also update θ and ω alternatively as presented in Algorithm 1. Note that our theoretical analysis also holds for the alternative gradient methods, as the ordinary differential equations of both simultaneous gradient methods and alternative gradient methods have the same Jacobians at the equilibrium point.

Throughout this paper we will use the notation $\nabla^{\top}(\cdot)$ to denote the row vector corresponding to the gradient that is being compute. First, we derive the Jacobian at equilibrium.

Theorem 8. For the dynamical system defined in Equation (39) and the updates defined in Equation (41), under Assumption 4, the Jacobian at an equilibrium point (θ^*, ω^*) is:

$$\mathbf{J} = \begin{bmatrix} -\nabla_{\theta}^2 V & -\nabla_{\omega} \nabla_{\theta} V \\ \nabla_{\theta} \nabla_{\omega} V & \nabla_{\omega}^2 V \end{bmatrix} = \begin{bmatrix} -f''(1) \mathbf{K}_{EE} & f''(1) \mathbf{K}_{EH} \\ -f''(1) \mathbf{K}_{EH}^{\top} & -f''(1) \mathbf{K}_{HH} \end{bmatrix}$$

where

$$\begin{aligned} \mathbf{K}_{EE} &:= (\mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})] - 2\mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x})] \mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})])|_{\theta^*} \\ \mathbf{K}_{EH} &:= (\mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x})] \mathbb{E}_{p(\mathbf{x})}[\nabla_{\omega}^{\top} H_{\omega}(\mathbf{x})])|_{(\theta^*, \omega^*)} \\ \mathbf{K}_{HH} &:= (\mathbb{E}_{p(\mathbf{x})}[\nabla_{\omega} H_{\omega}(\mathbf{x}) \nabla_{\omega}^{\top} H_{\omega}(\mathbf{x})])|_{\omega^*} \end{aligned}$$

Proof. First, let us derive the first- and second-order derivatives of $V(\theta, \omega)$ with respect to θ :

$$\begin{aligned} \nabla_{\theta} V(\theta, \omega) &= \int p(\mathbf{x}) A'(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x})) \nabla_{\theta} E_{\theta}(\mathbf{x}) d\mathbf{x} - \\ &\quad \int B(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x})) \nabla_{\theta} q_{\theta}(\mathbf{x}) d\mathbf{x} - \\ &\quad \int q_{\theta}(\mathbf{x}) B'(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x})) \nabla_{\theta} E_{\theta}(\mathbf{x}) d\mathbf{x} \end{aligned}$$

$$\begin{aligned}
 \nabla_{\theta}^2 V(\theta, \omega) = & \int p(\mathbf{x}) A''(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x})) \nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x}) d\mathbf{x} + \\
 & \int p(\mathbf{x}) A'(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x})) \nabla_{\theta}^2 E_{\theta}(\mathbf{x}) d\mathbf{x} - \\
 & \int B'(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x})) \nabla_{\theta} q_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x}) d\mathbf{x} - \\
 & \int B(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x})) \nabla_{\theta}^2 q_{\theta}(\mathbf{x}) d\mathbf{x} - \\
 & \int B'(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x})) \nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} q_{\theta}(\mathbf{x}) d\mathbf{x} - \\
 & \int q_{\theta}(\mathbf{x}) B''(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x})) \nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x}) d\mathbf{x} - \\
 & \int q_{\theta}(\mathbf{x}) B'(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x})) \nabla_{\theta}^2 E_{\theta}(\mathbf{x}) d\mathbf{x}
 \end{aligned}$$

At the equilibrium point (θ^*, ω^*) , we have:

$$\begin{aligned}
 \nabla_{\theta}^2 V(\theta, \omega)|_{(\theta^*, \omega^*)} = & \left(\int p(\mathbf{x}) A''(0) \nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x}) d\mathbf{x} + \right. \\
 & \int p(\mathbf{x}) A'(0) \nabla_{\theta}^2 E_{\theta}(\mathbf{x}) d\mathbf{x} - \\
 & \int B'(0) \nabla_{\theta} q_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x}) d\mathbf{x} - \\
 & \int B(0) \nabla_{\theta}^2 q_{\theta}(\mathbf{x}) d\mathbf{x} - \\
 & \int B'(0) \nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} q_{\theta}(\mathbf{x}) d\mathbf{x} - \\
 & \left. \int q_{\theta}(\mathbf{x}) B''(0) \nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x}) d\mathbf{x} - \right. \\
 & \left. \int q_{\theta}(\mathbf{x}) B'(0) \nabla_{\theta}^2 E_{\theta}(\mathbf{x}) d\mathbf{x} \right) \Big|_{\theta^*}
 \end{aligned}$$

With Theorem 7 and Assumption 4, as well as by definition of function $B(u)$ ($B(0) = 0$) we have:

$$\begin{aligned}
 \nabla_{\theta}^2 V(\theta, \omega)|_{(\theta^*, \omega^*)} = & \left(\int p(\mathbf{x}) A''(0) \nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x}) d\mathbf{x} - \right. \\
 & \int B'(0) \nabla_{\theta} q_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x}) d\mathbf{x} - \\
 & \int B'(0) \nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} q_{\theta}(\mathbf{x}) d\mathbf{x} - \\
 & \left. \int q_{\theta}(\mathbf{x}) B''(0) \nabla_{\theta} E_{\theta}(\mathbf{x}) \nabla_{\theta}^{\top} E_{\theta}(\mathbf{x}) d\mathbf{x} \right) \Big|_{\theta^*} \tag{42}
 \end{aligned}$$

Next, for a θ -parametrized energy based model $q_{\theta}(\mathbf{x}) = \frac{\exp(-E_{\theta}(\mathbf{x}))}{Z_{\theta}}$, we observe that

$$\nabla_{\theta} q_{\theta}(\mathbf{x}) = q_{\theta}(\mathbf{x}) \nabla_{\theta} \log q_{\theta}(\mathbf{x}) = q_{\theta}(\mathbf{x}) (-\nabla_{\theta} E_{\theta}(\mathbf{x}) - \nabla_{\theta} \log Z_{\theta}) \tag{43}$$

$$\nabla_{\theta} \log Z_{\theta} = \frac{\int \exp(-E_{\theta}(\mathbf{x})) (-\nabla_{\theta} E_{\theta}(\mathbf{x})) d\mathbf{x}}{Z_{\theta}} = - \int q_{\theta}(\mathbf{x}) \nabla_{\theta} E_{\theta}(\mathbf{x}) d\mathbf{x} \tag{44}$$

With Assumption 4 ($q_{\theta^*} = p$) and above observation, Equation (42) can be written as:

$$\begin{aligned}
 \nabla_{\theta}^2 V(\boldsymbol{\theta}, \boldsymbol{\omega})|_{(\theta^*, \omega^*)} &= \left(\int p(\mathbf{x})(A''(0) - B''(0))\nabla_{\theta} E_{\theta}(\mathbf{x})\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})d\mathbf{x} - \right. \\
 &\quad \left. \int B'(0)\nabla_{\theta} q_{\theta}(\mathbf{x})\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})d\mathbf{x} - \right. \\
 &\quad \left. \int B'(0)\nabla_{\theta} E_{\theta}(\mathbf{x})\nabla_{\theta}^{\top} q_{\theta}(\mathbf{x})d\mathbf{x} \right) \Big|_{\theta^*} \\
 &= \left((A''(0) - B''(0)) \int p(\mathbf{x})\nabla_{\theta} E_{\theta}(\mathbf{x})\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})d\mathbf{x} - \right. \\
 &\quad \left. \int B'(0)q_{\theta}(\mathbf{x})(-\nabla_{\theta} E_{\theta}(\mathbf{x}) - \nabla_{\theta} \log Z_{\theta})\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})d\mathbf{x} - \right. \\
 &\quad \left. \int B'(0)q_{\theta}(\mathbf{x})\nabla_{\theta} E_{\theta}(\mathbf{x})(-\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x}) - \nabla_{\theta}^{\top} \log Z_{\theta})d\mathbf{x} \right) \Big|_{\theta^*} \\
 &= \left((A''(0) - B''(0) + 2B'(0)) \int p(\mathbf{x})\nabla_{\theta} E_{\theta}(\mathbf{x})\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})d\mathbf{x} + \right. \\
 &\quad \left. B'(0)\nabla_{\theta} \log Z_{\theta} \cdot \int q_{\theta}(\mathbf{x})\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})d\mathbf{x} + \right. \\
 &\quad \left. B'(0) \int q_{\theta}(\mathbf{x})\nabla_{\theta} E_{\theta}(\mathbf{x})d\mathbf{x} \cdot \nabla_{\theta}^{\top} \log Z_{\theta} \right) \Big|_{\theta^*} \\
 &= \left((A''(0) - B''(0) + 2B'(0)) \int p(\mathbf{x})\nabla_{\theta} E_{\theta}(\mathbf{x})\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})d\mathbf{x} - \right. \\
 &\quad \left. 2B'(0)\nabla_{\theta} \log Z_{\theta} \nabla_{\theta}^{\top} \log Z_{\theta} \right) \Big|_{\theta^*}
 \end{aligned}$$

With Theorem 7, we have:

$$\begin{aligned}
 \nabla_{\theta}^2 V(\boldsymbol{\theta}, \boldsymbol{\omega})|_{(\theta^*, \omega^*)} &= f''(1) \left(\int p(\mathbf{x})\nabla_{\theta} E_{\theta}(\mathbf{x})\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})d\mathbf{x} - 2\nabla_{\theta} \log Z_{\theta} \nabla_{\theta}^{\top} \log Z_{\theta} \right) \Big|_{\theta^*} \\
 &= f''(1) \left(\int p(\mathbf{x})\nabla_{\theta} E_{\theta}(\mathbf{x})\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})d\mathbf{x} - 2 \int p(\mathbf{x})\nabla_{\theta} E_{\theta}(\mathbf{x})d\mathbf{x} \int p(\mathbf{x})\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})d\mathbf{x} \right) \Big|_{\theta^*} \\
 &= f''(1) (\mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x})\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})] - 2\mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta} E_{\theta}(\mathbf{x})]\mathbb{E}_{p(\mathbf{x})}[\nabla_{\theta}^{\top} E_{\theta}(\mathbf{x})]) \Big|_{\theta^*}
 \end{aligned}$$

Now let us derive the first- and second-order derivatives of $V(\boldsymbol{\theta}, \boldsymbol{\omega})$ with respect to $\boldsymbol{\omega}$:

$$\begin{aligned}
 \nabla_{\omega} V(\boldsymbol{\theta}, \boldsymbol{\omega}) &= \int p(\mathbf{x})A'(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x}))\nabla_{\omega} H_{\omega}(\mathbf{x})d\mathbf{x} - \int q_{\theta}(\mathbf{x})B'(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x}))\nabla_{\omega} H_{\omega}(\mathbf{x})d\mathbf{x} \\
 \nabla_{\omega}^2 V(\boldsymbol{\theta}, \boldsymbol{\omega}) &= \int p(\mathbf{x})A''(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x}))\nabla_{\omega} H_{\omega}(\mathbf{x})\nabla_{\omega}^{\top} H_{\omega}(\mathbf{x})d\mathbf{x} + \\
 &\quad \int p(\mathbf{x})A'(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x}))\nabla_{\omega}^2 H_{\omega}(\mathbf{x})d\mathbf{x} - \\
 &\quad \int q_{\theta}(\mathbf{x})B''(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x}))\nabla_{\omega} H_{\omega}(\mathbf{x})\nabla_{\omega}^{\top} H_{\omega}(\mathbf{x})d\mathbf{x} - \\
 &\quad \int q_{\theta}(\mathbf{x})B'(H_{\omega}(\mathbf{x}) + E_{\theta}(\mathbf{x}))\nabla_{\omega}^2 H_{\omega}(\mathbf{x})d\mathbf{x}
 \end{aligned}$$

At the equilibrium point $(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)$, under Assumption 4, we have:

$$\nabla_{\boldsymbol{\omega}}^2 V(\boldsymbol{\theta}, \boldsymbol{\omega})|_{(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)} = \left((A''(0) - B''(0)) \int p(\mathbf{x}) \nabla_{\boldsymbol{\omega}} H_{\boldsymbol{\omega}}(\mathbf{x}) \nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x}) d\mathbf{x} + (A'(0) - B'(0)) \int p(\mathbf{x}) \nabla_{\boldsymbol{\omega}}^2 H_{\boldsymbol{\omega}}(\mathbf{x}) d\mathbf{x} \right) \Big|_{\boldsymbol{\omega}^*}$$

With Theorem 7, we have:

$$\begin{aligned} \nabla_{\boldsymbol{\omega}}^2 V(\boldsymbol{\theta}, \boldsymbol{\omega})|_{(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)} &= -f''(1) \left(\int p(\mathbf{x}) \nabla_{\boldsymbol{\omega}} H_{\boldsymbol{\omega}}(\mathbf{x}) \nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x}) d\mathbf{x} \right) \Big|_{\boldsymbol{\omega}^*} \\ &= -f''(1) (\mathbb{E}_{p(\mathbf{x})} [\nabla_{\boldsymbol{\omega}} H_{\boldsymbol{\omega}}(\mathbf{x}) \nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x})]) \Big|_{\boldsymbol{\omega}^*} \end{aligned}$$

Finally, let us derive $\nabla_{\boldsymbol{\omega}} \nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta}, \boldsymbol{\omega})$:

$$\begin{aligned} \nabla_{\boldsymbol{\omega}} \nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta}, \boldsymbol{\omega}) &= \int p(\mathbf{x}) A''(H_{\boldsymbol{\omega}}(\mathbf{x}) + E_{\boldsymbol{\theta}}(\mathbf{x})) \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) \nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x}) d\mathbf{x} - \\ &\quad \int B'(H_{\boldsymbol{\omega}}(\mathbf{x}) + E_{\boldsymbol{\theta}}(\mathbf{x})) \nabla_{\boldsymbol{\theta}} q_{\boldsymbol{\theta}}(\mathbf{x}) \nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x}) d\mathbf{x} - \\ &\quad \int q_{\boldsymbol{\theta}}(\mathbf{x}) B''(H_{\boldsymbol{\omega}}(\mathbf{x}) + E_{\boldsymbol{\theta}}(\mathbf{x})) \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) \nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x}) d\mathbf{x} \end{aligned}$$

At the equilibrium point $(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)$, under Assumption 4, we have:

$$\begin{aligned} \nabla_{\boldsymbol{\omega}} \nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta}, \boldsymbol{\omega})|_{(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)} &= \left((A''(0) - B''(0)) \int p(\mathbf{x}) \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) \nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x}) d\mathbf{x} - \right. \\ &\quad \left. B'(0) \int q_{\boldsymbol{\theta}}(\mathbf{x}) (-\nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) - \nabla_{\boldsymbol{\theta}} \log Z_{\boldsymbol{\theta}}) \nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x}) d\mathbf{x} \right) \Big|_{(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)} \\ &= \left((A''(0) - B''(0) + B'(0)) \int p(\mathbf{x}) \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) \nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x}) d\mathbf{x} + \right. \\ &\quad \left. B'(0) \nabla_{\boldsymbol{\theta}} \log Z_{\boldsymbol{\theta}} \int p(\mathbf{x}) \nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x}) d\mathbf{x} \right) \Big|_{(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)} \\ &= \left((A''(0) - B''(0) + B'(0)) \int p(\mathbf{x}) \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) \nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x}) d\mathbf{x} - \right. \\ &\quad \left. B'(0) \int p(\mathbf{x}) \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) d\mathbf{x} \int p(\mathbf{x}) \nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x}) d\mathbf{x} \right) \Big|_{(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)} \end{aligned}$$

With Theorem 7, we have:

$$\nabla_{\boldsymbol{\omega}} \nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta}, \boldsymbol{\omega})|_{(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)} = -f''(1) (\mathbb{E}_{p(\mathbf{x})} [\nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x})] \mathbb{E}_{p(\mathbf{x})} [\nabla_{\boldsymbol{\omega}}^{\top} H_{\boldsymbol{\omega}}(\mathbf{x})]) \Big|_{(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)}$$

Similarly, for $\nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\omega}} V(\boldsymbol{\theta}, \boldsymbol{\omega})$, we have:

$$\nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\omega}} V(\boldsymbol{\theta}, \boldsymbol{\omega})|_{(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)} = -f''(1) (\mathbb{E}_{p(\mathbf{x})} [\nabla_{\boldsymbol{\omega}} H_{\boldsymbol{\omega}}(\mathbf{x})] \mathbb{E}_{p(\mathbf{x})} [\nabla_{\boldsymbol{\theta}}^{\top} E_{\boldsymbol{\theta}}(\mathbf{x})]) \Big|_{(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)}$$

□

Now we are ready to present the main theorem:

Theorem 9. *The dynamical system defined in Equation (39) and the updates defined in Equation (41) is locally exponentially stable with respect to an equilibrium point $(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)$ when the Assumptions 4, 5, 6 hold for $(\boldsymbol{\theta}^*, \boldsymbol{\omega}^*)$. Let $\lambda_{\max}(\cdot)$ and $\lambda_{\min}(\cdot)$ denote the largest and smallest eigenvalues of a non-zero positive semi-definite matrix. The rate of convergence is governed only by the eigenvalues λ of the Jacobian \mathbf{J} of the system at the equilibrium point, with a strictly negative real part upper bounded as:*

- When $\text{Im}(\lambda) = 0$,

$$\text{Re}(\lambda) < -f''(1) \frac{\lambda_{\min}(\mathbf{K}_{HH})\lambda_{\min}(\mathbf{K}_{EH}\mathbf{K}_{EH}^\top)}{\lambda_{\min}(\mathbf{K}_{HH})\lambda_{\max}(\mathbf{K}_{EE}, \mathbf{K}_{HH}) + \lambda_{\min}(\mathbf{K}_{EH}\mathbf{K}_{EH}^\top)} < 0$$

where $\lambda_{\max}(\mathbf{K}_{EE}, \mathbf{K}_{HH}) = \max(\lambda_{\max}(\mathbf{K}_{EE}), \lambda_{\max}(\mathbf{K}_{HH}))$.

- When $\text{Im}(\lambda) \neq 0$,

$$\text{Re}(\lambda) \leq -\frac{f''(1)}{2}(\lambda_{\min}(\mathbf{K}_{EE}) + \lambda_{\min}(\mathbf{K}_{HH})) < 0$$

Proof. In Theorem 8, for any *f*-divergences with strictly convex and differentiable generator function *f*, we derived the Jacobian of the system **J** under Assumption 4. With Assumptions 5 and 6, and the strict convexity of the function *f*, we know that $f''(1)\mathbf{K}_{EE}$ is positive semi-definite, $f''(1)\mathbf{K}_{HH}$ is positive definite (a full rank moment matrix with a positive multiplicative factor), and $f''(1)\mathbf{K}_{EH}^\top$ is full column rank. Therefore, with Theorem 6, we know that the Jacobian **J** is a Hurwitz matrix, (*i.e.*, all the eigenvalues of **J** have strictly negative real parts). Furthermore, with Theorem 6, we can obtain an upper bound of the real parts of the eigenvalues. Finally, with Theorem 4, we can conclude that the system is locally exponentially stable. \square

E. Implementation Details

E.1. Implementation of *f*-EBM Algorithm in PyTorch

```

1 def update_H(real_x, fake_x, model_h, model_e, optim_h, grad_exp, conjugate_grad_exp):
2     // Step 6-7 in Algorithm 1. Update the parameter of the variational function.
3     // - real_x and fake_x are samples from the data distribution and EBM respectively.
4     // - model_h is the neural network for the variational function $H_\omega$.
5     // - model_e is the neural network for the energy function $E_\theta$.
6     // - optim_h is the optimizer for model_h, e.g. torch.optim.SGD(model_h.parameters())
7     // - grad_exp and conjugate_grad_exp are functions defined by the used f-divergence.
8     real_e, fake_e = model_e(real_x), model_e(fake_x)
9     real_h, fake_h = model_h(real_x), model_h(fake_x)
10    loss_h = -(grad_exp(real_h+real_e) - conjugate_grad_exp(fake_h+fake_e)).mean()
11    optim_h.zero_grad()
12    loss_h.backward()
13    optim_h.step()
14
15 def update_E(real_x, fake_x, model_h, model_e, optim_e, grad_exp, conjugate_grad_exp):
16    // Step 8-9 in Algorithm 1. Update the parameter of the energy function.
17    // - optim_e is the optimizer for model_e, e.g. torch.optim.SGD(model_e.parameters())
18    real_e, fake_e = model_e(real_x), model_e(fake_x)
19    real_h, fake_h = model_h(real_x), model_h(fake_x)
20    loss_e = torch.mean(grad_exp(real_h+real_e)) + \
21              torch.mean(conjugate_grad_exp(fake_h+fake_e).detach() * fake_e) - \
22              torch.mean(conjugate_grad_exp(fake_h+fake_e)) - \
23              torch.mean(fake_e) * torch.mean(conjugate_grad_exp(fake_h+fake_e)).detach()
24    optim_e.zero_grad()
25    loss_e.backward()
26    optim_e.step()

```

Note that according to Lemma 3, technically we should use two independent batches of data to estimate the two expectations in the product $\mathbb{E}_{q_\theta(\mathbf{x})}[\nabla_\theta E_\theta(\mathbf{x})] \cdot \mathbb{E}_{q_\theta(\mathbf{x})}[f^*(f'(\exp(E_\theta(\mathbf{x}) + H_\omega(\mathbf{x}))))]$. Empirically we found that using the same set of samples also works well, since it is an asymptotically consistent estimator.

To implement *f*-EBM in Tensorflow (Abadi et al., 2016), the main change is to use `tf.stop_gradient(X)` to replace `X.detach()`. Functions `grad_exp` and `conjugate_grad_exp` correspond to $f'(\exp(u))$ and $f^*(f'(\exp(u)))$, which

can be derived based on the definitions of *f*-divergences (see Table 5 and Table 6 in (Nowozin et al., 2016) for reference). Some examples of $f'(\exp(u))$ and $f^*(f'(\exp(u)))$ can be found in Table 4.

Name	$D_f(P Q)$	$f'(\exp(u))$	$f^*(f'(\exp(u)))$
Kullback-Leibler	$\int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$	$1 + u$	$\exp(u)$
Reverse Kullback-Leibler	$\int q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x}$	$-\exp(-u)$	$-1 + u$
Pearson χ^2	$\int \frac{(q(\mathbf{x})-p(\mathbf{x}))^2}{p(\mathbf{x})} d\mathbf{x}$	$2 \exp(u) - 2$	$\exp(2u) - 1$
Neyman χ^2	$\int \frac{(p(\mathbf{x})-q(\mathbf{x}))^2}{q(\mathbf{x})} d\mathbf{x}$	$1 - \exp(-2u)$	$2 - 2 \exp(-u)$
Squared Hellinger	$\int \left(\sqrt{p(\mathbf{x})} - \sqrt{q(\mathbf{x})} \right)^2 d\mathbf{x}$	$1 - \exp(-\frac{u}{2})$	$\exp(\frac{u}{2}) - 1$
Jensen-Shannon	$\frac{1}{2} \int p(\mathbf{x}) \log \frac{2p(\mathbf{x})}{p(\mathbf{x})+q(\mathbf{x})} + q(\mathbf{x}) \log \frac{2q(\mathbf{x})}{p(\mathbf{x})+q(\mathbf{x})} d\mathbf{x}$	$\log(2) + u - \log(1 + \exp(u))$	$-\log(2) + \log(1 + \exp(u))$
α -divergence ($\alpha \notin \{0, 1\}$)	$\frac{1}{\alpha(\alpha-1)} \int \left(p(\mathbf{x}) \left[\left(\frac{q(\mathbf{x})}{p(\mathbf{x})} \right)^\alpha - 1 \right] - \alpha(q(\mathbf{x}) - p(\mathbf{x})) \right) d\mathbf{x}$	$\frac{1}{\alpha-1} (\exp((\alpha-1)u) - 1)$	$\frac{1}{\alpha} (\exp(\alpha u) - 1)$

Table 4. Some examples of *f*-divergences and corresponding $f'(\exp(u))$ and $f^*(f'(\exp(u)))$ functions.

E.2. Discussion on Differentiating Through Langevin Dynamics

In this section, we provide a more detailed discussion on gradient reparametrization that was initially introduced in Section 3.1.2. Specifically, we will discuss the possibility of directly extending the *f*-GANs framework by differentiating through the Langevin dynamics. As discussed before, we typically need hundreds of Langevin steps to produce a single sample while we cannot use a sample replay buffer to reduce the number of transition steps (because the initial distribution of the Markov chain cannot depend on the model parameters). During gradient backpropagation, we need the same number of backward steps, where each backward step further involves computing Hessian matrices that are proportional to the parameter and data dimension. This will lead to hundreds of times more memory consumption compared to only using Langevin dynamics for producing samples. More specifically, we provide the following implementation to differentiate through Langevin dynamics in PyTorch:

```

1 def train_energy(model_e, model_dis, optim_e, conjugate_fn, device='cuda'):
2     // Initialize Langevin dynamics with random uniform distribution.
3     fake_x = torch.rand(batch_size, channel_num, img_size, img_size, device=device)
4     fake_x.requires_grad = True
5     gaussian_noise = torch.randn(batch_size, 3, 32, 32, device=device)
6
7     for k in range(num_langevin_steps):
8         energy_value = model_e(fake_x)
9         fake_x_grad = torch.autograd.grad(energy_value.sum(), fake_x,
10                                         create_graph=True)[0]
11         gaussian_noise.normal_(0, gaussian_noise_std)
12         fake_x = fake_x - step_size * fake_x_grad + gaussian_noise
13         fake_x.data.clamp_(0, 1)
14
15     fake_dis = model_dis(fake_x)
16     energy_loss = - torch.mean(conjugate_fn(fake_dis))
17     optim_e.zero_grad()
18     energy_loss.backward()
19     optim_e.step()

```

Note that in Line 10, we need to set `create_graph=True` in order to compute the second-order derivatives when backpropagating through Langevin dynamics later, which will store all the computation graphs along the Langevin dynamics. By contrast, in *f*-EBMs we only use Langevin dynamics to produce samples and once we get the value of the gradient $\nabla_{\mathbf{x}} E_{\theta}(\mathbf{x})$, we can discard the computational graph immediately after each transition step (`create_graph=False`). As a result, with gradient reparametrization, the memory consumption grows linearly as the number of transitions. With modern GPU such as NVIDIA GeForce RTX 2080 Ti and the model architecture in Figure 15, we will run out of memory when the number of Langevin steps is larger than 5. Since gradient backpropagation through Langevin dynamics involves computing many Hessian matrices, this approach is also computationally less efficient compared to *f*-EBMs. In our experiments, we set `num_langevin_steps = 5` and we observed that this approach cannot produce reasonable images.

F. Additional Experimental Results for Fitting Univariate Mixture of Gaussians

F.1. Parameter Learning Results

Objective	μ^*	$\hat{\mu}$	σ^*	$\hat{\sigma}$
Contrastive Divergence	1.01065	1.01204	1.82895	1.82907
Kullback-Leibler	1.01065	1.01536	1.82895	1.83024
Reverse KL	1.58454	1.58523	1.63106	1.63453
Squared Hellinger	1.32024	1.32274	1.73089	1.74710
Jensen Shannon	1.30322	1.31669	1.76716	1.76041
Pearson χ^2	0.57581	0.56563	1.92172	1.93461
Neyman χ^2	1.83037	1.82676	1.51508	1.51598
α -divergence ($\alpha = -0.5$)	1.74642	1.74332	1.55569	1.55209
α -divergence ($\alpha = -1$)	1.82923	1.81979	1.51844	1.52513
α -divergence ($\alpha = 0.9$)	1.07056	1.07237	1.81091	1.81852

Table 5. Fitting a quadratic EBM to mixtures of Gaussians. μ^*, σ^* represent the desired optimal solution under a certain discrepancy measure, and $\hat{\mu}, \hat{\sigma}$ represent the learned parameters. The first row is for contrastive divergence (with KL divergence being the underlying objective). The other rows are for *f*-EBMs with various discrepancy measures as the training objectives.

F.2. Density Ratio Estimation Results

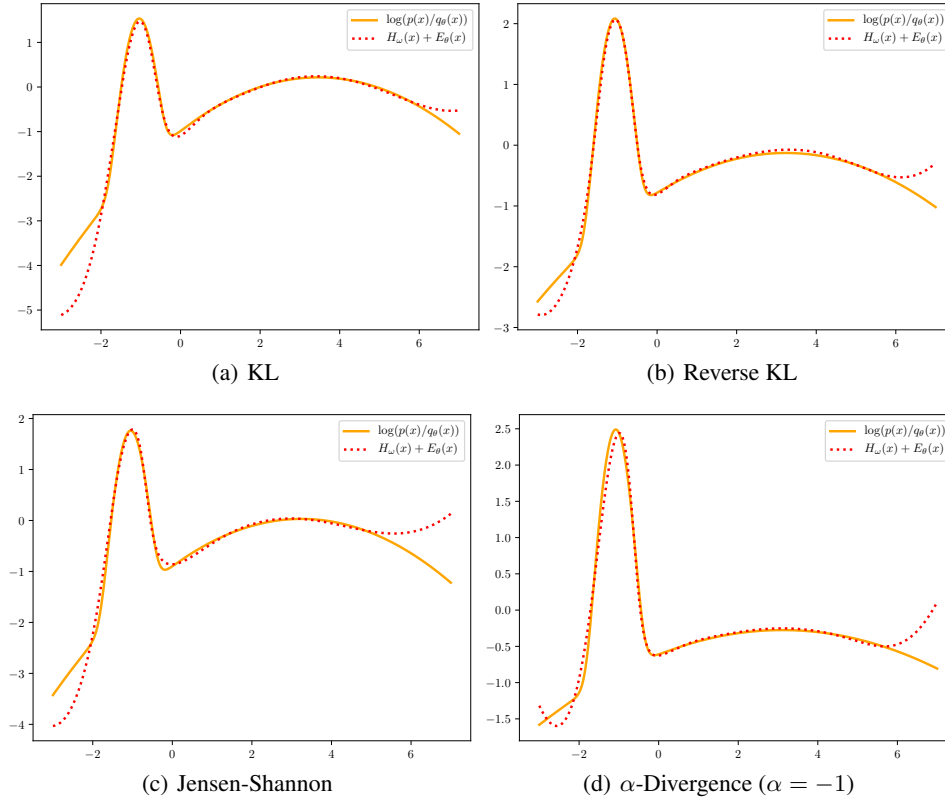


Figure 3. Density ratio estimation results for different *f*-divergences. The orange solid line represents the ground truth log-scale density ratio ($\log(p(x)/q_\theta(x))$) under a certain divergence; The red dashed line represents the estimated log-scale density ratio ($H_\omega(x) + E_\theta(x)$) learned by *f*-EBM. Note that the estimated density ratio is accurate in most areas except in the low density regimes (e.g. $(-\infty, -2]$ and $[6, +\infty)$) where very few training data comes from this region.

E.3. Optimization Trajectories of Single-Step *f*-EBM

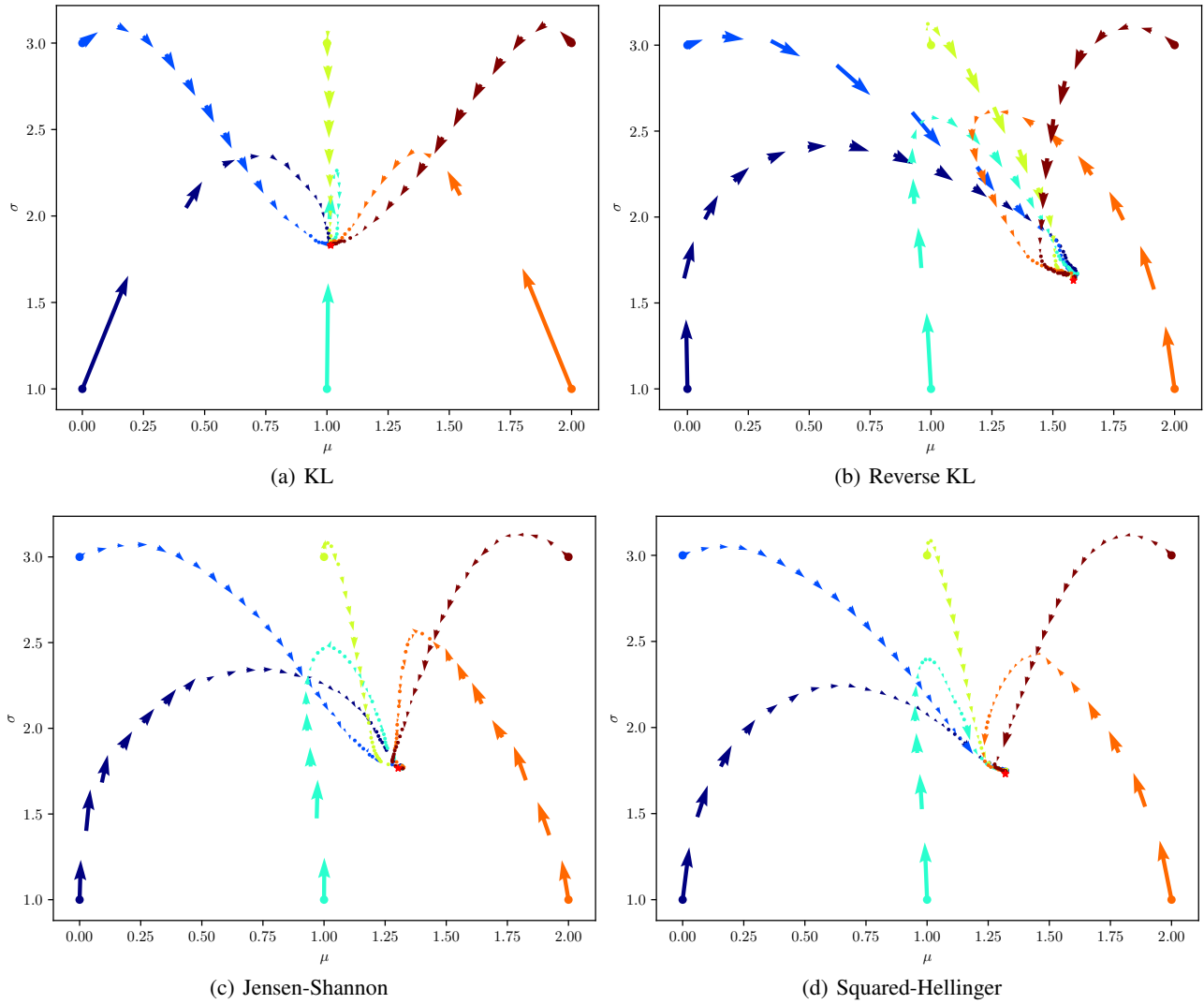


Figure 4. Convergence of Single-Step *f*-EBM algorithm for different *f*-divergences. Optimization trajectories in different colors start from different initializations. The length and direction of the arrows represent the scale and the direction of the gradient at a certain point. The red stars that the trajectories converge to represent the desired optimal solutions under corresponding *f*-divergences.

G. Additional Experimental Details for Modeling Natural Images

G.1. Samples from The Baseline Method in Section 3.1

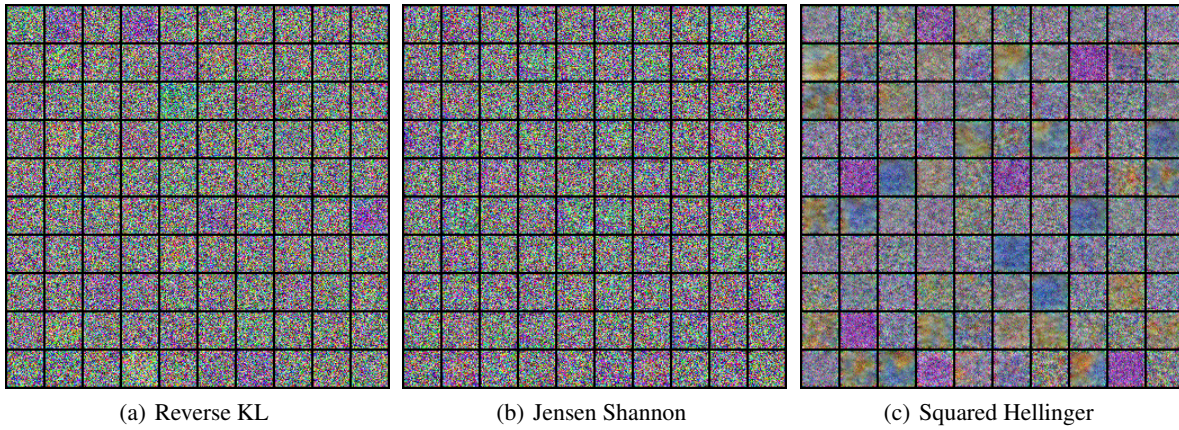


Figure 5. Uncurated samples from EBMs trained by the baseline approach described in Section 3.1 on CIFAR-10 dataset.

G.2. Uncurated CIFAR-10 Samples from *f*-EBM



Figure 6. Uncurated CIFAR-10 samples from *f*-EBM under the guidance of Jensen Shannon.

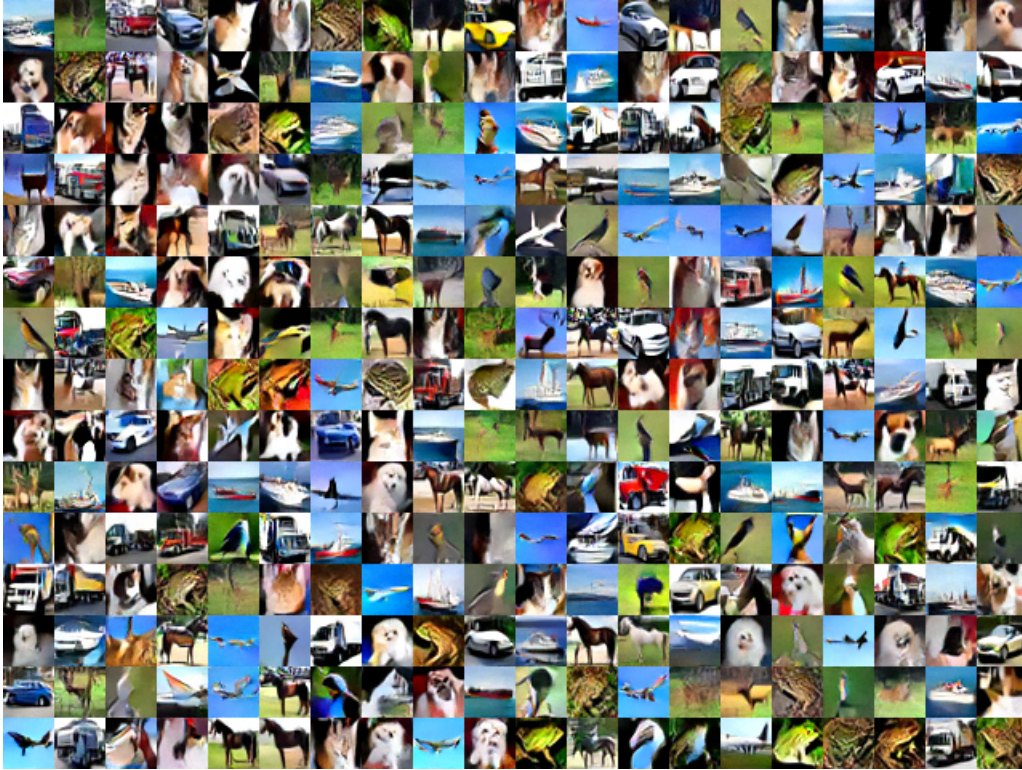
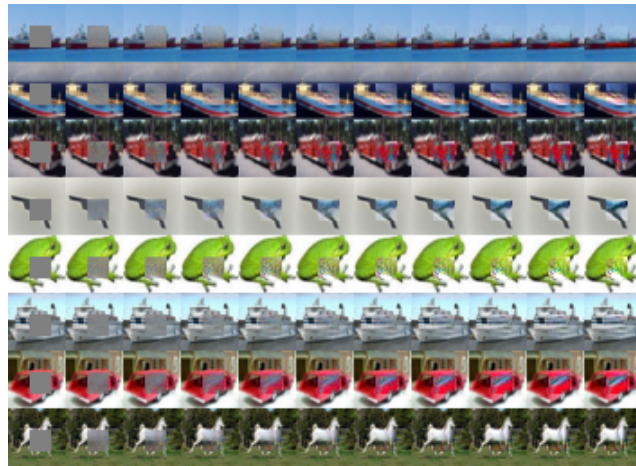


Figure 7. Uncurated CIFAR-10 samples from *f*-EBM under the guidance of Squared Hellinger.

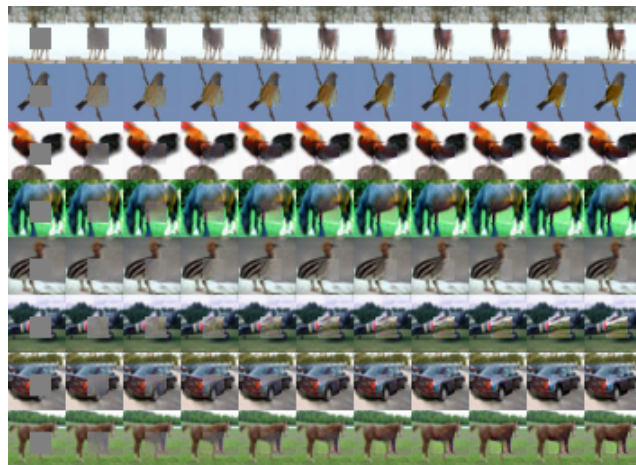


Figure 8. Uncurated CIFAR-10 samples from *f*-EBM under the guidance of Reverse KL.

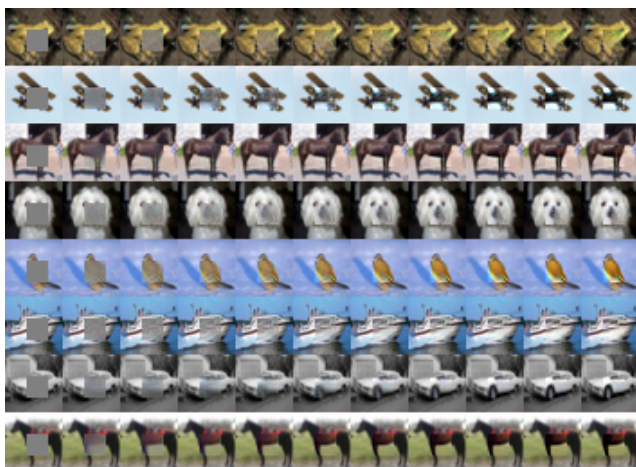
G.3. Image Inpainting Results



(a) Reverse KL



(b) Jensen Shannon



(c) Squared Hellinger

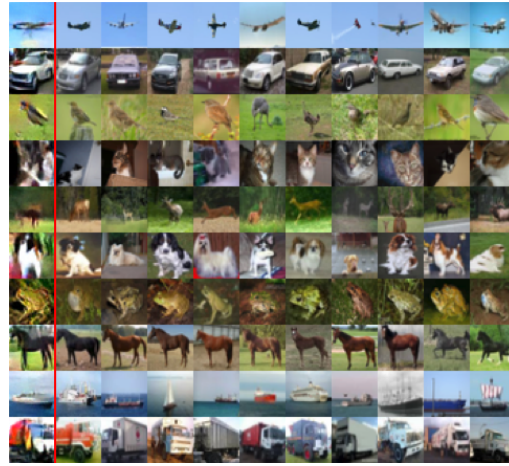
Figure 9. Image inpainting results for f -EBM trained with different f -divergences on CIFAR-10. We use Langevin dynamics sampling to restore the images which are corrupted by empty boxes.

G.4. Image Denoising Results

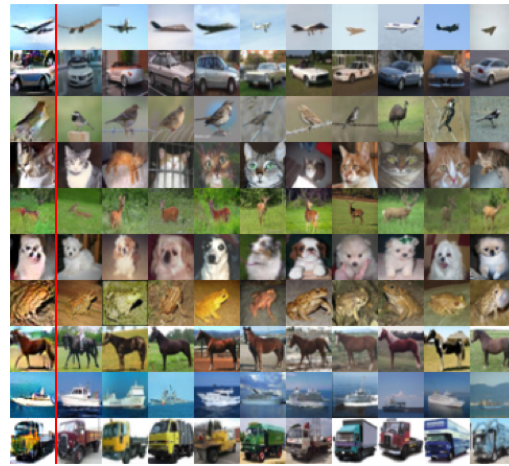


Figure 10. Image denoising results for f -EBM trained with different f -divergences on CIFAR-10. We apply 10% “salt and pepper” noise to the images in the test set and use Langevin dynamics sampling to restore the images.

G.5. Nearest Neighbor Images



(a) Reverse KL



(b) Jensen Shannon



(c) Squared Hellinger

Figure 11. Nearest neighbor images according to l_2 distance between images. Different rows are for different classes. In each row, the leftmost image (*i.e.*, on the left of the right vertical line) is generated by *f*-EBM, and the other images are nearest neighbors of the generated image in the training set.

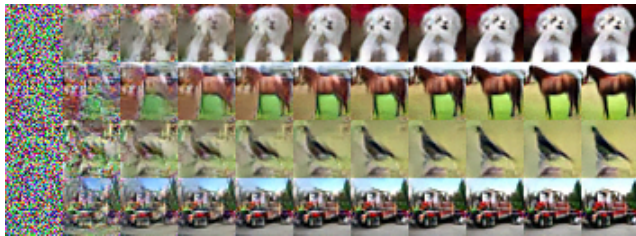
G.6. Intermediate Samples of Langevin Dynamics on CIFAR-10



(a) Reverse KL



(b) Jensen Shannon



(c) Squared Hellinger

Figure 12. Intermediate samples during Langevin dynamics sampling process for *f*-EBM.

G.7. Uncurated CelebA Samples from *f*-EBM



(a) Reverse KL



(b) Jensen Shannon



(c) Squared Hellinger

Figure 13. Uncurated CelebA samples from *f*-EBM.

G.8. Intermediate Samples of Langevin Dynamics on CelebA



(a) Reverse KL



(b) Jensen Shannon



(c) Squared Hellinger

Figure 14. Intermediate samples during Langevin dynamics sampling process for *f*-EBM.

G.9. Architectures and Training Hyperparameters

For all the experiments on natural images (conditional generation for CIFAR-10 and unconditional generation for CelebA), we use the residual network architecture (He et al., 2016) in Figure 15 (same as the conditional CIFAR-10 model in (Du & Mordatch, 2019)) to implement both the energy function and the variational function of *f*-EBM. For the contrastive divergence baseline, we also use the same model architecture for fair comparisons. We note that the model architecture is highly relevant to the model performance and we leave the investigation of better architectures to future works.

For the CelebA dataset, we first center-crop the images to 140×140 , then resize them to 32×32 . For both CelebA and CIFAR-10, we rescale the images to $[0, 1]$. Following (Du & Mordatch, 2019), we apply spectral normalization and L_2 regularization (on the outputs of the models) with coefficient 1.0 to improve the stability. We use 60 steps Langevin dynamics together with a sample replay buffer of size 10000 to produce samples in the training phase. In each Langevin step, we use a step size of 10.0 and a random noise with standard deviation of 0.005. We use Adam optimizer with $\beta_1 = 0.0, \beta_2 = 0.999$ and learning rate of 3×10^{-4} to optimize the parameters of both the energy function and the variational function. In each training iteration, we use a batch of 128 positive images and negative images. These training hyperparameters are used for both *f*-EBMs and the contrastive divergence baseline.

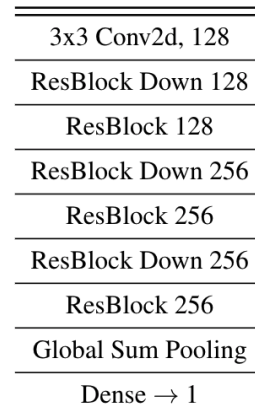


Figure 15. ResNet architecture for implementing the energy function and the variational function.

G.10. Time Complexity and Convergence Speed

In the experiments we use the same batch size, learning rate and MCMC step number as in CD (Du & Mordatch, 2019). Since we use single-step minimax optimization for both the energy function and the variational function, each iteration needs the same time for MCMC as (Du & Mordatch, 2019) and twice time for stochastic gradient descent parameter updates. Overall, the time complexity per iteration of *f*-EBM is comparable to that of CD (within a factor of 2). The convergence speed of *f*-EBM is also similar to CD in terms of number of iterations needed (CD: 75K iterations; *f*-EBM with Jensen Shannon: 50K iterations; *f*-EBM with Squared Hellinger: 80K iterations; *f*-EBM with Reverse KL: 70K iterations; *f*-EBM with KL: 75K iterations).