

Understanding Guidance Scale in Diffusion Models from a Geometric Perspective

Zhiyuan Zhan

*The University of Tokyo, Japan
RIKEN Center for AIP, Japan*

zhan@ms.k.u-tokyo.ac.jp

Liuzhuozheng Li

The University of Tokyo, Japan

liuzhuozheng-li@outlook.com

Masashi Sugiyama

*RIKEN Center for AIP, Japan
The University of Tokyo, Japan*

sugi@k.u-tokyo.ac.jp

Reviewed on OpenReview: <https://openreview.net/forum?id=nfHmL6g8G>

Abstract

Conditional diffusion models have become a leading approach for generating condition-consistent samples, such as class-specific images. In practice, the guidance scale is a key hyperparameter in conditional diffusion models, used to adjust the strength of the guidance term. While empirical studies have demonstrated that appropriately choosing the scale can significantly enhance generation quality, the theoretical understanding of its role remains limited. In this work, we analyze the probabilistic guidance term from a geometric view under the linear manifold assumption and, based on this analysis, construct a geometric guidance model that enables tractable theoretical study. To address regularity issues arising from multi-modal data, we introduce a mollification technique that ensures well-posed dynamics. Our theoretical results show that increasing the guidance scale improves alignment with the target data manifold, thereby enhancing generation performance. We further extend our framework to nonlinear manifolds, and empirical results on real-world datasets validate the effectiveness of the proposed model and are consistent with our theories.

1 Introduction

Diffusion models (Ho et al., 2020; Song et al., 2021a) have achieved state-of-the-art performance on generative tasks across various domains, including images (Dhariwal & Nichol, 2021; Rombach et al., 2022), text-to-image synthesis (Saharia et al., 2022), videos (Ho et al., 2022), and audio (Kong et al., 2021). As a result, their empirical success has led to increasing interest in understanding the theoretical foundations of diffusion models (De Bortoli et al., 2021; Lee et al., 2022; Chen et al., 2023c;a; Gao et al., 2025). In particular, under the manifold hypothesis (Bengio et al., 2013), the ability of diffusion models to output high-quality samples in high-dimensional spaces motivates researchers to investigate how these models can generate distributions supported on low-dimensional manifolds in high-dimensional ambient spaces (De Bortoli, 2022; Oko et al., 2023; Li & Yan, 2024; Wan et al., 2025).

Controlling diffusion models to generate conditional distributions is another active area of research. Based on the theoretical framework proposed by Song et al. (2021b), both classifier guidance and classifier-free guidance models (Dhariwal & Nichol, 2021; Ho & Salimans, 2022) apply a probabilistic guidance term—derived from Bayes’ rule—to guide the sampling process toward the target conditional distribution. These methods also introduce a scale to adjust the strength of the guidance, and they showed that the performance depends strongly on the choice of the guidance scale and an appropriate value can significantly improve generation

quality. Recent empirical studies further demonstrated the importance of the guidance scale in conditional generation tasks (Dinh et al., 2023; Sadat et al., 2024; 2025). However, the theoretical understanding of how the guidance scale affects the generation remains limited (Chidambaram et al., 2024; Wu et al., 2024).

In this work, we propose a new geometric guidance model to enable the theoretical analysis of the role of the guidance scale in conditional generation. A key challenge in studying the guidance scale in classifier(-free) models is the analytical complexity of the probabilistic guidance term. To address this, we replace the probabilistic guidance with a new geometric guidance term. Specifically, under the linear manifold hypothesis (Chung et al., 2022), we study the geometric property of the original probabilistic guidance term, building on an idea introduced by Chen et al. (2023b), and construct a linear geometric guidance term that plays the same role but more tractable for theoretical analysis.

As a next step, the analysis of the geometric guidance model requires certain regularity conditions on the score function, such as the Lipschitz continuity. However, because of the multi-modality of data distributions, these conditions generally fail to hold (Lee et al., 2022; Gao et al., 2025). To overcome this issue, we introduce a mollification technique inspired by mollifiers in mathematical analysis (Evans, 2018) to construct a surrogate score function that satisfies the required properties for our analysis.

Building on this, we construct a well-posed geometric guidance model through which we address two questions: (i) whether the model can recover the target data manifold, and (ii) what is the upper bound on the distance between the generated distribution and the target conditional distribution. Our results reveal the effects of the guidance scale: increasing the scale encourages the generated data to lie closer to the target manifold, and large guidance scales do not significantly increase an upper bound on the generation error.

Finally, for the nonlinear case and real-world data distributions, we extend our framework by constructing a nonlinear geometric guidance model. This model builds on the same principles as the linear case, with the theoretical foundation obtained by extending the results of Chung et al. (2022) to nonlinear data manifolds. Experimentally, we evaluate the nonlinear geometric guidance model on CIFAR-10 (Krizhevsky, 2009) and demonstrate its effectiveness for conditional generation. We also report how performance varies with the guidance scale, providing empirical evidence consistent with the behavior suggested by our linear analysis.

In summary, our contributions are:

1. We construct a new linear geometric guidance term to replace the original probabilistic guidance term by studying its geometric property under the linear manifold hypothesis.
2. To ensure the regularity of the unconditional score function, we apply a mollification technique to construct a surrogate score function, and build a well-posed geometric guidance model.
3. By analyzing the geometric guidance model, we uncover the role of the guidance scale: a large guidance scale encourages the generated data to lie closer to the target data manifold and does not significantly affect the upper bound of the generation error.
4. We propose a principled nonlinear geometric guidance model and evaluate it on CIFAR-10; the experiments demonstrate its effectiveness in conditional generation and illustrate guidance-scale effects beyond the linear setting.

The remainder of this paper is organized as follows. Section 2 reviews related work, and Section 3 summarizes the technical background on diffusion models. Section 4 introduces the construction of the geometric guidance term, and Section 5 presents the theoretical analysis of the geometric guidance model. Section 6 extends the model to nonlinear settings and reports experimental results. Section 7 concludes the paper and discusses limitations. Notation is summarized in Appendix A.

2 Related Works

Convergence analysis: A number of recent works have analyzed the convergence properties of diffusion models under various assumptions (De Bortoli et al., 2021; Lee et al., 2022; 2023; Chen et al., 2023c;a; Gao et al., 2025). De Bortoli et al. (2021) established total variation bounds under C^3 -regularity assumptions

on the score for the target distribution. Chen et al. (2023c) relaxed this requirement to Lipschitz continuity of the score function but for each intermediate density, which was further weakened in Chen et al. (2023a) to the Lipschitz continuity of the score only for the target density. Using functional inequalities, Lee et al. (2022; 2023) and Gao et al. (2025) have derived convergence guarantees under the assumption that the target density function is log-concave, with results in both total variation and Wasserstein distances. In contrast, our setting involves multi-modal target distributions for which log-concavity and smoothness assumptions do not hold (Lee et al., 2022). To address this, we introduce a technique that constructs a surrogate distribution satisfying the required regularity properties while closely approximating the original target.

Geometric structure: For real-world datasets, it is widely believed that high-dimensional data lie on a low-dimensional submanifold of the ambient space, a perspective known as the manifold hypothesis (Bengio et al., 2013). When generating such data distributions, deep generative models often encounter challenges such as the curse of dimensionality (Bronstein et al., 2021) and manifold overfitting (Loaiza-Ganem et al., 2022). However, the strong empirical performance suggests that diffusion models can avoid these issues. As a result, understanding the theoretical behavior of diffusion models under the manifold hypothesis has attracted increasing attention. For example, De Bortoli (2022) established a Wasserstein convergence bound assuming that the target distribution is supported on a compact set. Under the additional assumption that the target data manifold is linear, Oko et al. (2023) showed that diffusion models can avoid the curse of dimensionality by providing a Wasserstein bound that depends only on the intrinsic dimension. Chen et al. (2023b) further derived a total variation bound in terms of the intrinsic dimension based on a decomposition of the score function under the linear manifold assumption. Following this line of work, we further investigate the geometric structure of this decomposition to clarify the role of the score function in recovering the target data manifold, which in turn helps us construct the geometric guidance model.

Conditional generation: To control the generation (Song et al., 2021b), Dhariwal & Nichol (2021) and Ho & Salimans (2022) applied the probabilistic guidance term to generate conditional distributions. Following their works and based on the geometric structure of noisy data manifolds under the linear assumption of the target data manifold (Chung et al., 2022), Chung et al. (2022; 2024) and He et al. (2024) proposed using a new time-dependent guidance in conditional generation to constrain geometric structure of the generation process. From a different perspective, Song et al. (2023) and Bansal et al. (2023) constructed a time-independent guidance constructed by a loss function that is designed to enforce desired constraints on the generated data. Instead, our geometric guidance is constructed by studying the geometric property of the probabilistic guidance, with the goal of replacing its role in conditional generation.

To adjust the strength of guidance, Dhariwal & Nichol (2021) and Ho & Salimans (2022) also introduced a guidance scale, and their experiments showed that selecting an appropriate scale can significantly improve performance. However, there are limited works on theoretically analyzing the effects of the guidance scale in conditional generation. Chidambaram et al. (2024) studied one-dimensional case and showed that increasing the scale not only reduces diversity of generated distributions but also leads generated data to drift to the extreme points in the support of the conditional distribution. Wu et al. (2024) theoretically analyzed the influence of the guidance scale in the context of Gaussian mixture models, demonstrating that a large guidance scale diminishes distributional diversity while boosting classification confidence. Due to the analytical complexity of the probability guidance term, previous works have focused on special cases. Therefore, we propose a geometric guidance term that plays the same role as the probabilistic guidance but is more amenable to theoretical analysis of the guidance scale.

3 Background

3.1 Diffusion Model

Let $\mathbf{X} \sim \mathbb{P}_X \in \mathcal{P}(\mathbb{R}^D)$ denote the target data distribution. The forward process in denoising diffusion probabilistic models (DDPMs) (Ho et al., 2020) is governed by the stochastic differential equation (SDE)

$$d\mathbf{X}_t = -\frac{1}{2}\beta(t)\mathbf{X}_t dt + \sqrt{\beta(t)}d\mathbf{W}_t, \quad \forall t \in [0, T], \quad (1)$$

with the initial condition $\mathbf{X}_0 \sim \mathbb{P}_X$, where $(\mathbf{W}_t)_{t \geq 0}$ is a standard Brownian motion and $\beta: [0, T] \rightarrow (0, \infty)$ is smooth; see Song et al. (2021b). This SDE admits the following analytical solution:

$$\mathbf{X}_t \stackrel{d}{=} \sqrt{\alpha_t} \mathbf{X}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\xi}, \quad \forall t \in [0, T], \quad (2)$$

where $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$ a standard Gaussian, $\alpha_t := \exp\left(-\int_0^t \beta(s) ds\right)$, and “ $\stackrel{d}{=}$ ” means equal in distribution. The derivation is provided in Appendix B.1.

The reverse process of DDPMs aims to generate \mathbb{P}_X , which corresponds to the time-reversal process of (1). To this end, we need to consider the process

$$\mathbf{X}_t^{\leftarrow} := \mathbf{X}_{T-t}$$

and study its stochastic dynamics. As shown in Anderson (1982) and Haussmann & Pardoux (1986), the process $(\mathbf{X}_t^{\leftarrow})_{t \in [0, T]}$ satisfies the following SDE:

$$d\mathbf{X}_t^{\leftarrow} = \left(\frac{1}{2} \beta(T-t) \mathbf{X}_t^{\leftarrow} + \beta(T-t) \nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^{\leftarrow}) \right) dt + \sqrt{\beta(T-t)} d\bar{\mathbf{W}}_t, \quad (3)$$

where p_t is the density function of \mathbf{X}_t , and $(\bar{\mathbf{W}}_t)_{t \in [0, T]}$ is the Brownian motion in reverse time. A simplified proof can be found in Tang & Zhao (2024).

In practice, a neural network $\mathbf{s}_\theta(t, \cdot)$ with parameter θ is trained to estimate the score function $\nabla_{\mathbf{x}} \log p_t(\cdot)$ using the score matching method (Vincent, 2011). By substituting $\nabla_{\mathbf{x}} \log p_t$ with the estimator $\mathbf{s}_\theta(t, \cdot)$ in (3), experiments (Song & Ermon, 2019; Song et al., 2021b; Dhariwal & Nichol, 2021) showed that DDPMs achieve state-of-the-art performance in data generation tasks.

3.2 Probability Flow ODE

Instead of simulating the stochastic process (3), denoising diffusion implicit models (DDIMs) (Song et al., 2021a) employ a deterministic approach for generation, which corresponds to the following ordinary differential equation (ODE):

$$\frac{d}{dt} \mathbf{X}_t^{\leftarrow} = \frac{1}{2} \beta(T-t) (\mathbf{X}_t^{\leftarrow} + \nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^{\leftarrow})), \quad \forall t \in [0, T], \quad (4)$$

with the initial condition $\mathbf{X}_0^{\leftarrow} \sim p_T$, which is called the probability flow ODE. The evolution of the density functions of $\mathbf{X}_t^{\leftarrow}$ under this deterministic process is equivalent to that of the stochastic reverse process (3), as the continuity equation associated with the ODE coincides the Fokker–Planck equation corresponding to the SDE (1); see Song et al. (2021b) for details.

In this paper, we focus on the deterministic dynamics, as the Wasserstein distance used as the main metric makes analyzing the ODE formulation more convenient than the SDE. It naturally extends to the SDE via Itô’s formula (Gao et al., 2025). Following Chen et al. (2023c) and Chen et al. (2023a), we consider the Ornstein–Uhlenbeck process by setting $\beta(t) \equiv 2$ in Equation (1) for simplicity, where this constant choice is unimportant, as varying it merely rescales time.

3.3 Conditional Diffusion Model

When working with paired data $(\mathbf{X}, Y) \sim \mathbb{P}_{XY}$, the goal of conditional generation is to generate the conditional distribution $\mathbb{P}_{X|Y}(\cdot | Y)$. In Song et al. (2021b), diffusion models are directly applied to $\mathbb{P}_{X|Y}(\cdot | Y)$. Specifically, the forward process (1) is first run with the initial condition $\mathbf{X}_0 \sim \mathbb{P}_{X|Y}(\cdot | Y)$ to obtain the density functions p_t^y of \mathbf{X}_t . Then, the stochastic reverse process (3), or the deterministic process (4), is simulated to generate samples from $\mathbb{P}_{X|Y}(\cdot | Y)$.

Moreover, these intermediate densities p_t^y admit more explicit expressions. Suppose $(\mathbf{X}, Y) \sim \mathbb{P}_{XY}$ and we run the SDE (1) with initial condition $\mathbf{X}_0 \sim \mathbb{P}_X = \int \mathbb{P}_{XY}(\cdot, dy)$ to obtain \mathbf{X}_t . Let $p_t(\mathbf{x}_t, y)$ denote the joint density function of (\mathbf{X}_t, Y) . Then, it can be shown that

$$p_t(\mathbf{x}_t | y) = p_t^y(\mathbf{x}_t);$$

see Appendix B.2 for details.

Therefore, the score function for generating $\mathbb{P}_{X|Y}(\cdot | Y)$ can be decomposed as

$$\nabla_{\mathbf{x}} \log p_t^y(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_t(\mathbf{x} | y) = \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) + \nabla_{\mathbf{x}} \log p_t(y | \mathbf{x}), \quad (5)$$

where $p_t(\mathbf{x})$ is the marginal density of \mathbf{X}_t obtained by running (1) with the initial condition $\mathbf{X}_0 \sim \mathbb{P}_X$. This term can be estimated using standard methods from unconditional DDPMs. The remaining term, $\nabla_{\mathbf{x}} \log p_t(y | \mathbf{x})$, is known as the guidance term, and there are two main approaches for approximating it: classifier guidance and classifier-free guidance (Dhariwal & Nichol, 2021; Ho & Salimans, 2022). In classifier guidance, a time-dependent classifier is trained to approximate $p_t(y | \cdot)$ on all noisy data. In classifier-free guidance, a new neural network $\mathbf{s}_\theta(t, \mathbf{x}, y)$ is trained to estimate the conditional score $\nabla_{\mathbf{x}} \log p_t(\mathbf{x} | y)$, while $\mathbf{s}_\theta(t, \mathbf{x}, \emptyset)$ approximates the unconditional score $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$. The guidance term is then computed as

$$\nabla_{\mathbf{x}} \log p_t(y | \mathbf{x}) \approx \mathbf{s}_\theta(t, \mathbf{x}, y) - \mathbf{s}_\theta(t, \mathbf{x}, \emptyset).$$

In practice, a scaling parameter $\eta > 0$, known as the guidance scale, is typically introduced to control the strength of the guidance term (Dhariwal & Nichol, 2021). When using the deterministic dynamics (4), this modification is mathematically expressed as

$$\frac{d}{dt} \mathbf{X}_t^\leftarrow = \mathbf{X}_t^\leftarrow + \nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^\leftarrow) + \eta \nabla_{\mathbf{x}} \log p_{T-t}(y | \mathbf{X}_t^\leftarrow), \quad \forall t \in [0, T], \quad (6)$$

with the initial condition $\mathbf{X}_0^\leftarrow \sim p_T(\cdot | y)$.

As mentioned in Section 2, although setting $\eta \neq 1$ may seem counterintuitive from a theoretical perspective, empirical studies (Dhariwal & Nichol, 2021; Ho & Salimans, 2022) have shown that selecting an appropriate value of η can significantly improve performance. In particular, increasing the guidance scale η enhances the distinguishability of generated samples, but at the cost of reduced diversity (Ho & Salimans, 2022; Chidambaram et al., 2024; Wu et al., 2024). However, theoretical understanding of how the guidance scale η influences generation remains limited, due to the analytical complexity of the guidance term $\nabla_{\mathbf{x}} \log p_t(y | \mathbf{x})$ (Chidambaram et al., 2024; Wu et al., 2024).

Therefore, the main objective of this work is to provide a theoretical analysis of the guidance scale η , under the assumption that the target data concentrate on a low-dimensional linear subspace $\mathcal{M}_y \subset \mathbb{R}^D$, called the target data manifold, i.e., $\text{supp } \mathbb{P}_{X|Y}(\cdot | Y = y) \subset \mathcal{M}_y$. This analysis consists of two main steps:

- (i) First, we replace the probabilistic guidance term with a geometric guidance term in order to avoid the difficulty of handling $\nabla_{\mathbf{x}} \log p_t(y | \mathbf{x})$ (see Section 4).
- (ii) Second, we analyze the modified dynamics under the geometric guidance from two perspectives: (a) how η influences the recovery of the target data manifold, and (b) how it affects the distance between the generated distribution and the target distribution (see Section 5).

A central technical challenge in analyzing the geometric guidance dynamics is that $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ may fail to satisfy desirable properties, such as the L -Lipschitz continuity (and the log-concavity of $p_t(\mathbf{x})$), due to the fact that $p_t(\mathbf{x})$ arises from a diffusion process initialized with a multi-modal distribution (Lee et al., 2022; Gao et al., 2025). To address this issue, we introduce a novel technique inspired by mollification in mathematical analysis (Evans, 2018), which yields a surrogate distribution $p_t^\sigma(\mathbf{x})$ for which the geometric guidance dynamics is well-posed.

4 Geometric Guidance Model

In this section, our main objective is to construct a new guidance term to replace $\nabla_{\mathbf{x}} \log p_t(y | \mathbf{x})$ in Equation (6) from a geometric perspective. Specifically, the key idea is to understand the role that $\nabla_{\mathbf{x}} \log p_t(y | \mathbf{x})$ plays in recovering the target data manifold \mathcal{M}_y .

Note that $\nabla_{\mathbf{x}} \log p_t(y | \mathbf{x})$ appears as a component of $\nabla_{\mathbf{x}} \log p_t^y(\mathbf{x})$ by Equation (5). This motivates us to investigate the geometric interpretation of the score function $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ in the setting of unconditional

DDPMs (see Section 4.1). Based on Equation (5) and a basic property of $p_t(y \mid \mathbf{x})$, we then propose a replacement for $\nabla_{\mathbf{x}} \log p_t(y \mid \mathbf{x})$, which preserves its geometric role but is more tractable for theoretical analysis (see Section 4.2).

4.1 Geometric Interpretation of Score Function

To study the geometric properties of the score function $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$, we first examine the geometric structure of the noisy data manifolds that arise during the DDPM process. Chung et al. (2022) showed that, under the assumption that the target data lie on $\mathcal{M} \in \mathbb{R}^D$, a linear subspace of the ambient space \mathbb{R}^D with significantly lower dimension, the noisy data \mathbf{X}_t concentrate on a hypersurface, i.e., a $(D-1)$ -dimensional manifold embedded in \mathbb{R}^D , for any $t > 0$. We generalize this result (Chung et al., 2022, Proposition 1) in the following proposition; the proof is provided in Appendix C.1.

Proposition 1. Assume $\mathbf{Z} \sim \mathbb{P}^Z$ on \mathbb{R}^d , and $\mathbf{X} = A\mathbf{Z} \sim \mathbb{P}_X$ on \mathbb{R}^D for an $A \in \mathcal{O}^{D \times d}$, i.e., $A \in \mathbb{R}^{D \times d}$ and $A^\top A = \mathbf{I}_d$. Define

$$\mathcal{M}^t := \{\mathbf{x} \in \mathbb{R}^D : \|(\mathbf{I}_D - AA^\top)\mathbf{x}\| = r(t)\},$$

where $r(t) := \sqrt{(D-d)(1-\alpha_t)}$ and $\alpha_t = e^{-2t}$. Let \mathbf{X}_t be generated by the DDPM forward process (1) with the initial condition $\mathbf{X}_0 = \mathbf{X}$. If $d \ll D$, then \mathbf{X}_t concentrates on \mathcal{M}^t with high probability.

Based on this result, the next question is how the score function $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ contributes to recovering these noisy data manifolds \mathcal{M}^t during the reverse process (4).

Under the same assumptions as those in Proposition 1, Chen et al. (2023b) showed that

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = A \nabla_{\mathbf{z}} \log p_t^Z(\mathbf{z})|_{\mathbf{z}=A^\top \mathbf{x}} - \frac{1}{1-\alpha_t} (\mathbf{I}_D - AA^\top) \mathbf{x}, \quad (7)$$

where p_t^Z is the density associated with the forward process (1) initialized from p^Z . An alternative derivation of this formula, along with an analysis of its geometric properties, is provided in Appendix C.2.

Based on this orthogonal decomposition, we observe that the role of $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ can be understood as two components: (i) the first term serves as generating the distribution \mathbb{P}^Z in the latent space, and (ii) the second term controls the reconstruction of the noisy data manifolds \mathcal{M}^t in the ambient space. Informally, this decomposition can be summarized as

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = \text{Generate Latent Distribution} + \text{Recover Data Manifolds } \mathcal{M}^t.$$

We formalize this intuition in the following theorem; see the proof in Appendix C.1.

Theorem 2. Under the same setting as that in Proposition 1, let $\mathbf{X}_{t,\parallel}^\leftarrow = AA^\top \mathbf{X}_t^\leftarrow$ and $\mathbf{X}_{t,\perp}^\leftarrow = \mathbf{X}_t^\leftarrow - \mathbf{X}_{t,\parallel}^\leftarrow$, where $\mathbf{X}_t^\leftarrow = \mathbf{X}_{T-t}$.

(a) Let $\mathbf{X}_{t,\parallel}^\leftarrow = A\mathbf{Z}_t^\leftarrow$ with $\mathbf{Z}_t^\leftarrow = A^\top \mathbf{X}_t^\leftarrow$. Then \mathbf{Z}_t^\leftarrow satisfies

$$\frac{d}{dt} \mathbf{Z}_t^\leftarrow = \mathbf{Z}_t^\leftarrow + \nabla_{\mathbf{z}} \log p_{T-t}^Z(\mathbf{Z}_t^\leftarrow),$$

which implies that $\mathbf{Z}_t = A^\top \mathbf{X}_t = \mathbf{Z}_{T-t}^\leftarrow$ follows the forward process (1) initialized from p^Z .

(b) $\mathbf{X}_{t,\perp}^\leftarrow$ satisfies

$$\frac{d}{dt} \mathbf{X}_{t,\perp}^\leftarrow = \mathbf{X}_{t,\perp}^\leftarrow - \frac{1}{1-\alpha_{T-t}} \mathbf{X}_{t,\perp}^\leftarrow.$$

Moreover, $\|\mathbf{X}_{t_0,\perp}^\leftarrow\| = r(T-t_0)$ implies $\|\mathbf{X}_{t_0+\delta,\perp}^\leftarrow\| = r(T-t_0-\delta)$, where $r(t) = \sqrt{(D-d)(1-\alpha_t)}$.

In Theorem 2, statement (a) demonstrates that the parallel part $\nabla_{\mathbf{z}} \log p^Z(\mathbf{z})$ in the decomposition (7) is responsible for generating the target latent distribution p^Z via the reverse process of DDPMs, which has been thoroughly studied in Chen et al. (2023b). Meanwhile, statement (b) shows that, since

$$\|(\mathbf{I}_D - AA^\top) \mathbf{X}_t^\leftarrow\| = \|\mathbf{X}_{t,\perp}^\leftarrow\|,$$

the orthogonal part $(\mathbf{I}_D - AA^\top) \mathbf{x}$ plays a key role in guiding the recovery of the noisy data manifolds \mathcal{M}^t , which provides an insight for designing geometric guidance in conditional generation.

4.2 Geometric Guidance for Conditional Generation

Let us return to the conditional diffusion model. To apply the results from Section 4.1 in studying the role of $\nabla_{\mathbf{x}} \log p_t(y | \mathbf{x})$ in guidance, we first impose the linear assumption for the target data manifold.

We consider a two-class dataset $(\mathbf{X}, Y) \sim \mathbb{P}_{XY}$ on $\mathbb{R}^D \times \{1, 2\}$ for simplicity; the following analysis readily extends to the multi-class case. Let $\mathbb{P}(Y = 1) = w_1$ and $\mathbb{P}(Y = 2) = w_2$ so that

$$\mathbb{P}_X = w_1 \mathbb{P}_{X|Y}(\cdot | Y = 1) + w_2 \mathbb{P}_{X|Y}(\cdot | Y = 2).$$

The linear assumption states as follows.

Assumption I. For $i = 1, 2$, there exists a $\mathbf{Z}_i \sim p_i^Z$ on \mathbb{R}^{d_i} and an $A_i \in \mathcal{O}^{D \times d_i}$ such that

$$\mathbf{X}_i := A_i \mathbf{Z}_i \sim \mathbb{P}_{X|Y}(\cdot | Y = i),$$

and we further assume $A_1^\top A_2 = \mathbf{O}$.

Remark 1. For this assumption, we provide two remarks.

- (i) It basically means that the support $\text{supp } \mathbb{P}_{X|Y}(\cdot | Y = i) \subset \mathcal{M}_i := \text{Im } A_i$, the image of $\mathbf{x} \mapsto A\mathbf{x}$; in other words, $\mathbb{P}_{X|Y}(\cdot | Y = i)$ is supported on the linear space $\text{Im } A_i$. The definition of the support of a probability measure is provided in Appendix A.
- (ii) $A_1^\top A_2 = \mathbf{O}$ indicates $\mathcal{M}_1 \perp \mathcal{M}_2$. This orthogonality assumption is introduced to simplify the subsequent analysis, but it does not significantly affect our conclusions regarding the guidance scale; see Appendix E.1 for further discussion.

Next, we fix $Y = 1$ and our goal is to generate the conditional distribution, which needs to consider the geometric structure of the condition score function $\nabla_{\mathbf{x}} \log p_t(\mathbf{x} | y = 1)$. By combining the results in Section 4.1 with Equation (5), the conditional score function has two different types of decomposition:

$$\begin{aligned} \nabla_{\mathbf{x}} \log p_t(\mathbf{x} | y = 1) &= \text{Generate Latent Distribution} + \text{Recover Data Manifolds } \mathcal{M}_1^t \\ &= \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) + \nabla_{\mathbf{x}} \log p_t(y = 1 | \mathbf{x}). \end{aligned} \quad (8)$$

We will show that $\nabla_{\mathbf{x}} \log p_t(y = 1 | \mathbf{x})$ plays the role of recovering the data manifolds \mathcal{M}_1^t with respect to the first decomposition.

For the first decomposition in (8), based on Assumption I and Proposition 1, because the noisy data manifolds generated by the forward process starting from \mathcal{M}_1 are given by

$$\mathcal{M}_1^t = \{\mathbf{x} \in \mathbb{R}^D : \|(\mathbf{I}_D - A_1 A_1^\top) \mathbf{x}\| = r(t)\}, \quad (9)$$

the orthogonal part of $\nabla_{\mathbf{x}} \log p_t(\mathbf{x} | y = 1)$ in the first decomposition responsible for recovering \mathcal{M}_1^t is parallel to $(\mathbf{I}_D - A_1 A_1^\top) \mathbf{x}$ as shown in Section 4.1.

Intuitively, for the second decomposition in (8), since $p_t(y = 1 | \mathbf{x})$ acts as a classifier for \mathcal{M}_1^t , we have $p_t(y = 1 | \mathbf{x}) \approx 1$ for any $\mathbf{x} \in \mathcal{M}_1^t$, i.e., $\log p_t(y = 1 | \mathbf{x})$ is approximately constant on \mathcal{M}_1^t . Therefore, by Lemma I.1, $\nabla_{\mathbf{x}} \log p_t(y = 1 | \mathbf{x})$ is almost normal to \mathcal{M}_1^t ,

$$\nabla_{\mathbf{x}} \log p_t(y = 1 | \mathbf{x}) \approx -\eta(\mathbf{I}_D - A_1 A_1^\top) \mathbf{x}, \quad \text{for some } \eta > 0,$$

because $(\mathbf{I}_D - A_1 A_1^\top) \mathbf{x}$ is normal to \mathcal{M}_1^t by Lemma I.1. Rigorous details are provided in Appendix C.3.

Therefore, the guidance term $\nabla_{\mathbf{x}} \log p_t(y = 1 | \mathbf{x})$ partially contributes to the recovery of the data manifolds \mathcal{M}_1^t during the reverse process. Consequently, it can be replaced by $(\mathbf{I}_D - A_1 A_1^\top) \mathbf{x}$. Based on this insight, we propose the following geometric guidance model for conditional generation:

$$\frac{d}{dt} \mathbf{X}_t^\leftarrow = \mathbf{X}_t^\leftarrow + \nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^\leftarrow) - \eta P_1 \mathbf{X}_t^\leftarrow, \quad P_1 := \mathbf{I}_D - A_1 A_1^\top. \quad (10)$$

5 Main Results: Analysis of Geometric Guidance Model

In this section, we analyze the geometric guidance model (10) with the aim of uncovering the role of the guidance scale η . To understand its effects, we consider two related questions: whether the model can approximately estimate the target data manifold \mathcal{M}_1 (see Section 5.2), and how to quantify the distance between the generated and target distributions (see Section 5.3). These two problems serve as a lens through which we investigate the influence of η in conditional generation.

Before addressing these two questions, it is necessary to ensure the well-posedness of the ODE (10); that is, we must establish regularities of $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ such as its Lipschitz continuity and the log-concavity of $p_t(\mathbf{x})$, which requires careful analysis (see Section 5.1) because it is obtained from a multi-modal distribution \mathbb{P}_X .

5.1 Well-posedness of Geometric Guidance Model

In general, the Lipschitz continuity of $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ and the log-concavity of $p_t(\mathbf{x})$ induced by the DDPM forward process depend on properties of the initial distribution μ . A basic requirement is that μ admit a density $p(\mathbf{x})$. Log-concavity of $p(\mathbf{x})$ then implies log-concavity of $p_t(\mathbf{x})$ (Gao et al., 2025), and Lipschitz continuity of $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ implies the Lipschitz continuity of $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ (Chen et al., 2023a).

However, in our setting, it is clear that \mathbb{P}_X does not admit a density function. We therefore first deduce the necessary conditions on the latent distribution implied by \mathbb{P}_X ; see Sections 5.1.1 and 5.1.2. Second, the multi-modality of \mathbb{P}_X introduces irregularities in $p_t(\mathbf{x})$ (Lee et al., 2022), which we discuss in Section 5.1.3. By solving these two problems, we construct a surrogate $p_t^g(\mathbf{x})$ for use in the geometric guidance model (10), which is well-posed; see Section 5.1.4.

5.1.1 Problems in Latent Distribution

When μ does not admit a density function—for instance, when the support of μ lies on a lower-dimensional manifold in the ambient space—De Bortoli (2022) showed that the score function $\nabla_{\mathbf{x}} \log p_t$ is Lipschitz continuous under the assumption that $\text{supp } \mu$ is compact, i.e., closed and bounded. This setting aligns with our problem but guarantees only Lipschitz continuity. In contrast, we establish a stronger result in the following Proposition 3, which does not require the compactness, under the assumption that the target data manifold is linear. The proof is provided in Appendix D.1.

Proposition 3. *Let \mathbf{Z} be a random variable on \mathbb{R}^k with the density function p^Z , and let $B \in \mathbb{R}^{n \times k}$. Assume there are $m_0, \Lambda > 0$ such that*

$$-\nabla_{\mathbf{z}}^2 \log p^Z(\mathbf{z}) \succeq m_0 I_k, \quad \|B\|_{\text{op}}^2 \leq \Lambda,$$

and $\lambda := \lambda_{\min}(B^\top B) \geq 0$, the minimum of all eigenvalues of $B^\top B$. For $\alpha \in \mathbb{R}$ and $\beta > 0$, let

$$\mathbf{X} = \alpha B \mathbf{Z} + \beta \boldsymbol{\xi}, \quad \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I_n)$$

with the density function p_X on \mathbb{R}^n . We have

$$\|\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x})\|_{\text{op}} \leq L, \quad L := \frac{1}{\beta^2} + \frac{\alpha^2 \Lambda}{\beta^2(\alpha^2 \lambda + m_0 \beta^2)}.$$

Remark 2. A direct application of this proposition is that it extends the result of De Bortoli (2022) to a non-compact setting, under the additional assumption that the latent distribution is strongly log-concave, i.e., $-\nabla_{\mathbf{z}}^2 \log p^Z(\mathbf{z}) \succeq m_0 I_k$. If we are only concerned with the L -smoothness¹ of $\log p_X$, the log-concavity of p^Z can be relaxed to the L -smoothness; see Appendix F for details.

Corollary 4. *Using the same notations as in Proposition 3, we have*

$$\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x}) \preceq \left(\frac{\alpha^2 \Lambda}{\beta^2(\alpha^2 \lambda + m_0 \beta^2)} - \frac{1}{\beta^2} \right) I_n.$$

¹ L -smoothness of f and L -Lipschitz continuity of ∇f are equivalent for C^2 functions; we use them interchangeably.

Remark 3. Note that, if $\alpha^2\Lambda < \alpha^2\lambda + m_0\beta^2$, such as $\Lambda = \lambda$, then

$$-\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x}) \succeq m_x \mathbf{I}_n, \quad m_x := \frac{1}{\beta^2} \left(1 - \frac{\alpha^2\Lambda}{\alpha^2\lambda + m_0\beta^2} \right) > 0,$$

which implies that p_X is m_x -strongly log-concave. Therefore, it shows that the strong log-concavity of p^Z ensures not only the L -smoothness, but also the concavity of $\log p_X$.

Based on results in Proposition 3, even if \mathbb{P}_X does not admit a density function, the desired properties of the score function can still be guaranteed, provided that the latent distribution admits a density and satisfies strong log-concavity. However, in our setting, these two conditions are not satisfied:

- (i) For the latent distribution of \mathbb{P}_X , because $\mathbf{Z}_i \sim \mathbb{P}_i^Z$ on \mathbb{R}^{d_i} , we first lift them on \mathbb{R}^d with $d := d_1 + d_2$ by defining

$$\tilde{\mathbf{Z}}_1 = (\mathbf{I}_{d_1}, \mathbf{O}_{d_1 \times d_2})^\top \mathbf{Z}_1 \sim \tilde{\mathbb{P}}_1^Z, \quad \tilde{\mathbf{Z}}_2 = (\mathbf{O}_{d_2 \times d_1}, \mathbf{I}_{d_2})^\top \mathbf{Z}_2 \sim \tilde{\mathbb{P}}_2^Z.$$

Let $A = (A_1, A_2) \in \mathcal{O}^{D \times d}$. It follows that

$$A\tilde{\mathbf{Z}}_i = A_i \mathbf{Z}_i \sim \mathbb{P}_{X|Y}(\cdot | Y = i), \text{ i.e., } A_{\#} \tilde{\mathbb{P}}_i^Z = \mathbb{P}_{X|Y}(\cdot | Y = i).$$

Therefore, by Lemma H.7, if $\mathbf{Z} \sim \mathbb{P}^Z := w_1 \tilde{\mathbb{P}}_1^Z + w_2 \tilde{\mathbb{P}}_2^Z$, we have

$$\mathbf{X} = A\mathbf{Z} \sim \mathbb{P}_X = w_1 \mathbb{P}_{X|Y}(\cdot | Y = 1) + w_2 \mathbb{P}_{X|Y}(\cdot | Y = 2).$$

But the problem is that the latent distribution \mathbb{P}^Z does not admit a density function on \mathbb{R}^d .

- (ii) For log-concavity, even if the latent distribution admits a density function, it typically does not satisfy strong log-concavity due to its multi-modality (Lee et al., 2022).

Therefore, in the following, we first introduce a technique to address the log-concavity of the latent density (Sections 5.1.2 and 5.1.3), and then apply Proposition 3 to establish the desired properties of the score function (Section 5.1.4).

5.1.2 Mollification Technique

Mollification (Evans, 2018) is a standard technique in mathematical analysis to address non-smoothness of functions. When dealing with a non-smooth function f , the idea is to find a smooth kernel function k such that the convolution $g := f * k$, which is clearly smooth, is closed to f .

Following this idea, we choose a Gaussian distribution $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$ with some $\sigma > 0$ as the kernel, and consider its convolution with $\tilde{\mathbb{P}}^Z$; see Remark B.1 for the definition of convolution between measures. Let

$$\mathbb{P}_\sigma^Z := \tilde{\mathbb{P}}^Z * \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d) = w_1 \mathbb{P}_{1,\sigma}^Z + w_2 \mathbb{P}_{2,\sigma}^Z,$$

where $\mathbb{P}_{i,\sigma}^Z := \tilde{\mathbb{P}}_i^Z * \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$. Note that both \mathbb{P}_σ^Z and $\mathbb{P}_{i,\sigma}^Z$ admit density functions, denoted by p_σ^Z and $p_{i,\sigma}^Z$ respectively, and

$$p_\sigma^Z = w_1 p_{1,\sigma}^Z + w_2 p_{2,\sigma}^Z. \quad (11)$$

Moreover, by the definition of convolution, if $\mathbf{Z}_i \sim p_i^Z$, then

$$\mathbf{Z}_{1,\sigma} := (\mathbf{Z}_1, \mathbf{O})^\top + \sigma \boldsymbol{\zeta}_1 \sim p_{1,\sigma}^Z, \quad \mathbf{Z}_{2,\sigma} := (\mathbf{O}, \mathbf{Z}_2)^\top + \sigma \boldsymbol{\zeta}_2 \sim p_{2,\sigma}^Z, \quad (12)$$

for $\boldsymbol{\zeta}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ independent of \mathbf{Z}_i . Therefore, we obtain a smooth density p_σ^Z on the latent space \mathbb{R}^d .

Next, for \mathbb{P}_X , the following Proposition 5 addresses the question of whether sampling from p_σ^Z yields a $\mathbb{P}_X^\sigma := A_{\#} \mathbb{P}_\sigma^Z$ that is close \mathbb{P}_X . The proof is provided in Appendix D.1.

To measure the distance between probability measures, we use the 1-Wasserstein distance in this work for analytical convenience. For $\mu, \nu \in \mathcal{P}(\mathbb{R}^D)$, it is defined by

$$\mathcal{W}_1(\mu, \nu) := \inf \left\{ \int \|\mathbf{x} - \mathbf{y}\| d\gamma(\mathbf{x}, \mathbf{y}) : \gamma \in \Gamma(\mu, \nu) \right\} = \inf \{ \mathbb{E} [\|\mathbf{X} - \mathbf{Y}\|] : \mathbf{X} \sim \mu, \mathbf{Y} \sim \nu \},$$

where $\Gamma(\mu, \nu) := \{\gamma \in \mathcal{P}(\mathbb{R}^D \times \mathbb{R}^D) : \gamma(A \times \mathbb{R}^D) = \mu(A), \gamma(\mathbb{R}^D \times B) = \nu(B)\}$; see Chewi et al. (2024) for more details.

Proposition 5. *Using the above notation, if $\mathbf{Z}_\sigma \sim \mathbb{P}_\sigma^Z$, then for $\mathbf{X}^\sigma := A\mathbf{Z}_\sigma \sim \mathbb{P}_X^\sigma$, we have*

$$\mathbb{P}_X^\sigma = w_1 \mathbb{P}_{X|Y}^\sigma(\cdot | Y = 1) + w_2 \mathbb{P}_{X|Y}^\sigma(\cdot | Y = 2),$$

where $\mathbb{P}_{X|Y}^\sigma(\cdot | Y = i) := A_\# \mathbb{P}_{i,\sigma}^Z$ for $i = 1, 2$, and it follows that

$$\mathcal{W}_1(\mathbb{P}_X^\sigma, \mathbb{P}_X) \leq \sigma\sqrt{d}.$$

Therefore, the mollification technique provides a smooth latent density function p_σ^Z that induces a distribution \mathbb{P}_X^σ approximating \mathbb{P}_X .

5.1.3 Log-Concavity of Latent Density

In general, even if p_σ^Z is smooth, we cannot directly assume that it is strongly log-concave, as it is multi-modal by Equation (11). However, we can still assume that each of its components $p_{i,\sigma}^Z$ is strongly log-concave, which, in fact, follows from the assumption of strong log-concavity of the original latent density p_i^Z .

Assumption II. *Let p_i^Z be the density function of \mathbb{P}_i^Z defined on \mathbb{R}^{d_i} . There exists a large $m > 1$ such that*

$$-\nabla_z^2 \log p_i^Z(\mathbf{z}) \succeq m\mathbf{I}_{d_i},$$

i.e., p_i^Z is m -strongly log-concave for $i = 1, 2$.

Assumption II ensures the strong log-concavity of each component $p_{i,\sigma}^Z$, but it does not guarantee that the overall mixture p_σ^Z is strongly log-concave—this is a common difficulty in the case of multi-modal distributions. However, due to the mollification construction, the parameter σ can be freely chosen, which enables us to establish the strong log-concavity of p_σ^Z under the following assumption.

Assumption III. *For a chosen σ , we assume that*

$$M := \sup_{\mathbf{z}} \|\nabla_{\mathbf{z}} \log p_{1,\sigma}^Z(\mathbf{z}) - \nabla_{\mathbf{z}} \log p_{2,\sigma}^Z(\mathbf{z})\| < 2\sqrt{m-1}.$$

Remark 4. This assumption is novel but essential for addressing log-concavity in multi-modal settings. Characterizing the classes of p_i^Z that satisfy it is nontrivial. As a concrete example, if each $p_i^Z(\mathbf{z})$ is a Gaussian truncated to a compact, convex set, then compactness implies that the difference of $\nabla \log p_{i,\sigma}^Z$ is uniformly bounded by a quantity depending on σ . Therefore, with an appropriate choice of σ , Assumption III holds. The details, with a sufficient condition for Assumption III, are discussed in Appendix E.2.

Assumption III is required to obtain an upper bound on $\nabla^2 \log p$, even when the density p is multi-modal, as shown in the following lemma; see the proof in Appendix D.1.

Lemma 6. *Let p_1, p_2 be two probability density functions on \mathbb{R}^n such that $\nabla^2 \log p_i \preceq L_i \mathbf{I}_n$ for some constant $L_i \in \mathbb{R}$. Suppose that*

$$\sup_{\mathbf{x}} \|\nabla \log p_1(\mathbf{x}) - \nabla \log p_2(\mathbf{x})\| \leq M < \infty.$$

Then, for the mixture density $p = wp_1 + (1-w)p_2$ with $w \in (0, 1)$, it holds that

$$\nabla^2 \log p \preceq \left(\max\{L_1, L_2\} + \frac{1}{4}M^2 \right) \mathbf{I}_n.$$

By Lemma 6 and Proposition 3, the strong log-concavity of the multi-modal latent density function p_σ^Z can be guaranteed.

Theorem 7. *Under Assumptions II and III, if $\sigma^2 < (4m - M^2)/(M^2m)$, then p_σ^Z is strongly log-concave for $p_\sigma^Z = w_1 p_{1,\sigma}^Z + w_2 p_{2,\sigma}^Z$, i.e.,*

$$-\nabla_{\mathbf{z}}^2 \log p_\sigma^Z(\mathbf{z}) \succeq m_0^z \mathbf{I}_d, \quad m_0^z := \frac{4m - M^2(1 + m\sigma^2)}{4(1 + m\sigma^2)}.$$

Proof. Note that

$$\mathbf{Z}_{1,\sigma} = (\mathbf{I}_{d_1}, \mathbf{O})^\top \mathbf{Z}_1 + \sigma \boldsymbol{\zeta}_1 \sim p_{1,\sigma}^Z.$$

By Assumption II and Corollary 4, with the choices $B = (\mathbf{I}_{d_1}, \mathbf{O})^\top$, $m_0 = m$, $\alpha = 1$, and $\beta = \sigma$, we obtain

$$\nabla_{\mathbf{z}}^2 \log p_{1,\sigma}^Z(\mathbf{z}) \preceq \left(\frac{1}{\sigma^2(1+m\sigma^2)} - \frac{1}{\sigma^2} \right) \mathbf{I}_d.$$

For $p_{2,\sigma}^Z$, we similarly have

$$\nabla_{\mathbf{z}}^2 \log p_{2,\sigma}^Z(\mathbf{z}) \preceq \left(\frac{1}{\sigma^2(1+m\sigma^2)} - \frac{1}{\sigma^2} \right) \mathbf{I}_d.$$

Then, because $p_\sigma^Z = w_1 p_{1,\sigma}^Z + w_2 p_{2,\sigma}^Z$, it follows from Assumption III and Lemma 6 that

$$\nabla_{\mathbf{z}}^2 \log p_\sigma^Z(\mathbf{z}) \preceq \left(\frac{1}{\sigma^2(1+m\sigma^2)} - \frac{1}{\sigma^2} + \frac{1}{4} M^2 \right) \mathbf{I}_d = -m_0^z \mathbf{I}_d. \quad \square$$

5.1.4 Smoothness and Concavity

Before proceeding, let us recall that the latent distribution \mathbb{P}^Z of \mathbb{P}_X is not “good”. To address this, we construct a new distribution \mathbb{P}_X^σ whose latent distribution \mathbb{P}_σ^Z admits a “good” density function p_σ^Z , and which is close to \mathbb{P}_X . Consequently, instead of considering the score function associated with a DDPM initialized from \mathbb{P}_X , we consider a DDPM initialized from \mathbb{P}_X^σ , i.e.,

$$\mathbf{X}_t^\sigma = \sqrt{\alpha_t} A \mathbf{Z}_\sigma + \sqrt{1 - \alpha_t} \boldsymbol{\xi} \sim p_t^\sigma, \quad (13)$$

where $\mathbf{Z}_\sigma \sim p_\sigma^Z$ and $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$. We then modify the dynamics in (10) to define our final version of the geometric guidance model:

Definition 1 (Geometric Guidance Model). For any $t \in [0, T - \delta]$,

$$\frac{d}{dt} \tilde{\mathbf{X}}_t = \tilde{\mathbf{X}}_t + \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t) - \eta P_1 \tilde{\mathbf{X}}_t, \quad \tilde{\mathbf{X}}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D), \quad (*)$$

where $P_1 = \mathbf{I}_D - A_1 A_1^\top$.

Remark 5. (i) The initial condition is taken as $\mathcal{N}(\mathbf{0}, \mathbf{I}_D)$ instead of $p_T(\cdot | y)$ to reflect practical implementation settings. (ii) The time interval is chosen as $[0, T - \delta]$ for some $\delta > 0$ to avoid the singularity at time T .

Therefore, our main objective now becomes establishing the Lipschitz continuity of $\nabla_{\mathbf{x}} \log p_t^\sigma(\mathbf{x})$ and the log-concavity of $p_t^\sigma(\mathbf{x})$, which follows from the strong log-concavity of the latent density $p_\sigma^Z(\mathbf{z})$.

Theorem 8. Under Assumption I and the same settings as in Theorem 7, for the density function p_t^σ defined in Equation (13), we have

$$\|\nabla_{\mathbf{x}}^2 \log p_t^\sigma(\mathbf{x})\|_{\text{op}} \leq L_t, \quad L_t := \frac{2\alpha_t + (1 - \alpha_t)m_0^z}{(1 - \alpha_t)(\alpha_t + (1 - \alpha_t)m_0^z)},$$

and

$$-\nabla_{\mathbf{x}}^2 \log p_t^\sigma(\mathbf{x}) \succeq m_t \mathbf{I}_D, \quad m_t := \frac{m_0^z}{\alpha_t + (1 - \alpha_t)m_0^z}.$$

Proof. First, by Theorem 7, the latent density p_σ^Z is m_0^z -strongly log-concave. By the definition of p_t^σ , Proposition 3 implies that

$$\|\nabla_{\mathbf{x}}^2 \log p_t^\sigma(\mathbf{x})\|_{\text{op}} \leq \frac{2\alpha_t + (1 - \alpha_t)m_0^z}{(1 - \alpha_t)(\alpha_t + (1 - \alpha_t)m_0^z)},$$

with the choices $B = A$, $m_0 = m_0^z$, $\alpha = \sqrt{\alpha_t}$, and $\beta = \sqrt{1 - \alpha_t}$. This follows from the fact that $A^\top A = \mathbf{I}_d$ (Assumption I), which indicates $\|A\|_{\text{op}}^2 = 1$ and $\lambda_{\min}(A^\top A) = 1$.

For the log-concavity, Corollary 4 directly yields

$$-\nabla_{\mathbf{x}}^2 \log p_t^\sigma(\mathbf{x}) \succeq \left(\frac{1}{1 - \alpha_t} \left(1 - \frac{\alpha_t}{\alpha_t + (1 - \alpha_t)m_0^z} \right) \right) \mathbf{I}_D. \quad \square$$

Therefore, we have established the desired properties of p_t^σ , which ensure the well-posedness of the geometric guidance model (*). Moreover, from the definition of m_t in Theorem 8, we can derive a lower bound that will be useful in the following analysis; see Appendix D.1 for the proof.

Corollary 9. *There exists a small $\sigma > 0$ such that $m_0^z > 1$ and $m_I := \inf_{t \in (0, T]} m_t > 1$.*

5.2 Estimating Target Data Manifold

For the geometric guidance model (*), the first problem is whether it can estimate the target data manifold \mathcal{M}_1 . Specifically, we aim to show that the generated sample $\tilde{\mathbf{X}}_{T-\delta}$ approximately lies in \mathcal{M}_1 . Since $\mathcal{M}_1 = \text{Im } A_1$ is a linear subspace by Assumption I, it suffices to examine whether $\mathbb{E} [\|\tilde{\mathbf{Y}}_{T-\delta}\|] \approx 0$, where

$$\tilde{\mathbf{Y}}_t = P_1 \tilde{\mathbf{X}}_t, \quad P_1 = \mathbf{I}_D - A_1 A_1^\top.$$

Multiplying both sides of Equation (*) by P_1 , we obtain that $\tilde{\mathbf{Y}}_t$ satisfies the following dynamics:

$$\frac{d}{dt} \tilde{\mathbf{Y}}_t = \tilde{\mathbf{Y}}_t + P_1 \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t) - \eta \tilde{\mathbf{Y}}_t, \quad \tilde{\mathbf{Y}}_0 \sim \mathcal{N}(0, P_1), \quad (14)$$

for $t \in [0, T - \delta]$. By analyzing the dynamics (14), the following theorem provides a convergence rate of $\mathbb{E} [\|\tilde{\mathbf{Y}}_{T-\delta}\|] \rightarrow 0$ with respect to the guidance scale η .

Theorem 10. *Consider the dynamics (14) under Assumptions II and III. Then,*

$$\mathbb{E} [\|\tilde{\mathbf{Y}}_{T-\delta}\|] \leq \mathcal{O} \left(e^{-\eta} + \frac{1}{\eta} \right).$$

In particular, for any $\varepsilon > 0$, by choosing $\eta = \Theta(\max\{\log(1/\varepsilon), 1/\varepsilon\})$, $\mathbb{E} [\|\tilde{\mathbf{Y}}_{T-\delta}\|] < \varepsilon$.

Proof sketch. We provide a sketch of the proof here; the full proof is given in Appendix D.2.

The key idea is to derive a differential inequality for $\mathbb{E} [\|\tilde{\mathbf{Y}}_t\|]$. First, we have

$$\frac{d}{dt} \mathbb{E} [\|\tilde{\mathbf{Y}}_t\|] \leq (1 - \eta) \mathbb{E} [\|\tilde{\mathbf{Y}}_t\|] + \mathbb{E} [\|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t)\|]. \quad (15)$$

To bound $\mathbb{E} [\|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t)\|]$, the L_t -smoothness of $\log p_t^\sigma$ is required, which follows from Theorem 8. The smoothness implies that

$$\|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t)\| \leq L_S \|\tilde{\mathbf{X}}_t\| + C$$

for some constants L_S and C . Therefore, it suffices to bound $\mathbb{E} [\|\tilde{\mathbf{X}}_t\|]$. By deriving a differential inequality from Equation (*) and applying Grönwall's inequality (Lemma H.11), we obtain $\mathbb{E} [\|\tilde{\mathbf{X}}_t\|] \leq M_1$ for some constant M_1 , and thus $\mathbb{E} [\|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t)\|] \leq M_2$. Substituting this bound into (15) and applying Grönwall's inequality once more yields the desired result. \square

Remark 6. For this theorem, we provide two remarks.

- (i) Note that this result depends only on the L_t -smoothness of $\log p_t^\sigma$, and not on strong log-concavity. Therefore, Assumptions II and III can be relaxed; see further discussion in Appendix F.
- (ii) The universal guidance model,

$$\frac{d\mathbf{X}_t^\leftarrow}{dt} = \mathbf{X}_t^\leftarrow + \nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^\leftarrow) - \eta \nabla f(\mathbf{X}_t^\leftarrow),$$

was proposed by Bansal et al. (2023) to control the generation process such that the generated images match the prompt $g(\mathbf{X}_T^\leftarrow) \approx \mathbf{c}$. In their setting, $f(\mathbf{x}) = \ell(\mathbf{c}, g(\mathbf{x}))$ for some loss function ℓ . A similar idea used in the proof of Theorem 10 can be extended to theoretically analyze the universal guidance model. If the L -smoothness of $\log p_t$ holds (see Appendix F) and f is strongly convex,

$$\mathbb{E} [f(\mathbf{X}_T^\leftarrow)] \rightarrow \min f, \quad \text{as } \eta \rightarrow \infty;$$

see Appendix D.3 for more details.

Theorem 10 shows that the geometric guidance model can approximate the target data manifold. Specifically, as the guidance scale η increases, the generated data increasingly lie close to the target manifold. This result is consistent with empirical observations on both synthetic datasets (Wu et al., 2024; Chidambaram et al., 2024) and real-world datasets (Dhariwal & Nichol, 2021; Sadat et al., 2024; 2025), as well as with the theoretical results in the one-dimensional case studied by Chidambaram et al. (2024), which demonstrate that increasing η causes the generated data to move toward the extreme points in the support of the target conditional distribution.

5.3 Distance to Target Distribution

Let \tilde{p}_t be the density function of \tilde{X}_t in the geometric guidance model (*). The second question is how to measure the 1-Wasserstein distance between the generated density $\tilde{p}_{T-\delta}$ and the target conditional distribution $\mathbb{P}_{X|Y}(\cdot | Y = 1)$. Specifically, the goal is to provide an upper bound on $\mathcal{W}_1(\tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1))$.

First, we require an additional assumption: the boundedness of the first moment of each conditional distribution $\mathbb{P}_{X|Y}(\cdot | Y = i)$, which can be reduced to the same condition on the latent distribution p_i^Z .

Assumption IV. For $i = 1, 2$ and $\mathbf{Z}_i \sim p_i^Z$, $\mathbf{m}_i^Z := \mathbb{E}[\|\mathbf{Z}_i\|] < \infty$.

Theorem 11. Under Assumptions I, II, III, and IV, we obtain that

$$\mathcal{W}_1(\tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1)) \leq \mathcal{O}(e^{-T} + \delta^{1/2} + \sigma + \eta^{-1}) + \tilde{C}$$

for some constant \tilde{C} .

Proof sketch. The proof consists of two main steps:

- (i) Let $Q_1 = A_1 A_1^\top$ be the orthogonal projection onto $\mathcal{M}_1 = \text{Im } A_1$. By Theorem 10, we have

$$\mathcal{W}_1(\tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1)) \leq \mathcal{W}_1((Q_1)_\# \tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1)) + \mathcal{O}(e^{-T} + \eta^{-1}).$$

- (ii) For $\mathcal{W}_1((Q_1)_\# \tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1))$, it has

$$\mathcal{W}_1((Q_1)_\# \tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1)) \leq \mathcal{W}_1(\tilde{p}_{T-\delta}, \mathbb{P}_X) + \mathcal{W}_1((Q_1)_\# \mathbb{P}_X, \mathbb{P}_{X|Y}(\cdot | Y = 1)),$$

where the first term $\mathcal{W}_1(\tilde{p}_{T-\delta}, \mathbb{P}_X)$ can be bounded by comparing the geometric guidance model (*) with the unconditional reverse dynamics, and the second term is directly bounded by Lemma D.3.

The full proof is provided in Appendix D.4. \square

Remark 7. Since geometric guidance cannot carry as much information as probabilistic guidance due to its analytical simplicity, the error floor \tilde{C} does not vanish as $\eta = 1$, unlike in probabilistic guidance models; see further discussion in Remark D.1.

This result suggests that increasing the guidance scale does not harm the generating performance, which may appear counterintuitive and inconsistent with empirical observations. In practice, however, ODE dynamics are typically approximated using the Euler discretization (or the Euler–Maruyama scheme for SDEs), which introduces additional discretization error. In our setting, the Euler discretization error for the geometric guidance model (*) is bounded by $\mathcal{O}(h\eta^2)$, where h denotes the step size; see Appendix D.5 for details. Therefore, the performance degradation observed at large guidance scales arises not from the model formulation itself, but from the discretization algorithm. For example, Wu et al. (2024, Figure 3) showed that the large guidance scale would harm the modality of the original data, but this problem can be mitigated by reducing the discretization step size.

6 Nonlinear Extension

In this section, our main objective is to construct a nonlinear geometric guidance model suitable for real-world image datasets, and to evaluate its generation performance under varying guidance scales η .

The first challenge is to construct the geometric guidance term for image datasets, which may not lie in a linear subspace. To this end, we study the geometric structure of noisy data manifolds without assuming linearity of the target data manifold, by extending the result of Proposition 1 to the nonlinear case (see Section 6.1). Then, following the idea of Ross et al. (2024), we train functions $F_\theta^t: \mathbb{R}^D \rightarrow \mathbb{R}$ to model noisy data manifolds via $\mathcal{M}^t = (F_\theta^t)^{-1}(0)$ so that $\nabla_{\mathbf{x}} F_\theta^t$ can replace $(\mathbf{I}_D - AA^\top)\mathbf{x}$ to be the nonlinear geometric guidance term (see Section 6.2). Finally, we examine this nonlinear geometric guidance model on CIFAR-10 (Krizhevsky, 2009), and evaluate its performance under the different guidance scale (see Section 6.3).

6.1 Noisy Data Manifolds for Nonlinear Case

The geometric guidance term $(\mathbf{I}_D - AA^\top)\mathbf{x}$ is constructed based on the result in Proposition 1, which assumes that the target data manifold $\mathcal{M} = \text{Im } A$ is linear. However, for real-world image datasets, it may unrealistic to assume that the data lie in a linear subspace. Instead, it is more reasonable to assume that the target image data lie on a nonlinear manifold $\mathcal{M} \subset \mathbb{R}^D$ with intrinsic dimension $d \ll D$; see Appendix I for basic knowledge of manifolds. This assumption is known as the manifold hypothesis (Bengio et al., 2013), and it has been supported by both theoretical analyses (Fefferman et al., 2016) and empirical studies (Brown et al., 2022; Loaiza-Ganem et al., 2022).

To construct a new geometric guidance term, because of the nonlinearity of \mathcal{M} , we must extend the result of Proposition 1 to uncover the geometric structure of noisy data manifolds. Although the d -dimensional manifold $\mathcal{M} \subset \mathbb{R}^D$ is not assumed to be linear, we additionally require that it is locally isometric to \mathbb{R}^d . More precisely, we assume the existence of a C^∞ function $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$ such that $\text{Im } \phi = \mathcal{M}$ and ϕ is an isometry; that is, $J\phi^\top J\phi \equiv \mathbf{I}_d$. Then, by Lemma 12, we obtain an analogue of Proposition 1 in Theorem 13, which shows that the noisy data manifolds \mathcal{M}^t are hypersurfaces—i.e., $(D-1)$ -dimensional submanifolds of \mathbb{R}^D ; see the proofs in Appendix G.1.

Lemma 12. *Let $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$ be a C^∞ isometry such that $\mathcal{M} = \text{Im } \phi \subset \mathbb{R}^D$ is a d -dimensional submanifold. Then, there exists a C^∞ function $\phi^*: \mathbb{R}^D \rightarrow \mathbb{R}^d$ such that $\phi^* \circ \phi = \text{id}_{\mathbb{R}^d}$ and*

$$J\phi^*(\phi(\mathbf{z})) = J\phi(\mathbf{z})^\top, \quad \forall \mathbf{z} \in \mathbb{R}^d.$$

Remark 8. In fact, the isometry of ϕ implies that $\text{Im } \phi$ is a submanifold, because it is proper (i.e., the preimage of every compact set is compact) by the Hopf–Rinow theorem (Jost, 2008).

Theorem 13. *Let $\mathcal{M} \subset \mathbb{R}^D$ be a d -dimensional submanifold as defined in Lemma 12, and let \mathbb{P}_X on \mathbb{R}^D such that $\text{supp } \mathbb{P}_X \subset \mathcal{M}$. Let \mathbf{X}_t be generated by DDPM (1) initialized from \mathbb{P}_X . If $d \ll D$, then \mathbf{X}_t concentrates on a hypersurface $\mathcal{M}^t \subset \mathbb{R}^D$ with high probability, where*

$$\mathcal{M}^t := \{\mathbf{x}: f^t(\mathbf{x}) = r(t)\}, \quad r(t) = \sqrt{(D-d)(1-\alpha_t)},$$

for some C^∞ function $f^t: \mathbb{R}^D \rightarrow \mathbb{R}$.

6.2 Learning Geometric Guidance

For an image dataset $(\mathbf{X}, Y) \sim \mathbb{P}_{XY}$ with class label $Y \in \{1, 2, \dots, K\}$, we adopt the union of manifold hypothesis (Brown et al., 2022), that is,

$$\text{supp } \mathbb{P}_{X|Y}(\cdot | Y = y) \subset \mathcal{M}_y,$$

where $\mathcal{M}_y \subset \mathbb{R}^D$ is a d_y -dimensional submanifold. To apply Theorem 13, we further assume that, for each \mathcal{M}_y , there exists an isometry $\phi_y: \mathbb{R}^{d_y} \rightarrow \mathbb{R}^D$ such that $\text{Im } \phi_y = \mathcal{M}_y$. Then, the noisy data manifolds generated by the forward process initialized from \mathcal{M}_y are given by

$$\mathcal{M}_y^t := \{\mathbf{x} \in \mathbb{R}^D: f_y^t(\mathbf{x}) = r(t)\}, \quad r(t) = \sqrt{(D-d)(1-\alpha_t)},$$

for some function $f_y^t: \mathbb{R}^D \rightarrow \mathbb{R}$.

By adopting the same idea as in Section 4.2, for $\mathbf{x} \in \mathcal{M}_y^t$, the guidance term $\nabla_{\mathbf{x}} \log p_t(y | \mathbf{x})$ is approximately normal to \mathcal{M}_y^t at \mathbf{x} —that is, it is approximately parallel to $\nabla_{\mathbf{x}} f_y^t(\mathbf{x})$. Therefore, we construct the nonlinear

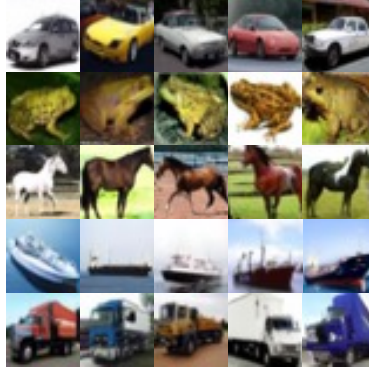


Figure 1: Images generated by GeGM on CIFAR-10

geometric guidance term as $\nabla_{\mathbf{x}} f_y^t(\mathbf{x})$ to replace the probabilistic guidance $\nabla_{\mathbf{x}} \log p_t(y | \mathbf{x})$ in the reverse process for conditional generation. The resulting nonlinear geometric guidance model (in deterministic form) is defined by

$$\frac{d}{dt} \mathbf{X}_t^{\leftarrow} = \mathbf{X}_t^{\leftarrow} + \nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^{\leftarrow}) - \eta \nabla_{\mathbf{x}} f_y^{T-t}(\mathbf{X}_t^{\leftarrow}), \quad (16)$$

where $\nabla_{\mathbf{x}} \log p_t$ is the score function of the unconditional DDPM initialized from \mathbb{P}_X .

To implement the nonlinear geometric guidance model, one must estimate both the score function and the nonlinear geometric guidance term. The score function $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ can be estimated using an unconditional diffusion model—specifically, by training a network $\mathbf{s}_{\theta}(t, \mathbf{x})$ via the score matching method (Vincent, 2011) on the unconditional data \mathbf{X} . The main task, then, is to estimate $\nabla_{\mathbf{x}} f_y^t(\mathbf{X}_t^{\leftarrow})$.

First, Theorem 10 shows that $\mathcal{M}_y^t = (f_y^t)^{-1}(r(t))$, so such function f_y^t is called a manifold-defining function in Ross et al. (2024). Following a similar idea, we train a network $F_{y,\theta}^t: \mathbb{R}^D \rightarrow \mathbb{R}$ to estimate $f_y^t - r(t)$, so $F_{y,\theta}^t$ needs to satisfy

$$F_{y,\theta}^t(\mathbf{x}) = 0, \text{ and } \nabla_{\mathbf{x}} F_{y,\theta}^t(\mathbf{x}) \neq \mathbf{0}, \quad \forall \mathbf{x} \in \mathcal{M}_y^t,$$

where the first condition follows directly from the definition of \mathcal{M}_y^t , and the second condition, called the rank condition, ensures $F_{y,\theta}^t$ a manifold-defining function, as guaranteed by the Constant Rank Theorem (Lemma I.2). Therefore, the loss function for training $F_{y,\theta}^t$ is designed as

$$\mathcal{L}_y^t(\theta) := \mathbb{E}_{\mathbf{X} \sim p_t(\cdot | y)} \left[|F_{y,\theta}^t(\mathbf{X})|^2 - \kappa \|\nabla_{\mathbf{x}} F_{y,\theta}^t(\mathbf{X})\|^2 \right], \quad (17)$$

where $\kappa > 0$ is chosen for controlling the strength of the rank condition. We simply set $\kappa = 1$.

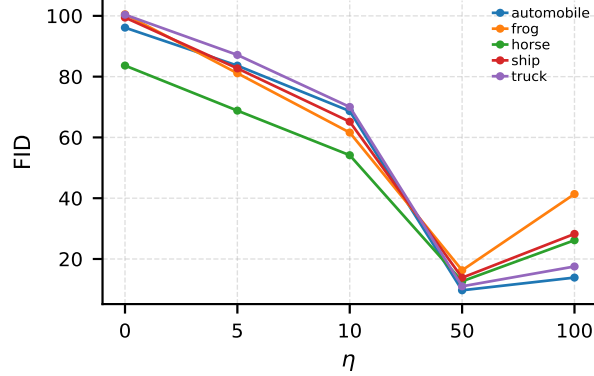
6.3 Experiments

Effectiveness of GeGM. We use the Fréchet Inception Distance (FID) (Heusel et al., 2017) as the metric for evaluating generation performance, because it can be regarded as a practical surrogate for the Wasserstein distance. We compare the FID of samples generated by the nonlinear geometric guidance model (16) (GeGM) with those generated by the classifier guidance model (CGM) (6). The results are reported in Table 1, where we present results for selected classes; the remaining classes are provided in Appendix G.2. Note that the guidance scales used for CGM and GeGM differ, since the norms of the probabilistic and geometric guidance terms are not comparable. For visualization, Figure 1 displays images generated by the nonlinear GeGM. These results demonstrate the effectiveness of the nonlinear GeGM in generating real-world images.

Performance vs. guidance scale. By applying the nonlinear GeGM (16), we evaluate how generation performance varies with the guidance scale η on selected classes from CIFAR-10; results for the remaining classes are provided in Appendix G.2. As shown in Figure 2, performance improves with increasing η within

Table 1: Comparison of FID on CIFAR-10

	Automobile	Frog	Horse	Ship	Truck
CGM ($\eta = 1$)	13.46	17.87	13.97	11.61	16.85
GeGM ($\eta = 50$)	9.70	16.28	12.65	13.84	11.02

Figure 2: FID v.s. guidance scale η of GeGM on selected classes of CIFAR-10

a reasonable range. Since FID serves as a practical approximation of the Wasserstein distance, this trend is consistent with Theorem 11, even in the nonlinear setting.

Remark 9. We emphasize that the observed trends are consistent with the spirit of Theorem 11 in nonlinear regimes, but they are not derived from it. Establishing nonlinear analogues of Theorems 10–11 will require additional analysis and is left for future work.

7 Conclusion

In this work, we studied the role of the guidance scale in conditional generation with diffusion models. To address the analytical intractability of the probabilistic guidance term, we introduced a geometric guidance model that enables theoretical analysis under the linear manifold hypothesis. To facilitate this analysis, we proposed a mollification technique to ensure the regularity of the score function in the presence of multimodality. Our results showed that increasing the guidance scale within a reasonable range can enhance generation performance, in line with empirical observations reported in prior studies. We further extended the model to nonlinear settings, and experiments on real-world datasets demonstrated the effectiveness of the geometric guidance model and provided additional evidence consistent with our theoretical findings.

Limitations: While the geometric guidance offers a more tractable alternative to probabilistic guidance, it comes with certain limitations. Notably, our analysis showed that the upper bound of the Wasserstein distance between the generated and target conditional distributions is bounded by a constant, regardless of the choice of the guidance scale. This implies that, unlike probabilistic guidance, which can approximate the target conditional distribution by setting the scale to 1, the geometric guidance does not guarantee convergence to the target distribution. This is a trade-off made for the sake of analytical tractability.

Although our experiments on the nonlinear extension partially supported the theoretical results, our current theoretical analysis is restricted to the linear manifold setting. In the nonlinear case, the geometric structure of the score function remains unclear. Regarding regularity of the score function, while Lipschitz continuity can be ensured under compactness assumptions, extending this to the non-compact setting remains an open problem. Furthermore, the log-concavity of the score function cannot be guaranteed, even in compact nonlinear cases.

Acknowledgments

We thank Ming Li and Luheng Wang for the helpful discussions. ZZ was supported by Institute for AI and Beyond at the University of Tokyo. MS was supported by JST ASPIRE Grant Number JPMJAP25B1. The authors also thank the anonymous reviewers for their careful reviews and insightful comments, which have been invaluable in improving both the clarity and rigor of this work.

References

- Brian D. O. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.
- Dominique Bakry, Ivan Gentil, and Michel Ledoux. *Analysis and Geometry of Markov Diffusion Operators*, volume 348. Springer Science & Business Media, 2013.
- Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 843–852, 2023.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, 2013.
- Sergey G. Bobkov and Michel Ledoux. From brunn–minkowski to brascamp–lieb and to logarithmic sobolev inequalities. *Geometric and Functional Analysis*, 10:1028–1052, 2000.
- Nicolas Boumal. *An introduction to optimization on smooth manifolds*. Cambridge University Press, 2023.
- Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021.
- Bradley C. A. Brown, Anthony L. Caterini, Brendan Leigh Ross, Jesse C. Cresswell, and Gabriel Loaiza-Ganem. Verifying the union of manifolds hypothesis for image data. In *International Conference on Learning Representations*, 2022.
- Hongrui Chen, Holden Lee, and Jianfeng Lu. Improved analysis of score-based generative modeling: User-friendly bounds under minimal smoothness assumptions. In *International Conference on Machine Learning*, pp. 4735–4763. PMLR, 2023a.
- Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estimation, and distribution recovery of diffusion models on low-dimensional data. In *International Conference on Machine Learning*, pp. 4672–4712. PMLR, 2023b.
- Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru Zhang. Sampling is as easy as learning the score: Theory for diffusion models with minimal data assumptions. In *International Conference on Learning Representations*, 2023c.
- Sinho Chewi, Jonathan Niles-Weed, and Philippe Rigollet. Statistical optimal transport. *arXiv preprint arXiv:2407.18163*, 2024.
- Muthu Chidambaram, Khashayar Gatmiry, Sitan Chen, Holden Lee, and Jianfeng Lu. What does guidance do? a fine-grained analysis in a simple setting. In *Advances in Neural Information Processing Systems*, 2024.
- Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Improving diffusion models for inverse problems using manifold constraints. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022.
- Hyungjin Chung, Suhyeon Lee, and Jong Chul Ye. Decomposed diffusion sampler for accelerating large-scale inverse problems. In *International Conference on Learning Representations*, 2024.

- Valentin De Bortoli. Convergence of denoising diffusion models under the manifold hypothesis. *Transactions on Machine Learning Research*, 2022.
- Valentin De Bortoli, James Thornton, Jeremy Heng, and Arnaud Doucet. Diffusion schrödinger bridge with applications to score-based generative modeling. *Advances in Neural Information Processing Systems*, 34: 17695–17709, 2021.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat GANs on image synthesis. In *Advances in Neural Information Processing Systems*, pp. 8780–8794, 2021.
- Anh-Dung Dinh, Daochang Liu, and Chang Xu. Rethinking conditional diffusion sampling with progressive guidance. *Advances in Neural Information Processing Systems*, 36:42285–42297, 2023.
- Lawrence C. Evans. *Measure Theory and Fine Properties of Functions*. Routledge, 2018.
- Charles Fefferman, Sanjoy Mitter, and Hariharan Narayanan. Testing the manifold hypothesis. *Journal of the American Mathematical Society*, 29(4):983–1049, 2016.
- Gerald B. Folland. *A Course in Abstract Harmonic Analysis*. CRC Press, 2016.
- Xuefeng Gao, Hoang M. Nguyen, and Lingjiong Zhu. Wasserstein convergence guarantees for a general class of score-based generative models. *Journal of Machine Learning Research*, 26(43):1–54, 2025.
- David Francis Griffiths and Desmond J Higham. *Numerical methods for ordinary differential equations: initial value problems*, volume 5. Springer, 2010.
- Ulrich G. Haussmann and Etienne Pardoux. Time reversal of diffusions. *The Annals of Probability*, pp. 1188–1205, 1986.
- Yutong He, Naoki Murata, Chieh-Hsin Lai, Yuhta Takida, Toshimitsu Uesaka, Dongjun Kim, Wei-Hsiang Liao, Yuki Mitsufuji, J. Zico Kolter, Ruslan Salakhutdinov, and Stefano Ermon. Manifold preserving guided diffusion. In *International Conference on Learning Representations*, 2024.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in Neural Information Processing Systems*, 30, 2017.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J. Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022.
- Jürgen Jost. *Riemannian Geometry and Geometric Analysis*, volume 42005. Springer, 2008.
- Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. DiffWave: A versatile diffusion model for audio synthesis. In *International Conference on Learning Representations*, 2021.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>, 2009. Technical Report, University of Toronto.
- Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pp. 1302–1338, 2000.
- Jean-François Le Gall. *Brownian Motion, Martingales, and Stochastic Calculus*, volume 274. Springer, 2016.
- Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence for score-based generative modeling with polynomial complexity. *Advances in Neural Information Processing Systems*, 35:22870–22882, 2022.

- Holden Lee, Jianfeng Lu, and Yixin Tan. Convergence of score-based generative modeling for general data distributions. In *International Conference on Algorithmic Learning Theory*, pp. 946–985. PMLR, 2023.
- John M. Lee. *Introduction to Smooth Manifolds*, volume 218. Springer Science & Business Media, 2012.
- John M. Lee. *Introduction to Riemannian Manifolds*. Graduate Texts in Mathematics. Springer International Publishing, 2019.
- Gen Li and Yuling Yan. Adapting to unknown low-dimensional structures in score-based diffusion models. *Advances in Neural Information Processing Systems*, 37:126297–126331, 2024.
- Gabriel Loaiza-Ganem, Brendan Leigh Ross, Jesse C. Cresswell, and Anthony L. Caterini. Diagnosing and fixing manifold overfitting in deep generative models. *Transactions on Machine Learning Research*, 2022.
- Gabriel Loaiza-Ganem, Brendan Leigh Ross, Rasa Hosseinzadeh, Anthony L. Caterini, and Jesse C. Cresswell. Deep generative models through the lens of the manifold hypothesis: A survey and new connections. *Transactions on Machine Learning Research*, 2024.
- James Raymond Munkres. *Topology*. Pearson, New York, NY, second edition, reissue edition, 2018.
- Kazusato Oko, Shunta Akiyama, and Taiji Suzuki. Diffusion models are minimax optimal distribution estimators. In *International Conference on Machine Learning*, pp. 26517–26582. PMLR, 2023.
- R. Tyrrell Rockafellar. *Convex Analysis*, volume 28. Princeton University Press, 1997.
- 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*, pp. 10684–10695, 2022.
- Brendan Leigh Ross, Gabriel Loaiza-Ganem, Anthony L. Caterini, and Jesse C. Cresswell. Neural implicit manifold learning for topology-aware density estimation. *Transactions on Machine Learning Research*, 2024.
- Seyedmorteza Sadat, Jakob Buhmann, Derek Bradley, Otmar Hilliges, and Romann M. Weber. CADs: Unleashing the diversity of diffusion models through condition-annealed sampling. In *International Conference on Learning Representations*, 2024.
- Seyedmorteza Sadat, Manuel Kansy, Otmar Hilliges, and Romann M. Weber. No training, no problem: Rethinking classifier-free guidance for diffusion models. In *International Conference on Learning Representations*, 2025.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021a.
- Jiaming Song, Qinsheng Zhang, Hongxu Yin, Morteza Mardani, Ming-Yu Liu, Jan Kautz, Yongxin Chen, and Arash Vahdat. Loss-guided diffusion models for plug-and-play controllable generation. In *International Conference on Machine Learning*, pp. 32483–32498. PMLR, 2023.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021b.

Wenpin Tang and Hanyang Zhao. Score-based diffusion models via stochastic differential equations: A technical tutorial. *arXiv preprint arXiv:2402.07487*, 2024.

Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural Computation*, 23(7):1661–1674, 2011.

Zhengchao Wan, Qingsong Wang, Gal Mishne, and Yusu Wang. Elucidating flow matching ODE dynamics via data geometry and denoisers. In *International Conference on Machine Learning*, 2025.

Yuchen Wu, Minshuo Chen, Zihao Li, Mengdi Wang, and Yuting Wei. Theoretical insights for diffusion guidance: A case study for gaussian mixture models. In *International Conference on Machine Learning*, pp. 53291–53327. PMLR, 2024.

A Notation

The symbols used throughout this paper are clarified below.

1. **Letters:** Unless otherwise specified, lowercase letters such as x and \mathbf{x} denote deterministic variables, while uppercase letters such as X and \mathbf{X} denote random variables. Scalars are typically represented by non-bold symbols such as x and Y , whereas vectors are denoted using bold symbols such as \mathbf{x} and \mathbf{X} . In particular, we use $\mathbf{I}_n \in \mathbb{R}^{n \times n}$ to denote the identity matrix and $\mathbf{0} \in \mathbb{R}^n$ to denote the zero vector.

2. **Linear Algebra:**

- (i) Let $\mathcal{O}^{m \times n} \subset \mathbb{R}^{m \times n}$ (with $m > n$) denote the set of matrices whose columns are orthonormal, i.e., those satisfying $A^\top A = \mathbf{I}_n$.
- (ii) For a vector $\mathbf{x} \in \mathbb{R}^n$, the notation $\|\mathbf{x}\|$ refers to the ℓ_2 -norm. For a matrix $A \in \mathbb{R}^{m \times n}$, the operator norm is defined as

$$\|A\|_{\text{op}} = \sup_{\|\mathbf{x}\|=1} \|A\mathbf{x}\| = \sqrt{\lambda_{\max}(A^\top A)},$$

where $\lambda_{\max}(\cdot)$ denotes the maximum eigenvalue.

- (iii) Let $A, B \in \mathbb{R}^{n \times n}$ be symmetric matrices, i.e., $A = A^\top$ and $B = B^\top$. We write $A \preceq B$ (or equivalently, $B \succeq A$) if $B - A$ is positive semi-definite, i.e.,

$$\mathbf{x}^\top (B - A) \mathbf{x} \geq 0, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

3. **Calculus:**

- (i) For a scalar-valued function $f: \mathbb{R}^n \rightarrow \mathbb{R}$, the gradient with respect to \mathbf{x} is denoted by $\nabla_{\mathbf{x}} f(\mathbf{x})$, and the Hessian matrix by $\nabla_{\mathbf{x}}^2 f(\mathbf{x})$.
- (ii) For a vector-valued function $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$, JF denotes the Jacobian matrix of F , and the second-order derivative $D^2 F$ is a bilinear map $D^2 F(\mathbf{x}): \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^m$ defined by

$$D^2 F(\mathbf{x})[\mathbf{v}, \mathbf{w}] = \frac{\partial^2}{\partial s \partial t} F(\mathbf{x} + s\mathbf{v} + t\mathbf{w}) = (\mathbf{v}^\top \nabla_{\mathbf{x}}^2 F^1(\mathbf{x}) \mathbf{w}, \dots, \mathbf{v}^\top \nabla_{\mathbf{x}}^2 F^m(\mathbf{x}) \mathbf{w})^\top,$$

where $F = (F^1, \dots, F^m)$. If each F^i has continuous derivative of order k , F is called C^k .

- (iii) For any set $U \subset \mathbb{R}^n$, the characteristic function $\chi_U: \mathbb{R}^n \rightarrow \mathbb{R}$ is defined by $\chi_U(\mathbf{x}) = 1$ if $\mathbf{x} \in U$, and $\chi_U(\mathbf{x}) = 0$ otherwise.
- (iv) For integrable functions $f, g: \mathbb{R}^n \rightarrow \mathbb{R}$, their convolution is denoted by

$$f * g(\mathbf{x}) = \int_{\mathbb{R}^n} f(\mathbf{y}) g(\mathbf{x} - \mathbf{y}) d\mathbf{y}.$$

- (v) For a function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, let $\text{Im } f = f(\mathbb{R}^n)$ denote the image of f . In particular, for a matrix $A \in \mathbb{R}^{m \times n}$, $\text{Im } A$ refers the image of the linear map $\mathbf{x} \mapsto A\mathbf{x}$.

4. Probability-related Symbols:

- (i) We fix the base probability space $(\Omega, \mathcal{F}, \mathbb{P})$, where Ω is the sample space, \mathcal{F} is a σ -algebra, and \mathbb{P} is a probability measure on \mathcal{F} .
- (ii) On \mathbb{R}^n , we typically work with the Borel σ -algebra $\mathcal{B}(\mathbb{R}^n)$, and let $\mathcal{P}(\mathbb{R}^n)$ denote the set of all probability measures defined on $\mathcal{B}(\mathbb{R}^n)$. Symbols such as μ and ν represent elements of $\mathcal{P}(\mathbb{R}^n)$. The integral with respect to a measure μ is denoted by $\int f(\mathbf{x})d\mu(\mathbf{x})$ or equivalently by $\int f(\mathbf{x})\mu(d\mathbf{x})$.
- (iii) For a measurable map $f: \Omega \rightarrow \mathbb{R}^n$, the push-forward measure of \mathbb{P} under f is denoted by $f_{\#}\mathbb{P}$, and is defined as

$$f_{\#}\mathbb{P}(U) = \mathbb{P}(f^{-1}(U)), \quad \forall U \in \mathcal{B}(\mathbb{R}^n).$$

- (iv) A random variable (or vector) $\mathbf{X}: \Omega \rightarrow \mathbb{R}^n$ is a measurable map. Its distribution, denoted by \mathbb{P}_X (or \mathbb{P}^X), is a probability measure on \mathbb{R}^n defined by $\mathbb{P}_X = \mathbf{X}_{\#}\mathbb{P}$. For some $\mu \in \mathcal{P}(\mathbb{R}^n)$, we say $\mathbf{X} \sim \mu$ if $\mu = \mathbb{P}_X$. Two random variables \mathbf{X} and \mathbf{Y} are said to be equal in distribution, denoted by $\mathbf{X} \stackrel{d}{=} \mathbf{Y}$, if $\mathbb{P}_X = \mathbb{P}_Y$.
- (v) For $\mathbf{X} \sim \mathbb{P}_X$, if \mathbb{P}_X is absolutely continuous with respect to the Lebesgue measure $d\mathbf{x}$, then by the Radon-Nikodym Theorem, there is a function p_X (or denoted by p^X) such that

$$\mathbb{P}_X(U) = \int_U p_X(\mathbf{x})d\mathbf{x}, \quad \forall U \in \mathcal{B}(\mathbb{R}^n),$$

and p_X is said the density function² of \mathbf{X} . For a measurable function $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$, if $\mathbf{X} \sim \mathbb{P}_X$, then $g(\mathbf{X}) \sim g_{\#}\mathbb{P}_X$. When \mathbb{P}_X admits a density p_X , the density of $g(\mathbf{X})$ is denoted by $g_{\#}p_X$. In particular, if $g(\mathbf{x}) = A\mathbf{x}$ for a matrix $A \in \mathbb{R}^{m \times n}$, $g_{\#}\mathbb{P}_X$ is also denoted by $A_{\#}\mathbb{P}_X$ for simplicity.

- (vi) For random variables $\mathbf{X}: \Omega \rightarrow \mathbb{R}^n$ and $Y: \Omega \rightarrow \mathbb{R}$, the joint distribution of $(\mathbf{X}, Y): \Omega \rightarrow \mathbb{R}^n \times \mathbb{R}$ is denoted by $\mathbb{P}_{XY} = (\mathbf{X}, Y)_{\#}\mathbb{P}$, a probability measure on $\mathbb{R}^n \times \mathbb{R}$. The conditional distribution $\mathbb{P}_{X|Y}(\cdot | Y)$ is defined as

$$\mathbb{P}_{X|Y}(U | Y) := \mathbb{P}(\mathbf{X} \in U | Y), \quad \forall U \in \mathcal{B}(\mathbb{R}^n),$$

which is a probability measure on \mathbb{R}^n .

- (vii) For a probability measure $\mu \in \mathcal{P}(\mathbb{R}^n)$, the support of μ is denoted by

$$\text{supp } \mu = \{\mathbf{x} \in \mathbb{R}^n: \mu(B_r(\mathbf{x})) > 0, \forall r > 0\}$$

where $B_r(\mathbf{x}) \subset \mathbb{R}^n$ denotes the open ball centered at \mathbf{x} with radius r . When μ admits a density function p ,

$$\text{supp } \mu = \overline{\{\mathbf{x} \in \mathbb{R}^n: p(\mathbf{x}) > 0\}}$$

B More Details in Background

B.1 Analytic Solution for DDPMs

To solve the SDE

$$d\mathbf{X}_t = -\frac{1}{2}\beta(t)\mathbf{X}_tdt + \sqrt{\beta(t)}d\mathbf{W}_t, \quad \forall t \in [0, T],$$

we multiply both sides by the integrating factor $e^{\frac{1}{2}\int_0^t \beta(s)ds}$. This gives

$$e^{\frac{1}{2}\int_0^t \beta(s)ds}d\mathbf{X}_t + \frac{1}{2}\beta(t)e^{\frac{1}{2}\int_0^t \beta(s)ds}\mathbf{X}_tdt = \sqrt{\beta(t)}e^{\frac{1}{2}\int_0^t \beta(s)ds}d\mathbf{W}_t,$$

which leads to

$$d\left(e^{\frac{1}{2}\int_0^t \beta(s)ds}\mathbf{X}_t\right) = \sqrt{\beta(t)}e^{\frac{1}{2}\int_0^t \beta(s)ds}d\mathbf{W}_t,$$

²When unambiguous, p_X is also occasionally referred to as the distribution.

by applying Itô's formula to $e^{\frac{1}{2} \int_0^t \beta(s) ds} \mathbf{X}_t$. Therefore, we obtain the solution

$$\mathbf{X}_t = \sqrt{\alpha_t} \mathbf{X}_0 + \boldsymbol{\xi}_t,$$

where $\alpha_t := \exp\left(-\int_0^t \beta(s) ds\right)$, and

$$\boldsymbol{\xi}_t := \int_0^t e^{-\frac{1}{2} \int_s^t \beta(r) dr} \sqrt{\beta(s)} d\mathbf{W}_s.$$

Since $(\mathbf{W}_t)_{t \geq 0}$ is a standard Brownian motion on \mathbb{R}^D , it follows that $\boldsymbol{\xi}_t \sim \mathcal{N}(\mathbf{0}, \sigma_{\xi_t}^2 \mathbf{I}_D)$. To compute $\sigma_{\xi_t}^2$, let $[\cdot, \cdot]$ denote the quadratic variation. Then

$$\begin{aligned} \sigma_{\xi_t}^2 &= \mathbb{E} [[\boldsymbol{\xi}, \boldsymbol{\xi}]_t] \\ &= \mathbb{E} \left[\int_0^t e^{-\int_s^t \beta(r) dr} \beta(s) d[\mathbf{W}, \mathbf{W}]_s \right] \\ &= \int_0^t e^{-\int_s^t \beta(r) dr} \beta(s) ds = 1 - \exp\left(-\int_0^t \beta(s) ds\right). \end{aligned}$$

(see Le Gall (2016) for details). As a result,

$$\boldsymbol{\xi}_t \stackrel{d}{=} \sqrt{1 - \alpha_t} \boldsymbol{\xi}, \quad \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D).$$

B.2 Density Functions in Conditional DDPMs

Proposition B.1. *Consider a joint data density function $p(\mathbf{x}, y)$ and the process governed by the SDE:*

$$d\mathbf{X}_t = -\frac{1}{2}\beta(t)\mathbf{X}_t dt + \sqrt{\beta(t)}d\mathbf{W}_t.$$

For the following two scenarios:

- (a) *Let $\mathbf{X} \sim p(\mathbf{x} \mid Y = y)$, and run the SDE for $\mathbf{X}_0 = \mathbf{X}$. Let $p_t^y(\mathbf{x}_t)$ be the distribution of \mathbf{X}_t ,*
- (b) *Let $(\mathbf{X}, Y) \sim p(\mathbf{x}, y)$, and run the SDE for $\mathbf{X}_0 = \mathbf{X}$. Let $p_t(\mathbf{x}_t, y)$ be the distribution of (\mathbf{X}_t, Y) ,*

Then, we have

$$p_t^y(\mathbf{x}_t) = p_t(\mathbf{x}_t \mid y).$$

Proof. As shown in Equation (2),

$$\mathbf{X}_t \stackrel{d}{=} \sqrt{\alpha_t} \mathbf{X}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\xi}, \quad \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D),$$

where $\alpha_t = \exp\left(-\int_0^t \beta(s) ds\right)$. Therefore, in the first case, we have

$$p_t^y(\mathbf{x}_t) = (\sqrt{\alpha_t})_{\#} p(\mathbf{x} \mid y) * \mathcal{N}(\mathbf{0}, (1 - \alpha_t) \mathbf{I}_D).$$

Moreover, by Lemma B.2, since $\boldsymbol{\xi}$ is independent of (\mathbf{X}_0, Y) , it follows that

$$\begin{aligned} p_t(\mathbf{x}_t \mid y) &= p(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\xi} \mid y) \\ &= p(\sqrt{\alpha_t} \mathbf{x} \mid y) * \mathcal{N}(\mathbf{0}, (1 - \alpha_t) \mathbf{I}_D) \\ &= (\sqrt{\alpha_t})_{\#} p(\mathbf{x} \mid y) * \mathcal{N}(\mathbf{0}, (1 - \alpha_t) \mathbf{I}_D). \end{aligned}$$

Consequently, we obtain:

$$p_t^y(\mathbf{x}_t) = p_t(\mathbf{x}_t \mid y). \quad \square$$

Lemma B.2. Consider three random variables, $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^n$, and $Z \in \mathbb{R}$. Let \mathbf{Y} be independent of paired (\mathbf{X}, Z) , and $\mathbf{W} = \mathbf{X} + \mathbf{Y}$. Then, we have

$$p_{W|Z}(\mathbf{w} | z) = (p_{X|Z}(\cdot | z) * p_Y(\cdot))(\mathbf{w}).$$

Or informally,

$$p_{XY|Z}(\mathbf{x} + \mathbf{y} | z) = p_{X|Z}(\mathbf{x} | z) * p_Y(\mathbf{y}).$$

Proof. Because \mathbf{Y} is independent of (\mathbf{X}, Z) ,

$$p_{XYZ}(\mathbf{x}, \mathbf{y}, z) = p_{XZ}(\mathbf{x}, z) p_Y(\mathbf{y}).$$

Let $D_w = \{(\mathbf{x}, \mathbf{y}) : \mathbf{x} + \mathbf{y} \leq \mathbf{w}\}$. Then, we have

$$\begin{aligned} \mathbb{P}(\mathbf{W} \leq \mathbf{w}, Z \leq z) &= \mathbb{P}(\mathbf{X} + \mathbf{Y} \leq \mathbf{w}, Z \leq z) \\ &= \int_0^z \left(\iint_{D_w} p_{XYZ}(\mathbf{x}, \mathbf{y}, z) d\mathbf{x} d\mathbf{y} \right) dz \\ &= \int_0^z \left(\iint_{D_w} p_{XZ}(\mathbf{x}, z) p_Y(\mathbf{y}) d\mathbf{x} d\mathbf{y} \right) dz \\ &= \int_0^z \int_0^{\mathbf{w}} (p_{XZ}(\cdot, z) * p_Y(\cdot))(\mathbf{s}) d\mathbf{s} dz, \end{aligned}$$

where $\mathbf{W} = (W_i)_i \leq \mathbf{w} = (w_i)_i$ means $W_i \leq w_i$ for all $i = 1, \dots, n$, and $\int_0^{\mathbf{w}} d\mathbf{s} = \int_0^{w_n} \dots \int_0^{w_1} ds_1 \dots ds_n$. It follows that

$$p_{WZ}(\mathbf{w}, z) = (p_{XZ}(\cdot, z) * p_Y(\cdot))(\mathbf{w}).$$

Therefore,

$$p_{W|Z}(\mathbf{w} | z) = \frac{p_{WZ}(\mathbf{w}, z)}{p_Z(z)} = \left(\frac{p_{XZ}(\cdot, z)}{p_Z(z)} * p_Y(\cdot) \right)(\mathbf{w}) = (p_{X|Z}(\cdot | z) * p_Y(\cdot))(\mathbf{w}). \quad \square$$

Remark B.1. In Lemma B.2, the existence of density functions is assumed, which also makes it necessary to assume the existence of the density for \mathbf{X}_0 in the proof of Proposition B.1. However, this condition is often not satisfied in practice. To address this limitation, consider the convolution of two probability measures $\mu, \nu \in \mathcal{P}(\mathbb{R}^n)$, defined by

$$\mu * \nu(U) := \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} \chi_U(\mathbf{x} + \mathbf{y}) d\mu(\mathbf{x}) d\nu(\mathbf{y}).$$

Note that $\mu * \nu$ is still a probability measure. Moreover, it follows that if $\mathbf{X} \sim \mu$ and $\mathbf{Y} \sim \nu$ with \mathbf{X} independent of \mathbf{Y} , then $\mathbf{X} + \mathbf{Y} \sim \mu * \nu$. Under this formulation, the conclusion of Lemma B.2 remains valid in the general case:

$$\mathbb{P}_{W|Z}(\cdot | Z) = \mathbb{P}_{X|Z}(\cdot | Z) * \mathbb{P}_{Y|Z}(\cdot | Z) = \mathbb{P}_{X|Z}(\cdot | Z) * \mathbb{P}_Y(\cdot),$$

where the first equality follows from the fact that independence of \mathbf{Y} and (\mathbf{X}, Z) implies that \mathbf{Y} is independent of \mathbf{X} conditional on Z , and the second equality holds because \mathbf{Y} is independent of Z due to its independence from the pair (\mathbf{X}, Z) . Therefore, by following a similar line of reasoning as in the proof of Proposition B.1—replacing statements about densities with statements about distributions—we can obtain the same result even when \mathbf{X}_0 does not admit a density function.

C More Details of Geometric Guidance

C.1 Omitted Poofs in Section 4

Proof of Proposition 1. Fix a time $t > 0$. By Equation (2),

$$\mathbf{X}_t = \sqrt{\alpha_t} A \mathbf{Z} + \sqrt{1 - \alpha_t} \boldsymbol{\xi},$$

for some $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$. It follows that

$$f(\mathbf{X}_t) := \|(\mathbf{I}_D - AA^\top)\mathbf{X}_t\| = \sqrt{1 - \alpha_t} \|(\mathbf{I}_D - AA^\top)\boldsymbol{\xi}\|.$$

Note that AA^\top is the orthogonal projection to $\text{Im } A$. Therefore, there exists a $U \in \mathcal{O}^{D \times D}$ such that

$$\mathbf{I}_D - AA^\top = U^\top \text{diag}(\underbrace{1, \dots, 1}_{D-d}, 0, \dots, 0)U.$$

Moreover, the orthogonality of U implies that $\boldsymbol{\nu} = (\nu_1, \dots, \nu_D)^\top = U\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$. Hence,

$$f(\mathbf{X}_t) = \sqrt{1 - \alpha_t} \|(\mathbf{I}_D - AA^\top)\boldsymbol{\xi}\| = \sqrt{1 - \alpha_t} (\nu_1^2 + \dots + \nu_{D-d}^2)^{\frac{1}{2}}.$$

For any $\varepsilon > 0$, by setting $\alpha = (D - d)\varepsilon$ in the Laurent-Massart bound (Lemma H.1), we obtain

$$\mathbb{P}\left(r(t)\sqrt{1 - 2\sqrt{\varepsilon}} \leq f(\mathbf{X}_t) \leq r(t)\sqrt{1 + 2\sqrt{\varepsilon} + 2\varepsilon}\right) \geq 1 - 2e^{-2(D-d)\varepsilon}.$$

Since $d \ll D$, we can choose ε sufficiently small such that $\delta = e^{-2(D-d)\varepsilon}$ is also sufficiently small. As a result, $\mathbb{P}(f(\mathbf{X}_t) \approx r(t)) \geq 1 - \delta$, i.e., \mathbf{X}_t concentrates on $\mathcal{M}^t = f^{-1}(r(t))$ with high probability. \square

Proof of Theorem 2. First, by applying the orthogonal decomposition of the score function in Equation (7), the deterministic reverse process (4) can be rewritten as

$$\frac{d}{dt}\mathbf{X}_t^\leftarrow = \mathbf{X}_t^\leftarrow + A\nabla_{\mathbf{z}} \log p_{T-t}^Z(A^\top \mathbf{X}_t^\leftarrow) - \frac{1}{1 - \alpha_{T-t}}(\mathbf{I}_D - AA^\top)\mathbf{X}_t^\leftarrow. \quad (18)$$

- (a) Because $A \in \mathcal{O}^{D \times d}$, we have $A^\top A = \mathbf{I}_d$ and $A^\top(\mathbf{I}_D - AA^\top) = \mathbf{0}$. Therefore, by multiplying A^\top on the both sides of (18),

$$\frac{d}{dt}\mathbf{Z}_t^\leftarrow = \mathbf{Z}_t^\leftarrow + \nabla_{\mathbf{z}} \log p_{T-t}^Z(\mathbf{Z}_t^\leftarrow),$$

for $\mathbf{Z}_t^\leftarrow = A^\top \mathbf{X}_t^\leftarrow$. Moreover, by the equivalence of the continuity equation of the Fokker-Planck equation (or by the statements in Appendix C.2), $\mathbf{Z}_t = \mathbf{Z}_{T-t}^\leftarrow$ satisfies the forward process of DDPMs starting from p^Z .

- (b) Similarly, by multiplying $\mathbf{I}_D - AA^\top$ on the both sides of (18),

$$\frac{d}{dt}\mathbf{X}_{t,\perp}^\leftarrow = \mathbf{X}_{t,\perp}^\leftarrow - \frac{1}{1 - \alpha_{T-t}}\mathbf{X}_{t,\perp}^\leftarrow = -\frac{\alpha_{T-t}}{1 - \alpha_{T-t}}\mathbf{X}_{t,\perp}^\leftarrow,$$

for $\mathbf{X}_{t,\perp}^\leftarrow = (\mathbf{I}_D - AA^\top)\mathbf{X}_t^\leftarrow$. Note that $\alpha_{T-t} = e^{-2(T-t)}$. Therefore, this equation has the analytical solution given by

$$\mathbf{X}_{t_0+\delta,\perp}^\leftarrow = \sqrt{\frac{1 - e^{-2(T-(t_0+\delta))}}{1 - e^{-2(T-t_0)}}}\mathbf{X}_{t_0,\perp}^\leftarrow.$$

When $\|\mathbf{X}_{t_0,\perp}^\leftarrow\| = \sqrt{(D-d)(1 - e^{-2(T-t_0)})}$, it follows that

$$\|\mathbf{X}_{t_0+\delta,\perp}^\leftarrow\| = \sqrt{\frac{1 - e^{-2(T-(t_0+\delta))}}{1 - e^{-2(T-t_0)}}}\|\mathbf{X}_{t_0,\perp}^\leftarrow\| = \sqrt{(D-d)(1 - e^{-2(T-(t_0+\delta))})}. \quad \square$$

C.2 Decomposition of Score Function

By Equation (2) and the assumption $\mathbf{X}_0 = A\mathbf{Z}$, we have

$$\begin{aligned}\mathbf{X}_t &= \sqrt{\alpha_t}\mathbf{X}_0 + \sqrt{1-\alpha_t}\boldsymbol{\xi} \\ &= \underbrace{\sqrt{\alpha_t}\mathbf{X}_0 + \sqrt{1-\alpha_t}Q\boldsymbol{\xi}}_{=: \mathbf{X}_{t,\parallel}} + \underbrace{\sqrt{1-\alpha_t}(\mathbf{I}_D - Q)\boldsymbol{\xi}}_{=: \mathbf{X}_{t,\perp}}\end{aligned}$$

for some $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$, where $Q = AA^\top$ is the orthogonal projection onto $\text{Im } A$.

We compute the covariance:

$$\begin{aligned}\text{Cov}(Q\boldsymbol{\xi}, (\mathbf{I}_D - Q)\boldsymbol{\xi}) &= \mathbb{E}[Q\boldsymbol{\xi} \cdot ((\mathbf{I}_D - Q)\boldsymbol{\xi})^\top] - \mathbb{E}[Q\boldsymbol{\xi}] \cdot \mathbb{E}[(\mathbf{I}_D - Q)\boldsymbol{\xi}]^\top \\ &= \mathbb{E}[Q\boldsymbol{\xi} \cdot ((\mathbf{I}_D - Q)\boldsymbol{\xi})^\top] = Q\mathbb{E}[\boldsymbol{\xi}\boldsymbol{\xi}^\top](\mathbf{I}_D - Q) \\ &= Q(\mathbf{I}_D - Q) = 0,\end{aligned}$$

which shows that $Q\boldsymbol{\xi}$ and $(\mathbf{I}_D - Q)\boldsymbol{\xi}$ are uncorrelated. Since both are Gaussian, they are independent. Hence, $\mathbf{X}_{t,\perp}$ is independent of $\sqrt{1-\alpha_t}Q\boldsymbol{\xi}$. Combined with the fact that $\boldsymbol{\xi}$ is independent of \mathbf{X}_0 , it follows that $\mathbf{X}_{t,\parallel}$ is independent of $\mathbf{X}_{t,\perp}$. By Lemma H.2, the density of \mathbf{X}_t admits the decomposition

$$p_t(\mathbf{x}) = p_{t,\parallel}(\mathbf{x}_\parallel)p_{t,\perp}(\mathbf{x}_\perp), \quad (19)$$

where $p_{t,\parallel}$ and $p_{t,\perp}$ are the densities of $\mathbf{X}_{t,\parallel}$ and $\mathbf{X}_{t,\perp}$ with respect to the canonical volume measures on $\text{Im } A$ and $(\text{Im } A)^\perp$, respectively. Here, $\mathbf{x}_\parallel = Q\mathbf{x}$ and $\mathbf{x}_\perp = \mathbf{x} - \mathbf{x}_\parallel$.

Next, let us analyze $p_{t,\parallel}$ and $p_{t,\perp}$, respectively.

- (i) For the parallel part, first define $\mathbf{Z}_t := A^\top \mathbf{X}_t$. Then, by multiplying A^\top on the both sides of Equation (1), we obtain

$$d\mathbf{Z}_t = -\mathbf{Z}_t dt + \sqrt{2}d\mathbf{B}_t,$$

where $(\mathbf{B}_t)_{t \geq 0} = (A^\top \mathbf{W}_t)_{t \geq 0}$ is a standard Brownian motion on \mathbb{R}^d by Lemma H.3. Therefore, the process $\mathbf{Z}_t \sim p_t^Z$ is governed by the DDPM dynamics initialized from p^Z . Since

$$\mathbf{X}_{t,\parallel} = Q\mathbf{X}_t = A\mathbf{Z}_t,$$

this shows that $\mathbf{X}_{t,\parallel}$ evolves as a diffusion process on the target data manifold $\mathcal{M} = \text{Im } A$.

Moreover, applying Lemma H.4 gives

$$p_{t,\parallel}(\mathbf{x}_\parallel) = A_\# p_t^Z(\mathbf{x}_\parallel) = p_t^Z(A^\top \mathbf{x}_\parallel) = p_t^Z(A^\top \mathbf{x}). \quad (20)$$

- (ii) For the orthogonal part, we have

$$\mathbf{X}_{t,\perp} = \sqrt{1-\alpha_t}P\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, (1-\alpha_t)P),$$

where $P = \mathbf{I}_D - Q$ is an orthogonal projection with rank $D-d$. So $P = B^\top B$ for some $B \in \mathcal{O}^{D \times (D-d)}$. It follows that $\mathbf{X}_{t,\perp}$ is a Gaussian on $\text{Im } B$, i.e., $\mathbf{X}_{t,\perp} = B\mathbf{W}$ for some $\mathbf{W} \sim \mathcal{N}(\mathbf{0}, (1-\alpha_t)\mathbf{I}_{D-d})$. Therefore, $\mathbf{X}_{t,\perp}$ is basically a $(D-d)$ -dimensional Gaussian. When $d \ll D$, as shown in the proof in Proposition 1,

$$\|(\mathbf{I}_D - AA^\top)\mathbf{X}_t\| = \|\mathbf{X}_{t,\perp}\| \approx r(t),$$

which implies that the orthogonal part $\mathbf{X}_{t,\perp}$ is responsible for the concentration of \mathbf{X}_t on \mathcal{M}^t and endows \mathbf{X}_t with its geometric structure. Furthermore, by Lemma H.5, p_t^\perp is approximately uniform on the sphere $\mathbb{S}^{(D-d)-1}(r(t))$. In other words, the density p_t , which is concentrated on the cylindrical-like surface \mathcal{M}^t , remains constant along radial directions and varies only in the longitudinal direction governed by $p_{t,\parallel}$ —a consequence of diffusion along the subspace $\text{Im } A$.

Moreover, applying Lemma H.4 again, we obtain

$$p_{t,\perp}(\mathbf{x}_\perp) = B_\# p^W(\mathbf{x}_\perp) = p^W(B^\top \mathbf{x}_\perp) = p^W(B^\top \mathbf{x}), \quad (21)$$

where

$$p^W(\mathbf{w}) = (2\pi(1 - \alpha_t))^{-\frac{D-d}{2}} \exp\left(-\frac{\|\mathbf{w}\|^2}{2(1 - \alpha_t)}\right).$$

Finally, for the decomposition, by combining (20) and (21) with (19), we get

$$\log p_t(\mathbf{x}) = \log p_t^Z(A^\top \mathbf{x}) + \log p^W(B^\top \mathbf{x}),$$

from which the orthogonal decomposition formula immediately follows:

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = A \nabla_{\mathbf{z}} \log p_t^Z(\mathbf{z})|_{\mathbf{z}=A^\top \mathbf{x}} - \frac{1}{1 - \alpha_t} (\mathbf{I}_D - P)\mathbf{x},$$

as originally derived via direct computation by Chen et al. (2023b).

For the geometric property, the randomness of \mathbf{X}_t arises from the diffusion process on the target data manifold $\mathcal{M} = \text{Im } A$, while the geometric structure of \mathbf{X}_t results from the concentration behavior of the orthogonal part.

C.3 Construction of Geometric Guidance

To clarify our intuition about $\nabla_{\mathbf{x}} \log p_t(y = 1 | \mathbf{x})$ “almost normal” to \mathcal{M}_1^t , we will show that there exists a small $\beta_t > 0$ such that

$$\|\nabla_{\mathbf{x}} \log p_t(y = 1 | \mathbf{x}) + \eta_t P_1 \mathbf{x}\| \leq \beta_t, \quad \forall \mathbf{x} \in \mathcal{M}_1^t,$$

for some scalar $\eta_t > 0$. But first, we need the following lemma.

Lemma C.1. *Let $\mathcal{M} \subset \mathbb{R}^D$ be a smooth manifold with dimension $D - 1$. Let $V \subset \mathbb{R}^D$ be a tubular neighborhood of \mathcal{M} with the orthogonal projection $\pi: V \rightarrow \mathcal{M}$. Let $f: V \rightarrow \mathbb{R}$ be a C^2 -function satisfying the following two conditions.*

$$(a) \quad \|\nabla_{\mathbf{x}}^2 f(\mathbf{x})\|_{\text{op}} \leq L.$$

$$(b) \quad f|_{\mathcal{M}} \text{ is } \beta\text{-Lipschitz with the induced distance of } \mathbb{R}^n \text{ on } \mathcal{M}.$$

Then for any $\mathbf{x} \in V$,

$$\|\nabla_{\mathbf{x}} f(\mathbf{x}) - \partial_n f(\pi(\mathbf{x}))n(\pi(\mathbf{x}))\| \leq \beta + L \text{dist}(\mathbf{x}, \mathcal{M}),$$

where $n: \mathcal{M} \rightarrow \mathbb{R}^n$ is a continuous unit normal vector field along \mathcal{M} , $\partial_n f = \langle \nabla f, n \rangle$ the derivative along n , and $\text{dist}(\mathbf{x}, \mathcal{M}) = \inf \{\|\mathbf{x} - \mathbf{y}\|: \mathbf{y} \in \mathcal{M}\}$ is the distance from \mathbf{x} to \mathcal{M} .

Proof. Let \mathcal{M} be equipped with the induced Riemannian structure of \mathbb{R}^n and ∇^M be the corresponding Levi-Civita connection. Because $\mathcal{M} \subset \mathbb{R}^D$ is a hypersurface, i.e., submanifold with dimension $D - 1$,

$$\nabla f = \nabla^M f + (\partial_n f)n, \quad (22)$$

see the details in Lee (2019, Chapter 8). Fix $\mathbf{x} \in V$ with $\mathbf{y} = \pi(\mathbf{x}) \in \mathcal{M}$. Note that

$$\text{dist}(\mathbf{x}, \mathcal{M}) = \|\mathbf{x} - \mathbf{y}\|, \quad (23)$$

by Lee (2019, Proposition 5.26 (c)). Writing

$$\|\nabla_{\mathbf{x}} f(\mathbf{x}) - \partial_n f(\mathbf{y})n(\mathbf{y})\| \leq \|\nabla_{\mathbf{x}} f(\mathbf{x}) - \nabla_{\mathbf{x}} f(\mathbf{y})\| + \|\nabla_{\mathbf{x}} f(\mathbf{y}) - \partial_n f(\mathbf{y})n(\mathbf{y})\|. \quad (24)$$

I. For the first term, by

$$\nabla_{\mathbf{x}} f(\mathbf{x}) - \nabla_{\mathbf{x}} f(\mathbf{y}) = \int_0^1 \nabla_{\mathbf{x}}^2 f(\mathbf{y} + s(\mathbf{x} - \mathbf{y})) (\mathbf{x} - \mathbf{y}) ds,$$

the fact that $\|\nabla_{\mathbf{x}}^2 f(\mathbf{x})\|_{\text{op}} \leq L$, and Equation (23), we have

$$\|\nabla_{\mathbf{x}} f(\mathbf{x}) - \nabla_{\mathbf{x}} f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\| = L \text{dist}(\mathbf{x}, \mathcal{M}). \quad (25)$$

II. For the second term, first, by (22),

$$\|\nabla_{\mathbf{x}} f(\mathbf{y}) - \partial_n f(\mathbf{y}) n(\mathbf{y})\| = \|\nabla^M f(\mathbf{y})\|.$$

By assumption, $f|_{\mathcal{M}}$ is β -Lipschitz with the induced distance of \mathbb{R}^n on \mathcal{M} , i.e.,

$$|f(\mathbf{y}_1) - f(\mathbf{y}_2)| \leq \beta d_{\mathcal{M}}(\mathbf{y}_1, \mathbf{y}_2),$$

where $d_{\mathcal{M}}$. It implies that

$$\|\nabla^M f(\mathbf{z})\| \leq \beta, \quad \forall \mathbf{z} \in \mathcal{M}, \quad (26)$$

see the details in Boumal (2023, Proposition 10.43).

Then combining the inequalities (25) and (26) with (24),

$$\|\nabla_{\mathbf{x}} f(\mathbf{x}) - \partial_n f(\mathbf{y}) n(\mathbf{y})\| \leq \beta + L \text{dist}(\mathbf{x}, \mathcal{M}). \quad \square$$

Let $f_t(\mathbf{x}) = \log p_t(y = 1 \mid \mathbf{x})$. It is natural to assume that f_t is C^2 on a tubular neighborhood V of \mathcal{M}_1^t , that $\|\nabla^2 f_t\|_{\text{op}} \leq L_t$ on V , and that f_t is β_t -Lipschitz continuous on \mathcal{M}_1^t . Then by Lemma C.1,

$$\|\nabla_{\mathbf{x}} f_t(\mathbf{x}) - \partial_n f_t(\pi(\mathbf{x})) n_t(\pi(\mathbf{x}))\| \leq \beta_t + L_t \text{dist}(\mathbf{x}, \mathcal{M}_1^t).$$

In particular, for any $\mathbf{x} \in \mathcal{M}_1^t$ and $\pi(\mathbf{x}) = \mathbf{x}$, we have

$$\|\nabla_{\mathbf{x}} f_t(\mathbf{x}) - \partial_n f_t(\mathbf{x}) n_t(\mathbf{x})\| \leq \beta_t.$$

Two questions remain: whether $\partial_n f_t(\mathbf{x}) n_t(\mathbf{x}) = -\eta_t P_1 \mathbf{x}$ for some scalar $\eta_t > 0$, and how to bound β_t .

For the first question, we can choose $n_t(\mathbf{x}) = P_1 \mathbf{x} / \|P_1 \mathbf{x}\|$ by the definition (9) of \mathcal{M}_1^t and Lemma I.1. So

$$\partial_n f_t(\mathbf{x}) n_t(\mathbf{x}) = -\eta_t P_1 \mathbf{x},$$

for

$$\eta_t = -\frac{\partial_n f_t(\mathbf{x})}{\|P_1 \mathbf{x}\|}.$$

Moreover, because $p_t(y = 1 \mid \mathbf{x})$ is the classifier for $(\mathbf{X}_t, y = 1)$ and such \mathbf{X}_t concentrates on \mathcal{M}_1^t by Proposition 1, $f_t(\mathbf{x}) = \log p_t(y = 1 \mid \mathbf{x})$ decreases when \mathbf{x} moves away from \mathcal{M}_1^t . So

$$\partial_n f_t(\mathbf{x}) < 0 \quad \Rightarrow \quad \eta_t > 0.$$

Next, to bound β_t , we introduce the following lemma.

Lemma C.2. *Let $p(y = k \mid \mathbf{x})$ be a softmax classifier with logits $g_k(\mathbf{x})$ for $k = 1, 2, \dots, K$, that is,*

$$p(y = k \mid \mathbf{x}) = \frac{\exp(g_k(\mathbf{x}))}{\sum_{k=1}^K \exp(g_k(\mathbf{x}))}.$$

Assume $\|\nabla_{\mathbf{x}} g_k(\mathbf{x})\| \leq L$ for all k, \mathbf{x} . Let \mathcal{M}_k be the region where points with label $y = k$ concentrate on. Assume classifier confidence

$$p(y = k \mid \mathbf{x}) > 1 - \varepsilon, \quad \forall \mathbf{x} \in \mathcal{M}_k.$$

Then

$$\|\nabla_{\mathbf{x}} \log p(y = k \mid \mathbf{x})\| \leq 2L\varepsilon, \quad \forall \mathbf{x} \in \mathcal{M}_k.$$

Proof. Fix k . Let $f(\mathbf{x}) = \log p(y = k \mid \mathbf{x})$.

$$\begin{aligned}\nabla_{\mathbf{x}} f(\mathbf{x}) &= \nabla_{\mathbf{x}} g_k(\mathbf{x}) - \sum_{j=1}^K p(y = j \mid \mathbf{x}) \nabla_{\mathbf{x}} g_j(\mathbf{x}) \\ &= \sum_{j=1}^K p(y = j \mid \mathbf{x}) (\nabla_{\mathbf{x}} g_j(\mathbf{x}) - \nabla_{\mathbf{x}} g_k(\mathbf{x})) \\ &= \sum_{j \neq k} p(y = j \mid \mathbf{x}) (\nabla_{\mathbf{x}} g_j(\mathbf{x}) - \nabla_{\mathbf{x}} g_k(\mathbf{x}))\end{aligned}$$

By assumption,

$$\|\nabla_{\mathbf{x}} g_j(\mathbf{x}) - \nabla_{\mathbf{x}} g_k(\mathbf{x})\| \leq \|\nabla_{\mathbf{x}} g_j(\mathbf{x})\| + \|\nabla_{\mathbf{x}} g_k(\mathbf{x})\| \leq 2L.$$

Therefore,

$$\begin{aligned}\|\nabla_{\mathbf{x}} f(\mathbf{x})\| &\leq \sum_{j \neq k} p(y = j \mid \mathbf{x}) \|\nabla_{\mathbf{x}} g_j(\mathbf{x}) - \nabla_{\mathbf{x}} g_k(\mathbf{x})\| \\ &\leq 2L \sum_{j \neq k} p(y = j \mid \mathbf{x}) = 2L(1 - p(y = k \mid \mathbf{x})).\end{aligned}$$

It implies that

$$\|\nabla_{\mathbf{x}} f(\mathbf{x})\| \leq 2L\varepsilon, \quad \forall \mathbf{x} \in \mathcal{M}_k,$$

by the assumption that classifier confidence $> 1 - \varepsilon$ on \mathcal{M}_k . □

Therefore, for all $p_t(y = 1 \mid \mathbf{x})$, we assume that they satisfy the conditions in Lemma C.2. Then if

$$p_t(y = 1 \mid \mathbf{x}) > 1 - \varepsilon_t, \quad \forall \mathbf{x} \in \mathcal{M}_1^t,$$

for a small ε_t , then

$$\|\nabla^M f_t(\mathbf{x})\| \leq \|\nabla_{\mathbf{x}} f_t(\mathbf{x})\| = \sqrt{\|\nabla^M f_t(\mathbf{x})\|^2 + |\partial_n f_t(\mathbf{x})|^2} \leq 2C\varepsilon_t, \quad \forall \mathbf{x} \in \mathcal{M}_1^t.$$

So $\beta_t \leq 2C\varepsilon_t$.

Combining above results, we have

$$\|\nabla_{\mathbf{x}} \log p_t(y = 1 \mid \mathbf{x}) + \eta_t P_1 \mathbf{x}\| \leq \beta_t, \quad \forall \mathbf{x} \in \mathcal{M}_1^t,$$

for some $\eta_t > 0$. Moreover, $\beta_t = \mathcal{O}(\varepsilon_t)$ for $p_t(y = 1 \mid \mathbf{x}) > 1 - \varepsilon_t$ on \mathcal{M}_1^t .

D More Details related to Main Results

D.1 Omitted Proofs in Section 5.1

Proof of Proposition 3. First, the Hessian is

$$\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x}) = \frac{\nabla_{\mathbf{x}}^2 p_X(\mathbf{x})}{p_X(\mathbf{x})} - \frac{\nabla_{\mathbf{x}} p_X(\mathbf{x}) \nabla_{\mathbf{x}} p_X(\mathbf{x})^\top}{p_X(\mathbf{x})^2}.$$

To express the above formula explicitly, by the definition, for any $\mathbf{x} \in \mathbb{R}^n$,

$$p_X(\mathbf{x}) = \int_{\mathbb{R}^k} K_z(\mathbf{x}) p^Z(\mathbf{z}) d\mathbf{z}, \quad K_z(\mathbf{x}) := (2\pi\beta^2)^{-\frac{n}{2}} \exp\left(-\frac{\|\mathbf{x} - \alpha B \mathbf{z}\|^2}{2\beta^2}\right),$$

and so

$$\nabla_{\mathbf{x}} K_z(\mathbf{x}) = \frac{\alpha B \mathbf{z} - \mathbf{x}}{\beta^2} K_z(\mathbf{x}),$$

$$\nabla_{\mathbf{x}}^2 K_z(\mathbf{x}) = -\frac{1}{\beta^2} K_z(\mathbf{x}) \mathbf{I}_n + \frac{(\mathbf{x} - \alpha B \mathbf{z})(\mathbf{x} - \alpha B \mathbf{z})^\top}{\beta^4} K_z(\mathbf{x}).$$

Let

$$d\mu_x(\mathbf{z}) = \frac{K_z(\mathbf{x}) p^Z(\mathbf{z})}{p^X(\mathbf{x})} d\mathbf{z}$$

be the posterior probability measure on \mathbb{R}^k . Then, for the first term

$$\frac{\nabla_{\mathbf{x}}^2 p_X(\mathbf{x})}{p_X(\mathbf{x})} = \frac{\int_{\mathbb{R}^k} \nabla_{\mathbf{x}}^2 K_z(\mathbf{x}) p^Z(\mathbf{z}) d\mathbf{z}}{p_X(\mathbf{x})} = -\frac{1}{\beta^2} \mathbf{I}_n + \frac{1}{\beta^4} \mathbb{E}_{\mathbf{Z} \sim \mu_x} [(\mathbf{x} - \alpha B \mathbf{Z})(\mathbf{x} - \alpha B \mathbf{Z})^\top],$$

and for the second term

$$\frac{\nabla_{\mathbf{x}} p_X(\mathbf{x}) \nabla_{\mathbf{x}} p_X(\mathbf{x})^\top}{p_X(\mathbf{x})^2} = \frac{1}{\beta^4} \mathbb{E}_{\mathbf{Z} \sim \mu_x} [\mathbf{x} - \alpha B \mathbf{Z}] \mathbb{E}_{\mathbf{Z} \sim \mu_x} [\mathbf{x} - \alpha B \mathbf{Z}]^\top.$$

Moreover, note that

$$\begin{aligned} & \mathbb{E}_{\mathbf{Z} \sim \mu_x} [(\mathbf{x} - \alpha B \mathbf{Z})(\mathbf{x} - \alpha B \mathbf{Z})^\top] - \mathbb{E}_{\mathbf{Z} \sim \mu_x} [\mathbf{x} - \alpha B \mathbf{Z}] \mathbb{E}_{\mathbf{Z} \sim \mu_x} [\mathbf{x} - \alpha B \mathbf{Z}]^\top \\ &= \text{Cov}_{\mathbf{Z} \sim \mu_x}(\mathbf{x} - \alpha B \mathbf{Z}) = \alpha^2 \text{Cov}_{\mu_x}(B \mathbf{Z}) = \alpha^2 B \text{Cov}_{\mu_x}(\mathbf{Z}) B^\top. \end{aligned}$$

Therefore, we get

$$\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x}) = \frac{\alpha^2}{\beta^4} B \text{Cov}_{\mu_x}(\mathbf{Z}) B^\top - \frac{1}{\beta^2} \mathbf{I}_n. \quad (27)$$

It follows that

$$\|\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x})\|_{\text{op}} \leq \frac{1}{\beta^2} + \frac{\alpha^2 \Lambda}{\beta^4} \|\text{Cov}_{\mu_x}(\mathbf{Z})\|_{\text{op}}. \quad (28)$$

It is sufficient to bound $\|\text{Cov}_{\mu_x}(\mathbf{Z})\|_{\text{op}}$. To do that, μ_x is required to satisfy the Poincaré Inequality. Let $p^Z(\mathbf{z}) = \exp(-V(\mathbf{z}))$ for some $V: \mathbb{R}^k \rightarrow \mathbb{R}$ and

$$U_x(\mathbf{z}) := \frac{\|\mathbf{x} - \alpha B \mathbf{z}\|^2}{2\beta^2} + V(\mathbf{z}),$$

which indicates that $d\mu_x(\mathbf{z}) = e^{-U_x(\mathbf{z})} d\mathbf{z} / \int e^{-U_x}$. Because $\nabla_{\mathbf{z}}^2 V(\mathbf{z}) = -\nabla_{\mathbf{z}}^2 \log p^Z(\mathbf{z}) \succeq m_0 \mathbf{I}_k$,

$$\nabla_{\mathbf{z}}^2 U_x(\mathbf{z}) = \frac{\alpha^2}{\beta^2} B^\top B + \nabla_{\mathbf{z}}^2 V(\mathbf{z}) \succeq \left(\frac{\alpha^2 \lambda}{\beta^2} + m_0 \right) \mathbf{I}_k.$$

Then, by Lemma H.6, μ_x satisfies the Poincaré Inequality with constant $m := \alpha^2 \lambda / \beta^2 + m_0$. Thus, for any C^1 function $f: \mathbb{R}^k \rightarrow \mathbb{R}$,

$$\text{Var}_{\mu_x}(f) \leq \frac{1}{m} \mathbb{E}_{\mu_x} [\|\nabla f\|^2].$$

For any $\mathbf{u} \in \mathbb{R}^n$, let $f_u(\mathbf{z}) = \langle \mathbf{u}, \mathbf{z} \rangle$ with $\nabla_{\mathbf{z}} f_u(\mathbf{z}) = \mathbf{u}$. The above inequality implies that

$$\mathbf{u}^\top \text{Cov}_{\mu_x}(\mathbf{Z}) \mathbf{u} = \text{Var}_{\mu_x}(f_u) \leq \frac{1}{m} \mathbb{E}_{\mu_x} [\|\nabla_{\mathbf{z}} f_u\|^2] \leq \frac{1}{m} \|\mathbf{u}\|^2,$$

for any $\mathbf{u} \in \mathbb{R}^k$. Therefore,

$$\|\text{Cov}_{\mu_x}(\mathbf{Z})\|_{\text{op}} \leq \frac{1}{m}. \quad (29)$$

Finally, by plugging inequality (29) into Equation (28), we get the result

$$\|\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x})\|_{\text{op}} \leq \frac{1}{\beta^2} + \frac{\alpha^2 \Lambda}{\beta^2 (\alpha^2 \lambda + m_0 \beta^2)}. \quad \square$$

Proof of Corollary 4. By Equation (29),

$$\|B \text{Cov}_{\mu_x}(\mathbf{Z})B^\top\|_{\text{op}} \leq \frac{\Lambda}{m} \Rightarrow B \text{Cov}_{\mu_x}(\mathbf{Z})B^\top \preceq \frac{\Lambda}{m} \mathbf{I}_n.$$

By combining this with Equation (27), we have

$$\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x}) \preceq \left(\frac{\alpha^2 \Lambda}{\beta^2(\alpha^2 \lambda + m_0 \beta^2)} - \frac{1}{\beta^2} \right) \mathbf{I}_n. \quad \square$$

Proof of Proposition 5. By Lemma H.7,

$$\mathbb{P}_X^\sigma = A_\# \mathbb{P}_\sigma^Z = w_1 A_\# \mathbb{P}_{1,\sigma}^Z + w_2 A_\# \mathbb{P}_{2,\sigma}^Z.$$

Moreover, because $\mathbf{Z}_{1,\sigma} = (\mathbf{Z}_1, 0)^\top + \sigma \boldsymbol{\zeta} \sim \mathbb{P}_{i,\sigma}^Z$ with $\boldsymbol{\zeta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$,

$$A \mathbf{Z}_{1,\sigma} = A_1 \mathbf{Z}_1 + \sigma A \boldsymbol{\zeta} \sim \mathbb{P}_{X|Y}^\sigma(\cdot | Y = 1).$$

Note that $A_1 \mathbf{Z}_1 \sim \mathbb{P}_{X|Y}(\cdot | Y = 1)$. Therefore,

$$\mathcal{W}_1(\mathbb{P}_{X|Y}^\sigma(\cdot | Y = 1), \mathbb{P}_{X|Y}(\cdot | Y = 1)) \leq \mathbb{E}[\|A \mathbf{Z}_{1,\sigma} - A_1 \mathbf{Z}_1\|] = \sigma \mathbb{E}[\|A \boldsymbol{\zeta}\|] \leq \sigma \sqrt{d},$$

where the final inequality is because $A \boldsymbol{\zeta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and Lemma H.8. Similarly, it can get

$$\mathcal{W}_1(\mathbb{P}_{X|Y}^\sigma(\cdot | Y = 2), \mathbb{P}_{X|Y}(\cdot | Y = 2)) \leq \sigma \sqrt{d}.$$

Combining these two inequality and by Lemma H.9, we have

$$\begin{aligned} \mathcal{W}_1(\mathbb{P}_X^\sigma, \mathbb{P}_X) &\leq w_1 \mathcal{W}_1(\mathbb{P}_{X|Y}^\sigma(\cdot | Y = 1), \mathbb{P}_{X|Y}(\cdot | Y = 1)) \\ &\quad + w_2 \mathcal{W}_1(\mathbb{P}_{X|Y}^\sigma(\cdot | Y = 2), \mathbb{P}_{X|Y}(\cdot | Y = 2)) \\ &\leq \sigma \sqrt{d}. \end{aligned} \quad \square$$

Proof of Lemma 6. Let

$$r_1(\mathbf{x}) := \frac{w p_1(\mathbf{x})}{p(\mathbf{x})}, \quad r_2(\mathbf{x}) := 1 - r_1(\mathbf{x}) = \frac{(1-w)p_2(\mathbf{x})}{p(\mathbf{x})}.$$

We have

$$\nabla \log p = \frac{w \nabla p_1 + (1-w) \nabla p_2}{p} = r_1 \nabla \log p_1 + r_2 \nabla \log p_2,$$

and

$$\nabla^2 \log p = r_1 \nabla^2 \log p_1 + r_2 \nabla^2 \log p_2 + \nabla r_1 (\nabla \log p_1 - \nabla \log p_2)^\top.$$

For $r_1 = w p_1 / p$,

$$\begin{aligned} \nabla r_1 &= w \frac{p \nabla p_1 - p_1 \nabla p}{p^2} \\ &= w \frac{(w p_1 + (1-w)p_2) \nabla p_1 - p_1 (w \nabla p_1 + (1-w) \nabla p_2)}{p^2} \\ &= \frac{w(1-w)}{p^2} (p_2 \nabla p_1 - p_1 \nabla p_2) \\ &= r_1 r_2 (\nabla \log p_1 - \nabla \log p_2). \end{aligned}$$

Therefore,

$$\nabla^2 \log p = r_1 \nabla^2 \log p_1 + r_2 \nabla^2 \log p_2 + r_1 r_2 (\nabla \log p_1 - \nabla \log p_2) (\nabla \log p_1 - \nabla \log p_2)^\top.$$

For the first two terms, by the assumption,

$$r_1 \nabla^2 \log p_1 + r_2 \nabla^2 \log p_2 \preceq r_1 L_1 \mathbf{I}_n + r_2 L_2 \mathbf{I}_n \preceq \max\{L_1, L_2\} \mathbf{I}_n.$$

For the third term, because $\sup_{\mathbf{x}} \|\nabla \log p_1(\mathbf{x}) - \nabla \log p_2(\mathbf{x})\| \leq M$,

$$\left\| (\nabla \log p_1 - \nabla \log p_2) (\nabla \log p_1 - \nabla \log p_2)^\top \right\|_{\text{op}} \leq M^2,$$

which implies that

$$(\nabla \log p_1 - \nabla \log p_2) (\nabla \log p_1 - \nabla \log p_2)^\top \preceq M^2 \mathbf{I}_n.$$

For the coefficients $r_1 r_2$, because $r_1, r_2 \in (0, 1)$, $r_1 r_2 \leq 1/4$. Combining these results, we have

$$\nabla^2 \log p \preceq \left(\max\{L_1, L_2\} + \frac{1}{4} M^2 \right) \mathbf{I}_n. \quad \square$$

Proof of Corollary 9. Because

$$m_0^z = m_0^z(\sigma) = \frac{m}{1 + m\sigma^2} - \frac{M^2}{4}$$

is decreasing in σ ,

$$m_0^z \leq m_0^z(0) = m - \frac{M^2}{4}.$$

With the Assumption III, we have

$$m - \frac{M^2}{4} > 1$$

Therefore, by choosing a small σ , we can also have $m_0^z > 1$. It follows that

$$m_t = \frac{m_0^z}{m_0^z + (1 - m_0^z)e^{-2t}}$$

is decreasing in t . So

$$m_I := \inf_{t \in (0, T]} m_t = m_T = \frac{m_0^z}{m_0^z + (1 - m_0^z)e^{-2T}} > 1. \quad \square$$

D.2 Proof of Theorem 10

Proof of Theorem 10. By differentiating $\|\tilde{\mathbf{Y}}_t\|^2$ from (14),

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \|\tilde{\mathbf{Y}}_t\|^2 &= \left\langle \tilde{\mathbf{Y}}_t, \frac{d}{dt} \tilde{\mathbf{Y}}_t \right\rangle \\ &= \langle \tilde{\mathbf{Y}}_t, \tilde{\mathbf{Y}}_t + P_1 \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t) - \eta \tilde{\mathbf{Y}}_t \rangle \\ &= (1 - \eta) \|\tilde{\mathbf{Y}}_t\|^2 + \langle \tilde{\mathbf{Y}}_t, P_1 \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t) \rangle \\ &\leq (1 - \eta) \|\tilde{\mathbf{Y}}_t\|^2 + \|\tilde{\mathbf{Y}}_t\| \|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t)\|. \end{aligned}$$

Therefore,

$$\frac{d}{dt} \|\tilde{\mathbf{Y}}_t\| \leq (1 - \eta) \|\tilde{\mathbf{Y}}_t\| + \|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t)\|.$$

Taking the expectation on the both sides yields

$$\frac{d}{dt} \mathbf{m}_t \leq (1 - \eta) \mathbf{m}_t + \mathbb{E} [\|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t)\|], \quad (30)$$

where $\mathbf{m}_t := \mathbb{E} [\|\tilde{\mathbf{Y}}_t\|]$. Therefore, the next step is to bound $\mathbb{E} [\|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t)\|]$.

Let

$$L_S := \sup_{t \in [\delta, T]} L_t, \quad C := \sup_{t \in [\delta, T]} \|\nabla_{\mathbf{x}} \log p_t^\sigma(\mathbf{0})\| < \infty, \quad (31)$$

where L_t is defined in Theorem 8. By the L_S -Lipschitz of $\nabla_{\mathbf{x}} \log p_t^\sigma$ (Theorem 8),

$$\begin{aligned} \|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t)\| &\leq \|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t) - \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\mathbf{0})\| + \|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\mathbf{0})\| \\ &\leq L_S \|\tilde{\mathbf{X}}_t\| + C \end{aligned} \quad (32)$$

For $\tilde{\mathbf{X}}_t$ in Equation (*), we have

$$\begin{aligned} \frac{d}{dt} \|\tilde{\mathbf{X}}_t\|^2 &= 2\|\tilde{\mathbf{X}}_t\|^2 + 2\langle \tilde{\mathbf{X}}_t, \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t) \rangle - 2\eta \langle \tilde{\mathbf{X}}_t, P_1 \tilde{\mathbf{X}}_t \rangle \\ &\leq 2\|\tilde{\mathbf{X}}_t\|^2 + 2\langle \tilde{\mathbf{X}}_t, \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t) \rangle \\ &\leq 2\|\tilde{\mathbf{X}}_t\|^2 + 2\|\tilde{\mathbf{X}}_t\| \|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t)\|, \end{aligned}$$

where the second inequality is because $\langle \tilde{\mathbf{X}}_t, P_1 \tilde{\mathbf{X}}_t \rangle \geq 0$. Combining this with (32),

$$\frac{d}{dt} \|\tilde{\mathbf{X}}_t\| \leq (1 + L_S) \|\tilde{\mathbf{X}}_t\| + C.$$

By taking the expectation on the both sides of above inequality, Grönwall's Inequality (Lemma H.11) implies

$$\mathbb{E} [\|\tilde{\mathbf{X}}_t\|] \leq \mathbb{E} [\|\tilde{\mathbf{X}}_0\|] e^{(1+L_S)t} + \frac{C}{1+L_S} (e^{(1+L_S)t} - 1). \quad (33)$$

Because $\tilde{\mathbf{X}}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$, $\mathbb{E} [\|\tilde{\mathbf{X}}_0\|] \leq \sqrt{D}$ (Lemma H.8). It follows that

$$\sup_{t \in [0, T-\delta]} \mathbb{E} [\|\tilde{\mathbf{X}}_t\|] \leq \sqrt{D} e^{(1+L_S)T} + \frac{C}{1+L_S} (e^{(1+L_S)T} - 1) =: M_1,$$

and (32) implies

$$\begin{aligned} \sup_{t \in [0, T-\delta]} \mathbb{E} [\|\nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\tilde{\mathbf{X}}_t)\|] &\leq \sup_{t \in [0, T-\delta]} L_S \mathbb{E} [\|\tilde{\mathbf{X}}_t\|] + C \\ &\leq L_S M_1 + C =: M_2. \end{aligned} \quad (34)$$

Then by substituting this into (30),

$$\frac{d}{dt} \mathbf{m}_t \leq -(\eta - 1) \mathbf{m}_t + M_2.$$

Because $\mathbf{m}_0 = \mathbb{E} [\|\tilde{\mathbf{Y}}_0\|] \leq \sqrt{D - d_1}$ by Lemma H.8, by applying Grönwall's Inequality again, we obtain

$$\mathbb{E} [\|\tilde{\mathbf{Y}}_t\|] = \mathbf{m}_t \leq \sqrt{D - d_1} e^{-(\eta-1)t} + \frac{M_2}{\eta-1} (1 - e^{-(\eta-1)t}) =: M_\eta(t), \quad (35)$$

which implies that

$$\mathbb{E} [\|\tilde{\mathbf{Y}}_{T-\delta}\|] \leq \sqrt{D - d_1} e^{-(\eta-1)(T-\delta)} + \frac{M_2}{\eta-1}.$$

For any $\varepsilon > 0$,

$$\frac{M_2}{\eta-1} \leq \frac{\varepsilon}{2} \Rightarrow \eta \geq \frac{2M_2}{\varepsilon} + 1,$$

and

$$\sqrt{D - d_1} e^{-(\eta-1)(T-\delta)} \leq \frac{\varepsilon}{2} \Rightarrow \eta \geq \frac{1}{T-\delta} \log \frac{2\sqrt{D - d_1}}{\varepsilon} + 1.$$

Therefore, for any $\varepsilon > 0$, by choosing

$$\eta \geq \max \left\{ \frac{2M_2}{\varepsilon}, \frac{1}{T-\delta} \log \frac{2\sqrt{D - d_1}}{\varepsilon} \right\} + 1,$$

we have

$$\mathbb{E} [\|\tilde{\mathbf{Y}}_{T-\delta}\|] \leq \varepsilon. \quad \square$$

D.3 Theoretical Analysis for Universal Guidance

Consider the universal guidance model

$$\frac{d\mathbf{X}_t^\leftarrow}{dt} = \mathbf{X}_t^\leftarrow + \nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^\leftarrow) - \eta \nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow), \quad \mathbf{X}_0^\leftarrow \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D), \quad (36)$$

for $t \in [0, T]$, where p_t is the density function in DDPMs.

Theorem D.1. *For the dynamics (36), assume that $\log p_t$ is L -smoothness, f is ρ -strongly convex, and $\mathbb{E}[f(\mathbf{X}_0^\leftarrow)] < \infty$. Then*

$$\mathbb{E}[f(\mathbf{X}_T^\leftarrow)] - f(\mathbf{x}_*) = \mathcal{O}\left(e^{-\eta} + \frac{1}{\eta}\right),$$

where \mathbf{x}_* is the unique minimizer of f .

Proof. By differentiating $f(\mathbf{X}_t^\leftarrow)$,

$$\begin{aligned} \frac{d}{dt} f(\mathbf{X}_t^\leftarrow) &= \left\langle \nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow), \frac{d}{dt} \mathbf{X}_t^\leftarrow \right\rangle \\ &= \langle \nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow), \mathbf{X}_t^\leftarrow + \nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^\leftarrow) - \eta \nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow) \rangle \\ &\leq \|\mathbf{X}_t^\leftarrow\| \|\nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow)\| + \|\nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^\leftarrow)\| \|\nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow)\| - \eta \|\nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow)\|^2. \end{aligned}$$

Let $C = \sup_{t \in [\delta, T]} \|\nabla_{\mathbf{x}} \log p_t(\mathbf{0})\| < \infty$. Then, the L -smoothness of $\log p_t$ implies that

$$\|\nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^\leftarrow)\| \leq L \|\mathbf{X}_t^\leftarrow\| + C.$$

Therefore, by $ab \leq (a^2 + b^2)/2$, we have

$$\begin{aligned} \frac{d}{dt} f(\mathbf{X}_t^\leftarrow) &\leq (1 + L) \|\mathbf{X}_t^\leftarrow\| \|\nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow)\| + C \|\nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow)\| - \eta \|\nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow)\|^2 \\ &\leq \frac{1+L}{2} \left(\|\mathbf{X}_t^\leftarrow\|^2 + \|\nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow)\|^2 \right) + \frac{1}{2} \left(C^2 + \|\nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow)\|^2 \right) - \eta \|\nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow)\|^2 \\ &= -\frac{1}{2}(\eta - 2 - L) \|\nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow)\|^2 + \frac{1+L}{2} \|\mathbf{X}_t^\leftarrow\|^2 + \frac{C^2}{2} \end{aligned}$$

Because f is ρ -strongly convex, by Lemma H.12, it satisfies the ρ -PL inequality,

$$\|\nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow)\|^2 \geq 2\rho (f(\mathbf{X}_t^\leftarrow) - f(\mathbf{x}_*)),$$

For $\eta > L + 2$, we obtain

$$\frac{d}{dt} f(\mathbf{X}_t^\leftarrow) \leq -\rho(\eta - 2 - L) (f(\mathbf{X}_t^\leftarrow) - f(\mathbf{x}_*)) + \frac{1+L}{2} \|\mathbf{X}_t^\leftarrow\|^2 + \frac{C^2}{2}.$$

Taking the expectation on the both sides yields that

$$\frac{d}{dt} \mathbb{E}[f(\mathbf{X}_t^\leftarrow)] \leq -\rho(\eta - 2 - L) (\mathbb{E}[f(\mathbf{X}_t^\leftarrow)] - f(\mathbf{x}_*)) + \frac{1+L}{2} \mathbb{E}[\|\mathbf{X}_t^\leftarrow\|^2] + \frac{C^2}{2}. \quad (37)$$

The next step is to bound $\mathbb{E}[\|\mathbf{X}_t^\leftarrow\|^2]$. Let $\mathbf{R}_t := \mathbf{X}_t^\leftarrow - \mathbf{x}_*$. Then

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \|\mathbf{R}_t\|^2 &= \langle \mathbf{R}_t, \mathbf{X}_t^\leftarrow + \nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^\leftarrow) - \eta \nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow) \rangle \\ &= \langle \mathbf{R}_t, \mathbf{X}_t^\leftarrow \rangle + \langle \mathbf{R}_t, \nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^\leftarrow) \rangle - \eta \langle \mathbf{R}_t, \nabla_{\mathbf{x}} f(\mathbf{X}_t^\leftarrow) \rangle. \end{aligned} \quad (38)$$

To obtain the desired inequality, we consider these three terms respectively. For the first term,

$$\langle \mathbf{R}_t, \mathbf{X}_t^\leftarrow \rangle = \|\mathbf{R}_t\|^2 + \langle \mathbf{R}_t, \mathbf{x}_* \rangle \leq \|\mathbf{R}_t\|^2 + \|\mathbf{x}_*\| \|\mathbf{R}_t\|. \quad (39)$$

Let $c = \|\nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{x}_*)\|$. By the L -smoothness of $\log p_t$, we have

$$\begin{aligned} \|\log p_{T-t}(\mathbf{X}_t^{\leftarrow})\| &\leq \|\log p_{T-t}(\mathbf{X}_t^{\leftarrow}) - \nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{x}_*)\| + \|\nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{x}_*)\| \\ &\leq L\|\mathbf{R}_t\| + c. \end{aligned}$$

Therefore, for the second term,

$$\begin{aligned} \langle \mathbf{R}_t, \nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^{\leftarrow}) \rangle &\leq \|\mathbf{R}_t\| \|\nabla_{\mathbf{x}} \log p_{T-t}(\mathbf{X}_t^{\leftarrow})\| \\ &\leq L\|\mathbf{R}_t\|^2 + c\|\mathbf{R}_t\|. \end{aligned} \tag{40}$$

For the third term, because f is ρ -strongly convex, $\nabla_{\mathbf{x}} f(\mathbf{x}_*) = 0$ and

$$\langle \mathbf{R}_t, \nabla_{\mathbf{x}} f(\mathbf{X}_t^{\leftarrow}) \rangle = \langle \mathbf{R}_t, \nabla_{\mathbf{x}} f(\mathbf{X}_t^{\leftarrow}) - \nabla_{\mathbf{x}} f(\mathbf{x}_*) \rangle \geq \rho\|\mathbf{R}_t\|^2. \tag{41}$$

Then, by combining (38) with (39) (40) (41), we have

$$\begin{aligned} \frac{d}{dt}\|\mathbf{R}_t\|^2 &\leq 2(L+1-\eta\rho)\|\mathbf{R}_t\|^2 + 2\tilde{c}\|\mathbf{R}_t\| \\ &\leq (2L+3-2\eta\rho)\|\mathbf{R}_t\|^2 + \tilde{c}^2. \end{aligned}$$

where $\tilde{c} = \|\mathbf{x}_*\| + c$. By taking the expectation on the both sides, Grönwall's Inequality (Lemma H.11) implies that

$$\mathbb{E}[\|\mathbf{R}_t\|^2] \leq \mathbb{E}[\|\mathbf{R}_0\|^2] e^{-(2\eta\rho-2L-3)t} + \frac{\tilde{c}^2}{2\eta\rho-2L-3} (1 - e^{-(2\eta\rho-2L-3)t})$$

By taking a sufficiently large η such that $2\eta\rho-2L-3 > \tilde{c}^2 > 0$, we have

$$\mathbb{E}[\|\mathbf{R}_t\|^2] \leq \mathbb{E}[\|\mathbf{R}_0\|^2] + 1$$

Note that $\mathbf{X}_0^{\leftarrow} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$, which implies that $\mathbb{E}[\|\mathbf{X}_0^{\leftarrow}\|^2] = D$. Therefore,

$$\mathbb{E}[\|\mathbf{R}_0\|^2] \leq \mathbb{E}[\|\mathbf{X}_0^{\leftarrow}\|^2] + \|\mathbf{x}_*\|^2 \leq D + \|\mathbf{x}_*\|^2,$$

and

$$\mathbb{E}[\|\mathbf{X}_t^{\leftarrow}\|^2] \leq \mathbb{E}[\|\mathbf{R}_t\|^2] + \|\mathbf{x}_*\|^2 \leq D + 2\|\mathbf{x}_*\|^2 + 1 =: M_3.$$

By substituting M_3 into (37), we obtain

$$\frac{d}{dt}\mathbb{E}[f(\mathbf{X}_t^{\leftarrow})] \leq -\rho(\eta-2-L)(\mathbb{E}[f(\mathbf{X}_t^{\leftarrow})] - f(\mathbf{x}_*)) + M_4$$

for $M_4 := ((1+L)M_3 + C^2)/2$. Then, by Grönwall's Inequality,

$$\mathbb{E}[f(\mathbf{X}_T^{\leftarrow})] - f(\mathbf{x}_*) \leq (\mathbb{E}[f(\mathbf{X}_0^{\leftarrow})] - f(\mathbf{x}_*)) e^{-\rho(\eta-2-L)T} + \frac{M_4}{\rho(\eta-2-L)},$$

which means that

$$\mathbb{E}[f(\mathbf{X}_T^{\leftarrow})] - f(\mathbf{x}_*) = \mathcal{O}\left(e^{-\eta} + \frac{1}{\eta}\right). \quad \square$$

D.4 Proof of Theorem 11

Proof of Theorem 11. The proof consists two main steps:

(i) Let $Q_1 = A_1 A_1^\top$. For any coupling $(\tilde{\mathbf{X}}, \mathbf{X}) \sim (\tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1))$, we have

$$\begin{aligned} \mathcal{W}_1(\tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1)) &\leq \mathbb{E}[\|\tilde{\mathbf{X}} - \mathbf{X}\|] \\ &= \mathbb{E}[\|Q_1 \tilde{\mathbf{X}} - Q_1 \mathbf{X}\|] + \mathbb{E}[\|P_1 \tilde{\mathbf{X}} - P_1 \mathbf{X}\|] \\ &= \mathbb{E}[\|Q_1 \tilde{\mathbf{X}} - \mathbf{X}\|] + \mathbb{E}[\|\tilde{\mathbf{Y}}_{T-\delta}\|], \end{aligned}$$

where the final equality holds because $\mathbf{X} \sim \mathbb{P}_{X|Y}(\cdot | Y = 1)$ implies that $Q_1 \mathbf{X} = \mathbf{X}$, and $\tilde{\mathbf{X}} \sim \tilde{p}_{T-\delta}$ implies that $P_1 \tilde{\mathbf{X}} = \tilde{\mathbf{Y}}_{T-\delta}$. And by (35),

$$\mathbb{E}[\|\tilde{\mathbf{Y}}_{T-\delta}\|] \leq M_\eta(T - \delta) = \mathcal{O}(e^{-T} + \eta^{-1}).$$

Let $(Q_1 \tilde{\mathbf{X}}, \mathbf{X})$ be chosen as the optimal coupling for $((Q_1)_\# \tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1))$, i.e.,

$$\mathcal{W}_1((Q_1)_\# \tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1)) = \mathbb{E}[\|Q_1 \tilde{\mathbf{X}} - \mathbf{X}\|].$$

Therefore, we have

$$\mathcal{W}_1(\tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1)) \leq \mathcal{W}_1((Q_1)_\# \tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1)) + \mathcal{O}(e^{-T} + \eta^{-1}). \quad (42)$$

(ii) For $\mathcal{W}_1((Q_1)_\# \tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1))$, by the triangular inequality,

$$\begin{aligned} \mathcal{W}_1((Q_1)_\# \tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1)) &\leq \mathcal{W}_1((Q_1)_\# \tilde{p}_{T-\delta}, (Q_1)_\# \mathbb{P}_X) \\ &\quad + \mathcal{W}_1((Q_1)_\# \mathbb{P}_X, \mathbb{P}_{X|Y}(\cdot | Y = 1)) \\ &\leq \mathcal{W}_1(\tilde{p}_{T-\delta}, \mathbb{P}_X) + \mathcal{W}_1((Q_1)_\# \mathbb{P}_X, \mathbb{P}_{X|Y}(\cdot | Y = 1)), \end{aligned} \quad (43)$$

where the final inequality is because Q_1 is an orthogonal projection (Lemma H.13).

By Lemma D.3, the second term in above inequality is bounded by

$$\mathcal{W}_1((Q_1)_\# \mathbb{P}_X, \mathbb{P}_{X|Y}(\cdot | Y = 1)) \leq \tilde{C}_1 \quad (44)$$

for some constant \tilde{C}_1 . For the first term, it can be divided into

$$\mathcal{W}_1(\tilde{p}_{T-\delta}, \mathbb{P}_X) \leq \mathcal{W}_1(\tilde{p}_{T-\delta}, \hat{p}_\delta) + \mathcal{W}_1(\hat{p}_\delta, p_\delta^\sigma) + \mathcal{W}_1(p_\delta^\sigma, \mathbb{P}_X^\sigma) + \mathcal{W}_1(\mathbb{P}_X^\sigma, \mathbb{P}_X), \quad (45)$$

where \hat{p}_t is defined in dynamics (51), p_t^σ is the density evolving in the DDPM initialized from \mathbb{P}_X^σ ; see (13), and \mathbb{P}_X^σ is defined in Proposition 5. For the four terms in (45):

(a) By Proposition D.5 and $m_I > 1$ (Corollary 9),

$$\mathcal{W}_1(\tilde{p}_{T-\delta}, \hat{p}_\delta^\sigma) \leq \mathcal{O}(e^{-T} + \eta^{-1}) + \tilde{C}_2 \quad (46)$$

for some constant \tilde{C}_2 .

(b) By Proposition D.4,

$$\mathcal{W}_1(\hat{p}_\delta, p_\delta^\sigma) \leq \mathcal{O}(e^{-T}). \quad (47)$$

(c) Note that

$$\mathbf{X}_\delta^\sigma = \sqrt{\alpha_\delta} \mathbf{A} \mathbf{Z} + \sqrt{1 - \alpha_\delta} \boldsymbol{\xi} \sim p_\delta^\sigma, \quad \alpha_\delta = e^{-2\delta}$$

for $\mathbf{Z} \sim p_\sigma^Z$. Moreover, $\mathbf{A} \mathbf{Z} \sim \mathbb{P}_X^\sigma$. Therefore,

$$\begin{aligned} \mathcal{W}_1(p_\delta^\sigma, \mathbb{P}_X^\sigma) &\leq \mathbb{E}[\|\mathbf{X}_\delta^\sigma - \mathbf{A} \mathbf{Z}\|] \\ &\leq \mathbb{E}[\|\mathbf{X}_\delta^\sigma - \sqrt{\alpha_\delta} \mathbf{A} \mathbf{Z}\|] + (1 - \sqrt{\alpha_\delta}) \mathbb{E}[\|\mathbf{A} \mathbf{Z}\|] \\ &= \sqrt{1 - \alpha_\delta} \mathbb{E}[\|\boldsymbol{\xi}\|] + (1 - \sqrt{\alpha_\delta}) \mathbb{E}_{\mathbf{Z} \sim p_\sigma^Z}[\|\mathbf{Z}\|] \\ &\leq \sqrt{2\delta D} + \delta \mathbf{m}_\sigma^Z, \end{aligned}$$

where $\mathbf{m}_\sigma^Z = \mathbb{E}_{\mathbf{Z} \sim p_\sigma^Z}[\|\mathbf{Z}\|] < \infty$ by Lemma D.2. It follows that

$$\mathcal{W}_1(p_\delta^\sigma, \mathbb{P}_X^\sigma) \leq \mathcal{O}(\delta^{1/2}). \quad (48)$$

(d) By Proposition 5,

$$\mathcal{W}_1(\mathbb{P}_X^\sigma, \mathbb{P}_X) \leq \mathcal{O}(\sigma). \quad (49)$$

Then, combining (46) (47) (48) (49) with (45), we have

$$\mathcal{W}_1(\tilde{p}_{T-\delta}, \mathbb{P}_X) \leq \mathcal{O}(e^{-T} + \delta^{1/2} + \sigma + \eta^{-1}) + \tilde{C}_2. \quad (50)$$

Combining (50) (44) with (43), it follows

$$\mathcal{W}_1((Q_1)_\# \tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1)) \leq \mathcal{O}(e^{-T} + \delta^{1/2} + \sigma + \eta^{-1}) + \tilde{C},$$

where $\tilde{C} = \tilde{C}_1 + \tilde{C}_2$. Therefore, substituting this in (42), we obtain

$$\mathcal{W}_1(\tilde{p}_{T-\delta}, \mathbb{P}_{X|Y}(\cdot | Y = 1)) \leq \mathcal{O}(e^{-T} + \delta^{1/2} + \sigma + \eta^{-1}) + \tilde{C}. \quad \square$$

Remark D.1. For the error floor \tilde{C} , we provide two further discussions.

(i) First, it follows from the above proof that $\tilde{C} = \tilde{C}_1 + \tilde{C}_2$, where

- \tilde{C}_1 is determined by (44) and Lemma D.3,

$$\tilde{C}_1 = w_2 \mathbf{m}_1^Z, \quad \mathbf{m}_1^Z := \mathbf{m}_1^Z = \mathbb{E}_{\mathbf{Z} \sim \mathbb{P}_1^Z}[\|\mathbf{Z}\|],$$

which is independent of the parameters T, δ, σ .

- \tilde{C}_2 is given by (46) and Proposition D.5,

$$\tilde{C}_2 = \frac{M_2}{m_I - 1},$$

where M_2 is defined in (34) and depends on $L_S = \sup_{t \in [0, T-\delta]} L_t$ and T , while $m_I = \inf_{t \in [0, T-\delta]} m_t$. Since L_t and m_t are specified by Theorem 8 through p_t^σ , \tilde{C}_2 depends implicitly on T, δ , and σ .

(ii) We believe the error floor is inherent to the geometric guidance model. Because of the analytical simplicity of the geometric guidance, it cannot provide as much information as the probability guidance term did. More precisely, in Appendix C.3, we show that

$$\|\nabla_{\mathbf{x}} \log p_t(y = 1 | \mathbf{x}) + \eta_t P_1 \mathbf{x}\| \leq \beta_t, \quad \forall \mathbf{x} \in \mathcal{M}_1^t,$$

for some scalar $\eta_t > 0$, and $\beta_t = \mathcal{O}(\varepsilon_t)$, when $p_t(y = 1 | \mathbf{x}) > 1 - \varepsilon_t$ for all $\mathbf{x} \in \mathcal{M}_1^t$. This shows that the probabilistic guidance $\nabla_{\mathbf{x}} \log p_t(y = 1 | \mathbf{x})$ is “almost parallel” to the geometric guidance $P_1 \mathbf{x}$, but the norm of the probabilistic guidance carries additional information that the geometric term cannot capture. This is a trade-off made for the sake of analytical tractability.

Lemma D.2. Let $\mathbf{Z}_i \sim p_i^Z$ for $i = 1, 2$. If $\mathbf{m}_i^Z = \mathbb{E}[\|\mathbf{Z}_i\|] < \infty$, then for p_σ^Z defined in (11) (12),

$$\mathbf{m}_\sigma^Z := \mathbb{E}_{\mathbf{Z} \sim p_\sigma^Z}[\|\mathbf{Z}\|] < \infty.$$

Proof. By the definition of (12),

$$\mathbb{E}[\|\mathbf{Z}_{i,\sigma}\|] \leq \mathbb{E}[\|\mathbf{Z}_i\|] + \sigma \mathbb{E}[\|\zeta_i\|] \leq \mathbf{m}_i^Z + \sigma \sqrt{d} < \infty$$

for $\mathbf{Z}_{i,\sigma} \sim p_{i,\sigma}^Z$, where the second inequality is by Lemma H.8. Then, by (11),

$$\begin{aligned} \mathbb{E}_{\mathbf{Z} \sim p_\sigma^Z}[\|\mathbf{Z}\|] &= \int_{\mathbb{R}^d} \mathbf{z} p_\sigma^Z(\mathbf{z}) d\mathbf{z} \\ &= w_1 \int_{\mathbb{R}^d} \mathbf{z} p_{1,\sigma}^Z(\mathbf{z}) d\mathbf{z} + w_2 \int_{\mathbb{R}^d} \mathbf{z} p_{2,\sigma}^Z(\mathbf{z}) d\mathbf{z} \\ &= w_1 \mathbb{E}[\|\mathbf{Z}_{1,\sigma}\|] + w_2 \mathbb{E}[\|\mathbf{Z}_{2,\sigma}\|] < \infty. \end{aligned} \quad \square$$

Lemma D.3. *For*

$$\mathbb{P}_X = w_1 \mathbb{P}_{X|Y}(\cdot | Y = 1) + w_2 \mathbb{P}_{X|Y}(\cdot | Y = 2)$$

under Assumption I,

$$\mathcal{W}_1((Q_1)_\# \mathbb{P}_X, \mathbb{P}_{X|Y}(\cdot | Y = 1)) \leq w_2 \mathbf{m}_1^Z,$$

where $Q_1 = A_1 A_1^\top$ and $\mathbf{m}_1^Z = \mathbb{E}_{\mathbf{Z} \sim \mathbb{P}_1^Z} [\|\mathbf{Z}\|]$.

Proof. First, by Lemma H.7,

$$(Q_1)_\# \mathbb{P}_X = (Q_1)_\# \mathbb{P}_{X|Y}(\cdot | Y = 1) + (Q_1)_\# \mathbb{P}_{X|Y}(\cdot | Y = 2)$$

For the two terms, if $\mathbf{X} \sim \mathbb{P}_{X|Y}(\cdot | Y = 1)$, then $Q_1 \mathbf{X} = \mathbf{X}$, which implies that

$$(Q_1)_\# \mathbb{P}_{X|Y}(\cdot | Y = 1) = \mathbb{P}_{X|Y}(\cdot | Y = 1)$$

On the other hand, $\mathbf{X} \sim \mathbb{P}_{X|Y}(\cdot | Y = 2)$ implies that $Q_1 \mathbf{X} = 0$ so that

$$(Q_1)_\# \mathbb{P}_{X|Y}(\cdot | Y = 2) = \delta_0,$$

the Dirichlet measure at 0. Therefore, by Lemma H.9,

$$\mathcal{W}_1((Q_1)_\# \mathbb{P}_X, \mathbb{P}_{X|Y}(\cdot | Y = 1)) \leq w_2 \mathcal{W}_1(\delta_0, \mathbb{P}_{X|Y}(\cdot | Y = 1)).$$

For any coupling $(\mathbf{D}, \mathbf{X}) \sim (\delta_0, \mathbb{P}_{X|Y}(\cdot | Y = 1))$,

$$\begin{aligned} \mathcal{W}_1(\delta_0, \mathbb{P}_{X|Y}(\cdot | Y = 1)) &\leq \mathbb{E} [\|\mathbf{D} - \mathbf{X}\|] \\ &= \mathbb{E} [\|\mathbf{X}\|] = \mathbb{E}_{\mathbf{Z} \sim \mathbb{P}_1^Z} [\|A_1 \mathbf{Z}\|] = \mathbb{E}_{\mathbf{Z} \sim \mathbb{P}_1^Z} [\|\mathbf{Z}\|], \end{aligned}$$

where the last two equalities are because $\mathbb{P}_{X|Y}(\cdot | Y = 1) = (A_1)_\# \mathbb{P}_1^Z$ and $A_1 \in \mathcal{O}^{D \times d_1}$ by Assumption I. \square

In the following, unless otherwise specified, we assume that Assumptions I, II, III, and IV hold.

Proposition D.4. *Let p_t^σ be defined in (13). Consider the following two dynamics:*

$$\frac{d\hat{\mathbf{X}}_t}{dt} = \hat{\mathbf{X}}_t + \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\hat{\mathbf{X}}_t), \quad \hat{\mathbf{X}}_0 \sim \mathcal{N}(0, \mathbf{I}_D) \quad (51)$$

with the notation $\hat{\mathbf{X}}_t \sim \hat{p}_{T-t}^\sigma$, and

$$\frac{d\bar{\mathbf{X}}_t}{dt} = \bar{\mathbf{X}}_t + \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\bar{\mathbf{X}}_t), \quad \bar{\mathbf{X}}_0 \sim p_T^\sigma,$$

where note that $\bar{\mathbf{X}}_t \sim p_{T-t}^\sigma$. For $\delta > 0$, we

$$\mathcal{W}_1(\hat{p}_\delta^\sigma, p_\delta^\sigma) \leq e^{-m_I(T-\delta)} \left(\mathbf{m}_\sigma^Z + \sqrt{D} \right),$$

where $\mathbf{m}_\sigma^Z = \mathbb{E}_{\mathbf{Z} \sim p_\sigma^Z} [\|\mathbf{Z}\|]$ and $m_I = \inf_{t \in [\delta, T]} m_t$ is defined in Theorem 8.

Proof. First, by the Theorem 8, p_{T-t}^σ is m_I -strong log-concavity for $t \in [0, T - \delta]$, which follows that

$$\begin{aligned} \left\langle \hat{\mathbf{X}}_t - \bar{\mathbf{X}}_t, \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\hat{\mathbf{X}}_t) - \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\bar{\mathbf{X}}_t) \right\rangle &= \left\langle \hat{\mathbf{X}}_t - \bar{\mathbf{X}}_t, \nabla_{\mathbf{x}}^2 \log p_{T-t}^\sigma(\mathbf{x}) \left(\hat{\mathbf{X}}_t - \bar{\mathbf{X}}_t \right) \right\rangle \\ &\leq -m_I \|\hat{\mathbf{X}}_t - \bar{\mathbf{X}}_t\|^2. \end{aligned}$$

Therefore, we have

$$\frac{d}{dt} \|\hat{\mathbf{X}}_t - \bar{\mathbf{X}}_t\|^2 = 2 \left\langle \hat{\mathbf{X}}_t - \bar{\mathbf{X}}_t, \frac{d}{dt} \left(\hat{\mathbf{X}}_t - \bar{\mathbf{X}}_t \right) \right\rangle$$

$$\begin{aligned}
&= 2\|\hat{\mathbf{X}}_t - \bar{\mathbf{X}}_t\|^2 + 2\left\langle \hat{\mathbf{X}}_t - \bar{\mathbf{X}}_t, \nabla_{\mathbf{x}} \log p_{T-t}^{\sigma}(\hat{\mathbf{X}}_t) - \nabla_{\mathbf{x}} \log p_{T-t}^{\sigma}(\bar{\mathbf{X}}_t) \right\rangle \\
&\leq -2(m_I - 1)\|\hat{\mathbf{X}}_t - \bar{\mathbf{X}}_t\|^2,
\end{aligned}$$

which indicates that

$$\frac{d}{dt}\|\hat{\mathbf{X}}_t - \bar{\mathbf{X}}_t\| \leq -(m_I - 1)\|\hat{\mathbf{X}}_t - \bar{\mathbf{X}}_t\|,$$

Then, by Grönwall's Inequality (Lemma H.11),

$$\|\hat{\mathbf{X}}_{T-\delta} - \bar{\mathbf{X}}_{T-\delta}\| \leq e^{-(m_I-1)(T-\delta)}\|\hat{\mathbf{X}}_0 - \bar{\mathbf{X}}_0\|.$$

Therefore, by the definition of Wasserstein distance,

$$\begin{aligned}
\mathcal{W}_1(\hat{p}_{\delta}^{\sigma}, p_{\delta}^{\sigma}) &\leq \mathbb{E} \left[\|\hat{\mathbf{X}}_{T-\delta} - \bar{\mathbf{X}}_{T-\delta}\| \right] \\
&\leq e^{-(m_I-1)(T-\delta)} \mathbb{E} \left[\|\hat{\mathbf{X}}_0 - \bar{\mathbf{X}}_0\| \right].
\end{aligned}$$

By choosing $(\hat{\mathbf{X}}_0, \bar{\mathbf{X}}_0)$ as the optimal coupling, we obtain that

$$\mathcal{W}_1(\hat{p}_{\delta}^{\sigma}, p_{\delta}^{\sigma}) \leq e^{-(m_I-1)(T-\delta)} \mathcal{W}_1(\mathcal{N}(0, \mathbf{I}_D), p_T^{\sigma}). \quad (52)$$

For the right hand side of (52), by the definition of p_t^{σ} in Equation (13), i.e.,

$$\mathbf{X}_t^{\sigma} = \sqrt{\alpha_t} A \mathbf{Z} + \sqrt{1 - \alpha_t} \boldsymbol{\xi} \sim p_t^{\sigma}$$

for $\mathbf{Z} \sim p_{\sigma}^Z$, $\boldsymbol{\xi} \sim \mathcal{N}(0, \mathbf{I}_D)$, and $\alpha_t = e^{-2t}$, we have

$$\mathcal{W}_1(p_t^{\sigma}, \mathcal{N}(0, (1 - \alpha_t)\mathbf{I}_D)) \leq \sqrt{\alpha_t} \mathbb{E}[\|\mathbf{A}\mathbf{Z}\|] = e^{-t} \mathbb{E}_{\mathbf{Z} \sim p_{\sigma}^Z}[\|\mathbf{Z}\|].$$

Moreover,

$$\mathcal{W}_1(\mathcal{N}(0, (1 - \alpha_T)\mathbf{I}_D), \mathcal{N}(0, \mathbf{I}_D)) \leq (1 - \sqrt{1 - \alpha_T}) \mathbb{E}[\|\boldsymbol{\xi}\|] \leq e^{-T} \sqrt{D}.$$

Therefore,

$$\begin{aligned}
\mathcal{W}_1(p_T^{\sigma}, \mathcal{N}(0, \mathbf{I}_D)) &\leq \mathcal{W}_1(p_T^{\sigma}, \mathcal{N}(0, (1 - \alpha_T)\mathbf{I}_D)) + \mathcal{W}_1(\mathcal{N}(0, (1 - \alpha_T)\mathbf{I}_D), \mathcal{N}(0, \mathbf{I}_D)) \\
&\leq e^{-T} \left(\mathfrak{m}_{\sigma}^Z + \sqrt{D} \right).
\end{aligned}$$

Substituting this in the inequality (52) implies that

$$\mathcal{W}_1(\hat{p}_{\delta}^{\sigma}, p_{\delta}^{\sigma}) \leq e^{-m_I(T-\delta)-\delta} \left(\mathfrak{m}_{\sigma}^Z + \sqrt{D} \right) \leq e^{-m_I(T-\delta)} \left(\mathfrak{m}_{\sigma}^Z + \sqrt{D} \right). \quad \square$$

Proposition D.5. Consider the geometric guidance model (*) and the dynamics (51), for the corresponding generated distribution \tilde{p}_t^{σ} and \hat{p}_t^{σ} , we have

$$\mathcal{W}_1(\tilde{p}_{T-\delta}^{\sigma}, \hat{p}_{\delta}^{\sigma}) \leq \frac{\eta \sqrt{D - d_1}}{\eta - m_I} e^{-(m_I-1)(T-\delta)} + \frac{\eta M_2}{(m_I - 1)(\eta - 1)},$$

where M_2 is the constant defined in (34).

Proof. By the m_I -strong log-concavity of p_t^{σ} (Theorem 8), we have

$$\begin{aligned}
\frac{1}{2} \frac{d}{dt} \|\hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t\|^2 &= \left\langle \hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t, \frac{d}{dt} (\hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t) \right\rangle \\
&= \left\langle \hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t, \hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t + \nabla_{\mathbf{x}} \log p_{T-t}^{\sigma}(\hat{\mathbf{X}}_t) - \nabla_{\mathbf{x}} \log p_{T-t}^{\sigma}(\tilde{\mathbf{X}}_t) + \eta P_1 \tilde{\mathbf{X}}_t \right\rangle \\
&\leq -(m_I - 1) \|\hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t\|^2 + \eta \|\hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t\| \|P_1 \tilde{\mathbf{X}}_t\|
\end{aligned}$$

Note that $P_1 \tilde{\mathbf{X}}_t = \tilde{\mathbf{Y}}_t$. It follows that

$$\frac{d}{dt} \|\hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t\| \leq -(m_I - 1) \|\hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t\| + \eta \|\tilde{\mathbf{Y}}_t\|.$$

Moreover, by (35), taking the expectation on the both sides yields

$$\frac{d}{dt} \mathbb{E} [\|\hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t\|] \leq -(m_I - 1) \mathbb{E} [\|\hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t\|] + \eta M_\eta(t).$$

Then, Grönwall's inequality implies that

$$\begin{aligned} \mathcal{W}_1(\tilde{p}_{T-\delta}, \hat{p}_\delta^\sigma) &\leq \mathbb{E} [\|\hat{\mathbf{X}}_{T-\delta} - \tilde{\mathbf{X}}_{T-\delta}\|] \\ &\leq \eta \int_0^{T-\delta} M_\eta(s) e^{-(m_I-1)(T-\delta-s)} ds =: I(\eta), \end{aligned}$$

when the initial coupling is chosen as $\hat{\mathbf{X}}_0 = \tilde{\mathbf{X}}_0 \sim \mathcal{N}(0, \mathbf{I}_D)$. For $I(\eta)$, by the definition of $M_\eta(t)$ in (35), we have

$$\begin{aligned} I(\eta) &\leq \eta \int_0^{T-\delta} \left(\sqrt{D-d_1} e^{-(\eta-1)s} + \frac{M_2}{\eta-1} \right) e^{-(m_I-1)(T-\delta-s)} ds \\ &= \frac{\eta \sqrt{D-d_1}}{\eta - m_I} e^{-(m_I-1)(T-\delta)} \left(1 - e^{-(\eta-m_I)(T-\delta)} \right) \\ &\quad + \frac{\eta M_2}{(m_I-1)(\eta-1)} \left(1 - e^{-(m_I-1)(T-\delta)} \right) \\ &\leq \frac{\eta \sqrt{D-d_1}}{\eta - m_I} e^{-(m_I-1)(T-\delta)} + \frac{\eta M_2}{(m_I-1)(\eta-1)}. \end{aligned} \quad \square$$

D.5 Discretization Error

To clarify why performance degrades in practice when η becomes too large, we analyze the discretization error of the geometric guidance model (*). In practice, ODEs are typically solved using the Euler method, while SDEs are solved using the Euler–Maruyama (EM) scheme. Since our model is formulated as a deterministic ODE in (*), we focus on the Euler approximation; the analysis for the corresponding SDE and the EM scheme is analogous.

More specifically, we partition the interval $[0, T - \delta]$ into N subintervals with step size $h = (T - \delta)/N$, and define $t_k = kh$ for $k = 0, 1, \dots, N$. The Euler scheme then constructs the sequence $\{\mathbf{X}_k^h\}_{k=0}^N$ via

$$\mathbf{X}_{k+1}^h = \mathbf{X}_k^h + h \left(\mathbf{X}_k^h + \nabla_{\mathbf{x}} \log p_{T-t_k}^\sigma(\mathbf{X}_k^h) - \eta P_1 \mathbf{X}_k^h \right), \quad \mathbf{X}_0^h \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D).$$

Let $\mathbf{X}_k^h \sim \tilde{p}_k^h$. Our goal is to bound the Wasserstein error $\mathcal{W}_1(\tilde{p}_{T-\delta}, \tilde{p}_N^h)$. Under the Lipschitz continuity of $\nabla_{\mathbf{x}} \log p_t^\sigma$, standard results yield $\mathcal{W}_1(\tilde{p}_{T-\delta}, \tilde{p}_N^h) \leq \mathcal{O}(h e^\eta)$ (Griffiths & Higham, 2010, Theorem 2.4). Because of Theorem 8, we not only have the L_S -smoothness

$$\|\nabla_{\mathbf{x}}^2 \log p_t^\sigma(\mathbf{x})\|_{\text{op}} \leq L_S, \quad L_S = \sup_{t \in [\delta, T]} L_t,$$

but also the m_I -strong log-concavity

$$-\nabla_{\mathbf{x}}^2 \log p_t^\sigma(\mathbf{x}) \succeq m_I \mathbf{I}_D, \quad m_I = \inf_{t \in [\delta, T]} m_t.$$

The additional strong log-concavity yields the improved bound

$$\mathcal{W}_1(\tilde{p}_{T-\delta}, \tilde{p}_N^h) \leq \mathcal{O}(h \eta^2).$$

Theorem D.6. Assume that $t \mapsto \nabla_{\mathbf{x}} \log p_t^\sigma(\mathbf{x})$ is C^1 for each \mathbf{x} , and there exist $A, B \geq 0$ such that

$$\|\partial_t \nabla_{\mathbf{x}} \log p_t^\sigma(\mathbf{x})\| \leq A + B\|\mathbf{x}\|.$$

For $\eta > 2$, if $h(\eta - 1) < 1$, we have

$$\mathcal{W}_1(\tilde{p}_{T-\delta}, \tilde{p}_N^h) \leq \mathcal{O}(h\eta^2).$$

Proof. Let

$$b(t, \mathbf{x}) = \mathbf{x} + \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\mathbf{x}) - \eta P_1 \mathbf{x}.$$

Let $\Phi_k(\mathbf{x}) =: \mathbf{x}_{t_{k+1}}$, where \mathbf{x}_t is the solution of ODE

$$\frac{d\mathbf{x}_t}{dt} = b(t, \mathbf{x}_t), \quad t \in [t_k, t_{k+1}], \quad (53)$$

with initial value $\mathbf{x}_{t_k} = \mathbf{x}$. By (*), we can see

$$\tilde{\mathbf{X}}_{t_{k+1}} = \Phi_k(\tilde{\mathbf{X}}_{t_k})$$

Moreover, define the Euler one-step map

$$\Psi_k(\mathbf{x}) = \mathbf{x} + hb(t_k, \mathbf{x}),$$

so that Euler scheme is

$$\mathbf{X}_{k+1}^h = \Psi_k(\mathbf{X}_k^h).$$

Therefore,

$$\begin{aligned} \mathbf{e}_{k+1} &:= \tilde{\mathbf{X}}_{t_{k+1}} - \mathbf{X}_{k+1}^h = \Phi_k(\tilde{\mathbf{X}}_{t_k}) - \Psi_k(\mathbf{X}_k^h) \\ &= (\Phi_k(\tilde{\mathbf{X}}_{t_k}) - \Phi_k(\mathbf{X}_k^h)) + (\Phi_k(\mathbf{X}_k^h) - \Psi_k(\mathbf{X}_k^h)). \end{aligned}$$

Next, we analyze these two terms respectively.

(i) By the m_I -strong log-concavity, we have

$$\langle \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq -m_I \|\mathbf{x} - \mathbf{y}\|^2.$$

Moreover, since P_1 is an orthogonal projection,

$$\langle (\mathbf{I}_D - \eta P_1)(\mathbf{x} - \mathbf{y}), \mathbf{x} - \mathbf{y} \rangle = \|\mathbf{x} - \mathbf{y}\|^2 - \eta \|P_1(\mathbf{x} - \mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2.$$

Therefore,

$$\langle b(t, \mathbf{x}) - b(t, \mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq -(m_I - 1) \|\mathbf{x} - \mathbf{y}\|^2$$

Let $\mathbf{x}_t, \mathbf{y}_t$ be the solution of (53) with the initial value $\mathbf{x}_{t_k} = \mathbf{x}$ and $\mathbf{y}_{t_k} = \mathbf{y}$. So we have

$$\frac{d}{dt} \|\mathbf{x}_t - \mathbf{y}_t\|^2 = 2 \langle b(t, \mathbf{x}_t) - b(t, \mathbf{y}_t), \mathbf{x}_t - \mathbf{y}_t \rangle \leq -2(m_I - 1) \|\mathbf{x}_t - \mathbf{y}_t\|^2.$$

Then by the Grönwall's Inequality (Lemma H.11),

$$\|\mathbf{x}_{t_{k+1}} - \mathbf{y}_{t_{k+1}}\| \leq e^{-(m_I - 1)h} \|\mathbf{x} - \mathbf{y}\|,$$

which implies that

$$\|\Phi_k(\mathbf{x}) - \Phi_k(\mathbf{y})\| \leq e^{-(m_I - 1)h} \|\mathbf{x} - \mathbf{y}\|. \quad (54)$$

(ii) Fix $\mathbf{x} \in \mathbb{R}^D$ and let \mathbf{x}_t be the solution of (53) with the initial value $\mathbf{x}_{t_k} = \mathbf{x}$. Note that by definition

$$\Phi_k(\mathbf{x}) = \mathbf{x} + \int_{t_k}^{t_{k+1}} b(t, \mathbf{x}_t) dt.$$

Therefore,

$$\begin{aligned} \Phi_k(\mathbf{x}) - \Psi_k(\mathbf{x}) &= \int_{t_k}^{t_{k+1}} (b(t, \mathbf{x}_t) - b(t_k, \mathbf{x})) dt \\ &= \int_{t_k}^{t_{k+1}} (b(t, \mathbf{x}_t) - b(t, \mathbf{x})) dt + \int_{t_k}^{t_{k+1}} (b(t, \mathbf{x}) - b(t_k, \mathbf{x})) dt. \end{aligned}$$

For above two terms, we analyze them respectively.

(a) Since $\|\mathbf{I}_D - \eta P_1\|_{\text{op}} = \eta - 1$ ($\eta > 2$) and $\nabla_{\mathbf{x}} \log p_t$ is L_S -Lipschitz continuous, $b(t, \cdot)$ is $K(\eta)$ -Lipschitz continuous for $K(\eta) = L_S + \eta - 1$. So

$$\int_{t_k}^{t_{k+1}} \|b(t, \mathbf{x}_t) - b(t, \mathbf{x})\| ds \leq K(\eta) \int_{t_k}^{t_{k+1}} \|\mathbf{x}_t - \mathbf{x}\| ds$$

Note that

$$\|\mathbf{x}_t - \mathbf{x}\| = \left\| \int_{t_k}^t b(s, \mathbf{x}_s) ds \right\| \leq \int_{t_k}^t \|b(s, \mathbf{x}_s)\| ds \leq (t - t_k) \sup_{s \in [t_k, t_{k+1}]} \|b(s, \mathbf{x}_s)\|.$$

Therefore,

$$\int_{t_k}^{t_{k+1}} \|b(t, \mathbf{x}_t) - b(t, \mathbf{x})\| ds \leq \frac{h^2}{2} K(\eta) \sup_{s \in [t_k, t_{k+1}]} \|b(s, \mathbf{x}_s)\|. \quad (55)$$

For the right hand side, using same notation as (31), let

$$C = \sup_{t \in [\delta, T]} \|\nabla_{\mathbf{x}} \log p_t^\sigma(\mathbf{0})\| < \infty.$$

It implies that $\|b(t, \mathbf{0})\| \leq C$ and so

$$\|b(s, \mathbf{x}_s)\| \leq \|b(s, \mathbf{x}_s) - b(s, \mathbf{0})\| + \|b(s, \mathbf{0})\| \leq C + K(\eta) \|\mathbf{x}_s\|.$$

Let $S = \sup_{s \in [t_k, t_{k+1}]} \|\mathbf{x}_s\| < \infty$. Using the similar idea as in the proof of Theorem 10 in Appendix D.2,

$$S = \sup_{s \in [t_k, t_{k+1}]} \|\mathbf{x}_s\| \leq C_1 \|\mathbf{x}\| + C_2$$

where $C_1 = \exp((1 + L_S)h)$ and $C_2 = C(\exp((1 + L_S)h) - 1)/(1 + L_S)$ as shown in (33) and they are independent of η . So

$$\sup_{s \in [t_k, t_{k+1}]} \|b(s, \mathbf{x}_s)\| \leq C + K(\eta) \sup_{s \in [t_k, t_{k+1}]} \|\mathbf{x}_s\| \leq C_1 K(\eta) \|\mathbf{x}\| + C_2 K(\eta) + C. \quad (56)$$

Combining (55) and (56), we have

$$\int_{t_k}^{t_{k+1}} \|b(t, \mathbf{x}_t) - b(t, \mathbf{x})\| dt \leq \frac{h^2}{2} (C_1 K(\eta)^2 \|\mathbf{x}\| + C_2 K(\eta)^2 + C K(\eta)). \quad (57)$$

(b) Since $b(t, \mathbf{x}) - b(t_k, \mathbf{x}) = \nabla_{\mathbf{x}} \log p_{T-t}^\sigma(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{T-t_k}^\sigma(\mathbf{x})$,

$$\|b(t, \mathbf{x}) - b(t_k, \mathbf{x})\| \leq \int_{t_k}^t \|\partial_t \nabla_{\mathbf{x}} \log p_{T-s}^\sigma(\mathbf{x})\| ds \leq (t - t_k) (A + B \|\mathbf{x}\|).$$

It implies that

$$\int_{t_k}^{t_{k+1}} \|b(t, \mathbf{x}) - b(t_k, \mathbf{x})\| dt \leq \frac{h^2}{2} (A + B \|\mathbf{x}\|). \quad (58)$$

Therefore, combining (57) and (58),

$$\begin{aligned} \|\Phi_k(\mathbf{x}) - \Psi_k(\mathbf{x})\| &\leq \int_{t_k}^{t_{k+1}} \|b(t, \mathbf{x}_t) - b(t, \mathbf{x})\| dt + \int_{t_k}^{t_{k+1}} \|b(t, \mathbf{x}) - b(t_k, \mathbf{x})\| dt \\ &\leq \frac{h^2}{2} (A + C_2 K(\eta)^2 + C K(\eta) + (B + C_1 K(\eta)^2) \|\mathbf{x}\|). \end{aligned} \quad (59)$$

Combining (54) and (59), by setting $\mathbf{x} = \mathbf{X}_k^h$ and $\mathbf{y} = \tilde{\mathbf{X}}_{t_k}$,

$$\begin{aligned} \|\mathbf{e}_{k+1}\| &\leq \|\Phi_k(\tilde{\mathbf{X}}_{t_k}) - \Phi_k(\mathbf{X}_k^h)\| + \|\Phi_k(\mathbf{X}_k^h) - \Psi_k(\mathbf{X}_k^h)\| \\ &\leq e^{-(m_I-1)h} \|\mathbf{e}_k\| + \frac{h^2}{2} (A + C_2 K(\eta)^2 + C K(\eta) + (B + C_1 K(\eta)^2) \|\mathbf{X}_k^h\|). \end{aligned}$$

Let $a_k = \mathbb{E}[\|\mathbf{e}_k\|]$. By the following Lemma D.7, because $h(\eta - 1) < 1$, $\mathbb{E}[\|\mathbf{X}_k^h\|] \leq M_e$. Taking the expectation of above inequality, we have

$$a_{k+1} \leq e^{-(m_I-1)h} a_k + \frac{h^2}{2} (A + C_2 K(\eta)^2 + C K(\eta) + (B + C_1 K(\eta)^2) M_e).$$

Therefore, by coupling $\tilde{\mathbf{X}}_0 = \mathbf{X}_0^h$, i.e., $a_0 = 0$, we have

$$\begin{aligned} a_N &= \mathbb{E}[\|\tilde{\mathbf{X}}_{T-\delta} - \mathbf{X}_N^h\|] \leq \frac{h^2}{2} (A + C_2 K(\eta)^2 + C K(\eta) + (B + C_1 K(\eta)^2) M_e) \sum_{k=0}^{N-1} e^{-(m_I-1)hk} \\ &\leq \frac{h}{2} (A + C_2 K(\eta)^2 + C K(\eta) + (B + C_1 K(\eta)^2) M_e) \frac{e^{(m_I-1)h}}{m_I - 1}. \end{aligned}$$

It follows that as $h \rightarrow 0$ and $\eta \rightarrow \infty$,

$$\mathcal{W}_1(\tilde{p}_{T-\delta}, \tilde{p}_N^h) \leq \mathbb{E}[\|\tilde{\mathbf{X}}_{T-\delta} - \mathbf{X}_N^h\|] \leq \mathcal{O}(h\eta^2). \quad \square$$

Lemma D.7. For $\eta > 2$, if $h(\eta - 1) < 1$, then

$$\sup_k \mathbb{E}[\|\mathbf{X}_k^h\|] \leq M_e,$$

where M_e is independent of η .

Proof. First, by construction,

$$\mathbf{X}_{k+1}^h = (\mathbf{I}_D + h(\mathbf{I}_D - \eta P_1)) \mathbf{X}_k^h + h \nabla_{\mathbf{x}} \log p_{T-t_k}^\sigma(\mathbf{X}_k^h).$$

Let $M_h = \mathbf{I}_D + h(\mathbf{I}_D - \eta P_1)$. Then because P_1 is an orthogonal projection, there are only two eigenvalues of M_h : for $\mathbf{x} \in \ker P_1$, $M_h \mathbf{x} = (1 + h)\mathbf{x}$, and for $\mathbf{x} \in \text{Im } P_1$, $M_h \mathbf{x} = (1 + h(1 - \eta))\mathbf{x}$. Because $h(\eta - 1) < 1$, $1 + h(1 - \eta) \in [0, 1]$. So

$$\|M_h\|_{\text{op}} = 1 + h.$$

Similarly, as shown in (32),

$$\|\nabla_{\mathbf{x}} \log p_{T-t_k}^\sigma(\mathbf{X}_k^h)\| \leq L_S \|\mathbf{X}_k^h\| + C.$$

Therefore,

$$\begin{aligned} \|\mathbf{X}_{k+1}^h\| &\leq \|M_h\|_{\text{op}} \|\mathbf{X}_k^h\| + h \|\nabla_{\mathbf{x}} \log p_{T-t_k}^\sigma(\mathbf{X}_k^h)\| \\ &\leq (1 + h(1 + L_S)) \|\mathbf{X}_k^h\| + Ch. \end{aligned}$$

Taking expectations on the both sides, we have

$$\mathbb{E}[\|\mathbf{X}_k^h\|] \leq (1 + h(1 + L_S))^k \mathbb{E}[\|\mathbf{X}_0^h\|] + Ch \sum_{j=0}^{k-1} (1 + h(1 + L_S))^j$$

$$\leq e^{(1+L_S)t_k} (\mathbb{E} [\|\mathbf{X}_0^h\|] + Ct_k).$$

Because $\mathbf{X}_0^h \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$, $\mathbb{E} [\|\mathbf{X}_0^h\|] \leq \sqrt{\mathbb{E} [\|\mathbf{X}_0^h\|^2]} = \sqrt{D}$. Therefore, if let

$$M_e := e^{(1+L_S)(T-\delta)} (\sqrt{D} + C(T-\delta)),$$

which is independent of η , then

$$\mathbb{E} [\|\mathbf{X}_k^h\|] \leq M_e. \quad \square$$

E Analysis for Assumptions

E.1 More Details for Orthogonality Assumption

For Assumption I, consider the case where $A_1^\top A_2 \neq \mathbf{O}$, i.e., \mathcal{M}_1 is not orthogonal to \mathcal{M}_2 . In this case, $A = (A_1, A_2) \notin \mathcal{O}^{D \times d}$, meaning that $A^\top A \neq \mathbf{I}_d$ and AA^\top is no longer an orthogonal projection. We claim that this relaxation does not affect our analysis regarding the guidance scale η . Based on our results, it is necessary to examine its influence from three perspective: the smoothness and concavity of $\log p_t^\sigma$ (Section 5.1), the estimation of the target manifold \mathcal{M}_1 (Section 5.2), and the distance between generated and target distributions (Section 5.3).

- (a) Smoothness and Convexity: First, the results on the strong log-concavity of the latent density in Theorem 7 are independent of the orthogonality of A . Therefore, to analyze the smoothness and concavity of $\log p_t^\sigma$, it suffices to revisit the proof of Theorem 8. Note that

$$\mathbf{X}_t^\sigma = \sqrt{\alpha_t} A \mathbf{Z}_\sigma + \sqrt{1 - \alpha_t} \boldsymbol{\xi} \sim p_t^\sigma.$$

By Proposition 3, Corollary 4, and the m_0^Z -strong log-concavity of the latent density p_σ^Z (Theorem 7), we obtain the following bounds:

$$\|\nabla_{\mathbf{x}}^2 \log p_t^\sigma(\mathbf{x})\|_{\text{op}} \leq L_t^A, \quad L_t^A := \frac{\alpha_t(\Lambda_A + \lambda_A) + (1 - \alpha_t)m_0^Z}{(1 - \alpha_t)(\alpha_t\lambda_A + m_0^Z(1 - \alpha_t))},$$

and

$$-\nabla_{\mathbf{x}}^2 \log p_t^\sigma(\mathbf{x}) \succeq m_t^Z \mathbf{I}_D, \quad m_t^Z := \frac{(1 - \alpha_t)m_0^Z - \alpha_t(\Lambda_A - \lambda_A)}{(1 - \alpha_t)(\alpha_t\lambda_A + m_0^Z(1 - \alpha_t))},$$

where

$$\Lambda_A = \|A\|_{\text{op}}^2 = \lambda_{\max}(A^\top A), \quad \lambda_A = \lambda_{\min}(A^\top A).$$

Because $A = (A_1, A_2)$ and A_i are orthogonal,

$$A^\top A = \begin{pmatrix} \mathbf{I}_{d_1} & C \\ C^\top & \mathbf{I}_{d_2} \end{pmatrix} = \mathbf{I}_d + \begin{pmatrix} \mathbf{O} & C \\ C^\top & \mathbf{O} \end{pmatrix}, \quad C := A_1^\top A_2.$$

Let $\sigma_{\max}(C)$ be the maximal singular value of C . Then, we have

$$1 - \sigma_{\max}(C) \leq \lambda_A \leq \Lambda_A \leq 1 + \sigma_{\max}(C).$$

Moreover, because $\|C\|_{\text{op}} \leq \|A_1\|_{\text{op}} \|A_2\|_{\text{op}} = 1$, $\sigma_{\max}(C) \leq 1$, which implies that $0 \leq \lambda_A \leq \Lambda_A \leq 2$.

For smoothness, it is clear that $0 < L_t^A < \infty$, so the non-orthogonality of A does not affect the L -smoothness of $\log p_t^\sigma$, except that the constant changes from L_t to L_t^A . However, for strong log-concavity, it requires $m_t^A > 1$ (Corollary 9), which holds if

$$t > \frac{1}{2} \log \frac{m_0^Z - \Lambda_A}{m_0^Z - \lambda_A}, \quad (60)$$

under the condition $m_0^Z > 2 \geq \Lambda_A$. This requires a modification of Assumption III, $M \leq 2\sqrt{m-2}$, for the same reason discussed in the proof of Corollary 9.

- (b) Estimating Target Manifold: Since Theorem 10 depends only on the L -smoothness of $\log p_t^\sigma$, and the geometric guidance model $(*)$ does not involve A , the result of Theorem 10 remains valid even when A is not orthogonal.
- (c) Distance to Target Distribution: Because the condition $m_t^A > 1$ requires inequality (60), one can set

$$\delta > \frac{1}{2} \log \frac{m_0^z - \Lambda_A}{m_0^z - \lambda_A},$$

and consider the geometric guidance model $(*)$ on interval $[0, T - \delta]$. With this adjustment, the results in Theorem 11 still hold, up to changes in certain constants. For instance, the non-orthogonality of A changes the bound on $\mathcal{W}_1((Q_1) \# \mathbb{P}_X, \mathbb{P}_{X|Y}(\cdot | Y = 1))$, specifically the constant \tilde{C} in Theorem 11.

E.2 More details of Assumption III

In the following, we demonstrate a family of distributions that satisfy both Assumption II and Assumption III. Consider the density function p_i^Z of the distribution \mathbb{P}_i^Z , given by the form:

$$p_i^Z(\mathbf{z}) = e^{-V_i(\mathbf{z})} \chi_{K_i}(\mathbf{z}),$$

where $K_i \subset \mathbb{R}^{d_i}$ is a convex and compact set, and

$$\nabla_{\mathbf{z}}^2 V_i(\mathbf{z}) \succeq m \mathbf{I}_{d_i}.$$

In other words, p_i^Z belongs to the class of strongly log-concave densities supported on convex and compact subsets of \mathbb{R}^{d_i} .

First, for such p_i^Z , strong log-concavity on a convex set does not perfectly align with Assumption II, which induces a question of whether this property can substitute for Assumption II in deriving the strong log-concavity of the mixture latent density p_σ^Z defined in Equation (11).

In the proof of Theorem 7, the strong log-concavity of p_σ^Z is inherited from that of the component densities $p_{i,\sigma}^Z$ defined in Equation (12), which are shown to be strongly log-concave via Corollary 4, under Assumption II. In other words, the key question is whether strong log-concavity on a convex set suffices to replace the strong log-concavity condition in Proposition 3, and thereby still allow us to deduce the conclusion of Corollary 4.

Proposition E.1. *Let \mathbf{Z} be a random variable on \mathbb{R}^k with the density function p^Z given by*

$$p^Z(\mathbf{z}) = e^{-V(\mathbf{z})} \chi_K(\mathbf{z}),$$

where K is a convex set. Let $B \in \mathbb{R}^{n \times k}$. Assume there are $m_0, \Lambda > 0$ such that

$$\nabla_{\mathbf{z}}^2 V(\mathbf{z}) \succeq m_0 \mathbf{I}_k, \quad \|B\|_{\text{op}}^2 \leq \Lambda,$$

and $\lambda := \lambda_{\min}(B^\top B) \geq 0$. For $\alpha \in \mathbb{R}$ and $\beta > 0$, let

$$\mathbf{X} = \alpha B \mathbf{Z} + \beta \boldsymbol{\xi}, \quad \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$$

with the density function p_X on \mathbb{R}^n . We have

$$\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x}) \preceq \left(\frac{\alpha^2 \Lambda}{\beta^2 (\alpha^2 \lambda + m_0 \beta^2)} - \frac{1}{\beta^2} \right) \mathbf{I}_n.$$

Proof. By the same calculation, we have

$$\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x}) = \frac{\alpha^2}{\beta^4} B \text{Cov}_{\mu_x}(\mathbf{Z}) B^\top - \frac{1}{\beta^2} \mathbf{I}_n,$$

Note in this case,

$$d\mu_x(\mathbf{z}) = \frac{e^{-U_x(\mathbf{z})} \chi_K(\mathbf{z}) d\mathbf{z}}{\int_K e^{-U}(\mathbf{y}) d\mathbf{y}},$$

where

$$U_x(\mathbf{z}) = \frac{1}{2\sigma^2} \|\mathbf{x} - B\mathbf{z}\|^2 + V(\mathbf{z}).$$

It follows that

$$\nabla_{\mathbf{z}}^2 U_x(\mathbf{z}) = \frac{\alpha^2}{\beta^2} B^\top B + \nabla_{\mathbf{z}}^2 V(\mathbf{z}) \succeq m \mathbf{I}_k, \quad m := \frac{\alpha^2 \lambda}{\beta^2} + m_0.$$

Instead of applying Lemma H.6, by using the Brascamp–Lieb Inequality on a convex set (Bobkov & Ledoux, 2000, Proposition 2.1), we still have

$$\text{Var}_{\mu_x}(f) \leq \frac{1}{m} \mathbb{E}_{\mu_x} \left[\|\nabla f\|^2 \right],$$

for any C^1 function $f: \mathbb{R}^k \rightarrow \mathbb{R}$, which also indicates that

$$\|\text{Cov}_{\mu_x}(\mathbf{Z})\|_{\text{op}} \leq \frac{1}{m}.$$

Then, the following proof is as same as the proof in Proposition 3 and in Corollary 4 so that we have the same result

$$\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x}) \preceq \left(\frac{\alpha^2 \Lambda}{\beta^2 (\alpha^2 \lambda + m_0 \beta^2)} - \frac{1}{\beta^2} \right) \mathbf{I}_n. \quad \square$$

Therefore, Proposition E.1 shows that, in our settings, Assumption II can be replaced Assumption II':

Assumption II'. For $i = 1, 2$, \mathbb{P}_i^Z admits the density function p_i^Z that has the form $p_i^Z(\mathbf{z}) = e^{-V_i(\mathbf{z})} \chi_{K_i}(\mathbf{z})$ such that $K_i \subset \mathbb{R}^{d_i}$ is a convex and compact set, and

$$\nabla_{\mathbf{z}}^2 V_i(\mathbf{z}) \succeq m \mathbf{I}_{d_i}.$$

Next, we verify if p_i^Z in such class can satisfy Assumption III.

Proposition E.2. For $i = 1, 2$, let p_i^Z satisfy Assumption II', and let $p_{i,\sigma}^Z$ defined by Equation (12). Fix a $\sigma > 0$, we have

$$\sup_{\mathbf{x}} \|\nabla_{\mathbf{x}} \log p_{1,\sigma}^Z(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{2,\sigma}^Z(\mathbf{x})\| \leq \frac{\sqrt{|K_1|^2 + |K_2|^2}}{\sigma^2},$$

where $|K_i| = \sup \{\|\mathbf{z}\| : \mathbf{z} \in K_i\}$.

Proof. First, by the definition (12),

$$p_{1,\sigma}^Z(\mathbf{z}) = (2\pi\sigma^2)^{-\frac{d}{2}} \int_{K_1} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{z} - (\mathbf{z}_1, 0)^\top\|^2\right) p_1^Z(\mathbf{z}_1) d\mathbf{z}_1,$$

Therefore,

$$\nabla_{\mathbf{z}} \log p_{1,\sigma}^Z(\mathbf{z}) = \frac{\nabla_{\mathbf{z}} p_{1,\sigma}^Z(\mathbf{z})}{p_{1,\sigma}^Z(\mathbf{z})} = \frac{-\frac{1}{\sigma^2} \int_{K_1} (\mathbf{z} - (\mathbf{z}_1, 0)^\top) \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{z} - (\mathbf{z}_1, 0)^\top\|^2\right) p_1^Z(\mathbf{z}_1) d\mathbf{z}_1}{\int_{K_1} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{z} - (\mathbf{z}_1, 0)^\top\|^2\right) p_1^Z(\mathbf{z}_1) d\mathbf{z}_1}.$$

It follows that

$$\nabla_{\mathbf{z}} \log p_{1,\sigma}^Z(\mathbf{z}) = \frac{1}{\sigma^2} ((m_1(\mathbf{z}), 0)^\top - \mathbf{z}),$$

where

$$m_1(\mathbf{z}) = \frac{\int_{K_1} \mathbf{z}_1 \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{z} - (\mathbf{z}_1, 0)^\top\|^2\right) p_1^Z(\mathbf{z}_1) d\mathbf{z}_1}{\int_{K_1} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{z} - (\mathbf{z}_1, 0)^\top\|^2\right) p_1^Z(\mathbf{z}_1) d\mathbf{z}_1}.$$

Similarly,

$$\nabla_{\mathbf{z}} \log p_{2,\sigma}^Z(\mathbf{z}) = \frac{1}{\sigma^2} ((0, m_2(\mathbf{z}))^\top - \mathbf{z}),$$

for

$$m_2(\mathbf{z}) = \frac{\int_{K_2} \mathbf{z}_2 \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{z} - (\mathbf{z}_2, 0)^\top\|^2\right) p_2^Z(\mathbf{z}_2) d\mathbf{z}_2}{\int_{K_2} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{z} - (\mathbf{z}_2, 0)^\top\|^2\right) p_2^Z(\mathbf{z}_2) d\mathbf{z}_2}.$$

Therefore,

$$\|\nabla_{\mathbf{z}} \log p_{1,\sigma}^Z(\mathbf{z}) - \nabla_{\mathbf{z}} \log p_{2,\sigma}^Z(\mathbf{z})\| = \frac{1}{\sigma^2} \|(m_1(\mathbf{z}), m_2(\mathbf{z}))^\top\| = \frac{1}{\sigma^2} \sqrt{\|m_1(\mathbf{z})\|^2 + \|m_2(\mathbf{z})\|^2}.$$

Note that

$$m_i(\mathbf{z}) = \mathbb{E}_{\mathbf{Z} \sim \mu_z^i}[\mathbf{Z}], \quad d\mu_z^i(\mathbf{z}_i) := \frac{\exp\left(-\frac{1}{2\sigma^2} \|\mathbf{z} - (\mathbf{z}_i, 0)^\top\|^2\right) p_i^Z(\mathbf{z}_i) d\mathbf{z}_i}{\int_{K_i} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{z} - (\mathbf{z}_i, 0)^\top\|^2\right) p_i^Z(\mathbf{z}_i) d\mathbf{z}_i}.$$

By the convexity of K_i and Lemma H.10, $m_i(\mathbf{z}) \in K_i$. Then, the boundedness of K_i implies

$$\sup_{\mathbf{x}} \|\nabla_{\mathbf{x}} \log p_{1,\sigma}^Z(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{2,\sigma}^Z(\mathbf{x})\| \leq \frac{\sqrt{|K_1|^2 + |K_2|^2}}{\sigma^2}. \quad \square$$

Sufficient conditions for Assumption III. If p_i^Z belongs to the class of distributions

$$\left\{ e^{-V(\mathbf{z})} \chi_K(\mathbf{z}) : \nabla^2 V \succeq m\mathbf{I}, K \text{ is compact and convex.} \right\}, \quad (61)$$

then $p_{i,\sigma}^Z$ given Equation (12) is strongly log-concave by Proposition E.1. Moreover, if we choose σ such that

$$M \leq \frac{\sqrt{|K_1|^2 + |K_2|^2}}{\sigma^2} \leq 2\sqrt{m-1} \quad \Leftrightarrow \quad \sigma^2 \geq \sqrt{\frac{|K_1|^2 + |K_2|^2}{4(m-1)}}, \quad (62)$$

Proposition E.2 shows that Assumption III is satisfied. Then the mixture latent distribution p_σ^Z given by Equation (11) is m_0^Z -strongly log-concave provided by Theorem 7, which further implies that p_t^σ in the geometric guidance model (*) satisfies:

$$\|\nabla_{\mathbf{x}}^2 \log p_t^\sigma(\mathbf{x})\|_{\text{op}} \leq L_t, \quad -\nabla_{\mathbf{x}}^2 \log p_t^\sigma(\mathbf{x}) \succeq m_t \mathbf{I}_D.$$

F Lipschitz Continuity of Score Function

If we only focus on the Lipschitz continuity of the score function $\nabla_{\mathbf{x}} \log p_t$, where p_t is obtained by a DDPM initialized from a distribution whose latent distribution admits a smooth density function p^Z , then the conditions in Proposition 3 can be relaxed. We consider two cases below.

The first case aligns with the setting considered in De Bortoli (2022), where $\text{supp } p^Z$ is assumed to be compact. We provide an alternative proof for this case, motivated by the argument used in the proof of Proposition 3.

Proposition F.1. *Let \mathbf{Z} be a random variable on \mathbb{R}^k with the density function p^Z , and let $\phi: \mathbb{R}^k \rightarrow \mathbb{R}^n$ be continuous. Assume $\text{supp } p^Z$ is compact. For $\alpha \in \mathbb{R}$ and $\beta > 0$, let*

$$\mathbf{X} = \alpha\phi(\mathbf{Z}) + \beta\xi, \quad \xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n),$$

with the density function p_X on \mathbb{R}^n . We have

$$\|\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x})\|_{\text{op}} \leq \frac{1}{\beta^2} + \frac{\alpha^2 R^2}{\beta^4},$$

for some constant $R > 0$.

Proof. By the similar calculation as in the proof of Proposition 3,

$$\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x}) = \frac{\alpha^2}{\beta^4} \text{Cov}_{\mu_x}(\phi(\mathbf{Z})) - \frac{1}{\beta^2} \mathbf{I}_n,$$

which follows that

$$\|\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x})\|_{\text{op}} \leq \frac{1}{\beta^2} + \frac{\alpha^2}{\beta^4} \|\text{Cov}_{\mu_x}(\phi(\mathbf{Z}))\|_{\text{op}}.$$

To bound the second term, first by the definition of μ_x , $\text{supp } \mu_x = \text{supp } p^Z$. Because $\text{supp } p^Z$ is compact and ϕ is continuous, $\phi(\text{supp } \mu_x)$ is compact, which means that there exists a $R > 0$ such that

$$\sup \{\|\phi(\mathbf{z})\| : \mathbf{z} \in \text{supp } \mu_x\} \leq R.$$

Then, we obtain that for any $\mathbf{u} \in \mathbb{R}^n$,

$$\mathbf{u}^\top \text{Cov}_{\mu_x}(\phi(\mathbf{Z})) \mathbf{u} = \text{Var}_{\mu_x}(\mathbf{u}^\top \phi(\mathbf{Z})) \leq \text{Var}_{\mu_x}(\|\mathbf{u}\| \|\phi(\mathbf{Z})\|) \leq R^2 \|\mathbf{u}\|^2,$$

which indicates that

$$\|\text{Cov}_{\mu_x}(\phi(\mathbf{Z}))\|_{\text{op}} \leq R^2.$$

Therefore, we have

$$\|\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x})\|_{\text{op}} \leq \frac{1}{\beta^2} + \frac{\alpha^2 R^2}{\beta^4}. \quad \square$$

Remark F.1. This proposition shows that when the latent density function has compact support, no additional conditions—such as log-concavity or L -smoothness—are required for the latent distribution. Moreover, under the compactness assumption, the results of Proposition 3 can be extended to the nonlinear case, as shown in De Bortoli (2022).

Next, in the non-compact case, Proposition 3 requires the strong log-concavity of the latent density p^Z , as it is used to establish not only the L -smoothness but also the concavity of $\log p_X$ (see Corollary 4). However, if we are only interested in the L -Lipschitz continuity of the score function, the assumption of concavity can be relaxed to the L_0 -smoothness of $\log p^Z$, i.e., $\|\nabla_{\mathbf{z}}^2 \log p^Z(\mathbf{z})\| \leq L_0$, or even to the weaker condition $\nabla_{\mathbf{z}}^2 \log p^Z(\mathbf{z}) \preceq L_0 \mathbf{I}_k$; see Proposition F.2 below.

Proposition F.2. *Let \mathbf{Z} be a random variable on \mathbb{R}^k with the density function p^Z and $B \in \mathbb{R}^{n \times k}$. Assume there are $L_0, \Lambda > 0$ such that*

$$\nabla_{\mathbf{z}}^2 \log p^Z(\mathbf{z}) \preceq L_0 \mathbf{I}_k, \quad \|B\|_{\text{op}}^2 \leq \Lambda,$$

and $\lambda := \lambda_{\min}(B^\top B) > 0$, the minimum of all eigenvalues of $B^\top B$. For $\alpha \in \mathbb{R}$ and $\beta > 0$, let

$$\mathbf{X} = \alpha B \mathbf{Z} + \beta \boldsymbol{\xi}, \quad \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n),$$

with the density function p_X on \mathbb{R}^n . If $\alpha^2 \lambda - L_0 \beta^2 > 0$, we have

$$\|\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x})\|_{\text{op}} \leq \frac{1}{\beta^2} + \frac{\alpha^2 \Lambda}{\beta^2(\alpha^2 \lambda - L_0 \beta^2)}$$

and

$$\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x}) \preceq \left(\frac{\alpha^2 \Lambda}{\beta^2(\alpha^2 \lambda - L_0 \beta^2)} - \frac{1}{\beta^2} \right) \mathbf{I}_n.$$

Proof. The main difference of this proof to the proof in Proposition 3 is how to bound $\|\text{Cov}_{\mu_x}(\mathbf{Z})\|_{\text{op}}$.

Note that

$$\nabla_{\mathbf{z}}^2 U_x(\mathbf{z}) = \frac{\alpha^2}{\beta^2} B^\top B + \nabla_{\mathbf{z}}^2 V(\mathbf{z}) \succeq \left(\frac{\alpha^2 \lambda}{\beta^2} - L_0 \right) \mathbf{I}_k,$$

because $-\nabla_{\mathbf{z}}^2 \log p^Z(\mathbf{z}) = \nabla_{\mathbf{z}}^2 V(\mathbf{z}) \succeq -L_0 \mathbf{I}_k$. When $\alpha^2 \lambda - L_0 \beta^2 > 0$, we similarly obtain

$$\|\text{Cov}_{\mu_x}(\mathbf{Z})\|_{\text{op}} \leq \frac{\beta^2}{\alpha^2 \lambda - L_0 \beta^2}.$$

Therefore,

$$\|\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x})\|_{\text{op}} \leq \frac{1}{\beta^2} + \frac{\alpha^2 \Lambda}{\beta^2(\alpha^2 \lambda - L_0 \beta^2)}.$$

On the other hand,

$$B \text{Cov}_{\mu_x}(\mathbf{Z}) B^\top \preceq \frac{\Lambda \beta^2}{\alpha^2 \lambda - L_0 \beta^2} \mathbf{I}_n,$$

which follows that

$$\nabla_{\mathbf{x}}^2 \log p_X(\mathbf{x}) \preceq \left(\frac{\alpha^2 \Lambda}{\beta^2(\alpha^2 \lambda - L_0 \beta^2)} - \frac{1}{\beta^2} \right) \mathbf{I}_n. \quad \square$$

G More Details for Nonlinear Extension

G.1 Omitted Proofs in Section 6

Proof of Lemma 12. The proof consists of the following three steps. First, let $\mathbf{z} \in \mathbb{R}^d$ be arbitrary.

- (i) Local construction: Because $\phi: \mathbb{R}^d \rightarrow \mathcal{M} \subset \mathbb{R}^D$ is an isometry, the columns of $J\phi(\mathbf{z})$ form an orthonormal basis for the tangent space $T_{\phi(\mathbf{z})}\mathcal{M}$. These vectors can be extended to an orthonormal basis of \mathbb{R}^D by adjoining

$$\{n_1(\mathbf{z}), n_2(\mathbf{z}), \dots, n_{D-d}(\mathbf{z})\},$$

where each n_i is a smooth normal vector fields along \mathcal{M} . For such n_i , one can define the Fermi coordinates map as

$$F: \mathbb{R}^d \times \mathbb{R}^{D-d} \longrightarrow \mathbb{R}^D, \quad F(\mathbf{z}, \mathbf{v}) = \phi(\mathbf{z}) + \sum_{i=1}^{D-d} v_i n_i(\mathbf{z}).$$

Then, by the Tubular Neighborhood Theorem (Theorem I.3), there exists a $\varepsilon: \mathbb{R}^d \rightarrow (0, \infty)$ such that for the set

$$V = \{(\mathbf{z}, \mathbf{v}) \in \mathcal{M} \times \mathbb{R}^{D-d} : \|\mathbf{v}\| < \varepsilon(\mathbf{z})\},$$

$F: V \rightarrow U = F(V)$ is a diffeomorphism, where the open set $U \subset \mathbb{R}^D$ is a neighborhood of \mathcal{M} , called a tubular neighborhood. Let $\pi: \mathbb{R}^d \times \mathbb{R}^{D-d} \rightarrow \mathbb{R}^d$ be the projection, i.e., $\pi(\mathbf{z}, \mathbf{v}) = \mathbf{z}$. Then, one can construct

$$\tilde{\phi}^*: U \longrightarrow \mathbb{R}^d, \quad \tilde{\phi}^*(\mathbf{x}) = \pi(F^{-1}(\mathbf{x}))$$

- (ii) Check conditions: First, because $\mathcal{M} \subset U$, and F is diffeomorphic from V to U with $F(\mathbf{z}, 0) = \phi(\mathbf{z})$,

$$\tilde{\phi}^*(\phi(\mathbf{z})) = \pi(F^{-1}(\phi(\mathbf{z}))) = \pi(\mathbf{z}, 0) = \mathbf{z}, \quad \forall \mathbf{z} \in \mathbb{R}^d.$$

For the derivative condition, by the definition of F , we have

$$JF(\mathbf{z}, 0) = (J_{\mathbf{z}}F(\mathbf{z}, 0), J_{\mathbf{v}}F(\mathbf{z}, 0)) = (J\phi(\mathbf{z}), \mathbf{n}(\mathbf{z})),$$

where $\mathbf{n} = (n_1(\mathbf{z}), \dots, n_{D-d}(\mathbf{z}))$. By $J\phi^\top J\phi = \mathbf{I}_d$, $JF(\mathbf{z}, 0)$ is orthogonal, which follows that

$$J(F^{-1})(F(\mathbf{z}, 0)) = JF(\mathbf{z}, 0)^{-1} = JF(\mathbf{z}, 0)^\top = \begin{pmatrix} J\phi(\mathbf{z})^\top \\ \mathbf{n}(\mathbf{z})^\top \end{pmatrix}.$$

On the other hand, F^{-1} can be written as $F^{-1}(\mathbf{x}) = (F_1(\mathbf{x}), F_2(\mathbf{x}))$, where $F_1 = \pi \circ F^{-1} = \tilde{\phi}^*$ on U . It implies that

$$J(F^{-1})(F(\mathbf{z}, 0)) = \begin{pmatrix} J\tilde{\phi}^*(\phi(\mathbf{z})) \\ JF_2(\phi(\mathbf{z})) \end{pmatrix}.$$

Therefore, $J\tilde{\phi}^*(\phi(\mathbf{z})) = J\phi(\mathbf{z})^\top$.

- (iii) Global construction: By the Urysohn Lemma (Munkres, 2018), there exists a smooth function $\chi: \mathbb{R}^D \rightarrow [0, 1]$ such that $\chi|_{\tilde{U}} \equiv 1$ and $\chi|_{\mathbb{R}^D \setminus U} \equiv 0$, where $\tilde{U} \subset U$ is an open neighborhood of \mathcal{M} . Let $h: \mathbb{R}^D \rightarrow \mathbb{R}^d$ be any smooth function—for instance, a constant function $h \equiv \mathbf{c}$. Define

$$\phi^*(\mathbf{x}) = \chi(\mathbf{x})\tilde{\phi}^*(\mathbf{x}) + (1 - \chi(\mathbf{x}))h(\mathbf{x}),$$

then the desired identities hold:

$$\phi^* \circ \phi = \text{id}_{\mathbb{R}^d}, \quad J\phi^*(\phi(\mathbf{z})) = J\phi(\mathbf{z})^\top. \quad \square$$

Proof of Theorem 13. Let $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$ be the isometry for defining $\mathcal{M} = \text{Im } \phi$. Because $\text{supp } \mathbb{P}_X \subset \mathcal{M}$, there exists a \mathbb{P}^Z defined on \mathbb{R}^d such that $\mathbf{X} = \phi(\mathbf{Z}) \sim \mathbb{P}_X$ when $\mathbf{Z} \sim \mathbb{P}_Z$. Let t be fixed in $(0, T]$. By (2),

$$\mathbf{X}_t = \sqrt{\alpha_t}\phi(\mathbf{Z}) + \sqrt{1 - \alpha_t}\boldsymbol{\xi}.$$

Define $F^t: \mathbb{R}^D \rightarrow \mathbb{R}^D$ by

$$F^t(\mathbf{x}) := \sqrt{\alpha_t}\phi \circ \phi^* \left(\frac{\mathbf{x}}{\sqrt{\alpha_t}} \right), \quad \alpha_t = e^{-2t},$$

where ϕ^* is defined in Lemma 12. Then we have

$$F^t(\mathbf{X}_t) = \sqrt{\alpha_t}\phi \circ \phi^* \left(\mathbf{X}_0 + \sqrt{\frac{1 - \alpha_t}{\alpha_t}}\boldsymbol{\xi} \right), \quad \mathbf{X}_0 := \phi(\mathbf{Z}).$$

Now consider the Taylor expansion of $\varphi := \phi \circ \phi^*$ at $\mathbf{X}_0 = \phi(\mathbf{Z})$, with integral remainder. We obtain

$$F^t(\mathbf{X}_t) = F^t(\mathbf{X}_0) + \sqrt{1 - \alpha_t}J\varphi(\mathbf{X}_0)\boldsymbol{\xi} + R(\boldsymbol{\xi}),$$

where $R(\boldsymbol{\xi})$ denotes the remainder term.

Next, we analyze the three terms on the right-hand side one by one. For the first term, because $\mathbf{X}_0 = \phi(\mathbf{Z}) \in \mathcal{M}$, $\mathbf{Z} = \phi^*(\mathbf{X}_0)$ by the definition of ϕ^* ; see the proof of Lemma 12. It implies that

$$F^t(\mathbf{X}_0) = \sqrt{\alpha_t}\phi \circ \phi^*(\mathbf{X}_0) = \sqrt{\alpha_t}\mathbf{X}_0.$$

For the second term, by Lemma 12,

$$J\varphi(\mathbf{X}_0) = J\phi(\mathbf{Z})J\phi^*(\phi(\mathbf{Z})) = J\phi(\mathbf{Z})J\phi(\mathbf{Z})^\top.$$

Moreover, because $J\phi^\top J\phi = \mathbf{I}_d$, $P := J\varphi(\mathbf{X}_0)$ is an orthogonal projection with rank d . For the third term,

$$R(\boldsymbol{\xi}) = \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \int_0^1 (1 - s) D^2\varphi \left(\mathbf{X}_0 + s\sqrt{(1 - \alpha_t)/\alpha_t}\boldsymbol{\xi} \right) [\boldsymbol{\xi}, \boldsymbol{\xi}] ds.$$

By the proof of Lemma 12, $\phi^* \equiv \mathbf{c}$ on $\mathbb{R}^D \setminus U$ for a tubular neighborhood U of \mathcal{M} , which means $J\phi^* = 0$ and $D^2\phi^* = 0$ on $\mathbb{R}^D \setminus U$. It follows that

$$D^2\varphi(\mathbf{x})[\mathbf{u}, \mathbf{v}] = D^2\phi(\phi^*(\mathbf{x})) [J\phi^*(\mathbf{x})\mathbf{u}, J\phi^*(\mathbf{x})\mathbf{v}] + J\phi(\phi^*(\mathbf{x}))(D^2\phi^*(\mathbf{x})[\mathbf{u}, \mathbf{v}]) = 0$$

for $\mathbf{x} \in \mathbb{R}^D \setminus U$. For a chosen δ , we can choose a tubular neighborhood U sufficiently thin such that $\mathbf{X}_0 + s\sqrt{(1 - \alpha_t)/\alpha_t}\boldsymbol{\xi} \notin U$ for $s > \delta$. Therefore, we have

$$R(\boldsymbol{\xi}) = \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \int_0^\delta (1 - s) D^2\varphi \left(\mathbf{X}_0 + s\sqrt{(1 - \alpha_t)/\alpha_t}\boldsymbol{\xi} \right) [\boldsymbol{\xi}, \boldsymbol{\xi}] ds.$$

Assume $D^2\varphi$ is bounded on U . Then, for any small $\varepsilon' > 0$, one can choose δ sufficiently small such that $\|R(\boldsymbol{\xi})\| \leq \varepsilon'$.

Combining these analyses, we obtain

$$\sqrt{1-\alpha_t}\|(\mathbf{I}_D - P)\boldsymbol{\xi}\| - \varepsilon' \leq \|\mathbf{X}_t - F^t(\mathbf{X}_t)\| \leq \sqrt{1-\alpha_t}\|(\mathbf{I}_D - P)\boldsymbol{\xi}\| + \varepsilon'. \quad (63)$$

Let $f^t(\mathbf{x}) := \|\mathbf{x} - F^t(\mathbf{x})\|$. Similarly as the proof of Proposition 1, by the Laurent-Massart bound (Lemma H.1), (63) implies that

$$\mathbb{P}\left(r(t)\sqrt{1-2\sqrt{\varepsilon}} - \varepsilon' \leq f^t(\mathbf{X}_t) \leq r(t)\sqrt{1+2\sqrt{\varepsilon}} + 2\varepsilon + \varepsilon'\right) \geq 1 - 2e^{-2(D-d)\varepsilon},$$

where $r(t) = \sqrt{(D-d)(1-\alpha_t)}$. Because $d \ll D$, one can choose small ε such that $\delta = e^{-2(D-d)\varepsilon}$ is also small enough. As a result, $\mathbb{P}(f^t(\mathbf{X}_t) \approx r(t)) \geq 1 - \delta$, i.e., \mathbf{X}_t concentrates on $\mathcal{M}^t = (f^t)^{-1}(r(t))$ with high probability. \square

G.2 More Results of Experiments

Comparison of FID. Table 2 serves as a complement to Table 1.

Table 2: Comparison of FID on CIFAR-10

	Airplane	Bird	Cat	Deer	Dog	Overall
CGM ($\eta = 1$)	17.95	21.69	20.34	19.24	23.62	4.07
GeGM ($\eta = 50$)	18.98	18.39	17.35	17.38	18.45	5.15

FID v.s. guidance scale on CIFAR-10. By sampling with the nonlinear GeGM (16), Figure 3 shows how the FID varies with the guidance scale η across all classes from CIFAR-10, which is consistent with the result of Theorem 11.

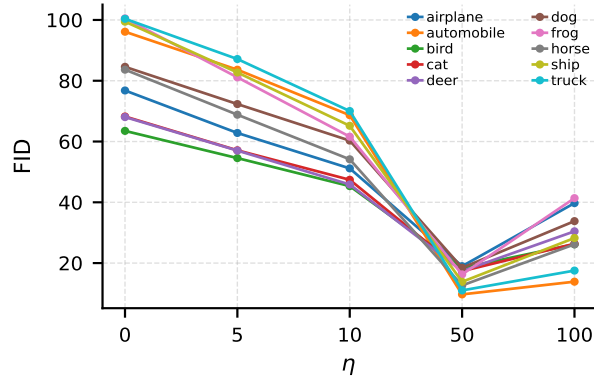


Figure 3: FID v.s. guidance scale η of GeGM on all classes of CIFAR-10

H Technical Lemmas

Lemma H.1 (Laurent & Massart (2000)). *Let X be a χ^2 -random variable with n degrees of freedom, i.e., $X = \sum_{i=1}^n \xi_i^2$ with $\xi_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$. Then, for any $\alpha > 0$, we have*

$$\begin{aligned} \mathbb{P}(X - n \geq 2\sqrt{n\alpha} + 2\alpha) &\leq e^{-\alpha}, \\ \mathbb{P}(X - n \leq -2\sqrt{n\alpha}) &\leq e^{-\alpha}. \end{aligned}$$

Lemma H.2. Let $\mathbb{R}^n = V \oplus V^\perp$ be an orthogonal decomposition of \mathbb{R}^n , where V is a linear subspace and V^\perp is its orthogonal complement. Let $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^n$ be random variables such that $\mathbf{X} \in V$, $\mathbf{Y} \in V^\perp$, and \mathbf{X} independent of \mathbf{Y} . Suppose that \mathbf{X} and \mathbf{Y} admit densities p_X and p_Y on V and V^\perp , respectively, with respect to the canonical volume measures on V and V^\perp . Then the density function of $\mathbf{Z} = \mathbf{X} + \mathbf{Y}$ is given by

$$p_Z(\mathbf{z}) = p_X(Q\mathbf{x})p_Y(Q^\perp\mathbf{x}),$$

where Q is the orthogonal projection onto V , and $Q^\perp = \mathbf{I}_n - Q$ is the orthogonal projection onto V^\perp .

Proof. Let m_V and m_{V^\perp} be the canonical volume measure on V and V^\perp , respectively. Define $\Phi: V \times V^\perp \rightarrow \mathbb{R}^n$ by $\phi(\mathbf{x}, \mathbf{y}) = \mathbf{x} + \mathbf{y}$. Clearly, Φ is an orthogonal linear map, which indicates $|\det J\Phi| = 1$, so

$$\Phi_\#(m_V \otimes m_{V^\perp}) = m_n,$$

where m_n is the Lebesgue measure on \mathbb{R}^n .

Let \mathbb{P}_X and \mathbb{P}_Y be the distributions of \mathbf{X} and \mathbf{Y} , respectively. Then $d\mathbb{P}_X = p_X dm_V$ and $d\mathbb{P}_Y = p_Y dm_{V^\perp}$. By the independence of \mathbf{X} and \mathbf{Y} , we have

$$d(\mathbb{P}_X \otimes \mathbb{P}_Y) = p_X(\mathbf{x})p_Y(\mathbf{y})d(m_V(\mathbf{x}) \otimes m_{V^\perp}(\mathbf{y})).$$

Since $\mathbf{Z} = \mathbf{X} + \mathbf{Y} = \Phi(\mathbf{X}, \mathbf{Y})$, it follows that $\mathbf{Z} \sim \mathbb{P}_Z = \Phi_\#(\mathbb{P}_X \otimes \mathbb{P}_Y)$, and thus

$$\begin{aligned} \mathbb{P}_Z(U) &= \int_{\mathbb{R}^n} \chi_U(\mathbf{z}) d\mathbb{P}_Z(\mathbf{z}) \\ &= \int_{V \times V^\perp} \chi_U(\mathbf{x} + \mathbf{y}) d(\mathbb{P}_X \otimes \mathbb{P}_Y) \\ &= \int_{V \times V^\perp} \chi_U(\mathbf{x} + \mathbf{y}) p_X(\mathbf{x}) p_Y(\mathbf{y}) d(m_V(\mathbf{x}) \otimes m_{V^\perp}(\mathbf{y})) \\ &= \int_{\mathbb{R}^n} \chi_U(\mathbf{z}) p_X(Q\mathbf{z}) p_Y(Q^\perp\mathbf{z}) d\Phi_\#(m_V(\mathbf{x}) \otimes m_{V^\perp}(\mathbf{y})) \\ &= \int_{\mathbb{R}^n} \chi_U(\mathbf{z}) p_X(Q\mathbf{z}) p_Y(Q^\perp\mathbf{z}) dm_n(\mathbf{z}). \end{aligned}$$

Therefore, we have

$$p_Z(\mathbf{z}) = p_X(Q\mathbf{x})p_Y(Q^\perp\mathbf{x}). \quad \square$$

Lemma H.3. Let $(\mathbf{W}_t)_{t \geq 0}$ be a standard Brownian motion on \mathbb{R}^m and $A \in \mathcal{O}^{m \times n}$. Let

$$\mathbf{B}_t := A^\top \mathbf{W}_t.$$

Then $(\mathbf{B}_t)_{t \geq 0}$ is a standard Brownian motion on \mathbb{R}^n .

Proof. The path continuity of $t \mapsto \mathbf{B}_t = A^\top \mathbf{W}_t$ follows directly from the path continuity of $t \mapsto \mathbf{W}_t$, as does the independence of increments. The initial condition $\mathbf{B}_0 = A^\top \mathbf{W}_0 = 0$ is immediate. Moreover, since $A \in \mathcal{O}^{m \times n}$, we have,

$$\mathbf{B}_t - \mathbf{B}_s = A^\top (\mathbf{W}_t - \mathbf{W}_s) \sim \mathcal{N}(0, (t-s)\mathbf{I}_m), \quad \forall t > s. \quad \square$$

Lemma H.4 (Jost (2008)). For a function $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$, if $g: \mathbb{R}^n \rightarrow \text{Im } g$ is a diffeomorphism, that is, both g and its inverse $g^{-1}: \text{Im } g \rightarrow \mathbb{R}^n$ are continuously differentiable, then $g_\# p_X$, the density function of $g_\# \mathbb{P}_X$ on $\text{Im } g$ with respect to the canonical volume measure on $\text{Im } g$, satisfies

$$g_\# p_X(\mathbf{y}) = p_X(\mathbf{x}) |\det (Jg(\mathbf{x}) Jg(\mathbf{x})^\top)|^{\frac{1}{2}}, \quad \mathbf{x} = g^{-1}(\mathbf{y}).$$

Moreover, when $g(\mathbf{x}) = A\mathbf{x}$ for an $A \in \mathcal{O}^{m \times n}$, $A_\# p_X(\mathbf{y}) = p_X(A^\top \mathbf{y})$.

Remark H.1. This result is essentially a general form of the change-of-variables formula, which has been widely used in the context of generative models on manifolds (see, e.g., Loaiza-Ganem et al. (2024)). To rigorously justify this result, some basic knowledge of Riemannian geometry is required. Since $g: \mathbb{R}^n \rightarrow \text{Im } g$ is a diffeomorphism, the image $\text{Im } g \subset \mathbb{R}^n$ is a submanifold. When $\text{Im } g$ is equipped with the canonical Riemannian structure induced from the ambient Euclidean space \mathbb{R}^n , the canonical volume measure on $\text{Im } g$ coincides with the Riemannian volume measure. Therefore, the relevant results from Jost (2008, Section 1.4) can be applied to establish the desired formula rigorously.

Lemma H.5. *Let $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ with large n . Then, with high probability, \mathbf{X} is approximately uniformly distributed on the sphere $\mathbb{S}^{n-1}(\sqrt{n})$, i.e., $\mathbf{X} \sim \text{Unif}(\mathbb{S}^{n-1}(\sqrt{n}))$.*

Proof. First, consider $\mathbf{Y} := \frac{\mathbf{X}}{\|\mathbf{X}\|}$. We first show that $\mathbf{Y} \sim \text{Unif}(\mathbb{S}^{n-1})$. Note that \mathbb{S}^{n-1} is a compact homogeneous space:

$$\mathbb{S}^{n-1} \cong \text{SO}(n)/\text{SO}(n-1),$$

where $\text{SO}(n) \subset \mathbb{R}^{n \times n}$ denotes the special orthogonal group. Consider the natural action of $\text{SO}(n)$ on \mathbb{S}^{n-1} given by $R: \mathbb{S}^{n-1} \rightarrow \mathbb{S}^{n-1}$, $\mathbf{x} \mapsto R\mathbf{x}$ for all $R \in \text{SO}(n)$. Then by the existence and uniqueness of Haar measure (Folland, 2016, Theorem 2.49), $\text{Unif}(\mathbb{S}^{n-1})$ is the unique rotation-invariant probability measure on \mathbb{S}^{n-1} . Therefore, it is sufficient to prove that the distribution of \mathbf{Y} is rotation-invariant, i.e., $\mathbf{Y} \stackrel{d}{=} R\mathbf{Y}$ for all $R \in \text{SO}(n)$.

Since $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ and $R \in \text{SO}(n)$, we have $R\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ and $\|R\mathbf{X}\| = \|\mathbf{X}\|$. Hence,

$$\mathbf{Y} = \frac{\mathbf{X}}{\|\mathbf{X}\|} \stackrel{d}{=} \frac{R\mathbf{X}}{\|R\mathbf{X}\|} = R\mathbf{Y},$$

which implies that $\mathbf{Y} \sim \text{Unif}(\mathbb{S}^{n-1})$. Similarly, by the uniqueness of the invariant measure,

$$\sqrt{n}\mathbf{Y} = \frac{\sqrt{n}}{\|\mathbf{X}\|}\mathbf{X} \sim \text{Unif}(\mathbb{S}^{n-1}(\sqrt{n})).$$

Moreover, as shown in the proof in Proposition 1, the Laurent-Massart Bound implies that $\|\mathbf{X}\| \approx \sqrt{n}$ with high probability when n is large. Therefore,

$$\mathbf{X} \approx \frac{\sqrt{n}}{\|\mathbf{X}\|}\mathbf{X} \sim \text{Unif}(\mathbb{S}^{n-1}(\sqrt{n})). \quad \square$$

Lemma H.6 (Corollary 4.8.2 of Bakry et al. (2013)). *Let $U: \mathbb{R}^n \rightarrow \mathbb{R}$ be C^2 such that $\nabla^2 U \succeq \rho \mathbf{I}_n$ for some $\rho > 0$. Then the probability measure*

$$d\mu(\mathbf{x}) = \frac{e^{-U(\mathbf{x})}}{\int e^{-U(\mathbf{y})} d\mathbf{y}} d\mathbf{x}$$

on \mathbb{R}^n satisfies the Poincaré Inequality with the constant $1/\rho$.

Lemma H.7. *Let $\mu, \nu \in \mathcal{P}(\mathbb{R}^n)$ be two probability measures, and let $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ be measurable. Then*

$$f_{\#}(w_1\mu + w_2\nu) = w_1 f_{\#}\mu + w_2 f_{\#}\nu,$$

for any $w_1, w_2 \in [0, 1]$ with $w_1 + w_2 = 1$.

Proof. For any $A \in \mathcal{B}(\mathbb{R}^m)$,

$$\begin{aligned} f_{\#}(w_1\mu + w_2\nu)(A) &= (w_1\mu + w_2\nu)(f^{-1}(A)) \\ &= w_1\mu(f^{-1}(A)) + w_2\nu(f^{-1}(A)) \\ &= w_1 f_{\#}\mu(A) + w_2 f_{\#}\nu(A). \end{aligned} \quad \square$$

Lemma H.8. Let μ be a probability measure on \mathbb{R}^n . Then, we have

$$\mathbb{E}_{\mathbf{X} \sim \mu} [\|\mathbf{X}\|] \leq \sqrt{\mathbb{E}_{\mathbf{X} \sim \mu} [\|\mathbf{X}\|^2]}.$$

In particular, if $\mu = \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$,

$$\mathbb{E} [\|\mathbf{X}\|] \leq \sqrt{n}.$$

Proof. Because μ is a probability measure, by the Hölder's Inequality,

$$\int_{\mathbb{R}^n} \|\mathbf{x}\| \cdot 1 d\mu(\mathbf{x}) \leq \left(\int_{\mathbb{R}^n} \|\mathbf{x}\|^2 d\mu(\mathbf{x}) \right)^{\frac{1}{2}} \left(\int_{\mathbb{R}^n} 1 d\mu(\mathbf{x}) \right)^{\frac{1}{2}},$$

that is, $\mathbb{E}_{\mathbf{X} \sim \mu} [\|\mathbf{X}\|] \leq \sqrt{\mathbb{E}_{\mathbf{X} \sim \mu} [\|\mathbf{X}\|^2]}$. In particular, if $\mu = \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$, $\mathbb{E}[\|\mathbf{X}\|^2] = n$. \square

Lemma H.9. Let $\mu_1, \mu_2, \nu_1, \nu_2$ be probability measures on \mathbb{R}^n , and let

$$\mu = w\mu_1 + (1-w)\mu_2, \quad \nu = w\nu_1 + (1-w)\nu_2, \quad w \in [0, 1].$$

Then, we have

$$\mathcal{W}_1(\mu, \nu) \leq w\mathcal{W}_1(\mu_1, \nu_1) + (1-w)\mathcal{W}_1(\mu_2, \nu_2).$$

Proof. By the existence of optimal coupling on \mathbb{R}^n (Chewi et al., 2024), there is a $\gamma_i \in \Gamma(\mu_i, \nu_i)$ for $i = 1, 2$ such that

$$\mathcal{W}_1(\mu_i, \nu_i) = \int_{\mathbb{R}^n \times \mathbb{R}^n} \|\mathbf{x} - \mathbf{y}\| d\gamma_i(\mathbf{x}, \mathbf{y}).$$

Let

$$\pi = w\gamma_1 + (1-w)\gamma_2.$$

Clearly, π is a probability measure on $\mathbb{R}^n \times \mathbb{R}^n$. Moreover, by definition,

$$\begin{aligned} \pi(A \times \mathbb{R}^n) &= w\gamma_1(A \times \mathbb{R}^n) + (1-w)\gamma_2(A \times \mathbb{R}^n) = w\mu_1(A) + (1-w)\mu_2(A) = \mu(A), \\ \pi(\mathbb{R}^n \times B) &= w\gamma_1(\mathbb{R}^n \times B) + (1-w)\gamma_2(\mathbb{R}^n \times B) = w\nu_1(B) + (1-w)\nu_2(B) = \nu(B), \end{aligned}$$

which means $\pi \in \Gamma(\mu, \nu)$. Therefore,

$$\begin{aligned} \mathcal{W}_1(\mu, \nu) &\leq \int_{\mathbb{R}^n \times \mathbb{R}^n} \|\mathbf{x} - \mathbf{y}\| d\pi(\mathbf{x}, \mathbf{y}) \\ &= w \int_{\mathbb{R}^n \times \mathbb{R}^n} \|\mathbf{x} - \mathbf{y}\| d\gamma_1(\mathbf{x}, \mathbf{y}) + (1-w) \int_{\mathbb{R}^n \times \mathbb{R}^n} \|\mathbf{x} - \mathbf{y}\| d\gamma_2(\mathbf{x}, \mathbf{y}) \\ &= w\mathcal{W}_1(\mu_1, \nu_1) + (1-w)\mathcal{W}_1(\mu_2, \nu_2). \end{aligned} \quad \square$$

Lemma H.10. Let $\mu \in \mathbb{R}^n$ be a probability measure such that its support K is closed and convex. Then

$$\mathbb{E}_{\mathbf{X} \sim \mu} [\mathbf{X}] \in K.$$

Proof. Suppose that $\mathbf{m} = \mathbb{E}_{\mathbf{X} \sim \mu} [\mathbf{X}] \notin K$. By the convexity and closedness of K , the strong separation theorem (Rockafellar, 1997) implies that there are $\mathbf{u} \in \mathbb{R}^n \setminus \{0\}$ and $c \in \mathbb{R}$ such that $\langle \mathbf{u}, \mathbf{m} \rangle > c$ and

$$\langle \mathbf{u}, \mathbf{x} \rangle \leq c, \quad \forall \mathbf{x} \in K.$$

Let $\mathbf{X} \sim \mu$. $\mathbf{X} \in K$ for almost everywhere and so

$$\langle \mathbf{u}, \mathbf{X} \rangle \leq c, \quad a.e..$$

Then taking the expectation on the both sides, we have

$$\langle \mathbf{u}, \mathbf{m} \rangle \leq c,$$

which induces a contradiction. \square

Lemma H.11 (Grönwall's Inequality). *If $u: [0, T] \rightarrow \mathbb{R}$ satisfies the linear ODE inequality as*

$$\frac{d}{dt}u(t) \leq a(t)u(t) + b(t),$$

then

$$u(t) \leq u(0)e^{\int_0^t a(r)dr} + \int_0^t b(s)e^{\int_s^t a(r)dr} ds.$$

Proof. Let $\Phi(t) = \exp\left(-\int_0^t a(s)ds\right)$. Then, $\Phi'(t) = -a(t)\Phi(t)$ and

$$\Phi(t) \frac{d}{dt}u(t) \leq \Phi(t)a(t)u(t) + \Phi(t)b(t) \Rightarrow \frac{d}{dt}(\Phi(t)u(t)) \leq \Phi(t)b(t).$$

By integrating on the both sides of above inequality, we have

$$u(t) \leq u(0)e^{\int_0^t a(r)dr} + \int_0^t b(s)e^{\int_s^t a(r)dr} ds. \quad \square$$

Lemma H.12. *If a C^1 function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is ρ -strongly convex, then it satisfies ρ -Polyak–Łojasiewicz (PL) inequality:*

$$\|\nabla_{\mathbf{x}} f(\mathbf{x})\|^2 \geq 2\rho(f(\mathbf{x}) - f(\mathbf{x}_*)),$$

where \mathbf{x}_ is the unique minimizer of f .*

Proof. Because f is ρ -strongly convex,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla_{\mathbf{x}} f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\rho}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

Minimizing the both sides with respect to \mathbf{y} , we obtain

$$f(\mathbf{x}_*) \geq f(\mathbf{x}) - \frac{1}{2\rho} \|\nabla_{\mathbf{x}} f(\mathbf{x})\|^2,$$

which is precisely the ρ -PL inequality. \square

Lemma H.13. *Let $f: \mathbb{R}^k \rightarrow \mathbb{R}^n$ be L -Lipschitz continuous. For two probability measures $\mu, \nu \in \mathcal{P}(\mathbb{R}^n)$,*

$$\mathcal{W}_1(f_{\#}\mu, f_{\#}\nu) \leq L\mathcal{W}_1(\mu, \nu).$$

Proof. Let (\mathbf{X}, \mathbf{Y}) be an optimal coupling for (μ, ν) , that is, $\mathbf{X} \sim \mu$, $\mathbf{Y} \sim \nu$, and $\mathcal{W}_1 = \mathbb{E}[\|\mathbf{X} - \mathbf{Y}\|]$. Besides, $f(\mathbf{X}) \sim f_{\#}\mu$ and $f(\mathbf{Y}) \sim f_{\#}\nu$. Then, by the Lipschitz continuity of f ,

$$\begin{aligned} \mathcal{W}_1(f_{\#}\mu, f_{\#}\nu) &\leq \mathbb{E}[\|f(\mathbf{X}) - f(\mathbf{Y})\|] \\ &\leq L\mathbb{E}[\|\mathbf{X} - \mathbf{Y}\|] \\ &= L\mathcal{W}_1(\mu, \nu). \end{aligned} \quad \square$$

I Preliminaries for Manifold

We provide only the minimal background on smooth manifolds necessary for this work. For a comprehensive treatment, we refer the reader to Lee (2012).

Definition I.1. A subset $\mathcal{M} \subset \mathbb{R}^n$ is called a m -dimensional (embedded) (sub)manifold of \mathbb{R}^n if there are a family open sets $\{U_{\alpha}\}_{\alpha \in \Gamma}$ in \mathbb{R}^n , a family of open sets $\{V_{\alpha}\}_{\alpha \in \Gamma}$ in \mathbb{R}^m , and a family of smooth (C^{∞}) maps $\{\phi_{\alpha}\}_{\alpha \in \Gamma}$ such that

$$\mathcal{M} \subset \bigcup_{\alpha \in \Gamma} U_{\alpha}, \text{ and } \phi_{\alpha}: V_{\alpha} \rightarrow U_{\alpha} \cap \mathcal{M}$$

is a diffeomorphism, i.e., $\phi_{\alpha}^{-1}: U_{\alpha} \cap \mathcal{M} \rightarrow V_{\alpha}$ is also smooth.

Each pair (ϕ_α, V_α) is called a chart, and $\{(\phi_\alpha, V_\alpha)\}_{\alpha \in \Gamma}$ is called an atlas of \mathcal{M} . In general, a single chart cannot cover the entire manifold \mathcal{M} . However, if \mathcal{M} is closed, then there exists a chart $\phi: V \rightarrow \mathcal{M}$ that can almost cover \mathcal{M} , in the sense that the volume measure of the set $\mathcal{M} \setminus \phi(V)$ is zero; see Lee (2019) for more details.

Definition I.2. Let $\mathcal{M} \subset \mathbb{R}^n$ be a m -dimensional manifold. For any $\mathbf{x} \in \mathcal{M}$, the tangent space, denoted $T_{\mathbf{x}}\mathcal{M}$, is a vector space defined as

$$T_{\mathbf{x}}\mathcal{M} := \{\gamma'(0) : \exists \varepsilon > 0, \gamma: (-\varepsilon, \varepsilon) \rightarrow \mathcal{M} \text{ smooth}, \gamma(0) = \mathbf{x}\}.$$

Lemma I.1. Let $\mathcal{M} \subset \mathbb{R}^n$ be a smooth submanifold. If a C^1 function $g: \mathbb{R}^n \rightarrow \mathbb{R}$ is constant on \mathcal{M} , then for any $\mathbf{x} \in \mathcal{M}$, $\nabla g(\mathbf{x})$ is normal to \mathcal{M} ; that is, $\nabla g(\mathbf{x}) \perp T_{\mathbf{x}}\mathcal{M}$.

Proof. For any $\mathbf{v} \in T_{\mathbf{x}}\mathcal{M}$, let $\gamma: [0, 1] \rightarrow \mathcal{M}$ be a smooth curve such that $\gamma(0) = \mathbf{x}$ and $\gamma'(0) = \mathbf{v}$. Then, because $g(\gamma(t)) \equiv c$,

$$0 = \left. \frac{d}{dt} \right|_{t=0} g(\gamma(t)) = \langle \nabla g(\gamma(0)), \gamma'(0) \rangle = \langle \nabla g(\mathbf{x}), \mathbf{v} \rangle$$

Therefore, $\nabla g(\mathbf{x}) \perp T_{\mathbf{x}}\mathcal{M}$. □

Theorem I.2 (Constant Rank Theorem (Lee, 2012)). Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^r$ be a smooth map and $\mathbf{c} \in \mathbb{R}^r$. Let

$$\mathcal{M} := \{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) = \mathbf{c}\}.$$

If $\text{rank } JF(\mathbf{x}) = r$ for any $\mathbf{x} \in \mathcal{M}$, then \mathcal{M} is a $(n - r)$ -dimensional manifold.

Theorem I.3 (Tubular Neighborhood Theorem (Lee, 2012)). Let $\mathcal{M} \subset \mathbb{R}^D$ be a d -dimensional submanifold. There is a smooth $\varepsilon: \mathcal{M} \rightarrow (0, \infty)$ such that for

$$V := \{(\mathbf{z}, \mathbf{v}) \in \mathcal{M} \times \mathbb{R}^{D-d} : \|\mathbf{v}\| < \varepsilon(\mathbf{z})\},$$

$F: V \rightarrow U = F(V)$ is a diffeomorphism and $U \subset \mathbb{R}^D$ is a neighborhood of \mathcal{M} .

Remark I.1. For a given tubular neighborhood V of \mathcal{M} , we also call $U = F(V) \subset \mathbb{R}^D$ its tubular neighborhood in \mathbb{R}^D . Moreover, we can define the corresponding orthogonal projection $\pi: U \rightarrow \mathcal{M}$ as

$$\pi(\mathbf{x}) = \pi_1(F^{-1}(\mathbf{x})),$$

where $\pi_1: V \rightarrow \mathcal{M}$ is $\pi_1(\mathbf{z}, \mathbf{v}) = \mathbf{z}$.