
Max-Sliced Score Matching

Safa Messaoud¹ Kuan-Chieh Wang² Alexander Schwing³

Abstract

Score matching methods have recently shown a lot of promise for estimating unnormalized densities. Scaling score matching to high dimensional spaces, is however challenging as it requires an accurate estimation of the trace of a Hessian. To improve upon sliced score matching, which circumvents computation of the Hessian by projecting the scores into random directions before comparing them, we study max-sliced score matching. Max-sliced score matching benefits from its use of the most informative projection directions rather than using randomly sampled ones. We prove that max-sliced score matching is consistent and asymptotically normal. We assess these models on estimation of densities and scores for implicit distributions in Variational and Wasserstein auto-encoders.

1. Introduction

Different methods have been proposed to learn unnormalized probabilistic models, including pseudo-likelihood models (Besag, 1975), contrastive divergence (Hinton, 2002; Tieleman & Hinton, 2009; Song & Ermon, 2019b) and noise-contrastive estimation (Gutmann & Hyvärinen, 2010). Notably, score matching methods (Vincent, 2011; Song & Ermon, 2019a; Jolicoeur-Martineau et al., 2020) which minimize the distance between the derivatives of the log-density functions, *i.e.*, the scores, of the data and model distributions have received increasing interest as they circumvent computation of the partition function.

However, score matching presents an additional challenge, as it requires computation of the trace of a Hessian, which involves a number of forward and backward propagations proportional to the data dimensionality (Vincent, 2011).

To overcome this issue, Kingma & LeCun (2010); Martens et al. (2012) propose to approximate the trace of the Hessian

via approximate back-propagation (Kingma & LeCun, 2010) or curvature propagation (Martens et al., 2012). However, these methods have no theoretical guarantees for the approximation error or are memory demanding even in modern back-propagation pipelines, as they either limit back-propagation to the diagonal of the Hessian (Kingma & LeCun, 2010) or require use of complex numbers (Martens et al., 2012).

More recently, denoising score matching (Vincent, 2011; Song & Ermon, 2019a; Jolicoeur-Martineau et al., 2020) was proposed to circumvent computation of the trace by perturbing the data points with a pre-specified noise distribution and using the score of the perturbed data as an estimator for the scores of the model. This procedure can however only recover the noise corrupted data and is very sensitive to the choice of the noise scales. Sliced score matching (Song et al., 2019) is another method that circumvents computation of the Hessian matrix. Instead of directly matching the high-dimensional scores, it matches their projections along random directions. Sliced score matching only involves Hessian-vector products, computed without explicitly evaluating the Hessian.

In this paper, we study max-sliced score matching, a variant of sliced score matching that finds the most informative directions for the projection of the scores. We show that those directions are the eigenvectors of the Hessian. We also show that using the Lanczos algorithm permits to compute those directions efficiently.

We assess the performance of max-sliced score matching on (1) density estimation, specifically for the class of Deep Kernel Exponential Families (DKEF) (Wenliang et al., 2019) and NICE flow models (Dinh et al., 2014), and (2) score estimation for implicit distributions in variational inference applications, specifically image generation using variational auto-encoders (VAE) (Kingma & Welling, 2014) and Wasserstein auto-encoders (WAE) (Tolstikhin et al., 2018) with implicit encoders.

2. Related Work

In the following, we review score matching in greater detail.

Score matching methods use surrogate objectives which have similar optima to the maximum likelihood objective.

¹University of Illinois at Urbana-Champaign ²University of Toronto ³University of Illinois at Urbana-Champaign. Correspondence to: Safa Messaoud <messaou2@illinois.edu>.

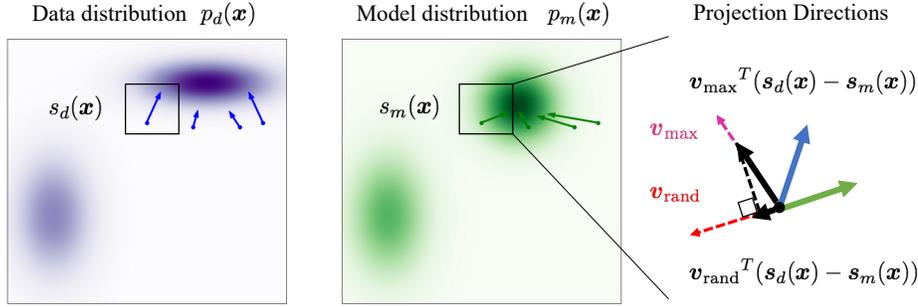


Figure 1. Visualization of the effect of slicing. $s_d(\mathbf{x})$, and $s_m(\mathbf{x})$ are the scores, *i.e.*, gradients of the log density of the data and model distributions respectively. Sliced score matching depends on projecting the difference in scores to a scalar. A random projection \mathbf{v}_{rand} could make the difference in scores appear smaller than it is. Our work aims to select the projection that maximizes the difference in scores (see \mathbf{v}_{max}). We find this to make learning more efficient.

Different variants of score matching have been used for different applications including image generation (Song et al., 2020; Song & Ermon, 2019a; Jolicoeur-Martineau et al., 2020), shape generation from point cloud data (Cai et al., 2020), waveform generation (Chen et al., 2020), and out-of-distribution image detection (Mahmood et al., 2020). Bao et al. (2020) extend score matching methods to learn general energy-based latent variable models (EBLVMs). Meng et al. (2020) propose a new divergence between distributions, which depends only on the derivatives of univariate log-conditionals (scores) of an autoregressive likelihood-based model.

Score matching (SM) (Hyvärinen, 2005) minimizes the Fisher divergence between the data distribution and the model, *i.e.*,

$$\min_{\theta} \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_d(\mathbf{x})} [(s_m(\mathbf{x}; \theta) - s_d(\mathbf{x}))^2], \quad (1)$$

with $s_d(\mathbf{x}) = \nabla_{\mathbf{x}} \log(p_d(\mathbf{x}))$ the score of the unknown data distribution and $s_m(\mathbf{x}; \theta) = \nabla_{\mathbf{x}} \log(p(\mathbf{x}; \theta))$ the score of the model distribution. By applying integration by parts, Hyvärinen (2005) shows that Eq. (1) is equivalent to

$$\min_{\theta} \mathbb{E}_{\mathbf{x} \sim p_d} [\text{tr}(\nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta)) + \frac{1}{2} \|s_m(\mathbf{x}; \theta)\|^2], \quad (2)$$

with $\text{tr}(\nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta))$ being the trace of the Hessian $\nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta)$. While this formulation avoids the computation of the partition function, it introduces a new challenge, *i.e.*, computing the trace of the Hessian, which scales quadratically with the input dimensions. This often makes the application of score matching slow and impractical for real applications with high-dimensional data. Different methods have been proposed to approximate the trace. Pang et al. (2020) propose to approximate the higher order directional derivative with finite differences (FD). This however trades accuracy for speed. Another approach, *Approximate*

Backpropagation (Kingma & LeCun, 2010) limits backpropagation to the diagonal of the Hessian instead of computing the full matrix. The method has no guarantees on the approximation errors. Similarly, Martens et al. (2012) propose *Curvature propagation (CP)*, an estimator that introduces noise for each node in the network. This however leads to high variance and requires operating on complex numbers in neural networks, which is memory demanding.

Denosing score matching (DSM) (Vincent, 2011; Song & Ermon, 2019a) is a different method for scaling score matching, that circumvents computation of the Hessian. For this, DSM applies the original score matching loss to noise-corrupted data, *i.e.*,

$$\min_{\theta} \frac{1}{2} \mathbb{E}_{q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})p_d(\mathbf{x})} [\|s_m(\tilde{\mathbf{x}}; \theta) - \nabla_{\tilde{\mathbf{x}}} \log q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})\|^2]. \quad (3)$$

Here, $q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})$ is a noise distribution. Song & Ermon (2019a) choose $q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})$ to be a Gaussian distribution $\mathcal{N}(\tilde{\mathbf{x}}|\mathbf{x}, \sigma^2 \mathbf{I})$. The noisy samples in this case are generated by injecting Gaussian noise to the original samples, *i.e.*, $\tilde{\mathbf{x}} = \mathbf{x} + \epsilon \sigma^2$, $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. DSM has the following limitations: (1) it can only recover the noise corrupted data distribution, (2) its performance is very sensitive to the choice of σ which is non trivial and often based on heuristics. Recently, new solutions for setting the parameters of the model have been proposed by Jolicoeur-Martineau et al. (2020). Additionally, Jolicoeur-Martineau et al. (2020) extend DSM with an adversarial loss. This led to better performance on image generation.

Sliced Score Matching (SSM) (Song et al., 2019) is another method aiming to scale the SM loss (Eq. (2)) by projecting $s_d(\mathbf{x})$ and $s_m(\mathbf{x}; \theta)$ onto some random direction \mathbf{v} and comparing their average difference along those random directions, *i.e.*,

$$\min_{\theta} \frac{1}{2} \mathbb{E}_{\mathbf{v} \sim p_{\mathbf{v}}} \mathbb{E}_{\mathbf{x} \sim p_d} [(\mathbf{v}^T s_m(\mathbf{x}; \theta) - \mathbf{v}^T s_d(\mathbf{x}))^2], \quad (4)$$

with p_v satisfying $\mathbb{E}_{p_v}[v v^T] = I$. Song et al. (2019) show, via integration by parts, that the program above is equivalent to

$$\min_{\theta} \mathbb{E}_{p_v} \mathbb{E}_{p_d} \left[\frac{1}{2} \|s_m(\mathbf{x}; \theta)\|^2 + \mathbf{v}^T \nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta) \mathbf{v} \right]. \quad (5)$$

Intuitively, the first term in the loss is lowest if \mathbf{x} is a stationary point of $\log p(\mathbf{x}; \theta)$. The second term pushes the Hessian $\nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta)$ toward negative definiteness, to ensure that $\log p(\mathbf{x}; \theta)$ is concave in \mathbf{x} , *i.e.*, \mathbf{x} is a maximum of $\log p(\mathbf{x}; \theta)$. Note that $\mathbf{v}^T \nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta) \mathbf{v}$ with $\mathbb{E}_{p_v}[v v^T]$ is the Hutchinson estimator (Hutchinson, 1989) of the trace of the Hessian $\nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta)$. The second order derivative term $\mathbf{v}^T \nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta) \mathbf{v}$ can be efficiently computed using two gradient operations for a single \mathbf{v} . If we average over M directions, $M + 1$ gradient operations are required. In contrast, estimating the trace requires d differentiations. Hence, $M < d$ should be ensured to achieve an advantage of the method. SSM’s main limitation is a high variance due to the small number of the projection directions. In this paper, we propose a new score matching objective, max-score matching (MSSM). It replaces the averaging over random projection directions with computation of the most informative direction.

3. Max-Sliced Score Matching (MSSM)

Given a dataset $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^{|\mathcal{D}|}$ of samples $\mathbf{x} \sim p_d(\mathbf{x})$ drawn from an unknown data distribution $p_d(\mathbf{x})$, the goal of density estimation is to find the parameters θ of a model distribution $p_m(\mathbf{x}; \theta)$ which best fits the samples \mathcal{D} .

Energy-based models (EBMs) are frequently used for this purpose, representing the model distribution via

$$p_m(\mathbf{x}; \theta) = \frac{\exp -E(\mathbf{x}; \theta)}{Z}, \quad Z = \int_{\mathbf{x}'} \exp -E(\mathbf{x}'; \theta) d\mathbf{x}',$$

where Z is referred to as the normalizing constant, or partition function.

There are two families of techniques for learning an EBM: by approximating the maximum likelihood gradient (see, *e.g.*, work by Du & Mordatch (2019) and references therein), or using score matching objectives. The first approach suffers from having to approximate the intractable partition function. Score matching methods circumvent this by only looking at the scores $s(\mathbf{x})$, *i.e.*, the gradient of the unnormalized density $s(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x})$. Hyvärinen (2005) proposed the Score Matching (SM) objective (Eq. (2)) for learning EBMs, but the Hessian term scales quadratically with the input dimensions. To reduce computation, Song et al. (2019) proposed Sliced Score Matching (SSM) which addresses

$$\min_{\theta} \frac{1}{2} \mathbb{E}_{\mathbf{v} \sim p_v} \mathbb{E}_{\mathbf{x} \sim p_d} [(\mathbf{v}^T s_m(\mathbf{x}; \theta) - \mathbf{v}^T s_d(\mathbf{x}))^2],$$

Algorithm 1 Max-Sliced Score Matching

Input: Update frequency F of \mathbf{V} ; Number of computed eigenvectors K of $\nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta)$; Dataset \mathcal{D} ; projection direction distribution $p_v(\mathbf{v})$; Model distribution $p_m(\mathbf{x}; \theta)$.
step $\leftarrow 0$

repeat

 Sample $\mathbf{x} \sim \mathcal{D}$

$s_m(\mathbf{x}; \theta) \leftarrow \nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \theta)$

$\mathcal{L}_{\text{MSSM}}(\mathbf{x}; \theta) \leftarrow \frac{1}{2} \|s_m(\mathbf{x}; \theta)\|^2$

if (step % F) == 0 **then**

 | $\mathbf{V} \leftarrow \text{eig}(\nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta))$

end

for $i \leftarrow 1$ to K **do**

 | $\mathcal{L}_{\text{MSSM}} \leftarrow \mathcal{L}_{\text{MSSM}} + \frac{\hat{p}_v(\mathbf{v}_i)}{K} \mathbf{v}_i^T \nabla_{\mathbf{x}} (s_m(\mathbf{x}; \theta) \mathbf{v}_i)$

end

 Update θ via Backpropagation

 step \leftarrow step + 1

Return θ

where $s_m(\cdot)$ and $s_d(\cdot)$ are the scores of the model and data distributions respectively. Intuitively, the SSM objective applies random projections to the scores in the Fisher divergence (Eq. (1)).

We observe the random nature of the projections \mathbf{v} to make learning inefficient. To address this we study a method which finds projections along directions that maximize the projected Fisher divergence via

$$\begin{aligned} \min_{\theta} \max_{\mathbf{V}} \quad & \mathbb{E}_{\mathbf{x} \sim p_d} [\|\mathbf{V}^T (s_d(\mathbf{x}) - s_m(\mathbf{x}; \theta))\|^2], \\ \text{s.t.} \quad & \mathbf{V} \mathbf{V}^T = I, \end{aligned} \quad (6)$$

where $\mathbf{V} \in \mathbb{R}^{d \times K}$, $K \leq d$, denotes the matrix with its i -th column $\mathbf{v}_i \in \mathbb{R}^d$ being a projection direction. We use $\|\cdot\|$ to denote the 2-norm of a vector. The constraint $\mathbf{V} \mathbf{V}^T = I$ is introduced to ensure that we make the most use of the space by requiring orthogonality. Beneficially, the constraint ensures boundedness of the loss (vectors \mathbf{v}_i can’t be scaled so that maximization reaches infinity). We refer to this objective as MSSM. In Fig. 1, we illustrate the scores $s_m(\mathbf{x}; \theta)$, $s_d(\mathbf{x})$ as well as the projection directions for a Gaussian mixture model.

We show that MSSM is consistent, *i.e.*, $p_m(\mathbf{x}; \theta)$ converges to $p_d(\mathbf{x})$ when Eq. (6) converges to 0. We sketch the proof in the following (see Appendix B for a full proof):

$$\begin{aligned} & \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} [\|\mathbf{V}^T (s_m(\mathbf{x}; \theta) - s_d(\mathbf{x}))\|^2] = 0 \\ \iff & \max_{\mathbf{V}} \|\mathbf{V}^T (s_m(\mathbf{x}; \theta) - s_d(\mathbf{x}))\|^2 = 0 \quad \forall \mathbf{x} \sim p_d(\mathbf{x}) \\ \iff & s_m(\mathbf{x}; \theta) = s_d(\mathbf{x}), \quad \forall \mathbf{x} \sim p_d(\mathbf{x}) \\ \iff & \log p_m(\mathbf{x}; \theta) = \log p_d(\mathbf{x}) + \text{Const}, \quad \forall \mathbf{x} \sim p_d(\mathbf{x}) \\ \iff & p_d(\mathbf{x}) = p_m(\mathbf{x}; \theta), \quad \forall \mathbf{x} \sim p_d(\mathbf{x}) \end{aligned}$$

Firstly, Eq. (6) can only be zero if for all real data the maximum projection of the difference between the model and groundtruth data scores are zero, forcing the model scores to match the data scores exactly. This is equivalent to matching the two normalized densities which concludes the proof.

Since $s_d(\mathbf{x})$ cannot be computed (the data distribution is unknown), we follow Hyvärinen (2005); Song et al. (2019) and use integration-by-parts to obtain an objective without $s_d(\mathbf{x})$ whose only difference from Eq. (6) is a constant which is independent of $s_m(\mathbf{x})$:

$$\begin{aligned} \min_{\theta} \quad & \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x}} \left[\sum_{i=1}^K \mathbf{v}_i^T \nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta) \mathbf{v}_i + \frac{1}{2} \|s_m(\mathbf{x}; \theta)\|^2 \right], \\ \text{s.t.} \quad & \mathbf{V}\mathbf{V}^T = \mathbf{I}. \end{aligned} \quad (7)$$

We defer the proof of equivalence to Appendix C.

We argue that MSSM leads to a better maximum likelihood estimation than SSM. We show in Appendix D, that MSSM allows a better approximation of the partition function, as d vectors are needed for estimating the trace exactly, whereas 2^d vectors are needed for the sum.

Optimization for MSSM: In the following we first discuss how to address the program in Eq. (7) before elaborating on efficiency.

Claim: The solution to the inner maximization problem in Eq. (7) is the eigenvector basis of the Hessian $\nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta)$.

Proof: The Lagrangian of the program above is $\mathcal{L}(\mathbf{V}; \theta) = \mathbb{E}_{\mathbf{x}} [\sum_{i=1}^d \mathbf{v}_i^T \nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta) \mathbf{v}_i + \frac{1}{2} \|s_m(\mathbf{x}; \theta)\|^2 - \sum_{i=1}^d \lambda_i (\|\mathbf{v}_i\|^2 - 1)]$ with $\mathbf{v}_i^T \mathbf{v}_j = 0, (\forall i \neq j)$. For a given \mathbf{v}_i , $\frac{\partial \mathcal{L}(\mathbf{v}_i; \theta)}{\partial \mathbf{v}_i} = 0$ is equivalent to $\nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta) \mathbf{v}_i = \lambda_i \mathbf{v}_i$. Hence, the direction \mathbf{v}_i that solves the linear system corresponds to the eigenvector of the symmetric Hessian $\nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta)$ with the eigenvalue λ_i . The condition $\mathbf{v}_i^T \mathbf{v}_j = 0$ is satisfied since the eigenvectors are orthogonal to each other which concludes the proof. \square

Unlike log-densities $\log p_d(\mathbf{x})$, which decrease as we move away from the data-manifold, the magnitude of the scores $\nabla_{\mathbf{x}} \log p_d(\mathbf{x})$ generally increases as we move away from the manifold. Hence, a good score estimation requires projection directions that encompass the entire space. The MSSM formulation uses a basis of d vectors to cover different directions in the d -dimensional space. As mentioned before, this is more efficient than sampling a large number of random directions to properly estimate the expectation in Eq. (4).

However, computing the eigenvectors using traditional eigenvector decomposition is very expensive and requires explicit computation and storage of the Hessian $\nabla_{\mathbf{x}} s_m(\mathbf{x}; \theta)$, which defies the purpose of applying slic-

Algorithm 2 Lanczos Algorithm for Score Matching

Input: Number of iterations K . The score $s_m(\mathbf{x}; \theta)$.
Initial $\mathbf{q}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

Choose $\mathbf{r} = \mathbf{q}_0$
 $\beta_0 = \|\mathbf{q}_0\|$

while $k = 1, \dots, K$ **do**

$$\begin{aligned} \mathbf{q}_k &= \frac{\mathbf{r}}{\beta_{k-1}} \\ \mathbf{r} &= \nabla_{\mathbf{x}} (s_m(\mathbf{x}; \theta)^T \mathbf{q}_k) \\ \mathbf{r} &= \mathbf{r} - \mathbf{q}_{k-1} \beta_{k-1} \\ \alpha_k &= \mathbf{q}_k^T \mathbf{r} \\ \mathbf{r} &= \mathbf{r} - \mathbf{q}_k \alpha_k \\ \beta_k &= \|\mathbf{r}\|^2 \end{aligned}$$

end

$$\mathbf{T}_K = \begin{bmatrix} \alpha_1 & \beta_1 & & & \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \beta_2 & \alpha_3 & & \\ & & & \ddots & \beta_{K-1} \\ & & & \beta_{K-1} & \alpha_K \end{bmatrix}$$

$$\begin{aligned} \mathbf{Q}_K &= [\mathbf{q}_0, \dots, \mathbf{q}_K] \\ \mathbf{S}_K &= \text{eig}(\mathbf{T}_K) \end{aligned}$$

Return $\mathbf{V}^* = \mathbf{Q}_K \mathbf{S}_K$

ing/projection to Score Matching. To address this concern we study the combination of three ideas: (1) use of the Lanczos algorithm (Lanczos, 1950) with K iterations to find the K most important directions $\mathbf{V} \in \mathbb{R}^{d \times K}$, while adapting the Lanczos algorithm to only operate on Hessian-vector products; (2) importance sampling following a distribution $p_{\mathbf{v}}(\mathbf{v})$ to approximate the sum in Eq. (7) which reduces the number of back-propagations; and (3) computing these vectors with a frequency F instead of at every iteration. Those ideas are summarized in Alg 1 and we discuss details next.

Lanczos Algorithm: The Lanczos algorithm (Lanczos, 1950) computes the spectrum of a symmetric matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ by first reducing it to a tridiagonal form $\mathbf{T}_K \in \mathbb{R}^{K \times K}$ and then computing the spectrum $\mathbf{S}_K \in \mathbb{R}^{K \times K}$ of that matrix instead. The computation of the spectrum of a tridiagonal matrix is very efficient. The algorithm works in steps $k \in \{1, \dots, K\}$ and progressively builds an adapted orthonormal basis $\mathbf{Q}_k = [\mathbf{q}_0, \dots, \mathbf{q}_k] \in \mathbb{R}^{d \times k}$, that satisfies at each iteration $\mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k = \mathbf{T}_k$. Each of the k iterations of the algorithm requires a single Hessian-vector multiplication, that we compute implicitly in two steps: (1) First, we compute $s_m(\mathbf{x}; \theta) = \nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \theta)$; then (2) the derivative of the vector product $\nabla_{\mathbf{x}} (\mathbf{v}^T s_m(\mathbf{x}; \theta))$. Overall, for K steps, $K + 1$ backpropagations are needed. The Lanczos algorithm is memory efficient as it only stores three vectors at every iteration k , i.e., \mathbf{q}_{k-1} , \mathbf{q}_k and \mathbf{q}_{k+1} . The eigenvectors \mathbf{V}^* of the matrix \mathbf{A} are then obtained via

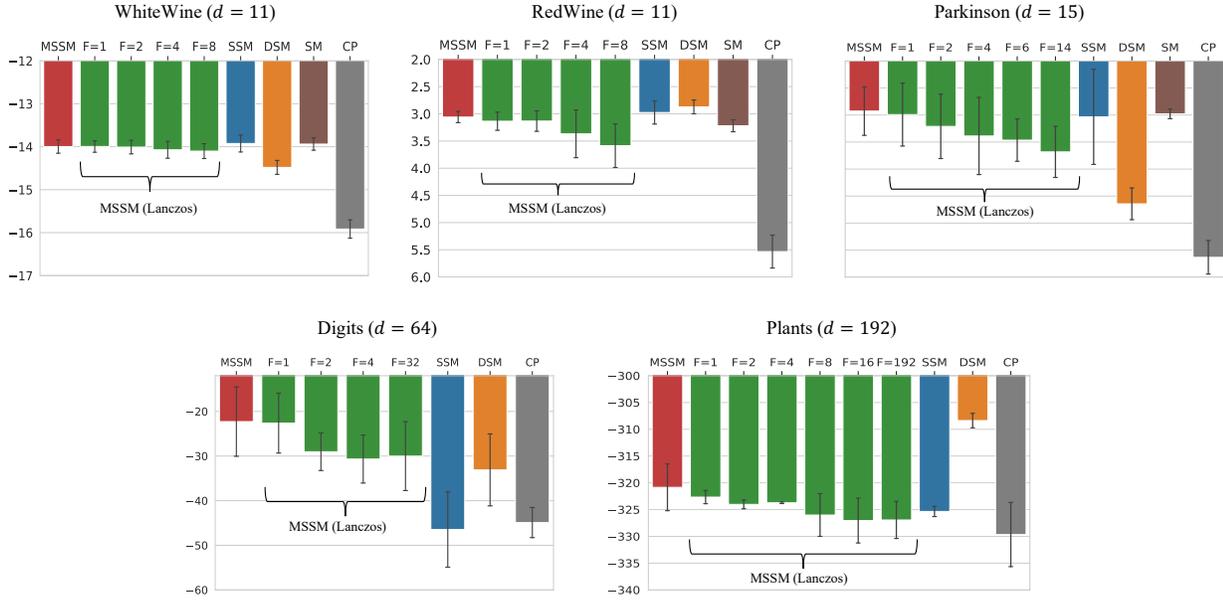


Figure 2. Log-likelihoods after training DKEF models on UCI datasets with different loss functions. Higher is better.

$\mathbf{V}^* = \mathbf{Q}_K \mathbf{S}_K \in \mathbb{R}^{d \times K}$, as illustrated in Alg. 2. Note that we use \mathbf{V}^* to denote the optimal directions obtained when using K steps of the Lanczos algorithm.

Sampling: After solving for the eigenvectors \mathbf{V}^* of $\nabla_{\mathbf{x}} s_m(\mathbf{x}; \boldsymbol{\theta})$ obtained from Alg. 2, we use importance sampling to avoid computing the sum over the projection directions in Eq. (7), as this requires K back-propagations. Instead we address the following problem,

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{\hat{p}_{\mathbf{v}}} \left[\frac{\mathbf{v}^T \nabla_{\mathbf{x}} s_m(\mathbf{x}; \boldsymbol{\theta}) \mathbf{v}}{\hat{p}_{\mathbf{v}}(\mathbf{v})} \right] + \frac{1}{2} \|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|^2. \quad (8)$$

We choose $\hat{p}_{\mathbf{v}}(\mathbf{v})$ to be the uniform distribution over the set of the computed eigenvectors \mathbf{V}^* . Note that the distribution $\hat{p}_{\mathbf{v}}$ differs from the distribution $p_{\mathbf{v}}$ over the space of all possible vectors \mathbf{v} which is used for sliced score matching in Eq. (4).

Overall the method requires $M \frac{K}{F}$ backpropagations. In practice, we chose $M = 1$. Here, M is the number of sampled vectors to approximate the expectation over \mathbf{v} in Eq. (8).

4. Experiments

We study the MSSM approach on estimating (1) densities (Sec. 4.1) and (2) scores for modeling implicit distributions (Sec. 4.2). Following Song et al. (2019), we set the number of vectors M to be 1 in Eq. (8).

4.1. Density Estimation

We consider two density estimation experiments: (1) estimating the parameters of an energy-based model, *i.e.*, ‘deep kernel exponential families’ (DKEF) and (2) learning a deepflow model.

Baselines: We use the following baselines: (1) Score matching (SM) (Hyvärinen, 2005), (2) Sliced score matching with reduced variance (SSM) and a multivariate Rademacher distribution $p_{\mathbf{v}}$ (Song et al., 2019), (3) Denoising score matching (DSM) (Vincent, 2011), (4) Approximate back-propagation (approx. BP) (Kingma & LeCun, 2010) and (5) Curvature propagation (CP) (Martens et al., 2012). We report results of MSSM with exact eigenvector computation and with an approximate one using the Lanczos algorithm.

Metric: As evaluation metric we use maximum likelihood which directly relates to the score matching loss, as explained in Appendix D.

1) Deep Kernel Exponential Families (DKEF): DKEF parameterizes the unnormalized log-likelihood via $\log p_f(\mathbf{x}) = \sum_{l=1}^L \alpha_l k(\mathbf{x}, \mathbf{z}_l) + \log q_0(\mathbf{x})$, where $q_0(\mathbf{x})$ is a fixed function, α_l are weights, and $k(\mathbf{x}, \mathbf{z}_l)$ is a mixture of R Gaussian kernels defined on the feature space ϕ_r of a deep net and evaluated at L different inducing points \mathbf{z}_l :

$$k(\mathbf{x}, \mathbf{z}_l) = \sum_{r=1}^R \rho_r \exp \left(\frac{-1}{2\sigma_r^2} \|\phi_r(\mathbf{x}) - \phi_r(\mathbf{z}_l)\|^2 \right). \quad (9)$$

The training parameters include the inducing points \mathbf{z}_l , the deep net ϕ_r , length scales σ_r and the non-negative mixture coefficients ρ_r . A closed-form expression for the weights α_l

Table 1. Log-likelihoods for NICE models trained on MNIST.

MLE	SSM	SSM-Staged	DSM ($\sigma = 0.1$)	DSM ($\sigma = 1.74$)	CP	Approx. BP	MSSM (Exact)	MSSM (Laczos, F=100)	MSSM (Laczos, F=200)	MSSM (Laczos, F=500)
-791	-3355	-1102	-4363	-8082	-1517	-2288	-1044	-1067	-1085	-1235

can be derived as a function of the other model parameters (see Appendix E for training details). We follow the same setup as Song et al. (2019). All models are trained with 15 different random seeds and training is stopped when the validation loss does not improve for 200 steps.

The obtained log-likelihoods on RedWine ($d = 11$), WhiteWine ($d = 11$), Parkinsons ($d = 15$), Digits ($d = 64$) and Plants ($d = 192$) (Dua & Graff, 2017) datasets are presented in Fig. 2. Note that the likelihoods are estimated using AIS (Neal, 2001), using a zero mean Gaussian proposal distribution $\mathcal{N}(0, 2I)$ and 1,000,000 samples. We observe that MSSM performance is comparable to SSM for data sets with small dimensionality, *i.e.*, RedWine, WhiteWine and Parkinsons. MSSM however scales better to the high-dimensional datasets, *i.e.*, Digits and Plants. Exact MSSM is consistently better than the Lanczos version. This is expected as we limit the number of iterations of the Lanczos algorithm to K , which may lead to smaller eigenvalues not converging. As we increase the frequency F of the eigenvector computation, MSSM (Laczos) performance decreases. This is expected as well, since the computed eigenvectors are no longer accurate and SSM ends-up sampling more diverse directions in the space. Our MSSM models, both exact and Lanczos-based with small update frequencies F , have lower variance than SSM for RedWine, WhiteWine, Parkinsons and Digits datasets. DSM performance strongly depends on the choice of the noise rate σ , which is hard to tune. Similarly to work by Song et al. (2019), we run a grid search over a range of values. CP performance is constantly worse as it injects noise to each node in the computation graph. This is not efficient for training large deep net models. Computing the exact score matching loss becomes already prohibitively slow for the digits and plants datasets. We omit the results for approx. BP, as the log-likelihoods were smaller than -10^6 for all datasets.

2) Deep Flow Models: Since Deepflow models enable access to a tractable maximum likelihood, we train a NICE (Dinh et al., 2014) model to better understand the relationship between MSSM, SSM and MLE. For this we consider MNIST data ($d = 728$) generation. Results are reported in Tab. 1. Since score matching is a weak form of MLE, as explained in Appendix D, we obtain better maximum likelihood when applying score matching loss on the intermediate layers of the flow model (SSM-staged). When trained with the SSM loss applied only to the last layer, the log-likelihood fluctuates and the best achieved performance

Table 2. Log-likelihoods on MNIST, estimated by AIS.

Latent Dim	VAE		WAE	
	8	32	8	32
ELBO	-96.87	-89.06	N/A	N/A
SSM	-95.50	-89.25	-98.24	-90.37
Stein	-96.71	-91.84	-99.05	-91.70
Spectral	-96.60	-94.67	-98.81	-92.55
MSSM (Exact)	-95.61	-89.26	-97.62	-89.53
MSSM (Laczos)	-95.75	-88.99	-98.25	-89.46

doesn’t correspond to the one at the end of training. Indeed, Song et al. (2019) report results from the best checkpoint. To further improve upon SSM-staged, we replace the SSM loss with MSSM at the output layer. We observe that the log-likelihood training curve has a monotonically increasing behavior. We improve upon SSM-staged at the end of the training when using MSSM (exact) or MSSM with eigenvectors computed every 100 or 200 iterations using the Lanczos algorithm.

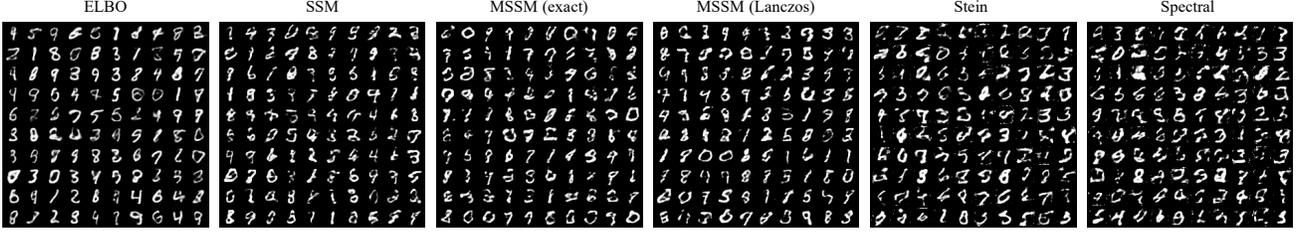
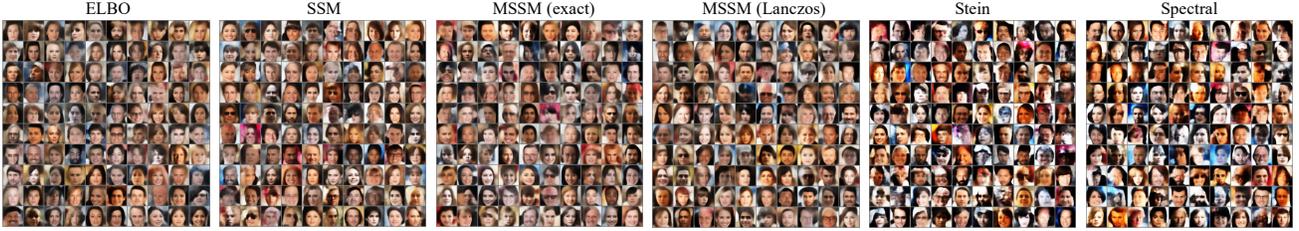
4.2. Score Estimation

We consider replacing simple explicit encoder distributions in Variational Auto-Encoders (VAE) and Wasserstein Auto-Encoders (WAE), with more expressive implicit ones that we train with an additional score matching loss. As baselines, we consider *Stein* (Li & Turner, 2018) and *Spectral* (Shi et al., 2018) score estimation techniques.

VAE with Implicit Encoders: A classical VAE is trained by minimizing the ELBO loss:

$$\min_{\phi, \theta} \mathbb{E}_{p_d(\mathbf{x})} \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [-\log p_\theta(\mathbf{x}|\mathbf{z}) - \log p(\mathbf{z}) + \log q_\phi(\mathbf{z}|\mathbf{x})]. \quad (10)$$

Here, $p(\mathbf{z})$ is the prior distribution, $q_\phi(\mathbf{z}|\mathbf{x})$ is the encoder distribution and $p_\theta(\mathbf{x}|\mathbf{z})$ is the decoder distribution. In order to compute a closed form of the loss, $q_\phi(\mathbf{z}|\mathbf{x})$ is usually chosen to be a simple explicit distribution (*e.g.*, a Gaussian). This obviously limits the expressivity of the model. Extending ELBO with a score matching loss permits to model arbitrarily complex distributions $q_\phi(\mathbf{z}|\mathbf{x})$ implicitly. This better captures the multi-modality and intricacies of common dataset. The latter is achieved by replacing the entropy term $-\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} \log q_\phi(\mathbf{z}|\mathbf{x})$ with the term $-\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\nabla_{\mathbf{z}} \log q_\phi(\mathbf{z}|\mathbf{x}) (\mathbf{z}|\mathbf{x})^T \mathbf{z}]$. It can be easily checked that both terms have the same derivative with


 Figure 3. VAE generated samples on MNIST for a latent dimension of 32 for z .

 Figure 4. VAE generated samples on CelebA for a latent dimension of 32 for z .

respect to ϕ , *i.e.*,

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z|\mathbf{x})} [\log q_{\phi}(z|\mathbf{x})] = \mathbb{E}_{q_{\phi}(z|\mathbf{x})} \left[\nabla_z \log q_{\phi}(z|\mathbf{x})^T \nabla_{\phi} z \right].$$

The gradient $\nabla_z \log q_{\phi}(z|\mathbf{x})$ is approximated with a deepnet $s_{\psi}(z|\mathbf{x})$ using a score matching loss, circumventing the computation of $q_{\phi}(z|\mathbf{x})$. The objective from Eq. (10) becomes:

$$\begin{aligned} & \min_{\theta, \phi, \psi} \max_{\mathbf{V}} \mathbb{E}_{p_d(\mathbf{x})} \mathbb{E}_{q_{\phi}(z|\mathbf{x})} [-\log p_{\theta}(\mathbf{x}|z) - \log p(z) + \\ & s_{\psi}(z|\mathbf{x})^T z + \sum_{i=1}^d \mathbf{v}_i^T \nabla_z s_{\psi}(z|\mathbf{x}) \mathbf{v}_i + \frac{1}{2} \|s_{\psi}(z|\mathbf{x})\|^2], \\ & \text{s.t. } \mathbf{V}\mathbf{V}^T = \mathbf{I}. \end{aligned}$$

We report negative log-likelihoods on MNIST, as estimated by AIS (Neal, 2001) with 1000 intermediate distributions, in Tab. 2 for latent dimensions of size 8 and 32. We present samples in Fig. 3. In Fig. 5, we present results for different Lanczos update frequencies. We observe MSSM to perform comparably to SSM. Both MSSM and SSM outperform the ELBO, which emphasizes the expressive power of an implicit encoder $q_{\phi}(z|\mathbf{x})$.

In Tab. 3, we evaluate sample quality on CelebA with FID scores (Heusel et al., 2017) computed after 10k, 40k, 70k and 100k iterations respectively. We observe that MSSM and its variants all converge faster, having low FID scores after 10k iterations already. In contrast, SSM does not perform as well early during training. Three MSSM variants outperform the ELBO loss. Samples are presented in Fig. 4.

WAE with Implicit Encoders: Unlike VAE which computes a KL-divergence between the encoder distribution $q_{\phi}(z|\mathbf{x})$ of each sample and the prior distribution $p(z)$,

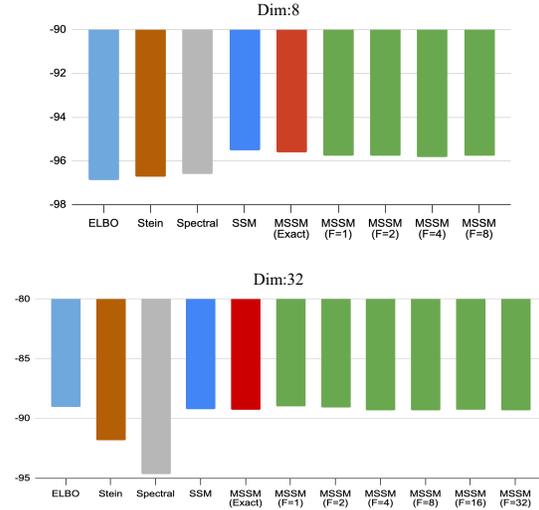


Figure 5. Log-likelihoods for VAE trained on MNIST.

WAE compares the distribution of the entire encoder space $q_{\phi}(z)$ to the prior $p(z)$. This alleviates the issue of erroneous reconstruction due to overlapping latent codes z originating from different samples \mathbf{x} . It also reduces the prior-hole problem.

Concretely, WAE minimizes the objective:

$$\min_{\theta, \phi} \mathbb{E}_{p_d(\mathbf{x})} \mathbb{E}_{q_{\phi}(z)} [c(\mathbf{x}, p_{\theta}(\mathbf{x}|z)) - \lambda \log p(z) - \lambda q_{\phi}(z)],$$

composed of the reconstruction loss $c(\mathbf{x}, p_{\theta}(\mathbf{x}|z))$ and the KL-divergence between the prior $p(z)$ and the encoder distribution $q_{\phi}(z|\mathbf{x})$. We model $q_{\phi}(z)$ implicitly. Following the VAE case we approximate $\nabla_z q_{\phi}(z)$ using a deepnet

Max-Sliced Score Matching

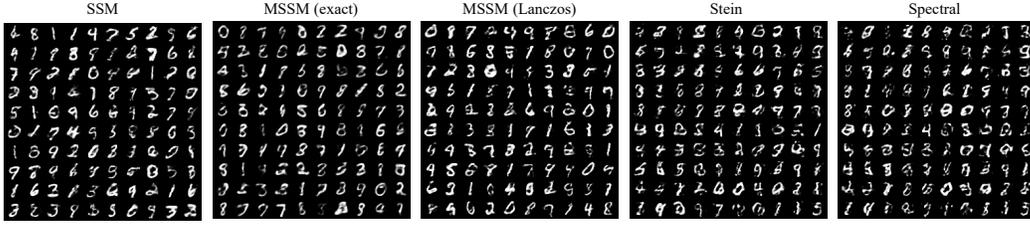


Figure 6. WAE generated samples on MNIST for a latent dimension of 32 for z .

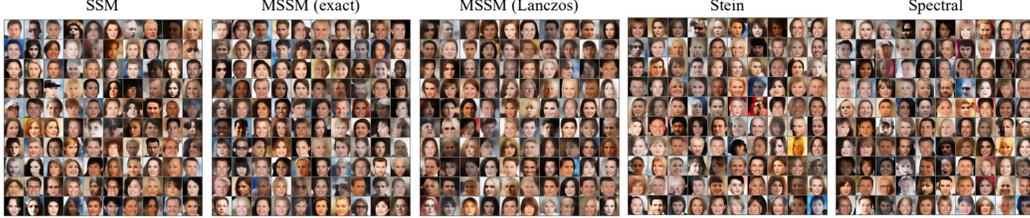


Figure 7. WAE generated samples on CelebA for a latent dimension of 32 for z .

Table 3. FID scores of different methods versus number of training iterations on the CelebA dataset.

		10k	40k	70k	100k
VAE	ELBO	96.20	73.70	69.42	66.32
	SSM	108.52	70.28	66.52	62.50
	Stein	126.60	118.87	120.51	126.76
	Spectral	131.90	125.04	128.36	133.93
	MSSM (exact)	100.09	69.17	64.37	60.68
	MSSM (Lanczos, F=1)	100.26	64.14	65.45	60.78
	MSSM (Lanczos, F=2)	98.38	65.34	64.91	62.46
	MSSM (Lanczos, F=8)	101.15	65.34	65.78	63.76
	MSSM (Lanczos, F=16)	100.65	64.71	62.68	63.33
	MSSM (Lanczos, F=32)	107.24	66.30	65.31	65.91
WAE	SSM	84.11	61.09	56.23	54.33
	Stein	82.93	63.46	58.53	57.61
	Spectral	82.30	62.47	58.03	55.96
	MSSM (exact)	87.12	60.93	56.17	53.98
	MSSM (Lanczos, F=1)	84.08	71.68	57.29	54.03
	MSSM (Lanczos, F=2)	82.33	71.00	60.66	54.23
	MSSM (Lanczos, F=8)	82.92	73.37	59.48	54.98
	MSSM (Lanczos, F=16)	82.21	73.48	59.62	58.49
	MSSM (Lanczos, F=32)	83.75	70.54	59.37	59.27

$s_\psi(z)$, resulting in the objective

$$\min_{\theta, \phi, \psi} \max_{\mathbf{V}} \mathbb{E}_{p_d(\mathbf{x})} \mathbb{E}_{q_\phi(z)} [c(\mathbf{x}, p_\theta(\mathbf{x}|z)) - \lambda \log p(z) - \lambda s_\psi(z)z + \mathbf{v}^T \nabla_z s_\psi(z)\mathbf{v} + \frac{1}{2} \|s_\psi(z)\|^2],$$

s.t. $\mathbf{V}\mathbf{V}^T = \mathbf{I}$.

On MNIST, MSSM achieves the best log-likelihoods for both latent dimensions across all baselines (Tab. 2, Fig. 8), as well as the best FID scores for CelebA data (Tab. 3). We visualize generated samples obtained from training on MNIST data in Fig. 6. Samples after training on CelebA are given in Fig. 7.

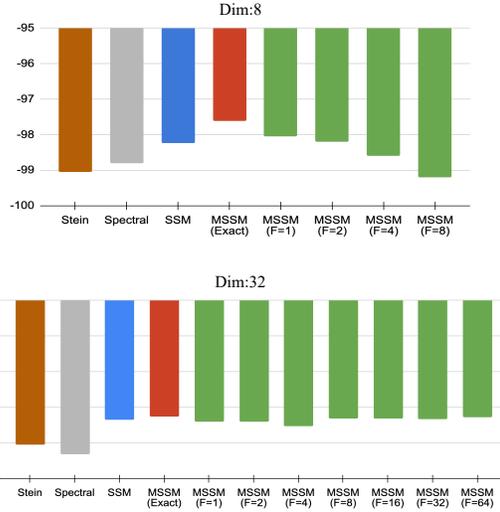


Figure 8. Log-likelihoods for WAE trained on MNIST.

5. Conclusion

We propose max-sliced score matching (MSSM) which improves upon sliced score matching (SSM) by finding the most informative projection directions for minimizing the difference between the data and model scores. We show that these directions correspond to the eigenvectors of the Hessian of the log-likelihood, which can be computed efficiently using the Lanczos algorithm, importance sampling to approximate the sum over the eigenvectors as well as a schedule for updating the eigenvectors at specific frequencies. We show improved performance on learning unnormalized statistical models and estimating scores for implicit distributions in VAEs and WAEs.

References

Bao, F., Li, C., Xu, T., Su, H., Zhu, J., and Zhang, B. Bi-level score matching for learning energy-based latent

- variable models. 2020.
- Besag, J. Statistical analysis of non-lattice data. *J. R. Stat. Soc.*, 1975.
- Cai, R., Yang, G., Averbuch-Elor, H., Hao, Z., Belongie, S., Snavely, N., and Hariharan, B. Learning gradient fields for shape generation. *arXiv preprint arXiv:2008.06520*, 2020.
- Chen, N., Zhang, Y., Zen, H., Weiss, R. J., Norouzi, M., and Chan, W. Wavegrad: Estimating gradients for waveform generation. *arXiv preprint arXiv:2009.00713*, 2020.
- Dinh, L., Krueger, D., and Bengio, Y. Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.
- Du, Y. and Mordatch, I. Implicit generation and generalization in energy-based models. *arXiv preprint arXiv:1903.08689*, 2019.
- Dua, D. and Graff, C. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- Gutmann, M. and Hyvärinen, A. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Procs. AISTATS*, 2010.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *arXiv preprint arXiv:1706.08500*, 2017.
- Hinton, G. E. Training products of experts by minimizing contrastive divergence. *Neural computation*, 2002.
- Hutchinson, M. F. A stochastic estimator of the trace of the influence matrix for laplacian smoothing splines. *Commun. Stat. Simul. Comput.*, 1989.
- Hyvärinen, A. Estimation of non-normalized statistical models by score matching. *JMLR*, 2005.
- Hyvarinen, A. Connections between score matching, contrastive divergence, and pseudolikelihood for continuous-valued variables. *IEEE Trans. Neural Netw. Learn. Syst.*, 2007.
- Jolicoeur-Martineau, A., Piché-Taillefer, R., Combes, R. T. d., and Mitliagkas, I. Adversarial score matching and improved sampling for image generation. 2020.
- Kingma, D. P. and LeCun, Y. Regularized estimation of image statistics by score matching. In *Procs. NeurIPS*, 2010.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. In *Procs. ICLR*, 2014.
- Lanczos, C. *An iteration method for the solution of the eigenvalue problem of linear differential and integral operators*. 1950.
- Li, Y. and Turner, R. E. Gradient estimators for implicit models. In *Procs. ICLR*, 2018.
- Mahmood, A., Oliva, J., and Styner, M. A. Multiscale score matching for out-of-distribution detection. *arXiv preprint arXiv:2010.13132*, 2020.
- Martens, J., Sutskever, I., and Swersky, K. Estimating the hessian by back-propagating curvature. In *Procs. ICML*, 2012.
- Meng, C., Yu, L., Song, Y., Song, J., and Ermon, S. Autoregressive score matching. 2020.
- Neal, R. M. Annealed importance sampling. *Statistics and computing*, 2001.
- Newey, W. K. and McFadden, D. Large sample estimation and hypothesis testing. *Handbook of econometrics*, 1994.
- Pang, T., Xu, T., Li, C., Song, Y., Ermon, S., and Zhu, J. Efficient learning of generative models via finite-difference score matching. 2020.
- Shi, J., Sun, S., and Zhu, J. A spectral approach to gradient estimation for implicit distributions. In *Procs. ICML*, 2018.
- Song, Y. and Ermon, S. Generative modeling by estimating gradients of the data distribution. In *Procs. NeurIPS*, 2019a.
- Song, Y. and Ermon, S. Generative modeling by estimating gradients of the data distribution. *arXiv preprint arXiv:1907.05600*, 2019b.
- Song, Y., Garg, S., Shi, J., and Ermon, S. Sliced score matching: A scalable approach to density and score estimation. *UAI*, 2019.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- Tieleman, T. and Hinton, G. Using fast weights to improve persistent contrastive divergence. In *Procs. ICML*, 2009.
- Tolstikhin, I., Bousquet, O., Gelly, S., and Schoelkopf, B. Wasserstein auto-encoders. In *Procs. ICLR*, 2018.
- Vincent, P. A connection between score matching and denoising autoencoders. *NC*, 2011.
- Wenliang, L., Sutherland, D., Strathmann, H., and Gretton, A. Learning deep kernels for exponential family densities. In *Procs. ICML*, 2019.

A. Appendix: Max-Sliced Score Matching

Sliced score matching (SSM) enabled scaling score matching (SM) to more complex problems and high-dimensional datasets. The main idea: instead of directly matching the high-dimensional scores, we match their projections along random directions, *i.e.*,

$$\min_{\theta} \mathbb{E}_{\mathbf{v} \sim p_{\mathbf{v}}} \mathbb{E}_{\mathbf{x} \sim p_d} [(\mathbf{v}^T \mathbf{s}_m(\mathbf{x}; \theta) - \mathbf{v}^T \mathbf{s}_d(\mathbf{x}))^2]. \quad (11)$$

Here, $p_d(\mathbf{x})$ and $p_m(\mathbf{x}; \theta)$ denote the data and the model distributions respectively. Moreover, $\mathbf{s}_d(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_d(\mathbf{x})$ and $\mathbf{s}_m(\mathbf{x}; \theta) = \nabla_{\mathbf{x}} \log p(\mathbf{x}; \theta)$ are the corresponding scores. Currently, high variance is the major drawback for the SSM loss. In order to address this, we propose Max-Sliced Score Matching (MSSM), a variant of SSM that chooses the most informative projection directions $\mathbf{V} \in \mathbb{R}^{d \times K}$ instead of averaging over random ones. We start our derivation by applying the max-slice to the Fisher Divergence, *i.e.*,

$$\begin{aligned} \min_{\theta} \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} [\|\mathbf{V}^T (\mathbf{s}_m(\mathbf{x}; \theta) - \mathbf{s}_d(\mathbf{x}))\|^2], \\ \text{s.t. } \mathbf{V}\mathbf{V}^T = \mathbf{I}. \end{aligned} \quad (12)$$

Since the score of the data distribution $\mathbf{s}_d(\mathbf{x})$ is unknown, our MSSM optimizes an equivalent loss that does not depend on that data score:

$$\begin{aligned} \min_{\theta} \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} \left[\sum_{i=1}^K \mathbf{v}_i^T \nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}; \theta) \mathbf{v}_i + \frac{1}{2} \|\mathbf{s}_m(\mathbf{x}; \theta)\|^2 \right], \\ \text{s.t. } \mathbf{V}\mathbf{V}^T = \mathbf{I}, \end{aligned} \quad (13)$$

where \mathbf{v}_i is the i -th column vector of matrix \mathbf{V} . Note, Eq. (12) and Eq. (13) are Eq. (6) and Eq. (7) in the main text, but we define them again here to avoid cross-referencing across documents.

In the following, we denote by $\mathcal{L}(\theta)$ the MSSM loss from Eq. (13), *i.e.*,

$$\mathcal{L}(\theta) = \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} [\|\mathbf{V}^T (\mathbf{s}_m(\mathbf{x}; \theta) - \mathbf{s}_d(\mathbf{x}))\|^2], \quad (14)$$

$J(\theta)$ the MSSM loss from Eq. (12), *i.e.*,

$$J(\theta) = \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} \left[\sum_{i=1}^K \mathbf{v}_i^T \nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}; \theta) \mathbf{v}_i + \frac{1}{2} \|\mathbf{s}_m(\mathbf{x}; \theta)\|^2 \right], \quad (15)$$

and by $J_{N,K}(\theta)$ the loss estimator resulting from empirically evaluating $J(\theta)$ on a batch size N using K vectors $\{\mathbf{v}_i\}_{i=1}^K$, *i.e.*,

$$J_{N,K}(\theta) = \max_{\mathbf{V}} \frac{1}{N} \sum_{j=1}^N \sum_{i=1}^K \mathbf{v}_i^T \nabla_{\mathbf{x}^{(j)}} \mathbf{s}_m(\mathbf{x}^{(j)}; \theta) \mathbf{v}_i + \frac{1}{2} \|\mathbf{s}_m(\mathbf{x}^{(j)}; \theta)\|^2. \quad (16)$$

Here, \mathbf{v}_i is the i -th column vector of \mathbf{V} . We also assume that $p_m(\mathbf{x}; \theta)$ is well specified, and that θ^* is the parameter of the data distribution, *i.e.* $p_d(\mathbf{x}) = p_m(\mathbf{x}; \theta^*)$. We denote by θ , $\hat{\theta}$ and $\hat{\theta}_{NK}$ the estimates resulting from minimizing $\mathcal{L}(\theta)$, $J(\theta)$ and $J_{NK}(\theta)$ under the constraint $\mathbf{V}\mathbf{V}^T = \mathbf{I}$, respectively, *i.e.*,

$$\theta = \operatorname{argmin}_{\theta} \mathcal{L}(\theta) \quad \text{s.t. } \mathbf{V}\mathbf{V}^T = \mathbf{I}, \quad (17)$$

$$\hat{\theta} = \operatorname{argmin}_{\theta} J(\theta) \quad \text{s.t. } \mathbf{V}\mathbf{V}^T = \mathbf{I}, \quad (18)$$

$$\hat{\theta}_{NK} = \operatorname{argmin}_{\theta} J_{NK}(\theta) \quad \text{s.t. } \mathbf{V}\mathbf{V}^T = \mathbf{I}. \quad (19)$$

We also define $f(\theta; \mathbf{x}, \mathbf{V}^*)$ as follows:

$$f(\theta; \mathbf{x}, \mathbf{V}^*) = \sum_{i=1}^K \mathbf{v}_i^{*T} \nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}; \theta) \mathbf{v}_i^* + \frac{1}{2} \|\mathbf{s}_m(\mathbf{x}; \theta)\|^2. \quad (20)$$

Here, \mathbf{V}^* is the optimal basis obtained from solving the inner optimization program given in Eq. (12). Hence, $J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim p_d}[f(\boldsymbol{\theta}; \mathbf{x}, \mathbf{V}^*)]$ and $J_{NK}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N [f(\boldsymbol{\theta}; \mathbf{x}^{(i)}, \mathbf{V}^*)]$.

In Appendix B, we prove consistency and asymptotic-normality of the MSSM objective. In Appendix C, we derive an equivalent version of the MSSM objective that does not depend on the unknown data score $s_d(\mathbf{x})$. We derive the relationship between MSSM loss and MLE loss in Appendix D. Further training details are presented in Appendix E.

B. Consistency and Asymptotic Normality of MSSM

In the following we show consistency and asymptotic normality. This is important to assess the convergence and asymptotic variance of the estimated parameters around the true parameters.

B.1. Consistency

We will show that $\hat{\boldsymbol{\theta}}_{NK}$ is consistent, *i.e.* as the sample size N increases, the sampling distribution of the estimator $\hat{\boldsymbol{\theta}}_{NK}$ becomes increasingly concentrated at the true parameter value $\boldsymbol{\theta}^*$. Formally, we will prove that for every $K \leq d$ with $K \in \mathbb{N}$,

$$\hat{\boldsymbol{\theta}}_{NK} \xrightarrow{P} \boldsymbol{\theta}^*, \quad N \rightarrow \infty.$$

Here, ' \xrightarrow{P} ' denotes convergence in probability, *i.e.* $\lim_{n \rightarrow \infty} P(|\hat{\boldsymbol{\theta}}_{NK} - \boldsymbol{\theta}^*| > \epsilon) = 0, \forall \epsilon > 0$. The proof follows two steps:

1. First, we prove that $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}^*$, in **Lemma 1** and **Theorem 1**.
2. Second, in **Theorem 2**, we prove that $\hat{\boldsymbol{\theta}}_{NK} \xrightarrow{P} \hat{\boldsymbol{\theta}}$ when $N \rightarrow \infty$. For this, we show that under the assumptions: (a) the parameter space Θ is compact, (b) $\|s_m(\mathbf{x}; \boldsymbol{\theta})\| < \infty$ and (c) $\|\nabla_{\mathbf{x}} \log s_m(\mathbf{x}; \boldsymbol{\theta})\|_{\infty} < \infty$, we have $|f(\boldsymbol{\theta}; \mathbf{x}, \mathbf{V}^*)| < \infty$. As a result, the conditions in Lemma 2.4 (Newey & McFadden, 1994) are satisfied. So, $\sup_{\boldsymbol{\theta}} |J_{NK}(\boldsymbol{\theta}) - J(\boldsymbol{\theta})| \xrightarrow{P} 0$. Thus, from Theorem 2.1 (Newey & McFadden, 1994), we conclude $\hat{\boldsymbol{\theta}}_{NK} \xrightarrow{P} \hat{\boldsymbol{\theta}}$.

Lemma 1 Assume that $p_m(\mathbf{x}; \boldsymbol{\theta})$ is well specified, *i.e.*, $p_d(\mathbf{x}) = p_m(\mathbf{x}; \boldsymbol{\theta}^*)$, and $p_m(\mathbf{x}; \boldsymbol{\theta}) \neq p_m(\mathbf{x}; \boldsymbol{\theta}^*)$ whenever $\boldsymbol{\theta} \neq \boldsymbol{\theta}^*$. In the following, we show that, when Eq. (12) converges to 0, $\boldsymbol{\theta}$ converges to $\boldsymbol{\theta}^*$, *i.e.*,

$$\max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} [\|\mathbf{V}^T s_m(\mathbf{x}; \boldsymbol{\theta}) - \mathbf{V}^T s_d(\mathbf{x})\|^2] = 0 \iff \boldsymbol{\theta} = \boldsymbol{\theta}^*.$$

Proof: The proof is based on the following sequence of equivalences:

$$\begin{aligned} & \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} [\|\mathbf{V}^T s_m(\mathbf{x}; \boldsymbol{\theta}) - \mathbf{V}^T s_d(\mathbf{x})\|^2] = 0, \quad \text{s.t. } \mathbf{V}\mathbf{V}^T = I \\ \iff & \max_{\mathbf{V}} (s_m(\mathbf{x}; \boldsymbol{\theta}) - s_d(\mathbf{x}))^T \mathbf{V}\mathbf{V}^T (s_m(\mathbf{x}; \boldsymbol{\theta}) - s_d(\mathbf{x})) = 0, \quad \forall \mathbf{x} \sim p_d(\mathbf{x}), \quad \text{s.t. } \mathbf{V}\mathbf{V}^T = I \\ \stackrel{(i)}{\iff} & \|s_m(\mathbf{x}; \boldsymbol{\theta}) - s_d(\mathbf{x})\|^2 = 0, \quad \forall \mathbf{x} \sim p_d(\mathbf{x}) \\ \iff & s_m(\mathbf{x}; \boldsymbol{\theta}) = s_d(\mathbf{x}), \quad \forall \mathbf{x} \sim p_d(\mathbf{x}) \\ \iff & \log p_m(\mathbf{x}; \boldsymbol{\theta}) = \log p_d(\mathbf{x}) + \text{Const.}, \quad \forall \mathbf{x} \sim p_d(\mathbf{x}) \\ \stackrel{(ii)}{\iff} & p_d(\mathbf{x}) = p_m(\mathbf{x}; \boldsymbol{\theta}), \quad \forall \mathbf{x} \sim p_d(\mathbf{x}) \end{aligned} \tag{21}$$

Note, (i) holds because $\mathbf{V}\mathbf{V}^T = I$. Further, (ii) holds as $p_m(\mathbf{x}, \boldsymbol{\theta})$ and $p_d(\mathbf{x})$ are normalized probability density functions. Finally, $\boldsymbol{\theta} = \boldsymbol{\theta}^*$ is a simple proof by contradiction. Suppose $\boldsymbol{\theta} \neq \boldsymbol{\theta}^*$ is true, then it implies $p_d(\mathbf{x}) \neq p_m(\mathbf{x}; \boldsymbol{\theta}^*)$ according to our assumption. This is not true, so $\boldsymbol{\theta} = \boldsymbol{\theta}^*$.

Theorem 1 Assume the results from Lemma 1. Given the equivalence Eq. (12) \iff Eq. (13) (Appendix C), the equality $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}^*$ holds.

Proof: In Lemma 1, we proved that $\boldsymbol{\theta} = \boldsymbol{\theta}^*$ holds. Besides, in Appendix C, we show that the programs in Eq. (12) and Eq. (13) are equivalent. This results in the equivalence $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}$. As a result, $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}^*$.

Theorem 2 Suppose that the parameter space Θ is compact, $\|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\| < \infty$, $\|\nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|_{\infty} < \infty$, and $f(\boldsymbol{\theta}; \mathbf{x}, \mathbf{V}^*)$ is continuous. Then, $\hat{\boldsymbol{\theta}}_{NK} \xrightarrow{p} \hat{\boldsymbol{\theta}}$.

Proof: We start by showing that $|f(\boldsymbol{\theta}; \mathbf{x}, \mathbf{V}^*)| < \infty$. Given that $\|\mathbf{v}_j\|^2 = 1 \forall j \in \{1, \dots, K\}$, then $|\mathbf{v}_j^{(k)} \mathbf{v}_j^{(l)}| < 1$, $\forall k, l \in \{1, \dots, d\}$. Since $\|\nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\| < \infty$, we conclude that $|\mathbf{v}_j^{*T} \nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) \mathbf{v}_j^*| < \infty$, $\forall j$. Besides, since $\|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|^2 < \infty$, we conclude that $|f(\boldsymbol{\theta}; \mathbf{x}, \mathbf{V}^*)| < \infty$. Hence, all the assumptions for Lemma 2.4 (Newey & McFadden, 1994) are satisfied. So, $\sup_{\boldsymbol{\theta}} |J_{NK}(\boldsymbol{\theta}) - J(\boldsymbol{\theta})| \xrightarrow{p} 0$. Thus, from Theorem 2.1 (Newey & McFadden, 1994), we conclude $\hat{\boldsymbol{\theta}}_{NK} \xrightarrow{p} \hat{\boldsymbol{\theta}}$.

B.2. Asymptotic Normality

We start by introducing some notation. We denote by \mathbf{H}_{ij} a matrix that depends on $p_m(\mathbf{x}; \boldsymbol{\theta})$ as follows:

$$\mathbf{H}_{ij} = \mathbb{E}_{p_d} [(\nabla_{\boldsymbol{\theta}} [\nabla_{\mathbf{x}}^2 \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_{ii} + \frac{1}{2} \nabla_{\boldsymbol{\theta}} ([\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_i)^2) (\nabla_{\boldsymbol{\theta}} [\nabla_{\mathbf{x}}^2 \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_{jj} + \frac{1}{2} \nabla_{\boldsymbol{\theta}} ([\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_j)^2)^T],$$

where, $[\nabla_{\mathbf{x}}^2 \log p_m(\mathbf{x}; \boldsymbol{\theta})]_{ii}$ is the ii -th entry of the matrix $\nabla_{\mathbf{x}}^2 p_m(\mathbf{x}; \boldsymbol{\theta})$ and $[\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta})]_i$ is the i -th entry of the vector $\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta})$. Also, we denote $\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*}$ and $J(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*}$ as $\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)$, and $J(\boldsymbol{\theta}^*)$, respectively.

In the following, we show that $\hat{\boldsymbol{\theta}}_{NK}$ is asymptotically normally distributed around the true parameter $\boldsymbol{\theta}^*$ when the batch size N becomes very large, *i.e.*,

$$\sqrt{N}(\hat{\boldsymbol{\theta}}_{NK} - \boldsymbol{\theta}^*) \xrightarrow{d} \mathcal{N}(0, \boldsymbol{\Sigma}).$$

Here ' \xrightarrow{d} ' denotes a convergence in distribution and $\boldsymbol{\Sigma}$ is the covariance matrix. Specifically, we prove the theorem:

Theorem 3 If (1) Θ is compact, (2) $\hat{\boldsymbol{\theta}}_{NK} \xrightarrow{p} \boldsymbol{\theta}^*$, (3) $\log p_m(\mathbf{x}; \boldsymbol{\theta})$ is twice continuously differentiable with respect to $\boldsymbol{\theta}$, (4) $J(\boldsymbol{\theta}^*)$ is invertible at $\boldsymbol{\theta}^*$, (5) $\|\nabla_{\boldsymbol{\theta}}^2 [\nabla_{\mathbf{x}}^2 \log p_m(\mathbf{x}; \boldsymbol{\theta})]_{ii}\|_{\infty} < \infty$, $\forall i \in \{1, \dots, d\}$ and (6) $\|\nabla_{\boldsymbol{\theta}}^2 \|\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta})\|^2\|_{\infty} < \infty$, then,

$$\sqrt{N}(\hat{\boldsymbol{\theta}}_{NK} - \boldsymbol{\theta}^*) \xrightarrow{d} \mathcal{N}\left(\mathbf{0}, [\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^*)]^{-1} \left(\sum_{i,j=1}^d \mathbf{H}_{ij} \right) [\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^*)]^{-1}\right).$$

The proof is based on Theorem 3.1 (Newey & McFadden, 1994). To apply it, we first prove:

1. $\sqrt{N} \nabla_{\boldsymbol{\theta}} J_{NK}(\boldsymbol{\theta}^*) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sum_{i,j=1}^d \mathbf{H}_{ij})$, in **Lemma 2**.
2. $\sup_{\boldsymbol{\theta}} |\nabla_{\boldsymbol{\theta}}^2 J_{NK}(\mathbf{x}; \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}}^2 J(\mathbf{x}; \boldsymbol{\theta})| \xrightarrow{p} \mathbf{0}$ in **Lemma 3**. Note that the absolute value is applied element wise to all the elements of the matrix.

Lemma 2 If $\log p_m(\mathbf{x}; \boldsymbol{\theta})$ is twice continuously differentiable with respect to $\boldsymbol{\theta}$, then

$$\sqrt{N} \nabla_{\boldsymbol{\theta}} J_{NK}(\boldsymbol{\theta}^*) \xrightarrow{d} \mathcal{N}\left(\mathbf{0}, \sum_{i,j=1}^d \mathbf{H}_{ij}\right).$$

Proof: The central limit theorem gives,

$$\sqrt{N} \sum_{i=1}^K \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^*; \mathbf{x}, \mathbf{V}^*) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \text{Var}[\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^*; \mathbf{x}, \mathbf{V}^*)]). \quad (22)$$

Note that $\mathbb{E}[\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}; \mathbf{x}, \mathbf{V})] = \mathbf{0}$ as $\boldsymbol{\theta}^*$ is the optimal parameter for $f(\boldsymbol{\theta}; \mathbf{x}, \mathbf{V}^*)$. We compute $\text{Var}[\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^*; \mathbf{x}, \mathbf{V}^*)]$ as

follows:

$$\begin{aligned}
 & \text{Var}[\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^*; \mathbf{x}, \mathbf{V}^*)] \\
 &= \mathbb{E}_{p_d}[\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^*; \mathbf{x}, \mathbf{V}^*) \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^*; \mathbf{x}, \mathbf{V}^*)^T] \\
 &= \mathbb{E}_{p_d}[(\sum_{i=1}^K \nabla_{\boldsymbol{\theta}} \mathbf{v}_i^T \nabla_{\mathbf{x}}^2 p_m(\mathbf{x}; \boldsymbol{\theta}^*) \mathbf{v}_i + \frac{1}{2} \nabla_{\boldsymbol{\theta}} \|\nabla_{\mathbf{x}} p_m(\mathbf{x}; \boldsymbol{\theta}^*)\|^2) (\sum_{j=1}^K \nabla_{\boldsymbol{\theta}} \mathbf{v}_j^T \nabla_{\mathbf{x}}^2 p_m(\mathbf{x}; \boldsymbol{\theta}^*) \mathbf{v}_j + \frac{1}{2} \nabla_{\boldsymbol{\theta}} \|\nabla_{\mathbf{x}} p_m(\mathbf{x}; \boldsymbol{\theta}^*)\|^2)^T] \\
 &= \mathbb{E}_{p_d}[(\sum_{i=1}^d \nabla_{\boldsymbol{\theta}} [\nabla_{\mathbf{x}}^2 \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_i + \frac{1}{2} \nabla_{\boldsymbol{\theta}} ([\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_i)^2) (\sum_{j=1}^d \nabla_{\boldsymbol{\theta}} [\nabla_{\mathbf{x}}^2 \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_j + \frac{1}{2} \nabla_{\boldsymbol{\theta}} ([\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_j)^2)^T] \\
 &= \sum_{i,j}^d H_{ij}
 \end{aligned}$$

Hence, $\sqrt{N} \nabla_{\boldsymbol{\theta}} J_{NK}(\boldsymbol{\theta}^*) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sum_{i,j=1}^K \mathbf{H}_{ij})$. \square

Lemma 3 Assuming Θ is compact and $p_m(\mathbf{x}; \boldsymbol{\theta})$ being twice continuously differentiable with respect to $\boldsymbol{\theta}$, then

$$\sup_{\boldsymbol{\theta}} |\nabla_{\boldsymbol{\theta}}^2 J_{NK}(\boldsymbol{\theta}^*) - \nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^*)| \xrightarrow{p} \mathbf{0}. \quad (23)$$

Proof: We write down the second order derivative of $f(\boldsymbol{\theta}; \mathbf{x}, \mathbf{V}^*)$ as:

$$\nabla_{\boldsymbol{\theta}}^2 f(\boldsymbol{\theta}^*; \mathbf{x}, \mathbf{V}^*) = \sum_{i=1}^K \nabla_{\boldsymbol{\theta}}^2 [\mathbf{v}_i^{*T} \nabla_{\mathbf{x}}^2 \log p_m(\mathbf{x}; \boldsymbol{\theta}^*) \mathbf{v}_i^*] + \frac{1}{2} \nabla_{\boldsymbol{\theta}}^2 \|\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)\|^2 \quad (24)$$

By assumption, we have $\|\nabla_{\boldsymbol{\theta}}^2 [\nabla_{\mathbf{x}}^2 \log p_m(\mathbf{x}; \boldsymbol{\theta})]_{ii}\|_{\infty} < \infty$ and $\|\nabla_{\boldsymbol{\theta}}^2 \|\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta})\|^2\|_{\infty} < \infty$. Hence, $\|\nabla_{\boldsymbol{\theta}}^2 f(\boldsymbol{\theta}^*; \mathbf{x}, \mathbf{V}^*)\| < \infty$. The assumptions for Lemma 2.4 (Newey & McFadden, 1994) are satisfied. So,

$$\sup_{\boldsymbol{\theta}} \left| \frac{1}{N} \sum_{i=1}^N \nabla_{\boldsymbol{\theta}}^2 f(\boldsymbol{\theta}^*; \mathbf{x}, \mathbf{V}^*) - \mathbb{E}_{\mathbf{x} \sim p_d} [\nabla_{\boldsymbol{\theta}}^2 f(\boldsymbol{\theta}^*; \mathbf{x}, \mathbf{V}^*)] \right| \xrightarrow{p} \mathbf{0}$$

which is equivalent to:

$$\sup_{\boldsymbol{\theta}} |\nabla_{\boldsymbol{\theta}}^2 J_{NK}(\boldsymbol{\theta}^*) - \nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^*)| \xrightarrow{p} \mathbf{0}. \square$$

Combining **Theorem 3** assumptions with **Lemma 2** and **Lemma 3** results, satisfies all the assumptions of Theorem 3.1 (Newey & McFadden, 1994). Hence,

$$\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \rightarrow \mathcal{N}\left(\mathbf{0}, (\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^*))^{-1} \left(\sum_{i,j=1}^d \mathbf{H}_{ij} \right) (\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^*))^{-1}\right). \quad (25)$$

Remark: As a result, the asymptotic variance of MSSM is the same as SM's asymptotic variance and smaller than the one for SSM. Proof are presented in Appendix B.4 by Song et al. (2019):

- SSM with $p_v \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$:

$$\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \rightarrow \mathcal{N}\left(\mathbf{0}, (\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^*; p_v))^{-1} \left(\sum_{i,j=1}^d H_{ij} + \frac{2}{M} \sum_{i=1}^d H_{ii} + \frac{2}{M} \sum_{\substack{i,j=1 \\ i \neq j}}^d W_{ij} \right) (\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^*; p_v))^{-1}\right) \quad (26)$$

Here,

$$\begin{aligned}
 W_{ij} = & \mathbb{E}_{p_d}[(\nabla_{\boldsymbol{\theta}} [\nabla_{\mathbf{x}}^2 \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_{ij} + \frac{1}{2} \nabla_{\boldsymbol{\theta}} ([\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_i [\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_j)) (\nabla_{\boldsymbol{\theta}} [\nabla_{\mathbf{x}}^2 \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_{ij} \\
 & + \frac{1}{2} \nabla_{\boldsymbol{\theta}} [\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_i [\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}^*)]_j)^T]. \quad (27)
 \end{aligned}$$

- **SSM** with $p_{\mathbf{v}}$ follow a Rademacher distribution:

$$\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \rightarrow \mathcal{N}\left(\mathbf{0}, (\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^*; p_{\mathbf{v}}))^{-1} \left(\sum_{i,j=1}^d H_{ij} + \frac{2}{M} \sum_{\substack{i,j=1 \\ i \neq j}}^d W_{ij} \right) (\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^*; p_{\mathbf{v}}))^{-1} \right) \quad (28)$$

- **SM** (Eq. (5)):

$$\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \rightarrow \mathcal{N}\left(\mathbf{0}, (\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^*; p_{\mathbf{v}}))^{-1} \left(\sum_{i,j=1}^d H_{ij} \right) (\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^*; p_{\mathbf{v}}))^{-1} \right) \quad (29)$$

C. Proof: MSSM objective (Eq. (12) \iff Eq. (13))

We assume that the model's score function $\mathbf{s}_m(\mathbf{x}, \boldsymbol{\theta})$ and the data score function $\mathbf{s}_d(\mathbf{x})$ are differentiable and satisfy $\mathbb{E}_{\mathbf{x} \sim p_d}[\|\mathbf{s}_d(\mathbf{x})\|^2] < \infty$ and $\mathbb{E}_{\mathbf{x} \sim p_d}[\|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|^2] < \infty$. Besides, we assume that the boundary condition holds, *i.e.*, $\lim_{x \rightarrow +\infty} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) p_d(\mathbf{x}) = 0, \forall \boldsymbol{\theta}$. In the following, we prove that the MSSM objective

$$\begin{aligned} \min_{\boldsymbol{\theta}} \max_{\mathbf{V}} \quad & \mathbb{E}_{\mathbf{x} \sim p_d}[\|\mathbf{V}^T(\mathbf{s}_d(\mathbf{x}) - \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}))\|^2], \\ \text{s.t.} \quad & \mathbf{V}\mathbf{V}^T = \mathbf{I}, \end{aligned} \quad (30)$$

is equivalent to:

$$\begin{aligned} \min_{\boldsymbol{\theta}} \max_{\mathbf{V}} \quad & \mathbb{E}_{\mathbf{x} \sim p_d}[\mathbf{V}^T \nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) \mathbf{V} + \frac{1}{2} \|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|^2], \\ \text{s.t.} \quad & \mathbf{V}\mathbf{V}^T = \mathbf{I}. \end{aligned} \quad (31)$$

Proof: We show how to derive a loss that doesn't depend on the unknown distribution $p_d(\mathbf{x})$ using: (1) the constraint $\mathbf{V}\mathbf{V}^T = \mathbf{I}$ and (2) Integration by parts (cf. Lemma 4 by Hyvärinen (2005)):

$$\begin{aligned} & \min_{\boldsymbol{\theta}} \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} \left[\frac{1}{2} \|\mathbf{V}^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) - \mathbf{V}^T \mathbf{s}_d(\mathbf{x})\|^2 \right], \quad \text{s.t.} \quad \mathbf{V}\mathbf{V}^T = \mathbf{I} \\ & \iff \min_{\boldsymbol{\theta}} \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} \left[\frac{1}{2} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})^T \mathbf{V}\mathbf{V}^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) + \mathbf{s}_d(\mathbf{x})^T \mathbf{V}\mathbf{V}^T \mathbf{s}_d(\mathbf{x}) - 2(\mathbf{V}^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}))^T (\mathbf{V}^T \mathbf{s}_d(\mathbf{x})) \right], \quad \text{s.t.} \quad \mathbf{V}\mathbf{V}^T = \mathbf{I} \\ & \stackrel{(i)}{\iff} \min_{\boldsymbol{\theta}} \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} \left[\frac{1}{2} \|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|^2 + \frac{1}{2} \|\mathbf{s}_d(\mathbf{x})\|^2 - (\mathbf{V}^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}))^T (\mathbf{V}^T \mathbf{s}_d(\mathbf{x})) \right], \quad \text{s.t.} \quad \mathbf{V}\mathbf{V}^T = \mathbf{I} \\ & \stackrel{(ii)}{\iff} \min_{\boldsymbol{\theta}} \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} \left[\frac{1}{2} \|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|^2 + \frac{1}{2} \|\mathbf{s}_d(\mathbf{x})\|^2 + \sum_{j=1}^K \mathbf{v}_j^T \nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) \mathbf{v}_j \right], \quad \text{s.t.} \quad \mathbf{V}\mathbf{V}^T = \mathbf{I} \\ & \stackrel{(iii)}{\iff} \min_{\boldsymbol{\theta}} \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} \left[\frac{1}{2} \|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|^2 + \sum_{j=1}^K \mathbf{v}_j^T \nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) \mathbf{v}_j \right], \quad \text{s.t.} \quad \mathbf{V}\mathbf{V}^T = \mathbf{I} \end{aligned}$$

Note, (i) is due to the fact that $\mathbf{V}\mathbf{V}^T = \mathbf{I}$. Further, (iii) is due to the fact that $\mathbf{s}_d(\mathbf{x})$ does not depend on either $\boldsymbol{\theta}$ or \mathbf{V} . Moreover, (ii) is due to the integration by parts trick, as we show in the following:

$$\begin{aligned} & \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} \left[\sum_{j=1}^K (\mathbf{v}_j^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})) (\mathbf{v}_j^T \mathbf{s}_d(\mathbf{x})) \right] \\ & \iff \max_{\mathbf{V}} \sum_{j=1}^K \int_{\mathbf{x}} p_d(\mathbf{x}) (\mathbf{v}_j^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})) (\mathbf{v}_j^T \mathbf{s}_d(\mathbf{x})) d\mathbf{x} \\ & \iff \max_{\mathbf{V}} \sum_{j=1}^K \int_{\mathbf{x}} p_d(\mathbf{x}) (\mathbf{v}_j^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})) (\mathbf{v}_j^T \nabla_{\mathbf{x}} \log p_d(\mathbf{x})) d\mathbf{x} \\ & \iff \max_{\mathbf{V}} \sum_{j=1}^K \int_{\mathbf{x}} (\mathbf{v}_j^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})) (\mathbf{v}_j^T \nabla_{\mathbf{x}} p_d(\mathbf{x})) d\mathbf{x} \end{aligned}$$

$$\begin{aligned}
 &\iff \max_{\mathbf{V}} \sum_{j=1}^K \int_{\mathbf{x}} (\mathbf{v}_j^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})) (\mathbf{v}_j^T \nabla_{\mathbf{x}} p_d(\mathbf{x})) d\mathbf{x} \\
 &\iff \max_{\mathbf{V}} \sum_{j=1}^K \sum_{i=1}^d \int_{\mathbf{x}} (\mathbf{v}_j^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})) (\mathbf{v}_j^{(i)} \frac{\partial p_d(\mathbf{x})}{\partial x_i}) d\mathbf{x} \\
 &\iff \max_{\mathbf{V}} \sum_{j=1}^K \sum_{i=1}^d \int_{\mathbf{x}} (\mathbf{v}_j^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})) (\mathbf{v}_j^{(i)} \frac{\partial p_d(\mathbf{x})}{\partial x_i}) d\mathbf{x} \\
 &\stackrel{(iv)}{\iff} \max_{\mathbf{V}} \sum_{j=1}^K \sum_{i=1}^d \left[(\mathbf{v}_j^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})) (\mathbf{v}_j^{(i)} p_d(\mathbf{x})) \right]_{-\infty}^{+\infty} - \int_{\mathbf{x}} \mathbf{v}_j^{(i)} p_d(\mathbf{x}) \mathbf{v}_j^T \frac{\partial \mathbf{s}_m(\mathbf{x})}{\partial x_i} d\mathbf{x} \\
 &\stackrel{(v)}{\iff} \max_{\mathbf{V}} \sum_{j=1}^K \sum_{i=1}^d - \int_{\mathbf{x}} \mathbf{v}_j^{(i)} p_d(\mathbf{x}) \mathbf{v}_j^T \frac{\partial \mathbf{s}_m(\mathbf{x})}{\partial x_i} d\mathbf{x} \\
 &\iff \max_{\mathbf{V}} - \sum_{j=1}^K \int_{\mathbf{x}} p_d(\mathbf{x}) \mathbf{v}_j^T \nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}) \mathbf{v}_j d\mathbf{x} \\
 &\iff \max_{\mathbf{V}} - \sum_{j=1}^K \mathbb{E}_{\mathbf{x} \sim p_d} [\mathbf{v}_j^T \nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}) \mathbf{v}_j]
 \end{aligned}$$

Note, (iv) is due to partial integration. Further, (v) is the boundary condition, *i.e.*, $p_d(\mathbf{x})$ vanishes at infinity. Hence,

$$\max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} \left[\sum_{j=1}^K (\mathbf{v}_j^T \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})) (\mathbf{v}_j^T \mathbf{s}_d(\mathbf{x})) \right] \iff \max_{\mathbf{V}} \mathbb{E}_{\mathbf{x} \sim p_d} \left[- \sum_{j=1}^K \mathbf{v}_j^T \nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) \mathbf{v}_j \right]. \quad (32)$$

D. Proof: MSSM relationship to MLE

As explained by Hyvarinen (2007), the objective: $\min_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x} \sim p_d} [\text{tr}(\nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})) + \frac{1}{2} \|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|^2]$, can be derived using the 2nd-order Taylor approximation of the contrastive divergence loss:

$$L_{\text{CD}}(\mathbf{x}; \boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim p_d(\mathbf{x})} [\log p_m(\mathbf{x}; \boldsymbol{\theta}) - \mathbb{E}_{\mathbf{x}' | \mathbf{x}} [\log p_m(\mathbf{x}'; \boldsymbol{\theta})]], \quad (33)$$

with a 1-step Langevin Monte-Carlo as the sampler from the model distribution $p_m(\mathbf{x}; \boldsymbol{\theta})$.

The Langevin Monte-Carlo sampler generates samples via:

$$\mathbf{x}' = \mathbf{x} + \frac{\mu}{2} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) + \sqrt{\mu} \boldsymbol{\sigma}, \quad \boldsymbol{\sigma} \sim \mathcal{N}(0, I). \quad (34)$$

The Taylor expansion gives:

$$\begin{aligned}
 \log p_m(\mathbf{x}'; \boldsymbol{\theta}) &= \log p_m(\mathbf{x}; \boldsymbol{\theta}) + \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})^T (\mathbf{x} - \mathbf{x}') + (\mathbf{x}' - \mathbf{x})^T \nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}) (\mathbf{x}' - \mathbf{x}) \\
 &= \log p_m(\mathbf{x}; \boldsymbol{\theta}) + \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})^T \left(\frac{\mu}{2} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) + \sqrt{\mu} \boldsymbol{\sigma} \right) + \left(\frac{\mu}{2} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) + \sqrt{\mu} \boldsymbol{\sigma} \right)^T \nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}) \left(\frac{\mu}{2} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) + \sqrt{\mu} \boldsymbol{\sigma} \right).
 \end{aligned}$$

Leveraging the fact that $\mathbb{E}[\boldsymbol{\sigma}] = 0$ and $\mathbb{E}[\boldsymbol{\sigma} \boldsymbol{\sigma}^T] = I$, we obtain:

$$\mathbb{E}_{\mathbf{x}' | \mathbf{x}} [\log p_m(\mathbf{x}'; \boldsymbol{\theta})] = \mathbb{E}_{\mathbf{x} \sim p_d(\mathbf{x})} [\log p_m(\mathbf{x}; \boldsymbol{\theta}) + \frac{\mu}{2} \|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|^2 + \mu \text{tr}(\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta})) + \frac{\mu^2}{4} \mathbf{s}_m^T(\mathbf{x}; \boldsymbol{\theta}) \nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}) \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})].$$

Hence, Eq. (33) becomes:

$$\begin{aligned}
 L_{\text{CD}}(\mathbf{x}; \boldsymbol{\theta}) &= - \mathbb{E}_{\mathbf{x} \sim p_d(\mathbf{x})} \left[\frac{\mu}{2} \|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|^2 + \mu \text{tr}(\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta})) + \frac{\mu^2}{4} \mathbf{s}_m^T(\mathbf{x}; \boldsymbol{\theta}) \nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta}) \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta}) \right] \\
 &= - \mathbb{E}_{\mathbf{x} \sim p_d(\mathbf{x})} \left[\frac{\mu}{2} \|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|^2 + \mu \text{tr}(\nabla_{\mathbf{x}} \log p_m(\mathbf{x}; \boldsymbol{\theta})) + \mathcal{O}(\mu^2) \right].
 \end{aligned}$$

SSM approximates the trace operator via the Hutchinson's trick, *i.e.*, the trace operator $\text{tr}(\nabla_{\mathbf{x}} s_m(\mathbf{x}; \boldsymbol{\theta}))$ is replaced with $\mathbb{E}_{\mathbf{v}}[\mathbf{v}^T \nabla_{\mathbf{x}} s_m(\mathbf{x}; \boldsymbol{\theta}) \mathbf{v}]$, such that $\mathbb{E}_{\mathbf{v}}[\mathbf{v} \mathbf{v}^T] = \mathbf{I}$. In case the projection directions follow a Rademacher distribution, *i.e.*, $\mathbf{v} \in \{-1, 1\}^d$, 2^d vectors are needed to exactly estimate the trace operator. However for MSSM, the equality $\text{tr}(\nabla_{\mathbf{x}} s_m(\mathbf{x}; \boldsymbol{\theta})) = \sum_{i=1}^d \mathbf{v}^T \nabla_{\mathbf{x}} s_m(\mathbf{x}; \boldsymbol{\theta}) \mathbf{v}$ holds when only considering d vectors.

E. Training Details

We closely follow the training setup of Song et al. (2019) for the different applications. In the following, we derive a closed form expression for the weight parameters $\boldsymbol{\alpha} = \{\alpha_l\}_{l=1}^L$ of the **Kernel Exponential Families** (DKEF) model, in case of the MSSM loss.

DKEF parameterizes the unnormalized log-likelihood via $\log p_f(\mathbf{x}) = \sum_{l=1}^L \alpha_l k(\mathbf{x}, \mathbf{z}_l) + \log q_0(\mathbf{x})$, where $q_0(\mathbf{x})$ is a fixed function, and $k(\mathbf{x}, \mathbf{z}_l)$ is a mixture of R Gaussian kernels defined on the feature space ϕ_r of a deep net and evaluated at L different inducing points \mathbf{z}_l :

$$k(\mathbf{x}, \mathbf{z}_l) = \sum_{r=1}^R \rho_r \exp\left(\frac{-1}{2\sigma_r^2} \|\phi_r(\mathbf{x}) - \phi_r(\mathbf{z}_l)\|^2\right). \quad (35)$$

The kernel parameters $\{\rho_r\}_{r=1}^R, \{\sigma_r\}_{r=1}^R, \{\phi_r\}_{r=1}^R$ and the inducing points $\{\mathbf{z}_l\}_{l=1}^L$ are learned via gradient descent. We prove that for fixed $\rho_r, \sigma_r, \phi_r, \mathbf{z}_l$ and \mathbf{V}^* , finding $\boldsymbol{\alpha}$ that minimizes the program,

$$\min_{\boldsymbol{\alpha}} \underbrace{\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \mathbf{v}_j^{*T} \nabla_{\mathbf{x}^{(i)}}^2 \log p_f(\mathbf{x}^{(i)}) \mathbf{v}_j^*}_{J_1} + \underbrace{\frac{1}{N} \sum_{i=1}^N \frac{1}{2} \|\nabla_{\mathbf{x}} \log p_f(\mathbf{x}^{(i)})\|^2 + \frac{1}{2} \lambda_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}\|^2}_{J_2}, \quad (36)$$

is equivalent to solving a linear system:

$$(\mathbf{G} + \lambda_{\boldsymbol{\alpha}} \mathbf{I}) \boldsymbol{\alpha} = -\mathbf{b} - \mathbf{c}, \quad (37)$$

with $\lambda_{\boldsymbol{\alpha}}$ being a regularization parameters, and $\mathbf{G} \in \mathbb{R}^{L \times L}$ and $\mathbf{b}, \mathbf{c} \in \mathbb{R}^L$ are defined as follows:

$$\begin{aligned} G_{l,l'} &= \frac{1}{N} \sum_{i=1}^N (\nabla_{\mathbf{x}^{(i)}} k(\mathbf{x}^{(i)}, \mathbf{z}_l))^T (\nabla_{\mathbf{x}^{(i)}} k(\mathbf{x}^{(i)}, \mathbf{z}_{l'})), \quad \forall l, l' \in [1, L] \\ b_l &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \mathbf{v}_j^{*T} \nabla_{\mathbf{x}^{(i)}}^2 k(\mathbf{x}^{(i)}, \mathbf{z}_l) \mathbf{v}_j^*, \quad \forall l \in [1, L] \\ c_l &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \nabla_{\mathbf{x}^{(i)}} q_0(\mathbf{x}^{(i)})^T k(\mathbf{x}^{(i)}, \mathbf{z}_l), \quad \forall l \in [1, L]. \end{aligned} \quad (38)$$

Note, N is the batch size and K is the number of eigenvectors.

Proof:

In the following we express the quadratic (J_1) and squared norm (J_2) terms of the loss in Eq. (36) as functions of $\boldsymbol{\alpha}, \mathbf{G}, \mathbf{b}$ and \mathbf{c} . Starting with J_1 :

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \mathbf{v}_j^{*T} \nabla_{\mathbf{x}}^2 \log p_f(\mathbf{x}) \mathbf{v}_j^* &= \frac{1}{N} \sum_{i=1}^N \sum_{l=1}^L \sum_{j=1}^K \alpha_l \mathbf{v}_j^{*T} \nabla_{\mathbf{x}^{(i)}} k(\mathbf{x}^{(i)}, \mathbf{z}_l) \mathbf{v}_j^* + \text{Const} \\ &= \sum_{l=1}^L \alpha_l b_l + \text{Const} \\ &= \boldsymbol{\alpha}^T \mathbf{b} + \text{Const} \end{aligned} \quad (39)$$

Next, we express J_2 as a function of $\boldsymbol{\alpha}, \mathbf{G}$ and \mathbf{c} :

$$\begin{aligned}
 & \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \|\nabla_{\mathbf{x}} \log p_f(\mathbf{x})\|^2 \\
 &= \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \left\| \sum_{l=1}^L \alpha_l \nabla_{\mathbf{x}^{(i)}} k(\mathbf{x}^{(i)}, \mathbf{z}_l) + \log q_0(\mathbf{x}^{(i)}) \right\|^2 \\
 &= \frac{1}{2N} \sum_{i=1}^N \left(\sum_{l,l'}^L \alpha_l \alpha_{l'} \nabla_{\mathbf{x}^{(i)}} k(\mathbf{x}^{(i)}, \mathbf{z}_l)^T \nabla_{\mathbf{x}^{(i)}} k(\mathbf{x}^{(i)}, \mathbf{z}_{l'}) + 2 \sum_{l=1}^L \alpha_l \nabla_{\mathbf{x}^{(i)}} k(\mathbf{x}^{(i)}, \mathbf{z}_l)^T \nabla_{\mathbf{x}^{(i)}} \log q_0(\mathbf{x}^{(i)}) \right) + \text{Const} \quad (40) \\
 &= \frac{1}{2} \sum_{l,l'}^L \alpha_l \alpha_{l'} G_{l,l'} + \sum_{l=1}^L \alpha_l c_l + \text{Const} \\
 &= \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{G} \boldsymbol{\alpha} + \boldsymbol{\alpha}^T \mathbf{c} + \text{Const}
 \end{aligned}$$

Hence, the optimization problem in Eq. (36) becomes:

$$\begin{aligned}
 \boldsymbol{\alpha} &= \underset{\boldsymbol{\alpha}}{\text{argmin}} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K [\mathbf{v}_j^{*T} \nabla_{\mathbf{x}^{(i)}}^2 \log p_f(\mathbf{x}^{(i)}) \mathbf{v}_j^* + \frac{1}{2} \|\nabla_{\mathbf{x}} \log p_f(\mathbf{x}^{(i)})\|^2] + \frac{1}{2} \lambda_{\alpha} \|\boldsymbol{\alpha}\|^2 \\
 &= \underset{\boldsymbol{\alpha}}{\text{argmin}} \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{G} \boldsymbol{\alpha} + \boldsymbol{\alpha}^T \mathbf{b} + \boldsymbol{\alpha}^T \mathbf{c} + \frac{1}{2} \lambda_{\alpha} \|\boldsymbol{\alpha}\|^2 \\
 &= -(\mathbf{G} + \lambda_{\alpha} \mathbf{I})^{-1} (\mathbf{b} + \mathbf{c}).
 \end{aligned} \quad (41)$$