

MULTI-SUBSPACE MULTI-MODAL MODELING FOR DIFFUSION MODELS: ESTIMATION, CONVERGENCE AND MIXTURE OF EXPERTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Recently, diffusion models have achieved a great performance with a small dataset of size n and a fast optimization process. Despite the impressive performance, the estimation error suffers from the curse of dimensionality $n^{-1/D}$, where D is the data dimension. Since images are usually a union of low-dimensional manifolds, current works model the data as a union of linear subspaces with Gaussian latent and achieve a $1/\sqrt{n}$ bound. Though this modeling reflects the multi-manifold property of data, the Gaussian latent can not capture the multi-modal property of the latent manifold. To bridge this gap, we propose the mixture subspace of low-rank mixture of Gaussian (MoLR-MoG) modeling, which models the target data as a union of K linear subspaces, and each subspace admits a mixture of Gaussian latent (n_k modals with dimension d_k). With this modeling, the corresponding score function naturally has a mixture of expert (MoE) structure, captures the multi-modal information, and contains nonlinear properties since each expert is a nonlinear latent MoG score. We first conduct real-world experiments to show that the generation results of MoE-latent MoG NN are much better than the results of MoE-latent Gaussian score. Furthermore, MoE-latent MoG NN achieves a comparable performance with MoE-latent Unet with $10\times$ parameters. These results indicate that the MoLR-MoG modeling is reasonable and suitable for real-world data. After that, based on such MoE-latent MoG score, we provide a $R^4 \sqrt{\sum_{k=1}^K n_k} \sqrt{\sum_{k=1}^K n_k d_k} / \sqrt{n}$ estimation error, which escapes the curse of dimensionality by using data structure. Finally, we study the optimization process and prove the convergence guarantee under the MoLR-MoG modeling. Combined with these results, under a setting close to real-world data, this work explains why diffusion models only require a small training sample and enjoy a fast optimization process to achieve a great performance.

1 INTRODUCTION

Recently, diffusion models have achieved impressive performance in many areas, such as 2D, 3D, and video generation (Rombach et al., 2022; Ho et al., 2022; Chen et al., 2023a; Ma et al., 2024; Liu et al., 2024). Due to the score matching technique, diffusion models enjoy a more stable training process and can achieve great performance with a small training dataset.

Despite the empirical success, the theoretical guarantee for the estimation and optimization error of the score matching process is lacking. For estimation error, current results suffer from the curse of dimensionality. More specifically, given training dataset $\{x^i\}_{i=1}^n$ with $x^i \in \mathbb{R}^D$, the estimation error of the score function achieve the minimax $n^{-s'/D}$ results for (conditional) diffusion models with deep ReLU NN and diffusion transformer, where s' is the smoothness parameter of the score function (Oko et al., 2023; Hu et al., 2024b;a; Fu et al., 2024). It is clear that this estimation error is heavily influenced by the external dimension D , which can not explain why diffusion models can generate great images with a small training dataset. Hence, a series of works studies estimation errors under specific target data structures and reduces the curse of dimensionality. There are two notable ways to model the target data: the multi-modal modeling and the low-dimensional modeling. For the multi-modal modeling, as the real-world target data is usually multi-modal, some works study the mixture of Gaussian (MOG) target data and improve the estimation error (Shah et al., 2023; Cui et al.,

2023; Chen et al., 2024b). When we delve deeper into the images and text data, a key feature is that the image and text data usually admit a low-dimensional structure (Pope et al., 2021; Brown et al., 2023; Kamkari et al., 2024). Hence, one notable way is to assume the data admits a low-dimensional structure. More specifically, some works assume the data admits a linear subspace $x = Az$, where $A \in \mathbb{R}^{D \times d}$ to convert data to the latent space and $z \in \mathbb{R}^d$ is a bounded support (Chen et al., 2023b; Yuan et al., 2023; Guo et al., 2024). Then, they reduce the estimation error to $n^{-2/d}$, which removes the dependence of D . However, as shown in Brown et al. (2023) and Kamkari et al. (2024), though the image dataset admits low dimension, it is a union of manifolds instead of one manifold. Inspired by this observation, Wang et al. (2024) model the image data as a union of linear subspaces, assume each subspace admits a low-dimensional Gaussian (mixture of low-rank Gaussians (MoLRG)), and achieve a $1/\sqrt{n}$ estimation error. Though the union of the linear subspace is closer to the real-world image dataset, the latent Gaussian assumption is far away from the low-dimensional multi-modal manifold Brown et al. (2023). Hence, the following two natural questions remain open:

Can we propose a modeling that reflects the multi-manifold multi-modal property of real-world data?

Can we escape the curse of dimensionality and enjoy a fast convergence rate based on this modeling?

In this work, for the first time, we propose and analyze the mixture of low-rank mixture of Gaussian (MoLR-MoG) distribution, which is more realistic than MoLRG since it captures the multi-modal property of real-world distribution and has a nonlinear score function. Based on this modeling, we first induce a MoE-latent nonlinear score function and conduct experiments to show that MoLR-MoG modeling is closer to the real-world data. After that, we simultaneously analyze the estimation and optimization error of diffusion models and explain why diffusion models achieve great performance.

1.1 OUR CONTRIBUTION

MoLR-MoG modeling and MoE Structure Nonlinear Score. We propose the MoLR-MoG modeling for the target data, which captures the multi low-dimensional manifold and multi-modal property of real-world data and naturally introduces the MoE-latent MoG score. Through the real-world experiments, we show that with this score, diffusion models can generate images that is comparable with the deep neural network MoE-latent Unet and only has $10\times$ smaller parameters. On the contrary, the MoE-latent Gaussian score induced by previous MoLRG modeling can only generate blurry images, which indicates MoLR-MoG is a suitable modeling for the real-world data.

Take Advantage of MoLR-MoG to Escape the Curse of Dimensionality. For the estimation error, we show that by taking advantage of the union of a low-dimensional linear subspace and the latent MoG property, diffusion models escape the curse of dimensionality. More specifically, we achieve the $R^4 \sqrt{\sum_{k=1}^K n_k} \sqrt{\sum_{k=1}^K n_k d_k} / \sqrt{n}$ estimation error, where R is the diameter of the target data, d_k is the latent dimension and n_k is the number of the modal in the k -the subspace. This result clearly shows the dependence on the number of linear subspaces, modal, and the latent dimensions R, d_k .

Strongly Convex Property and Convergence Guarantee. After directly analyzing the estimation error, we study how to optimize the highly non-convex score-matching objective function. Facing nonlinear latent MoG scores, we use the gradient descent (GD) algorithm to optimize the objective function. To obtain the convergence guarantee, we take advantage of the closed form of nonlinear MoG score and show that the landscape around the ground truth parameter is strongly convex. Then, with a great initialization area, we prove the convergence guarantee when considering MoLR-MoG.

2 RELATED WORK

Estimation Error Analysis for Diffusion Models. As shown in Section 1, a series of works Oko et al. (2023) study the general target data with a deep NN and achieve the minimax $n^{-s'/D}$ result. Then, some works analyze the general target data with a 2-layer wide NN and achieve $n^{-2/5}$ estimation error with $\exp(n)$ NN size (Li et al., 2023; Han et al., 2024). For the multi-modal modeling, some works study MoG data and improve the estimation error (Shah et al., 2023; Cui et al., 2023; Chen et al., 2024b). Except for the MoG modeling, Cole and Lu (2024) assume data is close to Gaussian and then prove the model escapes the curse of dimensionality. Mei and Wu (2023) analyze Ising models and prove that the term corresponds to n is $1/\sqrt{n}$. For the low-dimensional modeling, some works assume the target data admits a linear subspace (Chen et al., 2023b; Yuan et al., 2023). Chen et al. (2023b) assume data admit a linear subspace $x = Az$ with $z \in \mathbb{R}^d$ and achieve a $n^{-2/d}$. As the

image is a union of low-dimensional manifolds, Wang et al. (2024) models the target data as a union of linear subspaces with Gaussian latent and achieve $1/\sqrt{n}$ estimation error for each subspace.

Optimization Analysis for Diffusion Models. Since the score is highly nonlinear (except for Gaussian), only a few works analyze the optimization process, and most of them focus on the external dimensional space (Bruno et al., 2023; Cui and Zdeborová, 2023; Shah et al., 2023; Chen et al., 2024b; Li et al., 2023; Han et al., 2024). Since the score function of MoG has a nonlinear closed-form, a series of works design algorithms for diffusion models to learn the MoG (Bruno et al., 2023; Cui and Zdeborová, 2023; Shah et al., 2023; Chen et al., 2024b). For the general target data, Li et al. (2023) and Han et al. (2024) adopt a wide 2-layer ReLU NN to simplify the problem to a convex optimization. However, as discussed above, their NN has $\exp(n)$ size. For the latent space, only two works provide the optimization guarantee under the Gaussian latent (Yang et al., 2024a; Wang et al., 2024). Yang et al. (2024a) assume target data adopts a linear subspace with Gaussian latent and provide the closed-form minimizer. Wang et al. (2024) analyze the optimization process of each linear subspace separately, which is also reduced to the optimization for the Gaussian.

3 PRELIMINARIES

First, we introduce the basic knowledge and notation of diffusion models. Let p_0 be the data distribution. Given $x_0 \sim p_0 \in \mathbb{R}^D$, the forward process is defined by:

$$dx_t = f(t)x_t dt + g(t) dB_t,$$

where $\{B_t\}_{t \in [0, T]}$ is a D -dimensional Brownian motion, $f(t)$ is the coefficient of the drift term and $g(t)$ is the coefficient of the diffusion term. Let p_t be the density function of the forward process. After determining the forward process, the conditional distribution $p_t(x_t|x_0)$ has a closed-form

$$p_t(x_t|x_0) = \mathcal{N}(x_t; s_t x_0, s_t^2 \sigma_t^2 I_D),$$

where $s_t = \exp\left(\int_0^t f(\xi) d\xi\right)$, $\sigma_t = \sqrt{\int_0^t g^2(\xi)/s^2(\xi) d\xi}$. To generate samples from p_0 , diffusion models reverse the given forward process and obtain the following reverse process (Song et al., 2020):

$$dy_t = [f(t)y_t - g(t)^2 \nabla \log p_t(y_t)] dt + g(t) d\bar{B}_t, \quad y_0 \sim p_0$$

where \bar{B}_t is a reverse-time Brownian motion. A conceptual way to approximate the score function is to minimize the score matching (SM) objective function:

$$\min_{s_\theta \in \text{NN}} \mathcal{L}_{\text{SM}} = \int_\delta^T \mathbb{E}_{x_t \sim q_t} \|\nabla \log p_t(x_t) - s_\theta(x_t, t)\|_2^2 dt, \quad (1)$$

where NN is a given function class and $\delta > 0$ is the early stopping parameter to avoid a blow-up score. Since the ground truth score $\nabla \log p_t$ is unknown, this objective function can not be calculated. To avoid this problem, Vincent (2011) propose the denoised score matching (DSM) objective function:

$$\min_{s_\theta \in \text{NN}} \mathcal{L}_{\text{DSM}} = \int_\delta^T \mathbb{E}_{x_0 \sim q_0} \mathbb{E}_{x_t|x_0} \|\nabla \log p_t(x_t|x_0) - s_\theta(x_t, t)\|_2^2 dt.$$

As shown in Vincent (2011), the DSM and SM objective functions differ up to a constant independent of optimized parameters, which indicates these objective functions have the same landscape.

3.1 MIXTURE OF LOW-RANK MIXTURE OF GAUSSIAN (MoLR-MoG) MODELING

This part shows our MoLR-MoG modeling, which reflects the low-dimensional (Gong et al., 2019) and multi-modal property (Brown et al., 2023; Kamkari et al., 2024) of real-world data. More specifically, we assume the data distribution lives near a union of K linear subspaces rather than arbitrary manifolds. Concretely, for the k -th subspace of dimension d_k (represented by a **orthonormal basic** matrix $A_k^* \in \mathbb{R}^{D \times d_k}$ with orthonormal columns **for the k -th manifold**), we place a n_k -modal MoG within that subspace:

$$w_k(x) = \sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; A_k^* \mu_{k,l}^*, A_k^* \Sigma_{k,l}^* A_k^{*\top}),$$

where covariance $\Sigma_{k,l}^* = U_{k,l}^* U_{k,l}^{*\top}$, $l = 1, \dots, n_k$ with $U_{k,l}^* \in \mathbb{R}^{d_k \times d_{k,l}}$ ($d_{k,l} \leq d_k$) and $\mu_{k,l}^*$ is the mean of the l -th modal of the k -th subspace. **As shown in (Brown et al., 2023), the different manifold**

has different d_k and we do not require that d_k is exactly the same for each manifold. Then, the target distribution has the following form

$$p_0 = \sum_{k=1}^K \frac{1}{K} \sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; A_k^* \mu_{k,l}^*, A_k^* \Sigma_{k,l}^* A_k^{*\top}). \quad (2)$$

From the universal approximation perspective, by placing enough components and choosing parameters $\{\pi_{k,l}, \mu_{k,l}^*, \Sigma_{k,l}^*\}$, a MoG can approximate any smooth density arbitrarily well, which is more general than the Gaussian latent of Yang et al. (2024a) and Wang et al. (2024).

Nonlinear Mixture of Experts (MoE)-latent MoG score. Let $\gamma_t = s_t \sigma_t$, $\Sigma_{k,l,t,A} = s_t^2 A_k^* U_{k,l}^* U_{k,l}^{*\top} A_k^{*\top} + \gamma_t^2 I$ and $\delta_{k,l,t,A}(x) = x - s_t \mu_{k,l}^* - \frac{s_t^2}{s_t^2 + \gamma_t^2} A_k^* U_{k,l}^* U_{k,l}^{*\top} A_k^{*\top} (x - s_t \mu_{k,l}^* A_k^*)$. Under the MoLR-MoG modeling, the score function has the following form:

$$\nabla \log p_t(x) = -\frac{1}{\gamma_t^2} \frac{\sum_{k=1}^K \frac{1}{K} \sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}^* A_k^*, A_k^* \Sigma_{k,l,t,A} A_k^{*\top}) \delta_{k,l,t,A}(x)}{\sum_{k=1}^K \frac{1}{K} \sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}^* A_k^*, A_k^* \Sigma_{k,l,t,A} A_k^{*\top})},$$

This score function has a MoE structure, where each expert is the latent nonlinear MoG score. The linear encoder A_k first encodes images to the k -th manifold, and diffusion models run the denoising process. After that, the linear decoder A_k^\top decodes the denoised latent to the full-dimensional images. Since the estimation error introduced by the linear encoder and decoder has the order Dd_k^3/\sqrt{n} (Yang et al., 2024a) and is not the dominant term, we assume the linear encoder and decoder are perfectly learned and focus on the more difficult latent MoG diffusion part in this work. From the empirical part, this operation is similar to using the pretrained stable diffusion VAE and only training the diffusion models in the latent space. For the k -th low-dimensional manifold, the score function is

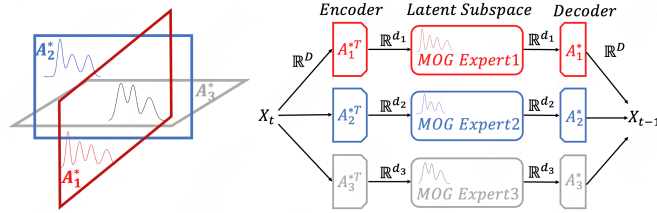
$$\nabla \log p_{t,k}(x^{\text{LD}}) = -\frac{1}{\gamma_t^2} \frac{\sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x^{\text{LD}}; s_t \mu_{k,l}^*, \Sigma_{k,l,t}^*) \delta_{k,l,t}(x^{\text{LD}})}{\sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}^*, \Sigma_{k,l,t}^*)}, \quad (3)$$

where $x^{\text{LD}} \in \mathbb{R}^{d_k}$ is a variable in the k -th low-dimensional subspace, $\Sigma_{k,l,t}^* = s_t^2 U_{k,l}^* U_{k,l}^{*\top} + \gamma_t^2 I$ and $\delta_{k,l,t}(x^{\text{LD}}) = x^{\text{LD}} - s_t \mu_{k,l}^* - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l}^* U_{k,l}^{*\top} (x^{\text{LD}} - s_t \mu_{k,l}^*)$. Let

$$s_k^*(x^{\text{LD}}, t) = \nabla \log p_{t,k}(x^{\text{LD}}), s^*(x^{\text{LD}}, t) = (s_1^*(x^{\text{LD}}, t), s_2^*(x^{\text{LD}}, t), \dots, s_K^*(x^{\text{LD}}, t)),$$

where the parameters are $\theta^* = \{\mu_{k,l}^*, U_{k,l}^*\}_{k=1, \dots, K}$. In this work, we want to learn the parameters of the ground truth score function. Hence, we construct a NN function class $s_\theta = (s_1(\cdot, \cdot), s_2(\cdot, \cdot), \dots, s_K(\cdot, \cdot))$ according to the above closed-form of MoE-latent MoG score. Let θ is the union of $\mu_{k,l}$ and $U_{k,l}$. Since we mainly focus on the estimation and optimization in the latent subspace, we omit the superscript LD of the latent subspace when there is no ambiguity.

We note that this modeling can capture the information of each low-dimensional manifold and the multi-modal property of each latent distribution. In the next section, through the real-world experiments, we show that the MoE-latent MoG score has a better performance compared with the MoE-latent Gaussian



(a) MoLR-MoG Modeling (b) MoE-nonlinear MoG Score

Figure 1: MoLR-MoG Modeling and Corresponding Nonlinear Score score induced by MoLRG modeling and compatible with the results of the MoE-latent Unet. In Section 5 and 6, we prove that by using the property of MoLR-MoG modeling, diffusion models can escape the curse of dimensionality and enjoy a fast convergence rate.

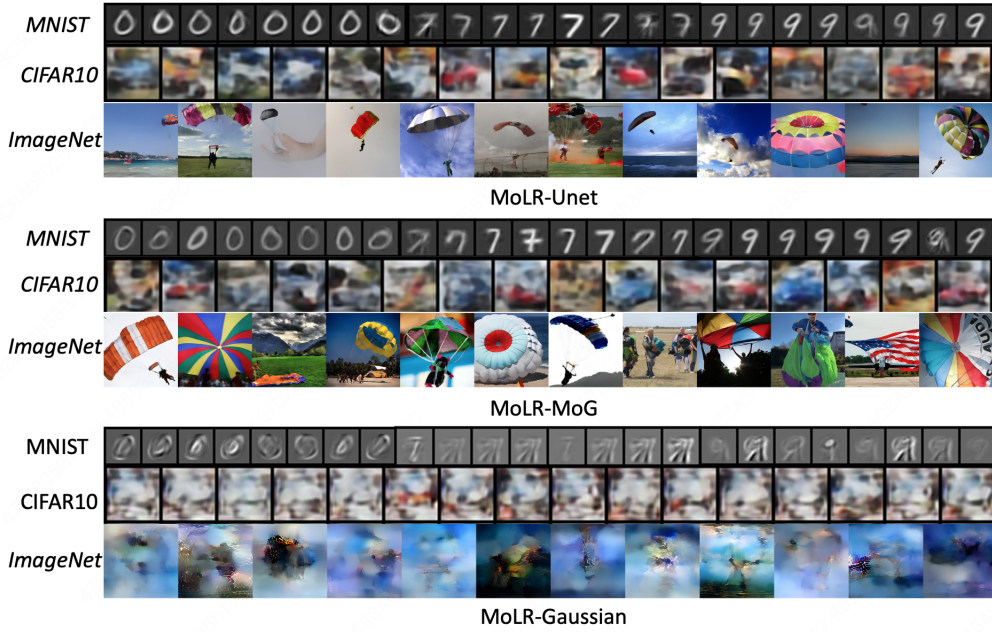


Figure 2: Results of Different Modeling on Real-world Data.

Remark 3.1 (Comparison with MoLRG modeling). Wang et al. (2024) provide the first multi-subspace modeling, which is an important and meaningful step. However, they assume a Gaussian latent with 0 mean, which can not capture the multi-modal property of real-world data. We also note that the MoLR-MoG modeling can not be viewed as MoLRG with $\sum_{k=1}^K n_k$ subspace since this modeling assumes there are $\sum_{k=1}^K n_k$ VAE, which is not reasonable in the real-world setting.

4 EXPERIMENTS FOR MOE-LATENT MOG SCORE

In this section, we conduct experiments using neural networks based on different modeling approaches (MoLR-MoG, MoLRG) as well as a general U-Net architecture. The goal is to demonstrate that MoLR-MoG provides a suitable modeling for real-world data, and that the MoE-latent MoG score is sufficient to generate images with clear semantic content. Specifically, we first show that training with MoLR-MoG yields significantly better results than the MoLRG model. Then, we show that the MoE-latent MoG network achieves performance comparable to that of the MoLR-U-Net, while using 10x fewer parameters for MNIST, CIFAR-10, ImageNet 256. (Figure 2)

Following Brown et al. (2023), we train 10 VAEs for each number in the MNIST, which represents our K low-dimensional manifold. In this part, we adopt nonlinear VAEs to achieve a good performance in real-world datasets. However, we still note that a series of theoretical works adopt linear subspaces, and our MoLR-MoG modeling with linear VAEs makes a step toward explaining the good performance of diffusion models. After obtaining these 10 VAE, we train diffusion models with different parametrized NNs. We adopt three different parameterizations: latent U-net, latent MoG NN, and latent Gaussian NN. For the latent MoG, we adopt the form of Eq. 3 with $n_k = 4, 8, 40$ in MNIST, CIFAR-10, and ImageNet256 for $k \in [K]$. For the latent Gaussian, we adopt the form of the closed-form score (Wang et al., 2024), which leads to a linear NN.

Discussion. From a qualitative perspective, as shown in Figure 2, the generation results with MoLRG modeling are difficult to distinguish specific numbers. On the contrary, the MoE-latent MoG can generate clean images comparable with the images generated by MoLR-Unet, which means this modeling captures the multi-modal property of each low-dimensional manifold. The training loss curve (Figure 3) shows that the loss of MoE-MoG NN is significantly smaller than the MoE-Gaussian and close to MoE-Uet, which indicates MoE-MoG NN efficiently approximates the ground-truth

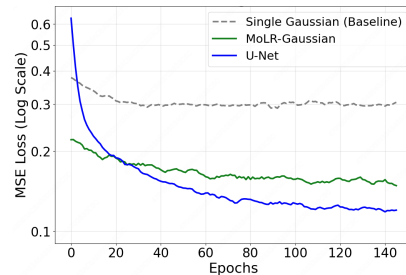


Figure 3: Loss Curve for CIFAR-10

score and supports our theoretical results. From a quantitative perspective, we calculate the CLIP score for the parachute class of ImageNet with text prompts "a photo of parachute". The Clip score for MoLR with Unet, MoG, and Gaussian NN is 0.304, 0.293, and 0.254, which indicates MoLR-MoG achieves almost comparable text-to-image alignment with MoE-Unet. Furthermore, the MoLR-MoG NN contains many fewer parameters compared to Unet since it uses the prior of latent MoG.

Discussion on Expert-Specific VAE. As shown in the score of MoLR-MoG, different from latent diffusion models with a single VAE, there are K VAEs to encode the input to the corresponding manifold. We note that this operation is important for MoLR-MoG with small MoG experts. As shown in Figure 4, with a unified VAE, the unified latent is complex, and a MoG expert can not learn a meaningful image with the target class. Hence, with a unified VAE, latent diffusion models require a large latent Unet. However, with an expert-specific VAE (for example, we fine-tune the pretrained VAE with the parachute class dataset), the latent manifold becomes simple, and latent MoG experts are enough to generate clear models, which also supports our theoretical modeling.

We note that these experiments aim to show that the MoLR-MoG modeling is reasonable instead of achieving the SOTA performance. It is possible to achieve great performance with a small-sized NN using MoLR-MoG modeling in the application. For large-scale datasets without labels, we can use a clustering algorithm to divide the data into different clusters. Then, we can train a VAE encoder, decoder, and latent MoG score for each cluster. For the VAE training, we do not require training the VAE from a sketch. We can LoRA fine-tune a VAE pretrained on large-scale datasets (for example, DC-AE (Chen et al., 2024a) for our ImageNet experiments) for each expert, which shares a pretrained VAE backbone and has a smaller model size. When generating images, we activate different VAE LoRA according to the clustering weight, which matches the spirit of MoE. We leave it as an interesting future work.



Figure 4: MoLR-MoG with Different VAE

5 ESCAPE THE CURSE OF DIMENSIONALITY WITH MoLR-MoG MODELING

This section shows that diffusion models can escape the curse of dimensionality by using MoLR-MoG properties. Before introducing our results, we first introduce the assumption on the target data.

Assumption 5.1. For $x \sim p_0$, we have that $\|x\|_2 \leq R$.

The bounded-support assumption is widely used in theoretical works (Chen et al., 2022; Yang et al., 2024a;b) and is naturally satisfied by image datasets. For a latent MoG, each component concentrates almost all mass within a few standard deviations of its mean, so by taking the most component means and variances, one can choose R large enough that $\|x\|_2 \leq R$ holds with high probability.

Since MoE-latent MoG score has a closed-form, we only need to learn the parameters $\mu_{k,l}$ and $U_{k,l}$ at a fixed time t . As a result, we consider the estimation error at a fixed time t . Let $\ell(\theta; x, t) = \|s_\theta(x, t) - s^*(x, t)\|_2^2$ be the per-sample squared error at time t . In this part, we study the estimation error with a limited training dataset $\{x_i\}_{i=1}^n$:

$$\left| \mathcal{L}(\theta) - \hat{\mathcal{L}}_n(\theta) \right|, \text{ with } \hat{\mathcal{L}}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\theta; x_i, t).$$

To obtain the estimation error, we first provide the Lipschitz constant for s_θ and the loss function by fully using the property of MoLR-MoG modeling and MoE-latent MoG score.

Lemma 5.2. [Lipschitz Continuity] Let L_{μ_l} and L_{U_k} be the Lipschitz constant w.r.t. s_θ . With MoLR-MoG modeling and Assumption 5.1, there is a constant

$$L \leq \sqrt{\sum_{i=1}^K n_k (L_{\mu_l}^2 + L_{U_k}^2)} = O\left((\sum_{k=1}^K n_k)^{\frac{1}{2}} C_w\right)$$

such that for any θ, θ' , $\|s_\theta(x, t) - s_{\theta'}(x, t)\|_2 \leq L \|\theta - \theta'\|_2$, where $C_w = \frac{(R + s_t B_\mu)^3 s_t^2}{\gamma_t^4}$, $B_\mu = \max_{k,l} \|\mu_{k,l}\|_2$. For s_θ and s^* , we have that $2\|s_\theta(x, t) - s^*(x, t)\|_2 \leq 2(R + s_t B_\mu)/\gamma_t^2 := L_t$.

Then, we obtain the Lipschitz constant $L' = L_l L$ for the whole loss function. With this Lipschitz property, the next step is to argue that fitting the network on n samples generalizes to the true population loss. We do so by controlling the Rademacher complexity of the loss class and then using a Bernstein concentration argument to obtain the following theorem.

Theorem 5.3. Denote by $\hat{\mathcal{L}}_n(\theta)$ the empirical loss on n i.i.d. samples and by $\mathcal{L}(\theta)$ its population counterpart. Then there exist constants C_1, C_2 such that with probability at least $1 - \delta$, for all $\theta \in \Theta$,

$$|\mathcal{L}(\theta) - \hat{\mathcal{L}}_n(\theta)| \leq O\left(C_1 \frac{(R + s_t B_\mu)^4 s_t^2 \sqrt{\sum_{k=1}^K n_k}}{\gamma_t^6} \sqrt{\frac{\sum_{k=1}^K n_k d_k}{n}} + C_2 \sqrt{\frac{\log(1/\delta)}{n}}\right).$$

where $C_1 = \max_{\theta \in \Theta} \|\theta_i - \theta_j\|_2$, $C_2 = \sigma \log 2$, $\sigma^2 = \sup_{\theta \in \Theta} \text{Var}[\ell(\theta; X, t)]$.

This result removes the exponential dependence on D with the number of latent subspace K , the latent dimension d_k , and the number of modalities n_k at each linear subspace, which reflects the key feature of the real-world data and escape the curse of dimensionality. The remaining question is why diffusion models enjoy a fast and stable optimization process. In the next part, we show that with MoLR-MoG modeling, the objective function is locally strongly convex and answer this question.

6 STRONGLY CONVEX PROPERTY AND CONVERGENCE GUARANTEE

In this part, by using the property of MoLR-MoG modeling, we derive explicit expressions for the *Jacobian* and *Hessian* of the objective function for 2-modal MoG latent and general MoG latent. Then, we establish conditions under which the resulting score-matching loss is locally strongly convex for each setting. Finally, we provide the convergence guarantee for the optimization.

6.1 2-MODAL LATENT MOG HESSIAN ANALYSIS AND OPTIMIZATION

In this section, we show that, under sufficient cluster separation, the Hessian matrix near θ^* simplifies to a block-diagonal form, yielding local strong convexity, which derives a linear convergence rate. As discussed in Section 3.1, following the real-world setting, we consider the optimization dynamic in the k -th latent subspace. While our modeling contains K encoders and decoders, facing an input image x , we can first determine which cluster image x belongs to, and then use the corresponding A_k to encode it into the corresponding latent space. Then, we only use data belonging to k clustering to train the k -th latent MoG score. This operation matches our experimental settings, and Wang et al. (2024) also adopts this operation. When considering the optimization problem, to simplify the calculation of the Hessian matrix, we set $d_{k,l} = 1$.

Similar to Shah et al. (2023), we start from a latent 2-modal MoG with the same covariance matrix Σ_k^* and $\mu_{k,1}^* = \mu_k^*, \mu_{k,2}^* = -\mu_k^*$, which leads to the following score:

$$\nabla \log p_{t,k}(x) = -\frac{1}{\gamma_t^2} \frac{\frac{1}{2} \mathcal{N}(x; s_t \mu_k^*, \Sigma_k^*) \delta'_k(x) + \frac{1}{2} \mathcal{N}(x; -s_t \mu_k^*, \Sigma_k^*) \epsilon_k(x)}{\frac{1}{2} \mathcal{N}(x; s_t \mu_k^*, \Sigma_k^*) + \frac{1}{2} \mathcal{N}(x; -s_t \mu_k^*, \Sigma_k^*)}, \quad (4)$$

where $\epsilon_k(x) = x - s_t \mu_k^* - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k^* U_k^{*\top} (x - s_t \mu_k^*)$, and $\delta'_k(x) = x + s_t \mu_k^* - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k^* U_k^{*\top} (x + s_t \mu_k^*)$. Before providing the convergence guarantee, we make an assumption on the 2-MoG latent distribution.

Assumption 6.1. [Separation within a cluster] Within each cluster k , the two symmetric peaks are well separated in the sense that $\|s_t \mu_k^* - (-s_t \mu_k^*)\| \geq \Delta_{\text{intra}}$, for some $\Delta_{\text{intra}} \gg \gamma_t$. Consequently, if a sample x is drawn from the “+” peak then its responsibility under the “−” peak satisfies

$$r_k^-(x) = \frac{\frac{1}{2} \mathcal{N}(x; -s_t \mu_k^*, \Sigma_k^*)}{\frac{1}{2} \mathcal{N}(x; s_t \mu_k^*, \Sigma_k^*) + \frac{1}{2} \mathcal{N}(x; -s_t \mu_k^*, \Sigma_k^*)} = O(e^{-\Delta_{\text{intra}}^2 / (2\gamma_t^2)}) \ll 1,$$

and symmetrically $r_k^+(x) \ll 1$ when x is drawn from the “−” peak.

The above assumption means that the separation of the two modals is sufficient. For each symmetric sub-peak, if the distance between them is relatively small, we can view them as having a mean of 0. Since they are the same distribution ($\mu = 0$ and $\Sigma = U_k U_k^\top + \gamma_t^2 I$), they are the same regardless of how they mix, which indicates that we can assume $r_k^+ \approx 1$ or $r_k^- \approx 1$. Moreover, in practice, if

raw data do not exhibit such clear gaps, one can always apply a simple linear embedding to magnify inter-mean distances relative to noise, thereby enforcing the same hard-assignment regime.

Since the ground truth score function has a closed-form under the MoLR-MoG modeling, we focus on the score matching objective function $\mathcal{L}_{\text{SM}}(\theta)$ instead of $\mathcal{L}_{\text{DSM}}(\theta)$ and abbreviate $\mathcal{L}_{\text{SM}}(\theta)$ as $\mathcal{L}(\theta)$. We note that $\mathcal{L}_{\text{SM}}(\theta)$ and $\mathcal{L}_{\text{DSM}}(\theta)$ are equivalent up to a constant independent of θ , which indicates the optimization landscape is the same. Furthermore, when considering the convergence guarantee under a 2-layer wide ReLU NN, Li et al. (2023) also adopt score matching objective \mathcal{L}_{SM} instead of \mathcal{L}_{DSM} . **Though calculating the bound of Jacobian $J_k^\mu(x) = \partial_{\mu_k} s_\theta$, $J_k^U(x)$ and the Hessian matrix w.r.t. \mathcal{L} , we provide the local strongly convexity parameters for the objective function.**

Lemma 6.2. [Local Strong Convexity] *Combining Lemma C.4 with continuity of $\nabla^2 \mathcal{L}$, there exist $\alpha > 0$ and neighborhood U of θ^* such that $\nabla^2 \mathcal{L}(\theta) \succeq \alpha I, \forall \theta \in \Theta$. If $\forall x \in \mathbb{R}^{d_k}, r_k^+(x) = 1$ or $r_k^-(x) = 1$ are strictly satisfied,*

$$\alpha = \min \left\{ \frac{s_t^2}{(s_t^2 + \gamma_t^2)^2}, \frac{4(U_k^\top \mu_k)^2 + \|U_k\|_2^2 \|\mu_k\|_2^2 - \|U_k\|_2 \|\mu_k\|_2 \sqrt{8(U_k^\top \mu_k)^2 + \|U_k\|_2^2 \|\mu_k\|_2^2}}{2} \right\}.$$

Theorem 6.3. [Local Linear Convergence] *Under Assumptions 5.1 and 6.1, if we take $\eta_m = \eta = 2/(\eta + L')$, and $\kappa = L'/\alpha$, then there exists a neighborhood U of θ^* such that*

$$\|\theta^{(m)} - \theta^*\|_2 \leq \left(\frac{\kappa-1}{\kappa+1} \right)^m \|\theta^{(0)} - \theta^*\|_2,$$

where m is the number of gradient descent iterations.

This result gives a lower bound on the convergence rate near θ^* . Due to its strongly convex property, the convergence rate is fast, which explains the fast and stable optimization process.

Proof Overview. Assumption 6.1 justifies the Jacobian simplification (Lemma C.2), which in turn yields the Hessian block structure (Lemma C.4). By Schur complement, this result gives local strong convexity (Lemma 6.2). Combining with the Lipschitz constant, we finish the proof.

6.2 GENERAL MOG LATENT HESSIAN ANALYSIS AND OPTIMIZATION

We now extend our analysis to the case where each subspace k carries an *asymmetric* Gaussian mixture (Equation 3). As before, we first state the key separation assumption and show that on each subspace, the individual Gaussian distributions in the mixture of Gaussian are highly separated from each other. Then, we simplify the Hessian and prove local convexity. Finally, we conclude a linear convergence rate based on the strongly convex and smooth property.

Assumption 6.4. [Highly Separated Gaussian] Consider the Gaussian mixture

$$p_k(x) = \sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; \mu_{k,l}, \Sigma_{k,l}), \quad r_{k,l}(x) := \frac{\pi_{k,l} \mathcal{N}(x; \mu_{k,l}, \Sigma_{k,l})}{\sum_{i=1}^{n_k} \pi_{k,i} \mathcal{N}(x; \mu_{k,i}, \Sigma_{k,i})}.$$

There exist constants $\varepsilon \ll 1$ and $\delta \ll 1$ such that when $x \sim p_k$ we have

$$\Pr_{x \sim p_k} \left(\exists l \in \{1, \dots, n_k\} \text{ with } r_{k,l}(x) \geq 1 - \varepsilon \right) \geq 1 - \delta.$$

Justification. With MoLR-MoG modeling, after adding diffusion noise of scale γ_t , each point x remains within $O(\gamma_t)$ of the subspace’s moment-matched center $\bar{\mu}_k$. Concretely, the subspace structure (or a preliminary projection onto principal components) ensures $\|x - \bar{\mu}_k\|_2 \leq \Delta = C\gamma_t$ with high probability, for some moderate constant C . Hence, any third-order Taylor term $\propto \|x - \bar{\mu}_k\|_2^3$ is $O(\gamma_t^3)$, which vanishes compared to the leading Hessian scale $O(\gamma_t^2)$. In the following corollary, we further show the approximation effect of equivalent Gaussians.

Corollary 6.5. Assume that $\|\mu_{k,i}^* - \mu_{k,j}^*\|_2 \leq \delta$, $\|U_{k,i}^* - U_{k,j}^*\|_2 \leq \epsilon$ and $\|x - \bar{\mu}_k^*\|_2 \leq \Delta$. We have

$$\|\log p(x) - \log \bar{p}(x)\|_2 = O(\epsilon + \delta\Delta + \Delta^3)$$

Remark 6.6 (Separated Gaussian simplification). For simplicity of description, we assume the individual Gaussian distributions in the mixture of Gaussians are highly separated. Actually, if there are n'_k Gaussians that are not separated from each other, we can employ clustering techniques to transform them into n_k mutually independent Gaussian distributions. The error caused by such an operation can be calculated using corollary 6.5. The core intuition is that the modals should not have much influence on each other. Hence, we can also use the idea of recursion to first cluster the general MoG into a 2-modal MoG latent. Then, we can use the analysis of Section 6.1 with Assumption 6.1.

Then, similar to the above section, we also calculate the Hessian matrix and show the local strong convex parameters. Finally, we provide the convergence guarantee for general MoLR-MoG modeling.

Lemma 6.7. [Eigenvalues of the Hessian] Assume Assumption 6.4, the Hessian at the k -th subspace is convex on a neighborhood of θ^* . If $\forall x \in \mathbb{R}^{d_k}$, $r_k^+(x) = 1$ or -1 are strictly satisfied, we have

$$\lambda_{\min}(H_{\mu_{k,l}\mu_{k,l}}) = \frac{\pi_{k,l}s_t^2}{(s_t^2 + \gamma_t^2)^2},$$

and $\lambda_{\min}(H_{U_{k,l}U_{k,l}})$ has the following form:

$$\left(\pi_{k,l} 4(U_{k,l}^\top \mu_{k,l})^2 + \|U_{k,l}\|_2^2 \|\mu_{k,l}\|_2^2 - \|U_{k,l}\|_2 \|\mu_{k,l}\|_2 \sqrt{8(U_{k,l}^\top \mu_{k,l})^2 + \|U_{k,l}\|_2^2 \|\mu_{k,l}\|_2^2} \right) / 2.$$

Lemma 6.8. [Local Strong Convexity] Assume Assumption 6.4, in a neighborhood of θ^* , $\nabla^2 \mathcal{L}(\theta) \succeq \alpha' I$, $\alpha' > 0$, $\forall \theta \in \Theta$. If $\forall x \in \mathbb{R}^{d_k}$, $\exists l \in [n_k]$, $r_{k,l}(x) = 1$ are strictly satisfied, $\alpha' = \min\{\lambda_1, \lambda_2\}$, where $\lambda_1 = \min_{l=1 \dots, n_k} \frac{c_{k,l}\gamma_t^4}{(s_t^2 + \gamma_t^2)^2}$, $\lambda_2 = \min_{l=1,2,\dots,n_k} \lambda_{\min}(H_{U_{k,l}U_{k,l}})$.

Thus, even without symmetry, equivalent Gaussians and sufficient subspace separation recover the same local convexity and linear convergence guarantees as in the asymmetric case. Similar to Theorem 6.3, under Assumption 6.4, we can obtain a convergence guarantee.

Remark 6.9 (Previous MoG Learning through Score Matching). Shah et al. (2023) and Chen et al. (2024b) consider MoG data and analyze the optimization process of diffusion models at the full space. However, these works aim to design a specific algorithm to learn the MoG distribution instead of using a standard optimization algorithm. On the contrary, by using the MoLR-MoG property to calculate the Hessian matrix, we adopt the GD algorithm and obtain the convergence guarantee.

Remark 6.10 (Initialization). Since the multi-modal GMM latent leads to a highly non-convex landscape, Theorem 6.3 and the corresponding asymmetric variant require the initialization to be around θ^* to guarantee local strong convexity and obtain a local convergence guarantee. As the MoLR-MoG is the first step to model the multi low-dimensional and multi-modal property, we leave the analysis of the global convergence guarantee as an interesting future work.

6.3 ANALYSIS WITHOUT HIGHLY SEPARATED CONDITION

In this part, we extend our analysis to latent MoG with overlap, which is closer to the real-world datasets. We define the pairwise overlap factor $\xi_{i,j}(x)$ between components i and j at the k -th manifold

$$\xi_{i,j}(x) \triangleq r_{k,i}(x)r_{k,j}(x).$$

and the maximum expected overlap for the manifold as: $\epsilon_{\text{overlap}} = \max_i \sum_{j \neq i} \mathbb{E}_- x \sim p_t[\xi_{i,j}(x)]$.

Without the high-separation assumption, our analysis proceeds in two steps. With the overlap factor $\epsilon_{\text{overlap}}$, we first examine the block-diagonal Hessian, deriving a refined lower bound α . Second, we analyze the full Hessian by treating off-diagonal interference as a perturbation bounded by the overlap factor. Applying Weyl's Inequality, we prove that the global matrix remains positive definite provided the perturbation (introduced by the overlap) is smaller than the effective diagonal curvature α , thus guaranteeing linear convergence.

Lemma 6.11 (Minimum Curvature for 2-Mode Mixture). Consider a mixture of two Gaussian components. Let $\epsilon_{\text{overlap}} = \sup_x r_k^+(x)r_k^-(x)$ denote the maximum pointwise overlap factor. The minimum eigenvalue of the ideal Hessian matrix, denoted as $\alpha_{2\text{-mode}}$, is bounded below by:

$$\alpha_{2\text{-mode}} \triangleq (1 - 4\epsilon_{\text{overlap}}) \min(\lambda_{\min}(H_{\mu_k\mu_k}), \lambda_{\min}(H_{U_kU_k})),$$

and

$$\lambda_{\min}(H) \geq \alpha_{2\text{-mode}} - C' \epsilon_{\text{overlap}} > 0,$$

where C' is defined in E.1.3.

Lemma 6.12 (Minimum Curvature for Multi-Modal). *Let $\epsilon_{k,l}^{\text{total}} = \sum_{j \neq l} \mathbb{E}[\xi_{j,l}(x)]$ represent the total probability mass leaking from the l -th component due to overlap. The minimum eigenvalue of the block-diagonal Hessian, denoted as $\alpha_{\text{Multi-Modal}}$, is determined by the component with the minimum effective mass:*

$$\alpha_{\text{Multi-Modal}} \triangleq \min_{l \in \{1, \dots, n_k\}} [(\pi_{k,l} - \epsilon_{k,l}^{\text{total}}) \min(\lambda_{\min}(H_{\mu_{k,l}, \mu_{k,l}}), \lambda_{\min}(H_{U_{k,l}, U_{k,l}}))] ,$$

and

$$\lambda_{\min}(H) \geq \alpha_{\text{Multi-Modal}} - \tilde{C} \cdot \epsilon_{\text{overlap}} ,$$

where \tilde{C} is defined in E.2.4.

For the Hessian to remain positive definite, the intrinsic weight of every cluster must exceed its total confusion with other clusters (i.e., $\pi_{k,l} > \epsilon_{k,l}^{\text{total}}$ for all l).

7 CONCLUSION

In this work, we provide a mixture of low-rank mixture of Gaussian (MoLR-MoG) modeling for target data, which reflects the low-dimensional and multi-modal property of real-world data. Through the real-world experiments, we first show that the MoLR-MoG is a suitable modeling for the real-world data. Then, we analyze the estimation error and optimization process under the MoLR-MoG modeling and explain why diffusion models can achieve great performance with a small training dataset and a fast optimization process.

For the estimation error, we show that with the MoLR-MoG modeling, the estimation error is $R^4 \sqrt{\sum_{k=1}^K n_k} \sqrt{\sum_{k=1}^K n_k d_k} / \sqrt{n}$, which means diffusion models can take fully use of the multi subspace, low-dimensional and multi-modal information to escape the curse of dimensionality. For the optimization process, we conducted a detailed analysis of the score-matching loss landscape. By formulating the exact score in both symmetric and asymmetric mixture settings, we derived explicit expressions for the parameter Jacobians and identified the dominant components under standard separation assumptions. Then, we prove that the population loss becomes strongly convex in a neighborhood of the ground truth score function, by estimating the Hessian and presenting lower bounds on both its minimal eigenvalue and the convergence rate. Then, we provide the local convergence guarantee for the score matching objective function, which explains the fast and stable training process of diffusion models.

Future work and limitation. Though we have extended the situation to multi-manifold MoG, how to extend the analysis to more general non-Gaussian sub-manifolds (e.g. heavy-tailed or multi-modal beyond second moments) by higher-order moment matching is still unknown. Meanwhile, we wish to design optimization algorithms or network architectures that explicitly leverage the block-diagonal Hessian structure for faster training. For example, we can perform a natural-gradient step separately in each block with a block-diagonal Hessian with decomposed data, which will accelerate the optimization process.

Ethics statement. Our work aims to deepen the understanding of the modeling of diffusion models and explain the success of diffusion models from a theoretical perspective. The MoLR-MoG modeling has the potential to achieve a great performance with fewer parameters. Hence, this work can be viewed as an important step in understanding diffusion models, and the societal impact is similar to general generative models (Mirsky and Lee, 2021).

Reproducibility statement. The detail and description of the real-world experiments are provided in Appendix F. We detail the model, hyperparameters and data.

REFERENCES

- Bradley CA Brown, Anthony L Caterini, Brendan Leigh Ross, Jesse C Cresswell, and Gabriel Loaiza-Ganem. Verifying the union of manifolds hypothesis for image data. In *ICLR*, 2023.
- Stefano Bruno, Ying Zhang, Dong-Young Lim, Ömer Deniz Akyildiz, and Sotirios Sabanis. On diffusion-based generative models and their error bounds: The log-concave case with full convergence estimates. *arXiv preprint arXiv:2311.13584*, 2023.
- Hong Chen, Xin Wang, Guanning Zeng, Yipeng Zhang, Yuwei Zhou, Feilin Han, and Wenwu Zhu. Video-dreamer: Customized multi-subject text-to-video generation with disen-mix finetuning. *arXiv preprint arXiv:2311.00990*, 2023a.
- Junyu Chen, Han Cai, Junsong Chen, Enze Xie, Shang Yang, Haotian Tang, Muiyang Li, Yao Lu, and Song Han. Deep compression autoencoder for efficient high-resolution diffusion models. *arXiv preprint arXiv:2410.10733*, 2024a.
- Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. *arXiv preprint arXiv:2302.07194*, 2023b.
- Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru R Zhang. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. *arXiv preprint arXiv:2209.11215*, 2022.
- Sitan Chen, Vasilis Kontonis, and Kulin Shah. Learning general gaussian mixtures with efficient score matching. *arXiv preprint arXiv:2404.18893*, 2024b.
- Frank Cole and Yulong Lu. Score-based generative models break the curse of dimensionality in learning a family of sub-gaussian probability distributions. *arXiv preprint arXiv:2402.08082*, 2024.
- Hugo Cui and Lenka Zdeborová. High-dimensional asymptotics of denoising autoencoders. *arXiv preprint arXiv:2305.11041*, 2023.
- Hugo Cui, Florent Krzakala, Eric Vanden-Eijnden, and Lenka Zdeborová. Analysis of learning a flow-based generative model from limited sample complexity. *arXiv preprint arXiv:2310.03575*, 2023.
- Hengyu Fu, Zhuoran Yang, Mengdi Wang, and Minshuo Chen. Unveil conditional diffusion models with classifier-free guidance: A sharp statistical theory. *arXiv preprint arXiv:2403.11968*, 2024.
- Sixue Gong, Vishnu Naresh Boddeti, and Anil K Jain. On the intrinsic dimensionality of image representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3987–3996, 2019.
- Yingqing Guo, Hui Yuan, Yukang Yang, Minshuo Chen, and Mengdi Wang. Gradient guidance for diffusion models: An optimization perspective. *arXiv preprint arXiv:2404.14743*, 2024.
- Yinbin Han, Meisam Razaviyayn, and Renyuan Xu. Neural network-based score estimation in diffusion models: Optimization and generalization. *arXiv preprint arXiv:2401.15604*, 2024.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022.
- Jerry Yao-Chieh Hu, Weimin Wu, Yi-Chen Lee, Yu-Chao Huang, Minshuo Chen, and Han Liu. On statistical rates of conditional diffusion transformers: Approximation, estimation and minimax optimality. *arXiv preprint arXiv:2411.17522*, 2024a.
- Jerry Yao-Chieh Hu, Weimin Wu, Zhuoru Li, Sophia Pi, Zhao Song, and Han Liu. On statistical rates and provably efficient criteria of latent diffusion transformers (dits). *Advances in Neural Information Processing Systems*, 37:31562–31628, 2024b.
- Hamid Kamkari, Brendan Ross, Rasa Hosseinzadeh, Jesse Cresswell, and Gabriel Loaiza-Ganem. A geometric view of data complexity: Efficient local intrinsic dimension estimation with diffusion models. *Advances in Neural Information Processing Systems*, 37:38307–38354, 2024.
- Puheng Li, Zhong Li, Huishuai Zhang, and Jiang Bian. On the generalization properties of diffusion models. *arXiv preprint arXiv:2311.01797*, 2023.

- Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng Chen, Chong Zeng, Jiayuan Gu, and Hao Su. One-2-3-45++: Fast single image to 3d objects with consistent multi-view generation and 3d diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10072–10083, 2024.
- Yue Ma, Yingqing He, Hongfa Wang, Andong Wang, Chenyang Qi, Chengfei Cai, Xiu Li, Zhifeng Li, Heung-Yeung Shum, Wei Liu, et al. Follow-your-click: Open-domain regional image animation via short prompts. *arXiv preprint arXiv:2403.08268*, 2024.
- Song Mei and Yuchen Wu. Deep networks as denoising algorithms: Sample-efficient learning of diffusion models in high-dimensional graphical models. *arXiv preprint arXiv:2309.11420*, 2023.
- Yisroel Mirsky and Wenke Lee. The creation and detection of deepfakes: A survey. *ACM Computing Surveys (CSUR)*, 54(1):1–41, 2021.
- Kazusato Oko, Shunta Akiyama, and Taiji Suzuki. Diffusion models are minimax optimal distribution estimators. *arXiv preprint arXiv:2303.01861*, 2023.
- Phillip Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. The intrinsic dimension of images and its impact on learning. *arXiv preprint arXiv:2104.08894*, 2021.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- Kulin Shah, Sitan Chen, and Adam Klivans. Learning mixtures of gaussians using the ddpm objective. *arXiv preprint arXiv:2307.01178*, 2023.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7): 1661–1674, 2011.
- Peng Wang, Huijie Zhang, Zekai Zhang, Siyi Chen, Yi Ma, and Qing Qu. Diffusion models learn low-dimensional distributions via subspace clustering. *arXiv preprint arXiv:2409.02426*, 2024.
- Ruofeng Yang, Bo Jiang, Cheng Chen, Ruinan Jin, Baoxiang Wang, and Shuai Li. Few-shot diffusion models escape the curse of dimensionality. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a.
- Ruofeng Yang, Zhijie Wang, Bo Jiang, and Shuai Li. Leveraging drift to improve sample complexity of variance exploding diffusion models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024b.
- Hui Yuan, Kaixuan Huang, Chengzhuo Ni, Minshuo Chen, and Mengdi Wang. Reward-directed conditional diffusion: Provable distribution estimation and reward improvement. *arXiv preprint arXiv:2307.07055*, 2023.

APPENDIX

A THE USE OF LARGE LANGUAGE MODELS (LLMs)

As this work mainly focus on the new modeling of diffusion models from a theoretical perspective, large language models were only used for minor language editing to check grammar. All ideas, new modelings, experiments, theoretical guarantee, discussion and writing decisions were made entirely by the authors.

B SCORE FUNCTION ERROR ESTIMATION

B.1 CALCULATE $\nabla \log p_t(x)$ AND DECOMPOSITION

Consider the k -th subspace

$$p_{t,k}(x) = \sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(\mu_{k,l}, \Sigma_{k,l})$$

where $\Sigma_{k,l} = s_t^2 U_{k,l} U_{k,l}^\top + \gamma_t^2 I$.

We know that

$$\begin{aligned} \Sigma_{k,l}^{-1} &= \frac{1}{\gamma_t^2} \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right), \\ \nabla p_{t,k}(x) &= \frac{1}{\gamma_t^2} \sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(\mu_{k,l}, \Sigma_{k,l}) \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right) (x - \mu_{k,l}), \end{aligned}$$

which indicates

$$\nabla \log p_{t,k}(x) = \frac{\nabla p_{t,k}(x)}{p_{t,k}(x)} = \frac{1}{\gamma_t^2} \frac{\sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(\mu_{k,l}, \Sigma_{k,l}) \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right) (x - \mu_{k,l})}{\sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(\mu_{k,l}, \Sigma_{k,l})}.$$

We want to learn the parameters of the score function:

$$s_k^*(x, t) = \nabla \log p_{t,k}(x),$$

where the parameters are $\{\mu_{k,l}^*, U_{k,l}^*\}, k = 1, \dots, K$.

And

$$s^*(x, t) = (s_1^*(x, t), s_2^*(x, t), \dots, s_K^*(x, t))$$

Define

$$R(s_k) = \mathbb{E} [\|s_k(x, t) - s_k^*(x, t)\|^2], \quad \hat{R}_n(s_k) = \frac{1}{n} \sum_{i=1}^n \|s_k(x_i, t_i) - s_k^*(x_i, t_i)\|^2$$

We have the following decomposition:

$$R(\hat{s}_{k, \hat{\theta}_n}) - \hat{R}_n(s_{k, \hat{\theta}_n}) = \underbrace{R(\hat{s}_{k, \hat{\theta}_n}) - \hat{R}(s_k^*)}_{\text{Estimation}} + \underbrace{\hat{R}(s_k^*) - \hat{R}(s_{k, \theta^*})}_{\text{Approximation}} + \underbrace{\hat{R}_n(s_{k, \theta^*}) - \hat{R}_n(\hat{s}_{k, \hat{\theta}_n})}_{\text{optimization}}$$

We can also obtain that

$$R(s) = \sum_{k=1}^K R(s_k)$$

Since *Estimation* and *Approximation* reflect the fitting ability of the network, we analyze the first term first. Then, in the next section, we analyze the optimization dynamic.

B.2 ESTIMATION

First, we show that f and loss function are Lipschitz. We will first prove that s_k is Lipschitz for $\forall k$, then we can know that s is Lipschitz.

Lemma B.1. [Lipschitz Continuity] Let L_{μ_l} and L_{U_k} be the Lipschitz constant w.r.t. s_θ . With MoLR-MoG modeling and Assumption 5.1, there is a constant

$$L \leq \sqrt{\sum_{i=1}^K n_k (L_{\mu_l}^2 + L_{U_k}^2)} = O\left((\sum_{k=1}^K n_k)^{\frac{1}{2}} C_w\right)$$

such that for any θ, θ' , $\|s_\theta(x, t) - s_{\theta'}(x, t)\|_2 \leq L \|\theta - \theta'\|_2$, where $C_w = \frac{(R+s_t B_\mu)^3 s_t^2}{\gamma_t^4}$, $B_\mu = \max_{k,l} \|\mu_{k,l}\|_2$. For s_θ and s^* , we have that $2\|s_\theta(x, t) - s^*(x, t)\|_2 \leq 2(R + s_t B_\mu)/\gamma_t^2 := L_t$.

Proof. Since we analyze the estimation error at a fixed time t , we ignore subscript t for $\Sigma_{k,l,t}$, $w_{k,t}$, $w_{l,k,t}$ and $\delta_{k,l,t}$ and define by

$$\begin{aligned}\Sigma_{k,l} &= s_t^2 U_{k,l} U_{k,l}^\top + \gamma_t^2 I \\ w_k(x) &= \sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}, \Sigma_{k,l}) \\ w_{k,l} &= \frac{1}{M} \pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}, \Sigma_{k,l}) \\ \delta_{k,l}(x) &= x + s_t \mu_{k,l} - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top (x + s_t \mu_{k,l}).\end{aligned}$$

Assume that $\|U_{k,l}\|_2 \leq B_U$, $\|\mu_{k,l}\|_2 \leq B_\mu$, $\max\{B_U, B_\mu\} = C$, and $\|x\|_2 \leq R$ for $\forall x \in X$.

For $\Sigma_{k,l}$, we know that

$$\Sigma_{k,l} = U_{k,l} U_{k,l}^\top + \gamma_t^2 I \succ \gamma_t^2 I \Rightarrow \lambda_{\min}(\Sigma_{k,l}) \geq \gamma_t^2 \Rightarrow \|\Sigma_{k,l}^{-1}\|_2 \leq \frac{1}{\gamma_t^2}.$$

To obtain the first L in this lemma, we need to bound $\left\| \frac{\partial s_{k,\theta}(x,t)}{\partial \mu_{k,l}} \right\|_2$ and $\left\| \frac{\partial s_{k,\theta}(x,t)}{\partial U_{k,l}} \right\|_2$.

The bound of $\left\| \frac{\partial s_{k,\theta}(x,t)}{\partial \mu_{k,l}} \right\|_2$. For the latent score of the k -th subspace, we have that

$$\begin{aligned}s_{k,\theta}(x, t) &= -\frac{1}{\gamma_t^2} \frac{\sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x)}{w_k(x)}, \\ \frac{\partial s_{k,\theta}(x, t)}{\partial \mu_{k,l}} &= -\frac{1}{\gamma_t^2} \frac{\sum_{l=1}^{n_k} (\frac{\partial w_{k,l}(x)}{\partial \mu_{k,l}} \delta_{k,l}(x) + \frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}} w_{k,l}(x)) w_k(x) - \frac{\partial w_k(x)}{\partial \mu_{k,l}} (\sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x))}{w_k^2(x)}, \\ \left\| \frac{\partial s_{k,\theta}(x, t)}{\partial \mu_{k,l}} \right\|_2 &\leq \frac{1}{\gamma_t^2} \left(\left\| \frac{\sum_{l=1}^{n_k} (\frac{\partial w_{k,l}(x)}{\partial \mu_{k,l}} \delta_{k,l}(x) + \frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}} w_{k,l}(x))}{w_k(x)} \right\|_2 + \left\| \frac{\frac{\partial w_k(x)}{\partial \mu_{k,l}} (\sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x))}{w_k^2(x)} \right\|_2 \right).\end{aligned}$$

To bound this term, we separately show that

(1) $w_k(x)$ has a lower bound.

(2) $w_{k,l}(x)$, $\delta_{k,l}(x)$, $\frac{\partial w_{k,l}(x)}{\partial \mu_{k,l}}$, $\frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}}$ have upper bounds.

(3) $\left\| \frac{\frac{\partial w_k(x)}{\partial \mu_{k,l}} \delta_{k,l}(x)}{w_k} \right\|_2$, $\left\| \frac{\sum_{l=1}^{n_k} \frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}} w_{k,l}(x)}{w_k(x)} \right\|_2$, $\left\| \frac{\frac{\partial w_k(x)}{\partial \mu_{k,l}} \sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x)}{w_k^2(x)} \right\|_2$ have upper bounds.

(1) $w_k(x)$ has a lower bound.

$$w_k(x) = \sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}, \Sigma_{k,l}), \text{ which is continuous.}$$

Since continuous function has maximum and minimum in a closed interval and $\|x\|_2 \leq R$, we can assume that $w_k(x) \geq m_w$. And for any x , $w_k(x) > 0$, so $m_w > 0$ holds.

(2) $w_{k,l}(x)$, $\delta_{k,l}(x)$, $\frac{\partial \delta_{k,l}(x)}{\partial \mu_k}$, $\frac{\partial w_{k,l}(x)}{\partial \mu_k}$ have upper bounds.

We already know that continuous function has maximum and minimum in a closed interval and $\|x\|_2 \leq R$. Thus, we can assume that $w_k(x) \leq M_{w_k}$. We also have that

$$w_k(x) \leq M_{w_k} \leq \sum_{l=1}^{n_k} \pi_{k,l} (2\pi)^{-\frac{n}{2}} |\Sigma_{k,l}|^{-\frac{1}{2}}.$$

For the second term, we have that

$$\delta_{k,l}(x) = x - s_t \mu_{k,l} - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top (x - s_t \mu_{k,l}) = \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right) (x - s_t \mu_{k,l}),$$

whose L_2 norm is bounded by

$$\left\| \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right) (x - s_t \mu_{k,l}) \right\|_2 \leq \|x - s_t \mu_{k,l}\|_2 \leq \|x\|_2 + \|s_t \mu_{k,l}\|_2 \leq R + s_t B_\mu.$$

Then, for the third term, we know that

$$\frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}} = -s_t + \frac{s_t^3}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top = -s_t \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right).$$

For the last term, we have we have the following expression

$$\frac{\partial w_{k,l}(x)}{\partial \mu_{k,l}} = -\frac{s_t}{2} \mathcal{N}(x; s_t \mu_{k,l}, \Sigma_{k,l}) \Sigma_{k,l}^{-1} (x - s_t \mu_{k,l}).$$

For term $\|\Sigma_{k,l}^{-1} (x - s_t \mu_{k,l})\|_2$, we have that

$$\|\Sigma_{k,l}^{-1} (x - s_t \mu_{k,l})\|_2 \leq \|\Sigma_{k,l}^{-1}\|_2 \|x - s_t \mu_{k,l}\|_2 = \frac{1}{\gamma_t^2} \|x - s_t \mu_{k,l}\|_2 \leq \frac{1}{\gamma_t^2} (R + \|s_t \mu_{k,l}\|_2),$$

which indicates

$$\begin{aligned} \left\| \frac{\partial w_{k,l}(x)}{\partial \mu_{k,l}} \right\|_2 &\leq s_t \mathcal{N}(x; s_t \mu_{k,l}, \Sigma_{k,l}) \frac{1}{\gamma_t^2} (R + \|s_t \mu_{k,l}\|_2) \leq s_t \mathcal{N}(x; s_t \mu_{k,l}, \Sigma_{k,l}) \frac{1}{\gamma_t^2} (R + s_t B_\mu) \\ \left\| \frac{\partial w_k(x)}{\partial \mu_{k,l}} \right\|_2 &\leq \sum_{l=1}^{n_k} s_t \mathcal{N}(x; s_t \mu_{k,l}, \Sigma_{k,l}) \frac{1}{\gamma_t^2} (R + s_t B_\mu). \end{aligned}$$

$$(3) \left\| \frac{\frac{\partial w_k(x)}{\partial \mu_{k,l}} \delta_{k,l}(x)}{w_k} \right\|_2, \left\| \frac{\sum_{l=1}^{n_k} \frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}} w_{k,l}(x)}{w_k(x)} \right\|_2, \left\| \frac{\frac{\partial w_k(x)}{\partial \mu_{k,l}} \sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x)}{w_k^2(x)} \right\|_2 \text{ have upper bounds.}$$

For the first two term,

$$\left\| \frac{\frac{\partial w_k(x)}{\partial \mu_{k,l}} \delta_{k,l}(x)}{w_k} \right\|_2 \leq \frac{s_t}{\gamma_t^2} (R + s_t B_\mu)^2,$$

and

$$\left\| \frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}} \right\|_2 = \text{Constant} \leq s_t, \left\| \frac{\sum_{l=1}^{n_k} \frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}} w_{k,l}(x)}{w_k(x)} \right\|_2 \leq s_t.$$

For the third term, we know that

$$\left\| \frac{\frac{\partial w_k(x)}{\partial \mu_{k,l}} \sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x)}{w_k^2(x)} \right\|_2 \leq \left\| \frac{s_t w_k^2(x) \frac{s_t}{\gamma_t^2} (R + s_t B_\mu)}{w_k^2(x)} \right\|_2 = \frac{s_t^2}{\gamma_t^2} (R + s_t B_\mu).$$

Combined with the above three, we obtain the bound for $\left\| \frac{\partial s_{k,\theta}(x,t)}{\partial \mu_{k,l}} \right\|_2$:

$$\begin{aligned} \left\| \frac{\partial s_{k,\theta}(x,t)}{\partial \mu_{k,l}} \right\|_2 &\leq \frac{1}{\gamma_t^2} \left(\left\| \frac{\sum_{l=1}^{n_k} \left(\frac{\partial w_{k,l}(x)}{\partial \mu_{k,l}} + \frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}} \right) \delta_{k,l}(x)}{w_k(x)} \right\|_2 + \left\| \frac{\frac{\partial w_k(x)}{\partial \mu_{k,l}} (\sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x))}{w_k^2(x)} \right\|_2 \right) \\ &\leq \frac{s_t^2}{\gamma_t^2} (R + s_t B_\mu)^2 + s_t + \frac{s_t}{\gamma_t^2} (R + s_t B_\mu) = O\left(\frac{s_t^2 (R + s_t B_\mu)^2}{\gamma_t^2}\right). \end{aligned}$$

The bound of $\left\| \frac{\partial s_{k,\theta}(x,t)}{\partial U_{k,l}} \right\|_2$. Now we compute the part about $U_{k,l}$. Through some simple algebra, we know that

$$\frac{\partial s_{k,\theta}(x,t)}{\partial U_{k,l}} = -\frac{1}{\gamma_t^2} \frac{\sum_{l=1}^{n_k} \left(\frac{\partial w_{k,l}(x)}{\partial U_{k,l}} \delta_{k,l}(x) + \frac{\partial \delta_{k,l}(x)}{\partial U_{k,l}} w_{k,l}(x) \right) w_k(x) - \frac{\partial w_k(x)}{\partial U_{k,l}} \left(\sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x) \right)}{w_k^2(x)}.$$

Then, we have the following inequality

$$\begin{aligned} \frac{\partial s_{k,\theta}(x,t)}{\partial U_{k,l}} &= -\frac{1}{\gamma_t^2} \frac{\sum_{l=1}^{n_k} \left(\frac{\partial w_{k,l}(x)}{\partial U_{k,l}} \delta_{k,l}(x) + \frac{\partial \delta_{k,l}(x)}{\partial U_{k,l}} w_{k,l}(x) \right) w_k(x) - \frac{\partial w_k(x)}{\partial U_{k,l}} \left(\sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x) \right)}{w_k^2(x)} \\ \left\| \frac{\partial s_{k,\theta}(x,t)}{\partial U_{k,l}} \right\|_2 &\leq \frac{1}{\gamma_t^2} \left(\left\| \frac{\sum_{l=1}^{n_k} \left(\frac{\partial w_{k,l}(x)}{\partial U_{k,l}} \delta_{k,l}(x) + \frac{\partial \delta_{k,l}(x)}{\partial U_{k,l}} w_{k,l}(x) \right)}{w_k(x)} \right\|_2 + \left\| \frac{\frac{\partial w_k(x)}{\partial U_{k,l}} * \left(\sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x) \right)}{w_k^2(x)} \right\|_2 \right). \end{aligned}$$

Similar with $\left\| \frac{\partial s_{k,\theta}(x,t)}{\partial \mu_{k,l}} \right\|_2$, we need to provide:

(1) The upper bound of $\frac{\partial w_{k,l}}{\partial U_{k,l}}$ and $\frac{\partial \delta_{k,l}}{\partial U_{k,l}}$,

(2) The upper bound of $\left\| \frac{\sum_{l=1}^{n_k} \left(\frac{\partial w_{k,l}(x)}{\partial U_{k,l}} \delta_{k,l}(x) + \frac{\partial \delta_{k,l}(x)}{\partial U_{k,l}} w_{k,l}(x) \right)}{w_k(x)} \right\|_2$ and $\left\| \frac{\frac{\partial w_k(x)}{\partial U_{k,l}} * \left(\sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x) \right)}{w_k^2(x)} \right\|_2$.

(1) The upper bound of $\frac{\partial w_{k,l}}{\partial U_{k,l}}$ and $\frac{\partial \delta_{k,l}}{\partial U_{k,l}}$.

For the first term, we have the following form

$$\begin{aligned} \frac{\partial w_{k,l}}{\partial U_{k,l}} &= \pi_{k,l} \frac{\partial \mathcal{N}(x; s_t \mu_{k,l}, \Sigma_{k,l})}{\partial U_{k,l}} \\ &= 2\pi_{k,l} s_t^2 [\mathcal{N}(x; s_t \mu_{k,l}, \Sigma_{k,l}) (\Sigma_k^{-1} (x - s_t \mu_{k,l}) (x - s_t \mu_{k,l})^\top \Sigma_{k,l}^{-1} - \Sigma_{k,l}^{-1})] U_{k,l}. \end{aligned}$$

Then, we know that

$$\begin{aligned} \left\| \frac{\partial w_{k,l}}{\partial U_{k,l}} \right\|_2 &\leq 2\pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}, \Sigma_{k,l}) s_t^2 \left(\frac{(R + s_t \|\mu_{k,l}\|_2)^2}{\gamma_t^4} + \frac{1}{\gamma_t^2} \right) \\ &\leq 2\pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}, \Sigma_{k,l}) s_t^2 \left(\frac{(R + s_t B_\mu)^2}{\gamma_t^4} + \frac{1}{\gamma_t^2} \right). \end{aligned}$$

For the second term, we have that

$$\frac{\partial \delta_{k,l}(x)}{\partial U_{k,l}} = -2 \frac{s_t^2}{s_t^2 + \gamma_t^2} (U_{k,l}^\top (x - s_t \mu_{k,l}) I + U_{k,l} (x - s_t \mu_{k,l})^\top),$$

which indicates

$$\begin{aligned} \left\| \frac{\partial \delta_{k,l}(x)}{\partial U_{k,l}} \right\|_2 &\leq 2 \frac{s_t^2}{s_t^2 + \gamma_t^2} (R + \|s_t \mu_{k,l}\|_2) \leq 2(R + \|s_t \mu_{k,l}\|_2) \\ &\leq 2(R + s_t B_\mu). \end{aligned}$$

(2) The upper bound of $\left\| \frac{\sum_{l=1}^{n_k} \left(\frac{\partial w_{k,l}(x)}{\partial U_{k,l}} \delta_{k,l}(x) + \frac{\partial \delta_{k,l}(x)}{\partial U_{k,l}} w_{k,l}(x) \right)}{w_k(x)} \right\|_2$ and $\left\| \frac{\frac{\partial w_k(x)}{\partial U_{k,l}} * \left(\sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x) \right)}{w_k^2(x)} \right\|_2$.

$$\left\| \frac{\sum_{l=1}^{n_k} \left(\frac{\partial w_{k,l}(x)}{\partial U_{k,l}} \delta_{k,l}(x) + \frac{\partial \delta_{k,l}(x)}{\partial U_{k,l}} w_{k,l}(x) \right)}{w_k(x)} \right\|_2 \leq s_t^2 \left(\frac{(R + s_t B_\mu)^3}{\gamma_t^4} + \frac{1}{\gamma_t^2} \right) + 2(R + s_t B_\mu)$$

We also have

$$\left\| \frac{\frac{\partial w_k(x)}{\partial U_{k,l}} \left(\sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x) \right)}{w_k^2(x)} \right\|_2 \leq s_t^2 \left(\frac{(R + s_t B_\mu)^2}{\gamma_t^4} + \frac{1}{\gamma_t^2} \right) (R + s_t B_\mu)$$

$$\begin{aligned} \left\| \frac{\partial s_{k,\theta}(x,t)}{\partial U_{k,l}} \right\|_2 &\leq s_t^2 \left(\frac{(R + s_t B_\mu)^2}{\gamma_t^4} + \frac{1}{\gamma_t^2} \right) + 2(R + s_t B_\mu) + s_t^2 \left(\frac{(R + s_t B_\mu)^2}{\gamma_t^4} + \frac{1}{\gamma_t^2} \right) (R + s_t B_\mu) \\ &= O \left(\frac{(R + s_t B_\mu)^3 s_t^2}{\gamma_t^4} \right). \end{aligned}$$

Therefore, $s_{\theta,k}$ is L_k -lipshiz, where

$$L_k \leq \sqrt{n_k(L_{\mu_{k,l}}^2 + L_{U_{k,l}}^2)} = O \left(n_k^{\frac{1}{2}} \frac{(R + s_t B_\mu)^3 s_t^2}{\gamma_t^4} \right).$$

Furthermore, we know that

$$\|s_\theta(x) - s_\theta(y)\|_2 = \left(\sum_{i=1}^K \left\| s_{\theta,i}(x^{(i)}) - s_{\theta,i}(y^{(i)}) \right\|^2 \right)^{\frac{1}{2}} \leq \left(\sum_{i=1}^K L_i \|(x^{(i)} - y^{(i)})\|_2^2 \right)^{\frac{1}{2}} \leq \sqrt{\sum_{i=1}^K L_i^2} \|x - y\|_2.$$

Thus,

$$L = \sqrt{\sum_{i=1}^k L_i^2} = O \left(\sqrt{\sum_{i=1}^k n_i^{\frac{1}{2}} \frac{(R + s_t B_\mu)^3 s_t^2}{\gamma_t^4}} \right).$$

After obtaining the Lipschitz constant for s_θ , we bound the gap between s_θ and s^* :

$$\nabla \log p_{t,k}(x) = -\frac{1}{\gamma_t^2} \frac{\sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}, s_t^2 U_{k,l}^* U_{k,l}^{*\top} + \gamma_t^2 I) \left(x - s_t \mu_{k,l} - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l}^* U_{k,l}^{*\top} (x - s_t \mu_{k,l}) \right)}{\sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}, s_t^2 U_{k,l}^* U_{k,l}^{*\top} + \gamma_t^2 I)}.$$

With the following bound

$$\|x - s_t \mu_{k,l} - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l}^* U_{k,l}^{*\top} (x - s_t \mu_{k,l})\|_2 \leq R + s_t B_\mu,$$

we have that

$$\|\nabla \log p_{t,k}(x)\|_2 \leq \frac{1}{\gamma_t^2} (R + s_t B_\mu), \text{ and } \|s_{k,\theta}(x)\|_2 \leq \frac{1}{\gamma_t^2} (R + s_t B_\mu),$$

which indicates

$$\|s_{k,\theta}(x) - \nabla \log p_{t,k}(x)\|_2 \leq \frac{2}{\gamma_t^2} (R + s_t B_\mu).$$

Hence, we obtain that

$$L_l \leq 2\|s_{k,\theta}(x) - \nabla \log p_{t,k}(x)\|_2 = O(R + s_t B_\mu).$$

■

Lemma B.2. [Rademacher Complexity] Let $\mathcal{F} = \{\ell(\theta; \cdot, \cdot) : \theta \in \Theta\}$ and suppose Θ has diameter R_Θ . Then the empirical Rademacher complexity satisfies

$$\hat{\mathfrak{R}}_n(\mathcal{F}) = O\left(L' \sqrt{\frac{p}{n}}\right).$$

Proof. Let function class $\mathcal{F} = \{s_\theta(x) : \theta = (\{\{\mu_{k,l}, U_{k,l}\}_{l=1}^{n_k}\}_{k=1}^K) \in \Theta\}$, where $\mu_{k,l} \in \mathbb{R}^d, U_{k,l} \in \mathbb{R}^d$

We know that the number of parameters

$$p = \sum_{k=1}^K n_k(d + d) = 2\sum_{k=1}^K n_k d_k.$$

And the covering number of the parameter space is

$$\mathcal{N}(\epsilon, \Theta, \|\cdot\|_2) \leq \left(\frac{C}{\epsilon}\right)^p$$

If f is L -lipschitz, we know that

$$\begin{aligned} \forall \theta_1, \theta_2 \in \Theta, \|f_{\theta_1} - f_{\theta_2}\|_{L_2(p)} &\leq L\|\theta_1 - \theta_2\|_2 \quad \text{and} \quad \forall \theta, \exists \theta_j, \text{ s.t. } \|\theta - \theta_j\|_2 \leq \frac{\epsilon}{L} \\ \Rightarrow \|f_{\theta} - f_{\theta_j}\|_{L_2(p)} &\leq L\|\theta - \theta_j\|_2 \leq \epsilon. \end{aligned}$$

Thus, assume that $\|\theta_i - \theta_j\|_2 \leq C_1$ for any $\theta_i, \theta_j \in \Theta$

$$\mathcal{N}(\epsilon, \Theta, \|\cdot\|_2) \leq \left(\frac{C_1}{\epsilon}\right)^p$$

$$\Rightarrow \mathcal{N}\left(\frac{\epsilon}{L}, \Theta, \|\cdot\|_2\right) \leq \left(\frac{C_1 L}{\epsilon}\right)^p$$

$$\Rightarrow \mathcal{N}(\epsilon, \mathcal{F}, \|\cdot\|_{L_2(p)}) \leq \mathcal{N}\left(\frac{\epsilon}{L}, \Theta, \|\cdot\|_2\right) \leq \left(\frac{C_1 L}{\epsilon}\right)^p \leq \left(\frac{C_1 L}{\epsilon}\right)^p, \quad \log \mathcal{N}\left(\frac{\epsilon}{L}, \mathcal{F}, \|\cdot\|_{L_2(p)}\right) \leq p \log\left(\frac{C_1 L}{\epsilon}\right).$$

We also know that $\text{diam}(\mathcal{F}) \leq L \text{diam}(\Theta) = C_1 L$, with Dudley integral, we have

$$\begin{aligned} \mathcal{R}_n(\mathcal{F}) &\leq \frac{12}{\sqrt{n}} \int_0^{\text{diam}(\mathcal{F})} \sqrt{\log N(\epsilon, \mathcal{F}, \|\cdot\|_{L_2(p)})} d\epsilon \\ &\leq \frac{12}{\sqrt{n}} \int_0^{C_1 L} \sqrt{p \log\left(\frac{C_1 L}{\epsilon}\right)} d\epsilon \\ &\leq \frac{12}{\sqrt{n}} \int_0^\infty p C L \sqrt{t} \exp(-t) dt = \frac{6\sqrt{\pi p}}{\sqrt{n}} C_1 L = O(C_1 L \sqrt{\frac{p}{n}}). \end{aligned}$$

We take the squared loss function.

$$\mathcal{R}_n(\mathcal{L}) \leq L_l \mathcal{R}_n(\mathcal{F}) = O(C_1 L_l L \sqrt{\frac{p}{n}}).$$

Theorem 5.3. Denote by $\hat{\mathcal{L}}_n(\theta)$ the empirical loss on n i.i.d. samples and by $\mathcal{L}(\theta)$ its population counterpart. Then there exist constants C_1, C_2 such that with probability at least $1 - \delta$, for all $\theta \in \Theta$,

$$|\mathcal{L}(\theta) - \hat{\mathcal{L}}_n(\theta)| \leq O\left(C_1 \frac{(R + s_t B_\mu)^4 s_t^2 \sqrt{\Sigma_{k=1}^K n_k}}{\gamma_t^6} \sqrt{\frac{\Sigma_{k=1}^K n_k d_k}{n}} + C_2 \sqrt{\frac{\log(1/\delta)}{n}}\right).$$

where $C_1 = \max_{\theta \in \Theta} \|\theta_i - \theta_j\|_2$, $C_2 = \sigma \log 2$, $\sigma^2 = \sup_{\theta \in \Theta} \text{Var}[\ell(\theta; X, t)]$.

Proof. Since

$$L_l \mathcal{R}_n(\mathcal{F}) = O(C_1 L_l L \sqrt{\frac{p}{n}}).$$

We have

$$\begin{aligned} \Delta &= \sup_{\theta \in \Theta} |\hat{\mathcal{L}}_n(\theta) - L(\theta)| = O(C_1 L_l L \sqrt{\frac{p}{n}}) \\ \Rightarrow \mathbb{E}[\Delta] &= O(C_1 L_l L \sqrt{\frac{p}{n}}). \end{aligned}$$

By Bernstein inequality, let $\sigma^2 = \sup_{\theta \in \Theta} \text{Var}[\ell(X; \theta)]$, we know that

$$\Pr(\sup_{\theta \in \Theta} |\hat{\mathcal{L}}_n(\theta) - L(\theta)| \geq \mathbb{E}[\Delta] + \epsilon) \leq 2 \exp\left(-\frac{n\epsilon^2}{2(\sigma^2 + L_l L C_1 \epsilon/3)}\right) \leq 2 \exp\left(-\frac{n\epsilon^2}{3\sigma^2}\right).$$

Let $2 \exp(-\frac{n\epsilon^2}{3\sigma^2}) < \delta$, we can obtain that

$$\Pr(\sup_{\theta \in \Theta} |\hat{\mathcal{L}}_n(\theta) - L(\theta)| \geq C_1 L L_l \sqrt{\frac{p}{n}} + C_2 \sqrt{\frac{\log(1/\delta)}{n}}) \leq \delta.$$

B.3 APPROXIMATION

Since our network can represent $\nabla \log p(x)$ strictly, we have

$$\text{Approximation Error} = 0$$

C 2-MODE MOG OPTIMIZATION

C.1 SETTING

In this section, we analyze

$$\nabla \log p_{t,k}(x) = \frac{\nabla p_{t,k}(x)}{p_{t,k}(x)} = -\frac{1}{\gamma_t^2} \frac{\frac{1}{2}\mathcal{N}(x; s_t\mu_k, s_t^2 U_k^* U_k^{*\top} + \gamma_t^2 I) \left(x - s_t\mu_k - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k^* U_k^{*\top} (x - s_t\mu_k) \right) + \frac{1}{2}\mathcal{N}(x; -s_t\mu_k, s_t^2 U_k^* U_k^{*\top} + \gamma_t^2 I) \left(x + s_t\mu_k - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k^* U_k^{*\top} (x + s_t\mu_k) \right)}{\frac{1}{2}\mathcal{N}(x; s_t\mu_k, s_t^2 U_k^* U_k^{*\top} + \gamma_t^2 I) + \frac{1}{2}\mathcal{N}(x; -s_t\mu_k, s_t^2 U_k^* U_k^{*\top} + \gamma_t^2 I)},$$

which can be reduced to

$$\nabla \log p_{t,k}(x) = -\frac{1}{\gamma_t^2} \frac{\frac{1}{2}\mathcal{N}(x; s_t\mu_k, \Sigma_k) \delta'_k(x) + \frac{1}{2}\mathcal{N}(x; -s_t\mu_k, \Sigma_k) \epsilon_k(x)}{\frac{1}{2}\mathcal{N}(x; s_t\mu_k, \Sigma_k) + \frac{1}{2}\mathcal{N}(x; -s_t\mu_k, \Sigma_k)}, \quad (5)$$

where $\epsilon_k(x) = x - s_t\mu_k - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k^* U_k^{*\top} (x - s_t\mu_k)$, and $\delta'_k(x) = x + s_t\mu_k - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k^* U_k^{*\top} (x + s_t\mu_k)$.

C.2 OPTIMIZATION

Assumption C.1. [Separation within a cluster] Within each cluster k , the two symmetric peaks are well separated in the sense that $\|s_t\mu_k^* - (-s_t\mu_k^*)\| \geq \Delta_{\text{intra}}$, for some $\Delta_{\text{intra}} \gg \gamma_t$. Consequently, if a sample x is drawn from the “+” peak then its responsibility under the “−” peak satisfies

$$r_k^-(x) = \frac{\frac{1}{2}\mathcal{N}(x; -s_t\mu_k^*, \Sigma_k^*)}{\frac{1}{2}\mathcal{N}(x; s_t\mu_k^*, \Sigma_k^*) + \frac{1}{2}\mathcal{N}(x; -s_t\mu_k^*, \Sigma_k^*)} = O(e^{-\Delta_{\text{intra}}^2/(2\gamma_t^2)}) \ll 1,$$

and symmetrically $r_k^+(x) \ll 1$ when x is drawn from the “−” peak.

In the following discussion, we assume that $x \in k$ -th manifold, which means that $w_i(x) = 0$ if $i \neq k$.

Lemma C.2. [Jacobian Simplification] Under Assumption 6.1, in a neighborhood of θ^* the first derivatives simplify to their “self-cluster” terms: $J_k^\mu(x) = \partial_{\mu_k} s_\theta \approx s_t(I - \alpha P_k)/\gamma_t^2$, and

$$J_k^U(x) \approx \frac{2s_t^2}{\gamma_t^2(s_t^2 + \gamma_t^2)} (r_k^-(x)(U_k^\top(x + s_t\mu_k)I + (x + s_t\mu_k)U_k^\top) + r_k^+(x)(U_k^\top(x - s_t\mu_k)I + U_k(x - s_t\mu_k)^\top)).$$

Proof.

$$\begin{aligned} J_k^\mu &= -\frac{1}{\gamma_t^2} \frac{\frac{\partial w_k^-(x)}{\partial \mu_k} \delta'_k(x) + \frac{\partial w_k^+(x)}{\partial \mu_k} \epsilon_k(x) + \frac{\partial \delta'_k(x)}{\partial \mu_k} w_k^-(x) + \frac{\partial \epsilon_k(x)}{\partial \mu_k} w_k^+(x) * \sum_{k=1}^K w_k(x) - \sum_{k=1}^K \frac{\partial w_k(x)}{\partial \mu_k} * \sum_{k=1}^K (w_k^-(x) \delta'_k(x) + w_k^+(x) \epsilon_k(x))}{w_k^2(x)} \\ &= \underbrace{\frac{w_k^-(x) \frac{\partial \delta'_k(x)}{\partial \mu_k} + w_k^+(x) \frac{\partial \epsilon_k(x)}{\partial \mu_k}}{\gamma_t^2 w_k(x)} - \frac{\frac{\partial w_k^-(x)}{\partial \mu_k} \delta'_k(x) + \frac{\partial w_k^+(x)}{\partial \mu_k} \epsilon_k(x)}{\gamma_t^2 w_k(x)}}_{\text{Term A}} \\ &\quad + \underbrace{\frac{\frac{\partial w_k(x)}{\partial \mu_k} (w_k^-(x) \delta'_k(x) + w_k^+(x) \epsilon_k(x))}{\gamma_t^2 w_k^2(x)}}_{\text{Term B}}. \end{aligned}$$

We will now prove that term B can be ignored compared to term A under our assumptions.

For term B, we have

$$\begin{aligned}
& \frac{\frac{\partial w_k(x)}{\partial \mu_k} (w_k^-(x) \delta'_k(x) + w_k^+(x) \epsilon_k(x))}{\gamma_t^2 w_k^2(x)} - \frac{\frac{\partial w_k^-(x)}{\partial \mu_k} \delta'_k(x) + \frac{\partial w_k^+(x)}{\partial \mu_k} \epsilon_k(x)}{\gamma_t^2 w_k(x)} \\
&= \frac{1}{\gamma_t^2 w_k^2(x)} \left(\frac{\partial w_k(x)}{\partial \mu_k} (w_k^-(x) \delta'_k(x) + w_k^+(x) \epsilon_k(x)) - w_k(x) \left(\frac{\partial w_k^-(x)}{\partial \mu_k} \delta'_k(x) + \frac{\partial w_k^+(x)}{\partial \mu_k} \epsilon_k(x) \right) \right) \\
&= \frac{1}{\gamma_t^2 w_k^2(x)} \left(\frac{\partial w_k^+(x)}{\partial \mu_k} w_k^-(x) \delta'_k(x) + \frac{\partial w_k^-(x)}{\partial \mu_k} w_k^+(x) \epsilon_k(x) - w_k^+(x) \frac{\partial w_k^-(x)}{\partial \mu_k} \delta'_k(x) - w_k^-(x) \frac{\partial w_k^+(x)}{\partial \mu_k} \epsilon_k(x) \right) \\
&= \frac{1}{\gamma_t^2 w_k^2(x)} \left(\frac{\partial w_k^+}{\partial \mu_k} w_k^- - \frac{\partial w_k^-}{\partial \mu_k} w_k^+ \right) (\epsilon_k(x) - \delta'_k(x)) \\
&= -\frac{2}{\gamma_t^2 w_k^2(x)} \left(\frac{\partial w_k^+}{\partial \mu_k} w_k^- - \frac{\partial w_k^-}{\partial \mu_k} w_k^+ \right) \left(I + \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k U_k^\top \right) s_t \mu_k \\
&= -\frac{4}{\gamma_t^2 w_k^2(x)} s_t^2 w_k^- w_k^+ \Sigma_k^{-1} x \left(I + \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k U_k^\top \right) \mu_k = O\left(\frac{r_k^+ r_k^-}{\gamma_t^4} s_t \|\mu_k\|_2 \|x\|_2\right).
\end{aligned}$$

And for term A, we have

$$\frac{w_k^-(x) \frac{\partial \delta'_k(x)}{\partial \mu_k} + w_k^+(x) \frac{\partial \epsilon_k(x)}{\partial \mu_k}}{\gamma_t^2 w_k(x)} = O\left(\frac{s_t \|\mu_k\|_2}{\gamma_t^2} |w_k^+ - w_k^-|\right).$$

Thus,

$$\frac{O\left(\frac{r_k^+ r_k^-}{\gamma_t^4} s_t \|\mu_k\|_2 \|x\|_2\right)}{O\left(\frac{s_t \|\mu_k\|_2}{\gamma_t^2} |w_k^+ - w_k^-|\right)} = O\left(\frac{r_k^+ r_k^- w_k \|x\|_2}{\gamma_t^2 |r_k^+ - r_k^-|}\right) = O\left(\frac{r_k^+ r_k^- w_k \|x\|_2}{\gamma_t^2}\right) \rightarrow 0.$$

$$\text{Thus, } J_k^\mu \approx -\frac{1}{\gamma_t^2} (r_k^+(x) \frac{\partial \delta'_k(x)}{\partial \mu_k} + r_k^-(x) \frac{\partial \epsilon_k(x)}{\partial \mu_k}) = -\frac{s_t}{\gamma_t^2} (r_k^+(x) - r_k^-(x)) \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k U_k^\top \right).$$

We will analyze J_k^U now.

$$\begin{aligned}
J_k^U &= -\frac{1}{\gamma_t^2} \frac{\left(\frac{\partial w_k^-(x)}{\partial U_k} \delta'_k(x) + \frac{\partial \delta'_k(x)}{\partial U_k} w_k^-(x) + \frac{\partial \epsilon_k(x)}{\partial U_k} w_k^+(x) + \frac{\partial w_k^+(x)}{\partial U_k} \epsilon_k(x) \right) w_k(x) - \frac{\partial w_k(x)}{\partial U_k} * (w_k^-(x) \delta'_k(x) + w_k^+(x) \epsilon_k(x))}{w_k^2(x)} \\
&= -\frac{1}{\gamma_t^2} \left(\frac{\frac{\partial \delta'_k(x)}{\partial U_k} w_k^-(x) + \frac{\partial \epsilon_k(x)}{\partial U_k} w_k^+(x)}{w_k(x)} \right. \\
&\quad \left. + \frac{\frac{\partial w_k^-(x)}{\partial U_k} \delta'_k(x) + \frac{\partial w_k^+(x)}{\partial U_k} \epsilon_k(x)}{w_k(x)} - \frac{\frac{\partial w_k(x)}{\partial U_k} (w_k^-(x) \delta'_k(x) + w_k^+(x) \epsilon_k(x))}{w_k^2(x)} \right).
\end{aligned}$$

By calculating, we have

$$\begin{aligned}
& \frac{\frac{\partial w_k^-(x)}{\partial U_k} \delta'_k(x) + \frac{\partial w_k^+(x)}{\partial U_k} \epsilon_k(x)}{w_k(x)} - \frac{\frac{\partial w_k(x)}{\partial U_k} * (w_k^-(x) \delta'_k(x) + w_k^+(x) \epsilon_k(x))}{w_k^2(x)} \\
&= \frac{1}{w_k^2(x)} \left((w_k(x) \left(\frac{\partial w_k^-(x)}{\partial U_k} \delta'_k(x) + \frac{\partial w_k^+(x)}{\partial U_k} \epsilon_k(x) \right) - \frac{\partial w_k(x)}{\partial U_k} (w_k^-(x) \delta'_k(x) + w_k^+(x) \epsilon_k(x))) \right) \\
&= \frac{1}{w_k^2(x)} \left(w_k(x) \left(\frac{\partial w_k^-(x)}{\partial U_k} \delta'_k(x) + \frac{\partial w_k^+(x)}{\partial U_k} \epsilon_k(x) \right) - \frac{\partial w_k(x)}{\partial U_k} (w_k^-(x) \delta'_k(x) + w_k^+(x) \epsilon_k(x)) \right) \\
&= \frac{1}{w_k^2(x)} \left(\frac{\partial w_k^+}{\partial U_k} w_k^- - \frac{\partial w_k^-}{\partial U_k} w_k^+ \right) (\epsilon_k(x) - \delta'_k(x)) \\
&= -\frac{2 s_t^3}{w_k^2(x)} \left[\mathcal{N}(x; s_t \mu_k, \Sigma) M^+(x) - \mathcal{N}(x; -s_t \mu_k, \Sigma) M^-(x) \right] U_k (I - \alpha U_k U_k^\top) \mu_k \\
&= O\left(r_k^+ r_k^- \frac{s_t^3}{\gamma_t^2 (s_t^2 + \gamma_t^2)}\right).
\end{aligned}$$

where $M^+(x) = \Sigma^{-1}(x - s_t\mu_k)(x - s_t\mu_k)^\top \Sigma^{-1} - \Sigma^{-1}$, $M^-(x) = \Sigma^{-1}(x + s_t\mu_k)(x + s_t\mu_k)^\top \Sigma^{-1} - \Sigma^{-1}$, $\alpha = \frac{s_t^2}{s_t^2 + \gamma_t^2}$.

We also know that

$$\frac{\sum_{k=1}^K (\frac{\partial \delta_k(x)}{\partial U_k} w_k^-(x) + \frac{\partial \epsilon_k(x)}{\partial U_k} w_k^+(x))}{\sum_{k=1}^K w_k(x)} = O\left(\frac{s_t^2 \|x\|_2}{s_t^2 + \gamma_t^2}\right) = O\left(\frac{s_t^3 \|\mu_k\|_2}{s_t^2 + \gamma_t^2}\right)$$

$$\frac{O(r_k^+ r_k^- \frac{s_t^3}{\gamma_t^2 (s_t^2 + \gamma_t^2)})}{O(\frac{s_t^3 \|\mu_k\|_2}{s_t^2 + \gamma_t^2})} \rightarrow 0.$$

Thus,

$$J_k^U \approx \frac{2s_t^2}{\gamma_t^2 (s_t^2 + \gamma_t^2)} (r_k^-(x)(U_k^\top (x + s_t\mu_k)I + (x + s_t\mu_k)U_k^\top) + r_k^+(x)(U_k^\top (x - s_t\mu_k)I + U_k(x - s_t\mu_k)^\top)).$$

Before we provide the simplification of Hessian, we first prove that for $a, b \in \mathbb{R}^n$ $M = a^\top b I_n + ba^\top$, MM^\top is positive-definite if and only if $b^\top a \neq 0$. At the same time, we provide the minimum eigenvalue of MM^\top , which will be used later.

Lemma C.3. *Let $a, b \in \mathbb{R}^n$ and $M = a^\top b I_n + ba^\top$. MM^\top is positive-definite if and only if $b^\top a \neq 0$.*

Moreover,

$$\lambda_{\min}(MM^\top) = \mu_2 = \frac{4(a^\top b)^2 + \|a\|_2^2 \|b\|_2^2 - \|a\|_2 \|b\|_2 \sqrt{8(a^\top b)^2 + \|a\|_2^2 \|b\|_2^2}}{2}.$$

Proof. Let $M = a^\top b I_n + ba^\top$, $c = a^\top b$. We know that $\forall x \in \mathbb{R}^n$,

$$\begin{aligned} x^\top MM^\top x &= (M^\top x)^\top (M^\top x) \\ &= \|M^\top x\|_2^2 \geq 0. \end{aligned}$$

Thus, MM^\top is semi-positive definite.

We can also have that

$$|M| = |a^\top b I_n + ba^\top| = c^n |I_n + \frac{1}{c} ba^\top| = 2c^n \geq 0,$$

where $c^n = 0$ if and only if $b^\top a = 0$.

The last equation holds because

$$|I_n + uv^\top| = 1 + v^\top u$$

Thus, $|MM^\top| > 0$, MM^\top is positive definite.

We can further get the eigenvalues of MM^\top .

Expanding gives the convenient representation

$$MM^\top = (a^\top b)^2 I_n + a^\top b (ba^\top + ab^\top) + a^\top a b b^\top. \quad (6)$$

$\forall x \in \mathbb{R}^n$, if $x^\top a = 0$ and $x^\top b = 0$, we have:

$$MM^\top x = (a^\top b)^2 x.$$

Thus, $(a^\top b)^2$ is an eigenvalue of M , and its eigenspace contains the orthogonal complement of $\text{span}\{a, b\}$. If a and b are linearly independent then $\dim(\text{span}\{a, b\}) = 2$, so the multiplicity of the eigenvalue α^2 is at least $n - 2$.

To find the remaining eigenvalues we restrict M to the subspace $\mathcal{S} := \text{span}\{a, b\}$. Assume first that a and b are linearly independent so that \mathcal{S} is two-dimensional.

Using equation 6, we can compute $\text{tr}(MM^\top)$, which is

$$\begin{aligned}\text{tr}(MM^\top) &= \text{tr}((a^\top b)^2 I_n + a^\top b(ba^\top + ab^\top) + a^\top abb^\top) \\ &= n(a^\top b)^2 + 2(a^\top b)^2 + \|a\|_2^2 \|b\|_2^2 \\ &= (n+2)(a^\top b)^2 + \|a\|_2^2 \|b\|_2^2.\end{aligned}$$

The second equation holds because of $\text{tr}(xy^\top) = \text{tr}(y^\top x) = y^\top x$.

We set the other two eigenvalues are μ_1 and μ_2 . Thus

$$\text{tr}(MM^\top) = \sum_{i=1}^n \lambda_i = (n-2)(a^\top b)^2 + \mu_1 + \mu_2 = (n+2)(a^\top b)^2 + \|a\|_2^2 \|b\|_2^2,$$

and

$$|MM^\top| = \prod_{i=1}^n \lambda_i = (a^\top b)^{2(n-2)} \mu_1 \mu_2 = 4(a^\top b)^{2n}.$$

So μ_1 and μ_2 are the two solutions of

$$x^2 - (4(a^\top b)^2 + \|a\|_2^2 \|b\|_2^2) x + 4(a^\top b)^{2n} = 0. \quad (7)$$

Solving equation 7, we have

$$\mu_1, \mu_2 = \frac{4(a^\top b)^2 + \|a\|_2^2 \|b\|_2^2 \pm \|a\|_2 \|b\|_2 \sqrt{8(a^\top b)^2 + \|a\|_2^2 \|b\|_2^2}}{2}.$$

Now we obtain all eigenvalues. Moreover, we can calculate the minimum of eigenvalues.

$$\lambda_{\min}(MM^\top) = \mu_2 = \frac{4(a^\top b)^2 + \|a\|_2^2 \|b\|_2^2 - \|a\|_2 \|b\|_2 \sqrt{8(a^\top b)^2 + \|a\|_2^2 \|b\|_2^2}}{2}.$$

■

Lemma C.4. [Eigenvalues of the Hessian blocks] Under the same conditions, H is convex. If $\forall x \in \mathbb{R}^{d_k}, r_k^+(x) = 1$ or $r_k^-(x) = 1$ are strictly satisfied, the eigenvalues of the Hessian at θ^* are

$$\lambda_{\min}(H_{\mu_k \mu_k}) = \frac{s_t^2}{(s_t^2 + \gamma_t^2)^2}, \text{ and}$$

$$\lambda_{\min}(H_{U_k U_k}) = \frac{4(U_k^\top \mu_k)^2 + \|U_k\|_2^2 \|\mu_k\|_2^2 - \|U_k\|_2 \|\mu_k\|_2 \sqrt{8(U_k^\top \mu_k)^2 + \|U_k\|_2^2 \|\mu_k\|_2^2}}{2}.$$

Proof. We first state the convexity of the loss function near the true value θ^* .

Let $\theta = \theta^* + \Delta\theta$

$$s_\theta(x, t) = s_{\theta^*}(x, t) + (\nabla_\theta s_\theta(x, t)|_{\theta^*})^\top [\Delta\theta] + O(\|\Delta\theta\|_2^2).$$

$$\begin{aligned}L(\theta) &= \mathbb{E}_{x \sim p_t(x)} [(s_\theta(x, t) - \nabla \log p_t(x))^\top (s_\theta(x, t) - \nabla \log p_t(x))] \\ &= \mathbb{E}_{x \sim p_t(x)} [(s_{\theta^*}(x, t) + (\nabla_\theta s_\theta(x, t)|_{\theta^*})^\top [\Delta\theta] + O(\|\Delta\theta\|_2^2) - \nabla \log p_t(x))^\top \\ &\quad (s_{\theta^*}(x, t) + (\nabla_\theta s_\theta(x, t)|_{\theta^*})^\top [\Delta\theta] + O(\|\Delta\theta\|_2^2) - \nabla \log p_t(x))] \\ &= \mathbb{E}_{x \sim p_t(x)} [((\nabla_\theta s_\theta(x, t)|_{\theta^*})^\top [\Delta\theta])^\top (\nabla_\theta s_\theta(x, t)|_{\theta^*} [\Delta\theta])] + O(\|\Delta\theta\|_2^3) \\ &= (\Delta\theta)^\top \mathbb{E}_{x \sim p_t(x)} [(\nabla_\theta s_\theta(x, t)|_{\theta^*}) (\nabla_\theta s_\theta(x, t)|_{\theta^*})^\top] \Delta\theta \\ &\triangleq (\Delta\theta)^\top H \Delta\theta.\end{aligned}$$

$$\frac{\partial^2 L(\theta)}{\partial \theta^2} = 2H.$$

We then analyze the convexity of $\mathbb{E}_{x \sim p_t(x)}[(\nabla_{\theta} s_{\theta}(x, t)|_{\theta^*})(\nabla_{\theta} s_{\theta}(x, t)|_{\theta^*})^{\top}] \triangleq H$. We can divide H into 4 parts: $H_{\mu\mu}, H_{UU}, H_{\mu U}$ and $H_{U\mu}$, where $H_{U\mu} = (H_{\mu U})^{\top}$.

Let $J_k^{\mu}|_{\theta} = \frac{\partial s_{\theta}}{\partial \mu_k}|_{\theta}$.

$$\begin{aligned} H &= \mathbb{E}_{x \sim p_t(x)}[(\nabla_{\theta} s_{\theta}(x, t)|_{\theta^*})(\nabla_{\theta} s_{\theta}(x, t)|_{\theta^*})^{\top}] \\ &= \mathbb{E}_{x \sim p_t(x)}[J_{\theta^*}(x, t)J_{\theta^*}(x, t)^{\top}]. \end{aligned}$$

Term $H_{\mu\mu}$

We will show that $H_{\mu_k\mu_k}$ is α -convex, where $\alpha > 0$.

$$H_{\mu_k\mu_k} = \mathbb{E}_{x \sim p_t(x)}[J_k^{\mu}J_k^{\mu\top}]$$

$$H_{\mu_k\mu_k} \approx \mathbb{E}_{x \sim p_t(x)}[J_k^{\mu}J_k^{\mu\top}] \approx \frac{s_t^2}{\gamma_t^4} \mathbb{E}_{x \sim p_t(x)}[(r_k^+(x) - r_k^-(x))^2](I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k U_k^{\top})^2.$$

Let $P_k = U_k U_k^{\top}$, $\alpha = \frac{s_t^2}{s_t^2 + \gamma_t^2}$,

$$(I - \alpha P_k)(I - \alpha P_k)^{\top} = (I - \alpha P_k)^2 = I - 2\alpha P_k + \alpha^2 P_k^2 = (I - \alpha P_k)^2.$$

We then prove that $\lambda_{\min}((I - \alpha P_k)^2) = (\frac{\gamma_t^2}{s_t^2 + \gamma_t^2})^2$.

First, we calculate the eigenvalue of P .

$$P^2 = P \Rightarrow \lambda_1 = 1, \lambda_2 = 0.$$

Then we take subspace $Col(P) = \{v; v = Px, x \in \mathcal{R}^D\}$ corresponding to λ_1 , and subspace $Ker(P) = \{v; Pv = 0, x \in \mathcal{R}^D\}$ corresponding to λ_2 .

If $w \in Col(P)$, $Pw = w$:

$$\begin{aligned} (I - \alpha P)w &= (1 - \alpha)w \\ (I - \alpha P)^2 w &= (1 - \alpha)^2 w \\ \Rightarrow \lambda'_1 &= (1 - \alpha)^2. \end{aligned}$$

If $w \in Ker(P)$, $Pw = 0$:

$$\begin{aligned} (I - \alpha P)w &= w \\ (I - \alpha P)^2 w &= w \\ \Rightarrow \lambda'_2 &= 1. \end{aligned}$$

$$\begin{aligned} H_{\mu\mu} &= \mathbb{E}[J_k^{\mu}(J_k^{\mu})^{\top}] \\ \lambda_{\min}(H_{\mu\mu}) &\approx \frac{s_t^2}{(s_t^2 + \gamma_t^2)^2}. \end{aligned}$$

Therefore, $\lambda_{\min}((I - \alpha P_k)^2) = (\frac{\gamma_t^2}{s_t^2 + \gamma_t^2})^2$. Hence, we have

$$\lambda_{\min}(H_{\mu_k\mu_k}) \geq \frac{c_k s_t^2}{(s_t^2 + \gamma_t^2)^2} \approx \frac{s_t^2}{(s_t^2 + \gamma_t^2)^2},$$

where $c_k = \mathbb{E}_{x \sim p_t(x)}[(r_k^+(x) - r_k^-(x))^2] \approx 1$.

Term $H_{U_k U_k}$

$$\begin{aligned} H_{U_k U_k} &\approx \mathbb{E}_{x \sim p_t(x)} [J_U^k J_U^{k\top}] \\ &\approx \frac{4s_t^4}{\gamma_t^4(s_t^2 + \gamma_t^2)^2} \mathbb{E}_{x \sim p_t(x)} [(U_k^\top(x + s_t \mu_k)I + (x + s_t \mu_k)U_k^\top)(U_k^\top(x + s_t \mu_k)I + (x + s_t \mu_k)U_k^\top)^\top] \\ &= \frac{4s_t^4}{\gamma_t^4(s_t^2 + \gamma_t^2)^2} (s_t^2 U_k^\top \mu_k \mu_k^\top U_k I + s_t^2 \mu_k^\top U_k (\mu_k U_k^\top + U_k \mu_k^\top) + \mu_k U_k^\top U_k \mu_k^\top + M(x)), \end{aligned}$$

where $M(x)$ is semi-positive for $\mathbb{E}_{x \sim p_t(x)}[x] = 0$.

Using lemma C.3, we can take $a = U_k$ and $b = \mu_k$ and obtain that

$H_{U_k U_k}$ is positive definite and

$$\lambda_{\min}(H_{U_k U_k}) = \frac{4(U_k^\top \mu_k)^2 + \|U_k\|_2^2 \|\mu_k\|_2^2 - \|U_k\|_2 \|\mu_k\|_2 \sqrt{8(U_k^\top \mu_k)^2 + \|U_k\|_2^2 \|\mu_k\|_2^2}}{2}.$$

Term $H_{\mu_k U_k}$ and Term $H_{U_k \mu_k}$

Since $H_{U_k \mu_k} = H_{\mu_k U_k}^\top$, we just analyze $H_{\mu_k U_k}$. We want to analyze the Hessian block

$$H_{\mu_k U_k} = \mathbb{E}_{x \sim p_t} [J_k^U(x) (J_k^\mu(x))^\top],$$

and show that under symmetric assumptions, this cross-term is zero.

The first-order derivative with respect to μ_k is approximately:

$$J_k^\mu(x) \approx -\frac{s_t}{\gamma_t^2} (r_k^+(x) - r_k^-(x)) (I - \alpha U_k U_k^\top), \quad \alpha = \frac{s_t^2}{s_t^2 + \gamma_t^2}.$$

The first-order derivative with respect to U_k is approximately:

$$J_k^U(x) \approx -\frac{1}{\gamma_t^2} \left[r_k^-(x) \frac{\partial \delta_k(x)}{\partial U_k} + r_k^+(x) \frac{\partial \epsilon_k(x)}{\partial U_k} \right],$$

with

$$\frac{\partial \delta_k(x)}{\partial U_k} = -2 \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k(x + s_t \mu_k), \quad \frac{\partial \epsilon_k(x)}{\partial U_k} = -2 \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k(x - s_t \mu_k).$$

combining terms:

$$J_k^U(x) = C \cdot U_k [r_k^-(x)(x + s_t \mu_k) + r_k^+(x)(x - s_t \mu_k)],$$

where $C = \frac{2s_t^2}{\gamma_t^2(s_t^2 + \gamma_t^2)}$. Assume that the underlying component distribution $p_k(x)$ is symmetric:

$$p_k(x) = p_k(-x),$$

and the weights satisfy:

$$r_k^+(-x) = r_k^-(x), \quad r_k^-(-x) = r_k^+(x).$$

Then we have:

(a) $J_k^\mu(x)$ is an odd function:

$$\begin{aligned} J_k^\mu(-x) &= -\frac{s_t}{\gamma_t^2} (r_k^+(-x) - r_k^-(-x)) (I - \alpha U_k U_k^\top) \\ &= -\frac{s_t}{\gamma_t^2} (r_k^-(x) - r_k^+(x)) (I - \alpha U_k U_k^\top) \\ &= -J_k^\mu(x). \end{aligned}$$

(b) $J_k^U(x)$ is an odd function:

$$\begin{aligned} J_k^U(-x) &= C U_k [r_k^-(x)(-x + s_t \mu_k) + r_k^+(x)(-x - s_t \mu_k)] \\ &= C U_k [r_k^+(x)(-x + s_t \mu_k) + r_k^-(x)(-x - s_t \mu_k)] \\ &= -C U_k [r_k^-(x)(x + s_t \mu_k) + r_k^+(x)(x - s_t \mu_k)] \\ &= -J_k^U(x). \end{aligned}$$

Now compute:

$$H_{\mu_k U_k} = \int J_k^U(x) (J_k^\mu(x))^\top p_k(x) dx.$$

Using symmetry:

$$= \int J_k^U(-x) (J_k^\mu(-x))^\top p_k(-x) dx = \int (-J_k^U(x)) (-J_k^\mu(x))^\top p_k(x) dx = H_{\mu_k U_k}.$$

Thus,

$$\begin{aligned} H_{\mu_k U_k} &= \mathbb{E}_{x \sim p_{data}} [J_k^\mu (J_k^U)^\top] = \mathbb{E}_{x \sim p_{data}} \left[\frac{2s_t^3}{\gamma_t^4(s_t^2 + \gamma_t^2)} (r_k^+(x) - r_k^-(x)) \left(1 - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k U_k^\top\right) \right. \\ &\quad \left. (r_k^-(x)(U_k^\top(x + s_t \mu_k)I + U_k(x + s_t \mu_k)^\top) + r_k^+(x)(U_k^\top(x - s_t \mu_k)I + U_k(x - s_t \mu_k)^\top)) \right]. \end{aligned}$$

$$\begin{aligned} \lambda_{H_{\mu\mu}} &= \mathbb{E}_{x \sim p_{data}} [(u^\top J_\mu^k)^2] \\ \lambda_{H_{UU}} &= \mathbb{E}_{x \sim p_{data}} [(u^\top J_U^k)^2] \\ \lambda_{H_{\mu U}} &= \mathbb{E}_{x \sim p_{data}} [(u^\top J_\mu^k)(u^\top J_U^k)] \leq \sqrt{\lambda_{H_{\mu\mu}} \lambda_{H_{UU}}}. \end{aligned}$$

■

Analyze H

$$H = \begin{pmatrix} H_{\mu_k \mu_k} & H_{\mu_k U_k} \\ H_{\mu_k U_k} & H_{U_k U_k} \end{pmatrix}$$

. If we can prove that $H_{\mu_k \mu_k} - H_{U_k \mu_k} H_{U_k U_k}^{-1} H_{U_k \mu_k}^\top$ is positive-definite, then H is positive-definite for Schur's Theorem.

$$\lambda_H \geq \lambda_S \geq \lambda_{H_{\mu_k \mu_k}} - \frac{r^2 \lambda_{H_{\mu_k \mu_k}} \lambda_{H_{U_k U_k}}}{\lambda_{H_{U_k U_k}}} = (1 - r^2) \lambda_{H_{\mu_k \mu_k}} \geq (1 - r^2) \frac{s_t^2}{(s_t^2 + \gamma_t^2)^2} > 0.$$

$$r = \max_{\|u\|=1, \|v\|=1} \frac{u^\top H_{\mu_k U_k} v}{\sqrt{u^\top H_{\mu_k \mu_k} u \cdot v^\top H_{U_k U_k} v}} \leq 1.$$

$r = 1$ if and only if $u^\top J_\mu^k = cv^\top J_U^k$, $c \neq 0$, which is almost impossible to happen.

More specially, if we assume that $\forall x \in \mathbb{R}^{d_k}, r_k^+ = 1$ or $r_k^- = 1$, for

$$\begin{aligned} H_{\mu_k U_k} &= \mathbb{E}_{x \sim p_{data}} [J_k^\mu (J_k^U)^\top] = \mathbb{E}_{x \sim p_{data}} \left[\frac{2s_t^3}{\gamma_t^4(s_t^2 + \gamma_t^2)} (r_k^+(x) - r_k^-(x)) \left(1 - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k U_k^\top\right) \right. \\ &\quad \left. (r_k^-(x)(U_k^\top(x + s_t \mu_k)I + U_k(x + s_t \mu_k)^\top) + r_k^+(x)(U_k^\top(x - s_t \mu_k)I + U_k(x - s_t \mu_k)^\top)) \right] \\ &= \mathbb{E}_{x \sim \mathcal{N}(s_t \mu_k, \Sigma_k)} \left[\frac{2s_t^3}{\gamma_t^4(s_t^2 + \gamma_t^2)} (r_k^+(x) - r_k^-(x)) \left(1 - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k U_k^\top\right) \right. \\ &\quad \left. (r_k^-(x)(U_k^\top(x + s_t \mu_k)I + U_k(x + s_t \mu_k)^\top) + r_k^+(x)(U_k^\top(x - s_t \mu_k)I + U_k(x - s_t \mu_k)^\top)) \right] \\ &= \mathbb{E}_{x \sim \mathcal{N}(s_t \mu_k, \Sigma_k)} \left[\frac{2s_t^3}{\gamma_t^4(s_t^2 + \gamma_t^2)} \left(1 - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k U_k^\top\right) + (U_k^\top(x - s_t \mu_k)I + U_k(x - s_t \mu_k)^\top) \right] \\ &= 0. \end{aligned}$$

We have $r = 0$,

$$\alpha = \min\left\{\frac{s_t^2}{(s_t^2 + \gamma_t^2)^2}, \frac{4(U_k^\top \mu_k)^2 + \|U_k\|_2^2 \|\mu_k\|_2^2 - \|U_k\|_2 \|\mu_k\|_2 \sqrt{8(U_k^\top \mu_k)^2 + \|U_k\|_2^2 \|\mu_k\|_2^2}}{2}\right\}.$$

Uttill now, We have shown that H is α -convex and L -lipschiz, where $\alpha = (1 - r^2)\lambda_{H_{\mu_k \mu_k}}$. And we can know that $L(\theta)$ is exponentially convergent.

Theorem C.5. *If we take $\eta_t = \eta = \frac{2}{\eta+L}$, and $\kappa = \frac{L}{\alpha}$, then*

$$\|\theta^t - \theta^*\|_2 \leq \left(\frac{\kappa - 1}{\kappa + 1}\right)^t \|\theta^{(0)} - \theta^*\|_2.$$

D K-MODE MOG OPTIMIZATION

D.1 SETTING

In this section, we analyze

$$\begin{aligned} \nabla \log p_{t,k}(x) &= \frac{\nabla p_{t,k}(x)}{p_{t,k}(x)} \\ &= -\frac{1}{\gamma_t^2} \frac{\sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}, s_t^2 U_{k,l}^* U_{k,l}^{*\top} + \gamma_t^2 I) \left(x - s_t \mu_{k,l} - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l}^* U_{k,l}^{*\top} (x - s_t \mu_{k,l})\right)}{\sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}, s_t^2 U_{k,l}^* U_{k,l}^{*\top} + \gamma_t^2 I)}. \end{aligned}$$

D.2 OPTIMIZATION

Assumption D.1. [Highly Separated Gaussian] Consider the Gaussian mixture

$$p_k(x) = \sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; \mu_{k,l}, \Sigma_{k,l}), \quad r_{k,l}(x) := \frac{\pi_{k,l} \mathcal{N}(x; \mu_{k,l}, \Sigma_{k,l})}{\sum_{i=1}^{n_k} \pi_{k,i} \mathcal{N}(x; \mu_{k,i}, \Sigma_{k,i})}.$$

There exist constants $\varepsilon \ll 1$ and $\delta \ll 1$ such that when $x \sim p_k$ we have

$$\Pr_{x \sim p_k} \left(\exists l \in \{1, \dots, n_k\} \text{ with } r_{k,l}(x) \geq 1 - \varepsilon \right) \geq 1 - \delta.$$

We assume that the gap between the subspaces is large, and the gap within the subspace is relatively small, and the equivalent Gaussian is used to replace the whole subspace.

Corollary D.2. *Assume that $\|\mu_{k,i}^* - \mu_{k,j}^*\|_2 \leq \delta$, $\|U_{k,i}^* - U_{k,j}^*\|_2 \leq \epsilon$ and $\|x - \bar{\mu}_k^*\|_2 \leq \Delta$. We have*

$$\|\log p(x) - \log \bar{p}(x)\|_2 = O(\epsilon + \delta \Delta + \Delta^3)$$

Proof. For k -th subspace, $w_k(x) = \sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; s_t \mu_{k,l}, \Sigma_{k,l})$, we take

$$\tilde{w}_k(x) = \mathcal{N}(x; \bar{\mu}_k, \bar{\Sigma}_k).$$

where

$$\begin{aligned} \mathbb{E}_{\tilde{w}_k}[x] &= \bar{\mu}_k = \mathbb{E}_{w_k}[x] = \sum_{l=1}^{n_k} \pi_{k,l} s_t \mu_{k,l} \\ \text{Cov}_{\tilde{w}_k}(x) &= \text{Cov}_{w_k}(x) = \mathbb{E}[(x - \bar{\mu}_k)(x - \bar{\mu}_k)^\top] = \sum_{l=1}^{n_k} \pi_{k,l} (\Sigma_{k,l} + s_t^2 \mu_{k,l} \mu_{k,l}^\top - s_t^2 \bar{\mu}_k \bar{\mu}_k^\top) \\ &\Rightarrow \bar{\Sigma}_k = \sum_{l=1}^{n_k} (\Sigma_{k,l} + s_t^2 \mu_{k,l} \mu_{k,l}^\top - s_t^2 \bar{\mu}_k \bar{\mu}_k^\top). \end{aligned}$$

We next show the order of the estimation under the condition that $\|\mu_{k,i} - \mu_{k,j}\|_2 \leq \delta$, $\|U_{k,i} - U_{k,j}\|_2 \leq \epsilon$ and $\|x - \bar{\mu}_k\|_2 \leq \Delta$. Using Taylor's Theorem and take $x_0 = \bar{\mu}_k$, we can obtain that

$$\begin{aligned} \log p(x) &= \log p(x_0) + (x - x_0)^\top \nabla \log p(x_0) + \frac{1}{2} (x - x_0)^\top \nabla^2 \log p(x_0) (x - x_0) + O(\|x - x_0\|^3) \\ \log \tilde{p}(x) &= \log \tilde{p}(x_0) + (x - x_0)^\top \nabla \log \tilde{p}(x_0) + \frac{1}{2} (x - x_0)^\top \nabla^2 \log \tilde{p}(x_0) (x - x_0) + O(\|x - x_0\|^3). \end{aligned}$$

$$\begin{aligned}
\log p(x_0) - \log \tilde{p}(x_0) &= \log \frac{\sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x_0; \mu_{k,l}, \Sigma_{k,l})}{\mathcal{N}(x_0; \bar{\mu}_k, \bar{\Sigma}_k)} \\
&= \log \left(\sum_{l=1}^{n_k} \pi_{k,l} \frac{1}{|\Sigma_{k,l}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\bar{\mu} - \mu_{k,l})^\top \Sigma_{k,l}^{-1}(\bar{\mu} - \mu_{k,l})\right) \right) + \frac{1}{2} \log |\bar{\Sigma}_k| \\
&= \log \left(\sum_{l=1}^{n_k} \pi_{k,l} \frac{1}{|\Sigma_{k,l}|^{\frac{1}{2}}} (1 + O(\delta^2)) \right) + \frac{1}{2} \log |\bar{\Sigma}_k| \\
&= \log \left(\sum_{l=1}^{n_k} \pi_{k,l} \frac{|\bar{\Sigma}_k|^{\frac{1}{2}}}{|\Sigma_{k,l}|^{\frac{1}{2}}} + O(\delta^2) \right) \\
&= O \left(\sum_{l=1}^{n_k} \pi_{k,l} \left(\frac{|\bar{\Sigma}_k|^{\frac{1}{2}}}{|\Sigma_{k,l}|^{\frac{1}{2}}} - 1 \right) + O(\delta^2) \right).
\end{aligned}$$

$$\|\log p(x_0) - \log \tilde{p}(x_0)\|_2 = O(\epsilon + \delta^2).$$

$$\begin{aligned}
\nabla \log p(x_0) - \nabla \log \tilde{p}(x_0) &= \nabla \log \sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x; \mu_{k,l}, \Sigma_{k,l})|_{x_0} \\
&= \frac{\sum_{l=1}^{n_k} \pi_{k,l} \mathcal{N}(x_0; \mu_{k,l}, \Sigma_{k,l}) (-\Sigma_{k,l}^{-1}(\bar{\mu} - \mu_{k,l}))}{p(x_0)}.
\end{aligned}$$

$$\|\nabla \log p(x_0) - \nabla \log \tilde{p}(x_0)\|_2 = O(\delta).$$

$$\begin{aligned}
\nabla^2 \log p(x_0) - \nabla^2 \log \tilde{p}(x_0) &= \frac{\nabla^2 p(x_0)}{p(x_0)} - \left(\frac{\nabla p(x_0)}{p(x_0)} \right) \left(\frac{\nabla p(x_0)}{p(x_0)} \right)^\top - \frac{\nabla^2 \tilde{p}(x_0)}{\tilde{p}(x_0)} \\
&= \left(\frac{\nabla^2 p(x_0)}{p(x_0)} - \frac{\nabla^2 \tilde{p}(x_0)}{\tilde{p}(x_0)} \right) - \left(\frac{\nabla p(x_0)}{p(x_0)} \right) \left(\frac{\nabla p(x_0)}{p(x_0)} \right)^\top.
\end{aligned}$$

$$\|\nabla^2 \log p(x_0) - \nabla^2 \log \tilde{p}(x_0)\|_2 = O(\epsilon^2 + \delta^2).$$

Thus, $\|\log p(x) - \log \tilde{p}(x)\|_2 = O(\epsilon + \delta\Delta + \Delta^3)$. ■

Lemma D.3. [Eigenvalues of the Hessian] Assume Assumption 6.4, the Hessian at the k -th subspace is convex on a neighborhood of θ^* . If $\forall x \in \mathbb{R}^{d_k}$, $r_k^+(x) = 1$ or -1 are strictly satisfied, we have

$$\lambda_{\min}(H_{\mu_{k,l}\mu_{k,l}}) = \frac{\pi_{k,l}s_t^2}{(s_t^2 + \gamma_t^2)^2},$$

and $\lambda_{\min}(H_{U_{k,l}U_{k,l}})$ has the following form:

$$\left(\pi_{k,l} 4(U_{k,l}^\top \mu_{k,l})^2 + \|U_{k,l}\|_2^2 \|\mu_{k,l}\|_2^2 - \|U_{k,l}\|_2 \|\mu_{k,l}\|_2 \sqrt{8(U_{k,l}^\top \mu_{k,l})^2 + \|U_{k,l}\|_2^2 \|\mu_{k,l}\|_2^2} \right) / 2.$$

Proof. According to the previous conclusion, we only need to calculate J_μ and J_U . With these assumptions and simplifications, similar to the symmetry case, we will prove that $J_{k,l}^\mu$ and $J_{k,l}^U$ have

dominant terms.

$$\begin{aligned}
J_{k,l}^\mu(x) &= -\frac{1}{\gamma_t^2} \frac{\partial s_\theta(x, t)}{\partial \mu_{k,l}} \\
&= -\frac{1}{\gamma_t^2} \frac{\sum_{l=1}^{n_k} \left(\frac{\partial w_{k,l}(x)}{\partial \mu_{k,l}} \delta_{k,l}(x) + \frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}} w_{k,l}(x) \right) w_k(x) - \frac{\partial w_k(x)}{\partial \mu_{k,l}} \sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x)}{w_k^2(x)} \\
&= -\frac{1}{\gamma_t^2} \left(\frac{\sum_{l=1}^{n_k} \frac{\partial w_{k,l}(x)}{\partial \mu_{k,l}} \delta_{k,l}(x)}{w_k(x)} + \frac{\sum_{l=1}^{n_k} \frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}} w_{k,l}(x)}{w_k(x)} - \frac{\frac{\partial w_k(x)}{\partial \mu_{k,l}} \sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x)}{w_k^2(x)} \right).
\end{aligned}$$

Let's go ahead and do the calculation.

$$\begin{aligned}
&\frac{\sum_{l=1}^{n_k} \frac{\partial w_{k,l}(x)}{\partial \mu_{k,l}} \delta_{k,l}(x)}{w_k(x)} - \frac{(\frac{\partial w_k(x)}{\partial \mu_{k,l}}) \sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x)}{w_k^2(x)} = \frac{\frac{\partial w_{k,l}(x)}{\partial \mu_{k,l}}}{w_k(x)} (\delta_{k,l}(x) - \bar{\delta}_k(x)) \\
&\frac{\sum_{l=1}^{n_k} \frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}} w_{k,l}(x)}{w_k(x)} \approx \frac{s_t}{\gamma_t^2} \sum_{l=1}^{n_k} r_{k,l}(x) \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right).
\end{aligned}$$

where $r_{k,l}(x) = \frac{\pi_{k,l} \mathcal{N}(x; \bar{\mu}_k, \bar{\Sigma}_k)}{\sum_{j=1}^K \mathcal{N}(x; \bar{\mu}_j, \bar{\Sigma}_j)}$.

Therefore, we can obtain that

$$\begin{aligned}
&\left\| \frac{\sum_{l=1}^{n_k} \frac{\partial w_{k,l}(x)}{\partial \mu_{k,l}} \delta_{k,l}(x)}{w_k(x)} - \frac{\frac{\partial w_k(x)}{\partial \mu_{k,l}} \sum_{l=1}^{n_k} (w_{k,l}(x) \delta_{k,l}(x))}{w_k^2(x)} \right\|_2 = O(\delta(R + s_t B_\mu) \frac{s_t^2}{\gamma_t^2}) \\
&\left\| \frac{\sum_{l=1}^{n_k} \frac{\partial \delta_{k,l}(x)}{\partial \mu_{k,l}} w_{k,l}(x)}{w_k(x)} \right\|_2 = O(s_t).
\end{aligned}$$

where $\delta \leq \|\mu_{k,i} - \mu_{k,j}\|_2 \ll 1$.

Thus, we have

$$J_{k,l}^\mu(x) = \frac{\partial s_\theta}{\partial \mu_{k,l}} \approx \frac{s_t}{\gamma_t^2} r_{k,l}(x) \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right).$$

$$\begin{aligned}
H_{\mu_k, l \mu_k, l} &= \mathbb{E}_{x \sim p_t} [J_{k,l}^\mu(x) J_{k,l}^\mu(x)^\top] \\
&= \frac{s_t^2}{\gamma_t^4} \mathbb{E}[r_{k,l}(x)^2] \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right) \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right)^\top.
\end{aligned}$$

For a given x , since we focus on the equivalent Gaussian distribution for each cluster, we have

$$H_{\mu_k \mu_k} \approx \text{diag}(\mathbb{E}[r_{k,1}^2] H_{\mu_k, 1 \mu_k, 1}, \mathbb{E}[r_{k,2}^2] H_{\mu_k, 2 \mu_k, 2}, \dots, \mathbb{E}[r_{k,n_k}^2] H_{\mu_k, n_k \mu_k, n_k}).$$

We first show that $\mathbb{E}[r_{k,l}^2] H_{\mu_k, l \mu_k, l}$ is positive-definite, then we will further show that $H_{\mu_k \mu_k}$ is positive-definite.

For $H_{\mu_k, l \mu_k, l}$, we know that

$$\begin{aligned}
\lambda_{\min}(H_{\mu_k, l \mu_k, l}) &= c_{k,l} \lambda_{\min}(J_{k,l}^\mu (J_{k,l}^\mu)^\top) \\
&= c_{k,l} \lambda_{\min}((I - \alpha P_k)^2) \\
&= \frac{c_{k,l} \gamma_t^4}{(s_t^2 + \gamma_t^2)^2},
\end{aligned}$$

where

$$c_{k,l} = \frac{s_t^2}{\gamma_t^4} \mathbb{E}[r_{k,l}^2] \approx \pi_{k,l} \frac{s_t^2}{\gamma_t^4}.$$

We know that for a block matrix $A = \text{diag}(A_1, A_2, \dots, A_k)$,

$$\lambda(A) = \cup_{i=1}^k \lambda(A_i).$$

Therefore,

$$\lambda_{\min}(H_{\mu_k \mu_k}) = \min_{l=1 \dots, n_k} \frac{c_{k,l} \gamma_t^4}{(s_t^2 + \gamma_t^2)^2}.$$

Thus, we take

$$\lambda_{H_{\mu_k \mu_k}} = \frac{c_{k,n_k} \gamma_t^4}{(s_t^2 + \gamma_t^2)^2}.$$

Similar to previous situation ,because

$$\frac{\left\| \frac{\sum_{l=1}^{n_k} \left(\frac{\partial \delta_{k,l}(x)}{\partial U_{k,l}} w_{k,l}(x) \right) (w_k(x) - \left(\frac{\partial w_k(x)}{\partial U_{k,l}} \right) \sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x))}{w_k^2(x)} \right\|_2}{\left\| \frac{\sum_{l=1}^{n_k} \frac{\partial \delta_{k,l}(x)}{\partial U_{k,l}} w_{k,l}(x)}{w_k(x)} \right\|_2} \rightarrow 0.$$

we can obtain that

$$\begin{aligned} J_{k,l}^U(x) &= -\frac{1}{\gamma_t^2} \frac{\sum_{l=1}^{n_k} \left(\frac{\partial w_{k,l}(x)}{\partial U_{k,l}} \delta_{k,l}(x) + w_{k,l}(x) \frac{\partial \delta_{k,l}(x)}{\partial U_{k,l}} \right) w_k(x) - \left(\frac{\partial w_k(x)}{\partial U_{k,l}} \right) \sum_{l=1}^{n_k} w_{k,l}(x) \delta_{k,l}(x)}{w_k^2(x)} \\ &= -\frac{1}{\gamma_t^2} \frac{\sum_{l=1}^{n_k} w_{k,l}(x) \frac{\partial \delta_{k,l}(x)}{\partial U_{k,l}}}{w_k(x)} \\ &\approx \frac{1}{\gamma_t^2} \frac{s_t^2}{s_t^2 + \gamma_t^2} r_{k,l}(x) \left[U_{k,l}(x - \mu_{k,l})^\top + (x - \mu_{k,l})^\top U_{k,l} I \right]. \end{aligned}$$

$$H_{U_k U_k} \approx \text{diag}(\mathbb{E}[r_{k,1}^2] H_{U_{k,1} U_{k,1}}, \mathbb{E}[r_{k,2}^2] H_{U_{k,2} U_{k,2}}, \dots, \mathbb{E}[r_{k,n_k}^2] H_{U_{k,n_k} U_{k,n_k}}).$$

$$\begin{aligned} H_{U_{k,l} U_{k,l}} &= \mathbb{E}[J_{k,l}^U(x) (J_{k,l}^U(x))^\top] \\ &= \mathbb{E}[(\frac{\alpha}{\gamma_t^2})^2 (U_{k,l}(x - \mu_{k,l})^\top (x - \mu_{k,l}) U_{k,l}^\top + U_{k,l}^\top (x - \mu_{k,l}) U_{k,l} (x - \mu_{k,l})^\top)] \\ &\quad + \mathbb{E}[(\frac{\alpha}{\gamma_t^2})^2 (U_{k,l}^\top (x - \mu_{k,l}) (x - \mu_{k,l}) U_{k,l}^\top + (U_{k,l}^\top (x - \mu_{k,l}))^2)]. \end{aligned}$$

Similar to our calculation in C.4, we can use C.3 to calculate the minimum eigenvalue of $H_{U_{k,l} U_{k,l}}$.

$H_{U_{k,l} U_{k,l}}$ is positive definite and

$$\lambda_{\min}(H_{U_{k,l} U_{k,l}}) = \frac{4(U_{k,l}^\top \mu_{k,l})^2 + \|U_{k,l}\|_2^2 \|\mu_{k,l}\|_2^2 - \|U_{k,l}\|_2 \|\mu_{k,l}\|_2 \sqrt{8(U_{k,l}^\top \mu_{k,l})^2 + \|U_{k,l}\|_2^2 \|\mu_{k,l}\|_2^2}}{2}.$$

Recall that

$$H_{U_k U_k} \approx \text{diag}(\mathbb{E}[r_{k,1}^2] H_{U_{k,1} U_{k,1}}, \mathbb{E}[r_{k,2}^2] H_{U_{k,2} U_{k,2}}, \dots, \mathbb{E}[r_{k,n_k}^2] H_{U_{k,n_k} U_{k,n_k}}).$$

and $\mathbb{E}[r_{k,l}^2] \approx \pi_{k,l}$, we can obtain the minimum eigenvalue of $H_{U_k U_k}$, which is

$$\min_{l=1,2,\dots,n_k} \pi_{k,l} \frac{4(U_{k,l}^\top \mu_{k,l})^2 + \|U_{k,l}\|_2^2 \|\mu_{k,l}\|_2^2 - \|U_{k,l}\|_2 \|\mu_{k,l}\|_2 \sqrt{8(U_{k,l}^\top \mu_{k,l})^2 + \|U_{k,l}\|_2^2 \|\mu_{k,l}\|_2^2}}{2}.$$

Lemma D.4. [Local Strong Convexity] Assume Assumption 6.4, in a neighborhood of θ^* , $\nabla^2 \mathcal{L}(\theta) \succeq \alpha' I$, $\alpha' > 0$, $\forall \theta \in \Theta$. If $\forall x \in \mathbb{R}^{d_k}$, $\exists l \in [n_k]$, $r_{k,l}(x) = 1$ are strictly satisfied, $\alpha' = \min\{\lambda_1, \lambda_2\}$, where $\lambda_1 = \min_{l=1 \dots, n_k} \frac{c_{k,l} \gamma_t^4}{(s_t^2 + \gamma_t^2)^2}$, $\lambda_2 = \min_{l=1,2,\dots,n_k} \lambda_{\min}(H_{U_{k,l} U_{k,l}})$.

Proof.

$$H_{\mu_k U_k} = \text{diag}(H_{\mu_{k,1} U_{k,1}}, H_{\mu_{k,2} U_{k,2}}, \dots, H_{\mu_{k,n_k} U_{k,n_k}}).$$

$$\|H_{\mu_k U_k}\| \leq \sqrt{\|H_{\mu_k \mu_k}\| \|H_{U_k U_k}\|} = O\left(\frac{s_t^3}{\gamma_t^2 (s_t^2 + \gamma_t^2)^2}\right).$$

$$H = \begin{pmatrix} \text{diag}(H_{\mu_{k,1} \mu_{k,1}}, \dots, H_{\mu_{k,n_k} \mu_{k,n_k}}) & \text{diag}(H_{\mu_{k,1} U_{k,1}}, \dots, H_{\mu_{k,n_k} U_{k,n_k}}) \\ \text{diag}(H_{\mu_{k,1} U_{k,1}}, \dots, H_{\mu_{k,n_k} U_{k,n_k}}) & \text{diag}(H_{U_{k,1} U_{k,1}}, \dots, H_{U_{k,n_k} U_{k,n_k}}) \end{pmatrix}.$$

Let

$$S = H_{\mu\mu} - H_{\mu U} H_{UU}^{-1} H_{U\mu}$$

we have

$$\lambda_H \geq \lambda_S \geq \lambda_{H_{\mu_k \mu_k}} - \frac{r^2 \lambda_{H_{\mu_k \mu_k}} \lambda_{H_{U_k U_k}}}{\lambda_{H_{U_k U_k}}} = (1 - r^2) \lambda_{H_{\mu_k \mu_k}} \geq (1 - r^2) \frac{s_t^2}{(s_t^2 + \gamma_t^2)^2} > 0.$$

$$r = \max_{\|u\|=1, \|v\|=1} \frac{u^\top H_{\mu_k U_k} v}{\sqrt{u^\top H_{\mu_k \mu_k} u \cdot v^\top H_{U_k U_k} v}} \leq 1.$$

$r = 1$ if and only if $u^\top J_\mu^k = cv^\top J_U^k$, $c \neq 0$, which is almost impossible to happen.

More specifically, if we assume that $\forall x \in \mathbb{R}^{d_k}$, $\exists l \in [n_k]$, $r_{k,l}(x) = 1$, we have

$$\begin{aligned} H_{\mu_{k,l} U_{k,l}} &= \mathbb{E}_{x \sim p_k} [J_{k,l}^U(x) (J_{k,l}^\mu(x))^\top] \\ &= \frac{1}{\gamma_t^4} \frac{s_t^3}{s_t^2 + \gamma_t^2} \mathbb{E}_{x \sim p_k} \left[r_{k,l}(x)^2 ((x - \mu_{k,l}) U_{k,l}^\top + (x - \mu_{k,l})^\top U_{k,l} I) \right] \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right) \\ &= \frac{1}{\gamma_t^4} \frac{s_t^3}{s_t^2 + \gamma_t^2} \mathbb{E}_{x \sim \pi_{k,l} \mathcal{N}_{k,l}} \left[r_{k,l}(x)^2 ((x - \mu_{k,l}) U_{k,l}^\top + (x - \mu_{k,l})^\top U_{k,l} I) \right] \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right) \\ &\approx 0 \end{aligned}$$

The second equation holds because $\forall x$, if $x \notin \mathcal{N}_{k,l}(\mu_{k,l}, \Sigma_{k,l})$, $r_{k,l}(x) = 0$. And the third equation holds because if $x \sim \mathcal{N}_{k,l}(\mu_{k,l}, \Sigma_{k,l})$, $\forall \text{Const } C$,

$$\mathbb{E}_{x \sim \pi_{k,l} \mathcal{N}_{k,l}} [C(x - \mu_{k,l})] = 0.$$

Thus, let α' be the minimum eigenvalue of H ,

$$\alpha' = \min\{\lambda_1, \lambda_2\}, \quad (8)$$

where

$$\lambda_1 = \min_{l=1 \dots, n_k} \frac{c_{k,l} \gamma_t^4}{(s_t^2 + \gamma_t^2)^2},$$

and

$$\lambda_2 = \min_{l=1,2,\dots,n_k} \pi_{k,l} \frac{4(U_{k,l}^\top \mu_{k,l})^2 + \|U_{k,l}\|_2^2 \|\mu_{k,l}\|_2^2 - \|U_{k,l}\|_2 \|\mu_{k,l}\|_2 \sqrt{8(U_{k,l}^\top \mu_{k,l})^2 + \|U_{k,l}\|_2^2 \|\mu_{k,l}\|_2^2}}{2}.$$

■

E EXTENSION TO MOG LATENT WITHOUT SEPARATION ASSUMPTION

E.1 2-MODE ANALYSIS

In this section, we relax the high separation assumption (where $r_k^+(x)r_k^-(x) \approx 0$). Instead, we treat the overlap between manifold components as a bounded perturbation to the ideal system. We aim to prove that the Hessian remains positive definite provided the overlap factor is sufficiently small.

E.1.1 DEFINITION OF OVERLAP FACTOR

We define the pointwise overlap factor $\xi_k(x)$ as the product of the assignment probabilities for the positive and negative components of the k -th manifold:

$$\xi_k(x) \triangleq r_k^+(x)r_k^-(x). \quad (9)$$

Since $r_k^+(x), r_k^-(x) \in [0, 1]$ and $r_k^+(x) + r_k^-(x) = 1$, the overlap factor is naturally bounded: $0 \leq \xi_k(x) \leq 0.25$.

We denote the maximum expected overlap magnitude as $\epsilon_{\text{overlap}}$:

$$\epsilon_{\text{overlap}} = \sup_{x \in \text{supp}(p_t)} \xi_k(x). \quad (10)$$

E.1.2 JACOBIAN ANALYSIS

We revisit the derivation of the Jacobian J_k^μ . In the original derivation, J_k^μ was decomposed into Term A (dominant term) and Term B (previously ignored):

$$J_k^\mu(x) = \underbrace{J_{\text{ideal}}^\mu(x)}_{\text{Term A}} + \underbrace{E^\mu(x)}_{\text{Term B}}.$$

When $\xi_k(x) \rightarrow 0$, we can recover the ideal Jacobian derived previously:

$$J_{\text{ideal}}^\mu(x) = -\frac{s_t}{\gamma_t^2}(r_k^+(x) - r_k^-(x)) \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k U_k^\top \right).$$

Term B contains the cross-product of weights, which is exactly our overlap factor $\xi_k(x)$. Specifically:

$$E^\mu(x) = -\frac{4s_t^2}{\gamma_t^2 w_k^2(x)} \cdot \xi_k(x) \cdot \Sigma_k^{-1} x \left(I + \frac{s_t^2}{s_t^2 + \gamma_t^2} U_k U_k^\top \right) \mu_k.$$

We can bound the norm of this error term. Since terms like $\frac{x}{w_k(x)}$ and projection matrices are bounded within the support, there exists a constant C_1 such that:

$$\|E^\mu(x)\|_2 \leq C_1 \cdot \xi_k(x). \quad (11)$$

Similarly, for the Jacobian with respect to U_k , we can decompose it into an ideal part and an error part proportional to the overlap:

$$J_k^U(x) = J_{\text{ideal}}^U(x) + E^U(x), \quad \text{where } \|E^U(x)\|_F \leq C_2 \cdot \xi_k(x).$$

E.1.3 HESSIAN ANALYSIS

The Hessian matrix H is defined as the expected outer product of the Jacobians:

$$H = \mathbb{E}_{x \sim p_t(x)} [J(x)J(x)^\top].$$

Let $J(x) = J_{\text{ideal}}(x) + E(x)$. Substituting this into the Hessian definition:

$$\begin{aligned} H &= \mathbb{E} [(J_{\text{ideal}} + E)(J_{\text{ideal}} + E)^\top] \\ &= \underbrace{\mathbb{E}[J_{\text{ideal}} J_{\text{ideal}}^\top]}_{H_{\text{ideal}}} + \underbrace{\mathbb{E}[J_{\text{ideal}} E^\top + E J_{\text{ideal}}^\top + E E^\top]}_{\Delta H}. \end{aligned}$$

Here, H_{ideal} is the Hessian matrix under the high separation assumption and ΔH is the perturbation matrix induced by the overlap.

From the previous proof, we established that H_{ideal} is block-diagonal (or has negligible off-diagonals due to symmetry) and positive definite. Let $\alpha > 0$ be its minimum eigenvalue:

$$\begin{aligned} \lambda_{\min}(H_{\text{ideal}}) &\approx \mathbb{E}[(r_k^+(x) - r_k^-(x))^2] \min(\lambda_{\min}(H_{\mu_k \mu_k}), \lambda_{\min}(H_{U_k U_k})) \\ &= \mathbb{E}[(1 - 4\xi_k(x))] \min(\lambda_{\min}(H_{\mu_k \mu_k}), \lambda_{\min}(H_{U_k U_k})) \\ &\geq (1 - 4\epsilon_{\text{overlap}}) \min(\lambda_{\min}(H_{\mu_k \mu_k}), \lambda_{\min}(H_{U_k U_k})) \triangleq \alpha. \end{aligned}$$

We apply the Triangle Inequality and Cauchy-Schwarz inequality to bound the spectral norm of ΔH :

$$\begin{aligned} \|\Delta H\|_2 &\leq 2\|\mathbb{E}[J_{\text{ideal}} E^\top]\|_2 + \|\mathbb{E}[E E^\top]\|_2 \\ &\leq 2\sqrt{\mathbb{E}[\|J_{\text{ideal}}\|^2] \mathbb{E}[\|E\|^2]} + \mathbb{E}[\|E\|^2]. \end{aligned}$$

Since $\|E^\mu(x)\| \leq C_1 \cdot \xi_k(x)$ and $\|E^U(x)\| \leq C_2 \cdot \xi_k(x)$, the perturbation norm is dominated by the overlap factor:

$$J_k^U(x) = J_{\text{ideal}}^U(x) + E^U(x), \quad \text{where } \|E^U(x)\|_F \leq C_2 \cdot \xi_k(x).$$

The Hessian perturbation matrix is given by $\Delta H \approx \mathbb{E}[J_{\text{ideal}} E^\top + E J_{\text{ideal}}^\top]$. To bound its spectral norm $\|\Delta H\|_2$, we define the signal bounds

$$S_\mu \triangleq \sup_x \|J_{\text{ideal}}^\mu(x)\|_2 \approx \frac{s_t}{\gamma_t^2}$$

and

$$S_U \triangleq \sup_x \|J_{\text{ideal}}^U(x)\|_2 \approx \frac{s_t R^2}{\gamma_t^2}.$$

We can define the composite perturbation constant C' as:

$$C' = 2(S_\mu + S_U)(C_1 + C_2).$$

And thus,

$$\|\Delta H\|_2 \leq C' \cdot \epsilon_{\text{overlap}}.$$

E.1.4 POSITIVE DEFINITENESS VIA WEYL'S INEQUALITY

We now use Matrix Perturbation Theory to prove the convexity of the actual loss landscape. With Weyl's Inequality for Hermitian Matrices, we have: Let $H = H_{\text{ideal}} + \Delta H$. The eigenvalues of H are bounded by:

$$\lambda_{\min}(H) \geq \lambda_{\min}(H_{\text{ideal}}) - \|\Delta H\|_2. \quad (12)$$

Substituting our bounds:

$$\lambda_{\min}(H) \geq \alpha - C' \cdot \epsilon_{\text{overlap}}. \quad (13)$$

Condition for Convexity: For the Hessian H to remain positive definite (ensuring strong convexity), we require:

$$\alpha - C' \cdot \epsilon_{\text{overlap}} > 0 \implies \epsilon_{\text{overlap}} < \frac{\alpha}{C'}. \quad (14)$$

This physically implies that as long as the manifolds are not excessively overlapping, the loss function remains locally strongly convex.

E.1.5 CONVERGENCE ANALYSIS

Based on the perturbation analysis, we state the revised convergence theorem.

Theorem E.1 (Linear Convergence under Bounded Overlap). *Let $L(\theta)$ be the loss function. Assume the overlap factor satisfies $\epsilon_{\text{overlap}} < \frac{\alpha}{C'}$. Then, the Hessian H at θ^* is positive definite with minimum eigenvalue:*

$$\lambda_{\min}(H) \geq \alpha_{\text{eff}} = \alpha - C' \epsilon_{\text{overlap}} > 0.$$

Consequently, gradient descent with step size η converges linearly:

$$\|\theta^t - \theta^*\|_2 \leq \left(\frac{\kappa_{\text{eff}} - 1}{\kappa_{\text{eff}} + 1} \right)^t \|\theta^{(0)} - \theta^*\|_2,$$

where the effective condition number is degraded by the overlap:

$$\kappa_{\text{eff}} = \frac{L}{\alpha - C' \epsilon_{\text{overlap}}}.$$

Proof. The proof follows directly from the strong convexity of $L(\theta)$ established by Weyl's inequality. As $\epsilon_{\text{overlap}} \rightarrow 0$, we recover the ideal convergence rate. ■

E.2 MULTI-MODAL ANALYSIS

In this section, we analyze the convergence properties for the K-Mode Mixture of Gaussians model. We explicitly model the **overlap** between Gaussian components as a perturbation.

E.2.1 THE OVERLAP FACTOR

We formally define the **Pairwise Overlap Factor** $\xi_{i,j}(x)$ between two components i and j :

$$\xi_{i,j}(x) \triangleq r_{k,i}(x) r_{k,j}(x). \quad (15)$$

And we define the **Maximum Expected Overlap** $\epsilon_{\text{overlap}}$ for the manifold as:

$$\epsilon_{\text{overlap}} = \max_i \sum_{j \neq i} \mathbb{E}_{x \sim p_t} [\xi_{i,j}(x)]. \quad (16)$$

This scalar $\epsilon_{\text{overlap}}$ quantifies the deviation from the ideal high separation regime. If components are perfectly separated, $\xi_{i,j} \rightarrow 0$ and $\epsilon_{\text{overlap}} \rightarrow 0$.

E.2.2 JACOBIAN DERIVATION

We need to compute the Jacobian of the score matching error vector $s_\theta(x, t) - \nabla \log p_t(x)$ with respect to the parameter $\mu_{k,l}$. Let $J_l^\mu(x) = \frac{\partial}{\partial \mu_{k,l}} \nabla \log p_{t,k}(x)$.

Similarly, we decompose the Jacobian for the l -th component into a **Signal Term** (Self) and a **Noise Term** (Interference).

$$J_\mu^l(x) = \underbrace{J_{\mu,\text{ideal}}^l(x)}_{\text{Signal}} + \underbrace{E_{\mu,\text{cross}}^l(x)}_{\text{Noise}}.$$

This term arises when we ignore the change in weights of other clusters ($j \neq l$). It dominates when $r_{k,l} \approx 1$:

$$J_{\mu,\text{ideal}}^l(x) \approx -\frac{s_t}{\gamma_t^2} r_{k,l}(x) \left(I - \frac{s_t^2}{s_t^2 + \gamma_t^2} U_{k,l} U_{k,l}^\top \right).$$

This term captures the gradient leaking into other clusters due to overlap:

$$E_{\mu,\text{cross}}^l(x) = \sum_{j=1}^{n_k} C'_1(x) \cdot \underbrace{r_{k,j}(x) r_{k,l}(x)}_{\xi_{j,l}(x)}, \quad (17)$$

where $C'_1(x)$ collects bounded vector terms. The norm of the error term is strictly bounded by the overlap:

$$\|E_{\mu,\text{cross}}^l(x)\|_2 \leq C'_1 \sum_{j \neq l} \xi_{j,l}(x).$$

For the Jacobian with respect to U_k , we have Similar derivation.

$$\|E_{U,\text{cross}}^l(x)\|_2 \leq C'_2 \sum_{j \neq l} \xi_{j,l}(x).$$

E.2.3 HESSIAN BLOCK STRUCTURE

The Hessian H for the parameters $\mu = [\mu_{k,1}, \dots, \mu_{k,n_k}]$ is a block matrix composed of $n_k \times n_k$ blocks, where each block is $D \times D$.

$$H_{\mu\mu} = \begin{pmatrix} H_{1,1} & H_{1,2} & \cdots & H_{1,n_k} \\ H_{2,1} & H_{2,2} & \cdots & H_{2,n_k} \\ \vdots & \vdots & \ddots & \vdots \\ H_{n_k,1} & H_{n_k,2} & \cdots & H_{n_k,n_k} \end{pmatrix}.$$

The (i, j) -th block is defined as:

$$H_{i,j} = \mathbb{E}_x[J_i^\mu(x)(J_j^\mu(x))^\top].$$

For diagonal blocks ($i = j = l$), the curvature is strictly determined by the expectation of the squared weights $\mathbb{E}[r_{k,l}(x)^2]$. Crucially, overlap causes **signal attenuation**, as the weight $r_{k,l}(x)$ drops below 1 in transition regions.

Using the identity $r_{k,l}(x)^2 = r_{k,l}(x)(1 - \sum_{j \neq l} r_{k,j}(x))$, we derive the exact expectation:

$$\begin{aligned} \mathbb{E}[r_{k,l}(x)^2] &= \mathbb{E}[r_{k,l}(x)] - \sum_{j \neq l} \mathbb{E}[r_{k,l}(x) r_{k,j}(x)] \\ &= \pi_{k,l} - \sum_{j \neq l} \mathbb{E}[\xi_{j,l}(x)] \\ &= \pi_{k,l} - \epsilon_{k,l}^{\text{total}}. \end{aligned}$$

Thus, we lower-bound the diagonal curvature by accounting for the total overlap mass $\epsilon_{k,l}^{\text{total}}$ leaking from cluster l :

$$H_{l,l} \approx \mathbb{E}[(J_l^{\text{ideal}})(J_l^{\text{ideal}})^\top] \succeq \lambda_{\text{diag},l} \cdot I,$$

where the effective base curvature is:

$$\lambda_{\text{diag},l} = (\pi_{k,l} - \epsilon_{k,l}^{\text{total}}) \min(\lambda_{\min}(H_{\mu_{k,l}\mu_{k,l}}), \lambda_{\min}(H_{U_{k,l}U_{k,l}}))$$

Here, the term $(\pi_{k,l} - \epsilon_{k,l}^{\text{total}})$ represents the effective probability mass contributing to convexity. This formulation explicitly shows that smaller clusters (small $\pi_{k,l}$) are significantly more vulnerable to instability, as the effective mass can vanish if the overlap $\epsilon_{k,l}^{\text{total}}$ becomes comparable to the cluster size $\pi_{k,l}$.

For $i \neq j$, the block $H_{i,j}$ represents the interference.

$$H_{i,j} \approx \mathbb{E}_x[J_i^{\text{ideal}}(J_j^{\text{ideal}})^\top] \propto \mathbb{E}[r_{k,i}(x)r_{k,j}(x)].$$

E.2.4 PERTURBATION ANALYSIS

We write the full Hessian as a sum of a block-diagonal matrix and a perturbation matrix:

$$H_{\mu\mu} = H_{\text{diag}} + \Delta H_{\text{overlap}}.$$

For the minimum eigenvalue of H_{diag} ,

$$\lambda_{\min}(H_{\text{diag}}) = \min_l \lambda_{\min}(H_{l,l}) = \min_l \lambda_{\text{diag},l} \triangleq \lambda_{\text{base}}.$$

For **Spectral Norm of $\Delta H_{\text{overlap}}$** , by Weyl's Inequality, the minimum eigenvalue of the full Hessian is:

$$\lambda_{\min}(H) \geq \lambda_{\min}(H_{\text{diag}}) - \|\Delta H_{\text{overlap}}\|_2.$$

and

$$\Delta H_{\text{overlap}} \leq \tilde{C} \cdot \mathbb{E}[\xi_{i,j}(x)], \quad (18)$$

where

$$\tilde{C} = 2(S_\mu C'_1 + S_U C'_2)$$

Substituting the bounds:

$$\lambda_{\min}(H) \geq \lambda_{\text{base}} - \tilde{C} \cdot \epsilon_{\text{overlap}}.$$

Therefore, H is positive definite **if and only if**:

$$\epsilon_{\text{overlap}} < \frac{\lambda_{\text{base}}}{\tilde{C}}.$$

Interpretation: The optimization landscape is locally strictly convex provided the overlap between clusters is smaller than the intrinsic curvature of the individual Gaussians.

E.2.5 FULL CONVERGENCE THEOREM

Combining the analysis of μ and the similar decoupling argument for U (using Schur complements to handle $H_{\mu U}$ terms which are also $O(\epsilon)$), we arrive at the final result.

Theorem E.2. *Let $\mathcal{L}(\theta)$ be the score matching loss. Assume the maximum expected overlap $\epsilon_{\text{overlap}}$ satisfies the condition $\epsilon_{\text{overlap}} < \tau$ for some threshold $\tau \propto \lambda_{\text{base}}$. Then the Hessian $H(\theta^*)$ is strictly positive definite.*

Linear Convergence: Gradient descent with step size η converges as:

$$\|\theta^{(t)} - \theta^*\|_2 \leq \rho^t \|\theta^{(0)} - \theta^*\|_2,$$

where the convergence rate $\rho < 1$ is determined by the effective condition number:

$$\kappa_{\text{eff}} = \frac{L}{\lambda_{\text{base}} - \tilde{C}\epsilon_{\text{overlap}}}.$$

This proves that the High Separation Assumption is not a binary requirement, but rather a continuum. The algorithm is robust to finite overlap, with the convergence rate degrading gracefully as the overlap increases.

Remark E.3. It is important to note that physically, $\epsilon_{\text{overlap}}$ will not be arbitrarily large.

F THE DETAIL OF THE REAL-WORLD EXPERIMENTS

In the part, we provide the detail of the experiments, including dataset and training pipeline. We use MNIST and CIFAR-10 as the datasets, and we adopt the mixture Gaussian distribution as the prior distribution in both cases.

For MNIST, our model consists of MLP-based encoder and decoder networks, each with a single hidden layer of 256 dimensions. The model is trained with the AdamW optimizer at a learning rate of 0.0005. We train 10 VAEs with the numbers 1 to 10 as the ten clusters.

On CIFAR-10, we implement a 3-layer RNN encoder and decoder for CIFAR-10. The encoder hidden dimensions are [64, 128, 256], and the decoder's are [256, 128, 64]. And we train 10 VAEs for each of the ten clusters based on the classification by category. Each layer in both networks stacks 3 recurrent blocks. The model is trained with the AdamW optimizer at a learning rate of 0.0001.

Our experiment was conducted on RTX4090.