
Benign Overfitting Does Not Occur in Diffusion Models

Anonymous Authors¹

Abstract

Benign overfitting and double descent have come to shape our understanding of generalization in deep learning, painting a consistent picture: overfitting is not only compatible with good generalization but can actively benefit it. Since diffusion models share much of the machinery of standard deep learning, it is natural to assume that they also exhibit these properties. In this work, we show that this assumption is largely incorrect. We first establish fundamental impossibility results, showing that overfitting and good generalization cannot occur simultaneously except when the sample size grows exponentially with data dimension. We then identify a key difference between regression and score matching in a simplified setting: regression benefits from an alignment between the kernel of the empirical covariance and the target, whereas no such alignment exists in score matching, making overfitting irreparably harmful. We further examine mechanisms that prevent overfitting, identifying implicit regularization arising from time-smoothness in the score function, and early stopping in training. We support our theoretical findings with high-dimensional experiments on U-Net architectures in image generation settings. Our results reveal that generalization is governed by mechanisms distinct from those of classical settings, motivating new theory for diffusion models.

1. Introduction

Diffusion models form a popular class of generative models (Ho et al., 2020; Song et al., 2021; Sohl-Dickstein et al., 2015). Their key insight is to reframe generation as the reversal of a noising process which is characterized by a score function, whose estimation turns out to be equivalent

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the FoGen Workshop at ICML 2026. Do not distribute.

to a denoising problem that can be solved by training a neural network on a simple regression-like objective (Song & Ermon, 2019; Vincent, 2011; Hyvärinen, 2005). Despite their conceptual simplicity, they are capable of achieving state-of-the-art performance in a wide range of applications (Dhariwal & Nichol, 2021; Esser et al., 2024; Rombach et al., 2022; Saharia et al., 2022; Watson et al., 2023; Zhang et al., 2024), and appear in many of the most widely deployed generative AI systems (Yang et al., 2023).

The widespread success of diffusion models naturally raises questions about the mechanisms underpinning their strengths, motivating a growing body of theoretical work on diffusion models (Oko et al., 2023; Azangulov et al., 2024; Benton et al., 2024; Conforti et al., 2025). Of particular interest is the question of generalization (Li et al., 2023; Bonnaire et al., 2025): understanding when a score network trained on finite data faithfully represents the underlying distribution, rather than memorizing the training set. This question is further sharpened by practical concerns around data privacy (Carlini et al., 2023), making a principled understanding of the generalization-memorization tradeoff both scientifically and practically important.

Preceding these developments, a rich theory was developed to understand the generalization properties of overparameterized neural networks in *supervised learning*, which pose a challenge to classical statistics based on the bias-variance tradeoff (Shalev-Shwartz & Ben-David, 2014). Indeed, deep neural networks have demonstrated the ability to generalize well while achieving zero training loss and even fit random labels on regression or classification tasks (Zhang et al., 2017). This is theoretically captured as *benign overfitting*, where exact interpolation of noisy data need not harm test performance (Bartlett et al., 2020). It has been shown in regression settings that, under certain conditions in the overparameterized setting, a form of self-induced regularization occurs that maintains generalization.

The relationship between model complexity and generalization is captured by the double descent phenomenon (Belkin et al., 2019), where the test error follows a non-monotonic curve: it decreases, peaks near the interpolation threshold, and then decreases again as model size grows (see Figure 1, left). This key behavior has been empirically studied in high-dimensional problems (Nakkiran et al., 2019) and theo-

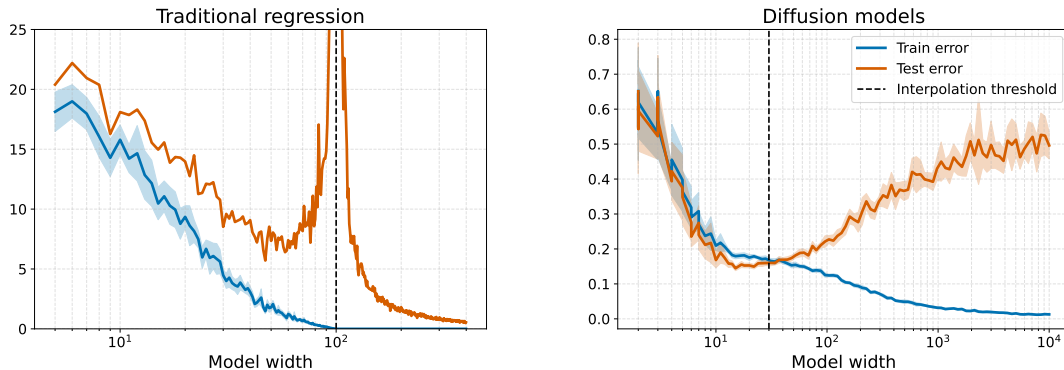


Figure 1. (Left) 2-layer random feature network **regression**, displaying the classical double descent curve. (Right) 2-layer random feature network **diffusion**, showing the absence of benign overfitting. Details on the setup and hyperparameters for these experiments is available in Section C.1.

retical grounding has been provided through the analysis of random feature models (Mei & Montanari, 2022). This allows a precise characterization of risk across regimes using tools from high-dimensional statistics and random matrix theory (Couillet & Liao, 2022). The takeaway from these theories is a now widely internalized picture: in highly over-parameterized models trained on regression-like objectives, fitting the training data perfectly is not only compatible with generalization but, in the right regime, actively beneficial.

One of the benefits of the diffusion model framework is that its construction manages to reuse much of the existing machinery of deep learning: diffusion models are built from existing architectures, trained with standard optimization algorithms (e.g., ADAM (Kingma & Ba, 2015)), and fit by minimizing a regression-like objective (Vincent, 2011; Song et al., 2021). As a result, it is tempting to conclude that the same mechanisms of benign overfitting and double descent are at work, and that the previously mentioned deep learning phenomena carry over more or less intact. *In this work, we show that such an assumption is largely wrong.* In several recent works, it has been hinted that diffusion models may be incapable of benign overfitting (Farghly et al., 2025; Dupuis et al., 2025; Merger & Goldt, 2025). Furthermore, the related facts that exactly minimizing the train error or over-training leads to memorization of the train set has been explored in previous works (Pidstrigach, 2022; Bonnaire et al., 2025). However, to our knowledge, the claim that benign overfitting and double descent do not occur in diffusion models is yet to receive rigorous empirical or theoretical investigation.

For diffusion models, we show that in most practical settings, it is fundamentally impossible have both the train and test loss simultaneously small, invalidating benign overfitting. We further show that instead of a double descent curve, the excess risk forms a U-shaped curve with increasing complexity, akin to the classical bias-variance tradeoff. This is illustrated in the random feature network experiment in

Figure 1 where instead of a double descent curve, the excess risk forms a U-shape that tapers off as the width of the network grows large. In Figure 2, we see how this phenomenon effects high-dimensional image generation problems based on a U-Net network (Ho et al., 2020; Ronneberger et al., 2015), showing that as the number of U-Net features are increased, the overfit diffusion model goes from low quality generation, to generalization, and then towards exactly memorizing the train set. Through a combination of fundamental theoretical results, simplified models and experiments we explore this phenomenon, the mechanisms underpinning it and how diffusion models generalize in spite of it.

Concretely, we present the following contributions:

- In Section 3, we derive fundamental negative results showing that the population and empirical score matching loss cannot be simultaneously small in practical settings, unless the sample size is extremely large, regardless of the score network architecture.
- In Section 4, through the analysis of a linear random feature model, we provide a precise understanding of the interplay between overfitting, dimension, sample size, noise scale, and model complexity. By comparing with a classical regression setup, we illustrate the key differences that prevent benign overfitting to occur in score matching.
- In Section 5, we push further the analysis of Section 4 by identifying key components of the score matching framework that can act as implicit regularization and help prevent overfitting in practice. More precisely, we identify the role of the time-smoothness of the score estimator and discuss the role of early stopping during training.
- We validate our theoretical findings by experiments conducted in high-dimensional settings.

All proofs are available in the appendix. We will make our implementation available upon publication.

Notations. For probability measures μ and ν in \mathbb{R}^d , if μ is absolutely continuous with respect to ν , denoted by $\mu \ll \nu$, we denote $d\mu/d\nu$ its density. If $\mu \ll \nu$, the Kullback-Leibler (KL) divergence is $\text{KL}(\mu||\nu) = \int \log(d\mu/d\nu)d\mu$, the relative Fisher information is $\mathcal{J}(\mu||\nu) = \int \|\nabla \log(d\mu/d\nu)\|^2 d\mu$, the Fisher information is $\mathcal{J}(\nu) = \int \|\nabla \log(d\nu/dx)\|^2 d\nu$. For $a, b \in \mathbb{R}$, $a \wedge b = \min(a, b)$, $a \vee b = \max(a, b)$, and $a_+ = a \vee 0$.

2. Background

By learning to reverse a noising process, diffusion models approximate a distribution ν on \mathbb{R}^d from a dataset sampled from ν . We briefly review this construction and fix the notation used in the paper.

Noising process. The noising process is given by an SDE of Ornstein–Uhlenbeck-type in \mathbb{R}^d ,

$$dX_t = -\kappa X_t dt + \sqrt{2}dB_t, \quad X_0 \sim \nu, \quad (1)$$

with $\kappa \geq 0$ and B_t is a standard d -dimensional Brownian motion. We assume that $\mathbb{E}_{X \sim \nu}[\|X\|^2] < \infty$. It is well-known that $X_t|X_0 \sim N(\alpha_t X_0, \sigma_t^2 I_d)$, where the mean and variance are given by

$$\alpha_t = e^{-\kappa t}, \quad \sigma_t^2 = \kappa^{-1}(1 - \alpha_t^2), \quad (2)$$

with the convention that $\sigma_t^2 = 2t$ when $\kappa = 0$. Fixing a time horizon $T > 0$ and writing $Y_t := X_{T-t}$, the classical result of (Haussmann & Pardoux, 1986; Anderson, 1982) identifies the law of the time-reversed process as the solution of

$$dY_t = \kappa Y_t dt + 2\nabla \log p_{T-t}(Y_t) dt + \sqrt{2}dB_t, \quad Y_0 \sim p_T,$$

where p_t denotes the density of X_t . Thus, if we could simulate Y_t starting from $Y_0 \sim p_T$, we would obtain samples $Y_T \sim \nu$. The only unknown quantity is the *score function* $\nabla \log p_t$, which must be learned from data. Once we have an approximate score function $\hat{s} : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}^d$, we generate samples by numerically simulating

$$d\hat{Y}_t = \kappa \hat{Y}_t dt + 2\hat{s}(\hat{Y}_t, t) dt + \sqrt{2}dB_t, \quad Y_0 \sim q_0.$$

For $\kappa > 0$, we set $q_0 = N(0, \kappa^{-1}I_d)$ as it approximates p_T for large T . For $\kappa = 0$, we set $q_0 = N(0, 2TI_d)$. We denote by q_t the probability density of \hat{Y}_t for $t \in [0, T]$.

Score matching. The score function is approximated using a deep neural network. Ideally, the goal is to minimize the (*population*) *explicit score matching (ESM) loss*,

$$\begin{aligned} \mathcal{L}_{\text{ESM}}(s, t) &:= \mathbb{E}[\|s(t, X_t) - \nabla \log p_t(X_t)\|^2], \\ \mathcal{L}_{\text{ESM}}(s, \varpi) &:= \int \mathcal{L}_{\text{ESM}}(s, t) d\varpi(t), \end{aligned}$$

where ϖ is some non-negative measure on $[0, T]$. While the time-weighting ϖ can be chosen freely, a canonical choice is the *evidence lower bound (ELBO) weighting*, $\varpi_{\text{ELBO}}^\epsilon(dt) = \mathbb{1}_{[\epsilon, T]}(t)$, for $\epsilon > 0$. This weighting makes the explicit score matching loss becomes an upper bound for the sample KL divergence, $\text{KL}(p_\epsilon||q_{T-\epsilon}) \leq \mathcal{L}_{\text{ESM}}(s, \varpi_{\text{ELBO}}^\epsilon) + \text{KL}(p_T||q_0)$. While the ESM loss depends on the intractable $\nabla \log p_t$, a tractable alternative is the *denoising score matching (DSM) loss*,

$$\mathcal{L}_{\text{DSM}}(s, \varpi) := \int \mathbb{E}[\ell_t(s, X_0)] d\varpi(t),$$

$$\ell_t(s, X_0) := \mathbb{E}[\|s(t, X_t) - \nabla \log p_{t|0}(X_t|X_0)\|^2 | X_0],$$

where $p_{t|0}(\cdot|\cdot)$ is the density of X_t given X_0 . This loss is equal to the ESM loss up to a constant, while using $\nabla \log p_{t|0}$, which is known and turns the learning task into a denoising problem (Vincent, 2011; Hyvärinen, 2005). It is then estimated using a dataset $S = \{Z_i\}_{i=1}^n \sim \nu^{\otimes n}$, leading to the *empirical DSM loss*,

$$\widehat{\mathcal{L}}_{\text{DSM}}(s, \varpi) := \frac{1}{n} \sum_{i=1}^n \int \ell_t(s, Z_i) d\varpi(t), \quad n \geq 1. \quad (3)$$

We also define the empirical ESM loss, related to the above by an additive data-dependent constant

$$\widehat{\mathcal{L}}_{\text{ESM}}(s, t) := \mathbb{E}[\|s(t, \hat{X}_t) - \nabla \log \hat{p}_t(\hat{X}_t)\|^2 | S],$$

$$\widehat{\mathcal{L}}_{\text{ESM}}(s, \varpi) := \int \widehat{\mathcal{L}}_{\text{ESM}}(s, t) d\varpi(t),$$

where $d\hat{X}_t = -\kappa \hat{X}_t dt + \sqrt{2}dB_t$ with $\hat{X}_0 \sim n^{-1} \sum_{i=1}^n \delta_{Z_i}$ and for $t > 0$, \hat{p}_t is the density of \hat{X}_t .

3. Fundamental limitations of overfitting in diffusion models

Benign overfitting occurs when good generalization performance is achieved despite interpolating training data. In this section, we derive impossibility results that show that this cannot happen in diffusion models without an exceedingly large sample size. Unlike Sections 4 and 5, the results presented here hold regardless of the score network architecture. In this sense, they identify fundamental limitations of the score matching framework itself, rather than of any particular class of models.

3.1. Warm-up: A Fisher information lower bound

We begin with the following simple lemma, which is a generic lower bound on the sum of the population and empirical ESM losses. The proof can be found in Section B.1.1.

Lemma 3.1. *Let $t \in (0, T]$, for any measurable score function $s : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}^d$, we have*

$$\mathbb{E}[\widehat{\mathcal{L}}_{\text{ESM}}(s, t) + \mathcal{L}_{\text{ESM}}(s, t)] \geq \frac{1}{2} \mathbb{E}[\mathcal{J}(\hat{p}_t||p_t)]. \quad (4)$$

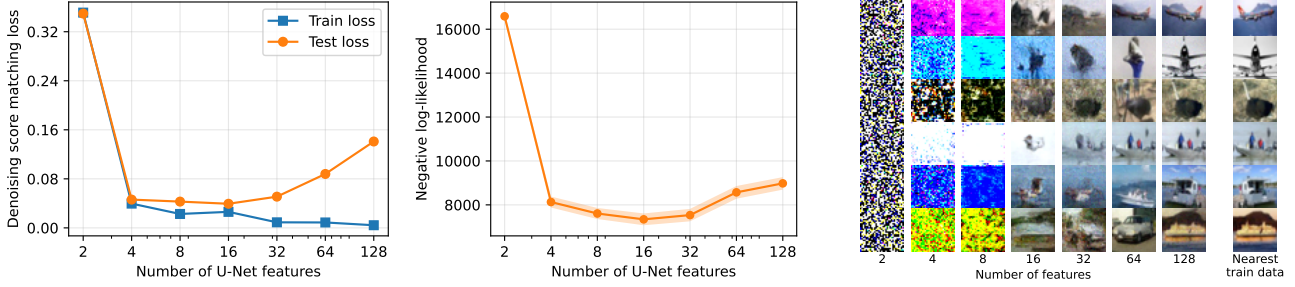


Figure 2. We train a DDPM U-Net model on a subset of CIFAR10 to convergence. We observe train and error for different configurations of the model, varying the number of features (NF) in the U-Net. (Left) Test and train error at convergence varying NF and thus, varying the number of parameters in the model. (Center) The same setting but for the population negative log-likelihood. (Right) Generated samples for different NF values along side the closest train image for NF = 8. Details of the experiment are given in Appendix C.2.

In particular, if $\mathcal{J}(\nu) < +\infty$, we have that $\inf_s \mathbb{E}[\widehat{\mathcal{L}}_{\text{ESM}}(s, t) + \mathcal{L}_{\text{ESM}}(s, t)] \rightarrow +\infty$ as $t \rightarrow 0^+$.

This shows that for both the empirical and population score matching losses to be small, we must have the empirical and population measures \hat{p}_t and p_t close also. The measure of closeness here is the relative Fisher information, which under additional concentration assumptions, can be shown to upper bound the KL divergence and, hence, the total variation distance (see Section B.1.2). Under different conditions on ν , lower bounds on a variant of the Fisher information were derived in (Ye et al., 2026); see Section 6 for details.

Lemma 3.1 shows that for both the empirical and population ESM losses to be small, and therefore for benign overfitting to occur, we require \hat{p}_t to be close to p_t . As \hat{p}_t is initialized at the empirical distribution $n^{-1} \sum_{i=1}^n \delta_{Z_i}$ where $(Z_1, \dots, Z_n) \sim \nu^{\otimes n}$, this suggests that benign overfitting only happens when the sample size is large enough. We make this argument precise in the next subsection.

3.2. A quantitative lower bound under time-integration

We now refine the previous result by obtaining lower bounds on the amount of data required for the empirical and population losses to be simultaneously small. The lower bound is a function of the effective dimension of the data distribution ν which we capture through the lower Rényi dimension, d' , which we formally define in Section A.1. In cases where ν is supported on a compact manifold and under mild conditions (see Section A.1), d' is exactly the manifold dimension.

In the next result, Lemma 3.1 is also extended by considering the full time-integrated loss as opposed to an individual time-slice, requiring only that the weighting ϖ has positive non-decreasing density on an interval. The bound is generic and holds for any random score function \hat{s} , that can be taken as a function of training data as well as additional sources of randomness.

Theorem 3.1. *Suppose that ϖ is a positive measure supported in $[\epsilon, T]$ with $\epsilon > 0$, with non-decreasing density*

$w : [\epsilon, T] \rightarrow \mathbb{R}_+$, that ν has finite second moments, and both T and d are sufficiently large. Suppose that the random score function \hat{s} , satisfies

$$\mathbb{E}_S[\widehat{\mathcal{L}}_{\text{ESM}}(\hat{s}, \varpi)] \leq \epsilon,$$

for some $\epsilon \leq \sup_t w(t)/32$. Then, if either of the following conditions hold:

$$\begin{aligned} \mathbb{E}_S[\mathcal{L}_{\text{ESM}}(\hat{s}, \varpi)] \leq \epsilon, \text{ or } \mathbb{E}_S[\mathcal{L}_{\text{ESM}}(\hat{s}, \varpi_{\text{ELBO}}^\epsilon)] \leq \frac{1}{32}, \\ \text{or } \text{KL}(p_\delta || q_{T-\delta}) \leq \frac{1}{32}, \end{aligned} \quad (5)$$

and with $\delta := w^{-1}(32\epsilon)$ sufficiently small, it must follow that,

$$n \geq (8\sqrt{d}\sigma_\delta)^{-\frac{d'}{2}},$$

where we define $w^{-1}(r) = \inf\{t \in \text{supp}(\varpi) : w(t) \geq r\}$ and d' is the lower Rényi dimension of ν .

The assumption on the measure ϖ , encompasses several standard choices of time-weighting within the continuous training framework (Song et al., 2021). In the following remark, we consider several of these, obtaining closed-form bounds.

Remark 3.1. Under the ELBO weighting $\varpi = \varpi_{\text{ELBO}}^\epsilon$, the sample size lower bound becomes $n \geq (8\sigma_\epsilon\sqrt{d})^{-d'/2}$ for ϵ sufficiently small. Under the variance weighting $w(t) = \sigma_t^2$, commonly used for training (Ho et al., 2020), the analogous bound is $n \geq (2(32\epsilon \vee \sigma_\epsilon)\sqrt{d})^{-d'/2}$ for $\epsilon, \sigma_\epsilon$ sufficiently small. Most strikingly, when $\epsilon = 0$ and $\epsilon \leq \frac{1}{32} \limsup_{t \rightarrow 0^+} w(t)$, there is no sample size at which benign overfitting can occur: $\mathbb{E}[\widehat{\mathcal{L}}_{\text{ESM}}(\hat{s}, \varpi)] \leq \epsilon$ necessarily implies $\mathbb{E}[\mathcal{L}_{\text{ESM}}(\hat{s}, \varpi)] > \epsilon$, $\mathbb{E}[\mathcal{L}_{\text{ESM}}(\hat{s}, \varpi_{\text{ELBO}}^0)] > 1/32$, and $\text{KL}(p_\delta || q_{T-\delta}) > 1/32$.

We highlight that the lower bound is increasing exponentially in the intrinsic dimension and depends on the minimum support of ϵ , supporting both the idea that manifold distributed data is beneficial to diffusion models and that

early stopping in the reverse process can benefit generalization (Azangulov et al., 2024; Farghly et al., 2026; Bortoli, 2022; Farghly et al., 2025). Since many popular implementations of diffusion models take $\epsilon \in \{10^{-3}, 10^{-5}\}$, the expectation that the amount of data scales with $\epsilon^{-d'/2}$ is unreasonably strong. Thus, in most practical settings, diffusion models will not exhibit the property of benign overfitting.

The proof, presented in Section B.1.3, uses a similar approach as Lemma 3.1, showing that for empirical and population losses to be small, \hat{p}_t and p_t must be close. Instead of the Fisher information, we derive a connection to the total variation via an argument based on Girsanov's theorem (Øksendal, 2003).

4. Fine-grained analysis on linear RFNNs

Section 3 establishes, at a general level, that benign overfitting cannot occur in diffusion models in any reasonable setting. However, it does not reveal how overfitting manifests, nor how it depends on the interplay between model size, data dimension, and noise scale. To get a finer understanding, we now analyze a concrete and tractable model class: two-layer linear random feature networks (Rahimi & Recht, 2007), which is similar with the models of (George & Macris, 2026; George et al., 2025; Bonnaire et al., 2025) with linear activation. Despite the simplicity of the linear setting, benign overfitting and double descent have been studied in it (Belkin et al., 2019; Hastie et al., 2022), suggesting that it is rich enough to apprehend overfitting in diffusion models.

Setup. We study 2-layer linear score networks of the form $s_{W,A} : x \mapsto AWx$, where $W \in \mathbb{R}^{p \times d}$ is fixed and $A \in \mathbb{R}^{d \times p}$ is a trainable parameter. We fix a time $t \in (0, T]$ and consider the minimization of the empirical risk defined by $A \mapsto \sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,A}, t)$. We denote by Σ the covariance matrix of the data distribution ν and define the empirical covariance matrix $\widehat{\Sigma} := n^{-1} \sum_{i=1}^n Z_i Z_i^T$ along with $\Sigma_t := \alpha_t^2 \Sigma + \sigma_t^2 \text{I}_d$ and $\widehat{\Sigma}_t := \alpha_t^2 \widehat{\Sigma} + \sigma_t^2 \text{I}_d$.

In the following lemma, we characterize the DSM and ESM loss associated with the empirical risk minimization problem. The proof can be found in Section B.2.

Lemma 4.1. *Assume that W has full rank. Let $\widehat{A} \in \arg \min \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,A}, t)$. Let $\Omega_t := \widehat{\Sigma}_t^{-1}$ if $p \geq d$ and $\Omega_t := W^T (W \widehat{\Sigma}_t W^T)^{-1} W$ if $p < d$. Then $\widehat{A}W = -\Omega_t$ and*

$$\begin{aligned} \sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,\widehat{A}}, t) &= d - \sigma_t^2 \text{Tr}(\Omega_t), \\ \mathcal{L}_{\text{ESM}}(s_{W,\widehat{A}}, t) &= \|(\Omega_t - \Sigma_t^{-1}) \Sigma_t^{1/2}\|_{\text{F}}^2. \end{aligned}$$

For $p \geq d$, the ESM and DSM loss at \widehat{A} are independent of p . By the proof of Lemma 4.1, this is because the empiri-

cal risk minimization (ERM) problem becomes equivalent to the ERM associated with a d -dimensional linear model (independent of p), similar to (Merger & Goldt, 2025). See Section B.2 for details. We specify our setup further in the next assumption.

Assumption 4.1. $W \in \mathbb{R}^{p \times d}$ has $\text{N}(0, 1)$ i.i.d. entries, $A \in \mathbb{R}^{d \times p}$ and $\nu = \text{N}(0, \beta \text{I}_d)$, with $\beta > 0$.

The next proposition is a precise characterization of overfitting in a proportional asymptotic regime.

Proposition 4.1 (Overfitting in linear networks). *Under Assumption 4.1, we consider the limit $d, n, p \rightarrow \infty$ such that $p/n \rightarrow \psi_p$ and $d/n \rightarrow \psi_d$. When $\psi_d, \psi_p > 1$, we have, almost-surely,*

$$\begin{aligned} \frac{1}{d} \mathcal{L}_{\text{ESM}}(s_{W,\widehat{A}}, t) &\rightarrow \\ &\frac{1}{\alpha_t^2 \beta + \sigma_t^2} \left(\frac{\psi_p \wedge \psi_d - 1}{\psi_d} \frac{\alpha_t^4 \beta^2}{\sigma_t^4} + \left(1 - \frac{\psi_p}{\psi_d}\right)_+ \right) + \mathcal{O}\left(\frac{1}{\psi_d}\right), \end{aligned} \quad (6)$$

where $\widehat{A} \in \arg \min \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,A}, t)$. For the empirical risk, we have

$$\frac{\sigma_t^2}{d} \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,\widehat{A}}, t) \rightarrow \frac{1}{\psi_d} + \left(1 - \frac{\psi_p}{\psi_d}\right)_+ + \mathcal{O}\left(\frac{1}{\psi_d(\psi_p \wedge \psi_d)}\right).$$

In this result, we report the empirical DSM loss scaled by σ_t^2 as it is the quantity minimized in practice. The population ESM loss is also given in Equation (6), as it is the quantity that is directly related to the ELBO. When $\psi_p > \psi_d$, we observe that the limit of the empirical risk is of order $\mathcal{O}(\psi_d^{-1})$ and can, therefore, become arbitrarily small and interpolate training data when the data dimension is large. In the same regime, we observe that the population score matching loss is of order σ_t^{-4} and therefore explodes as $t \rightarrow 0^+$. This can also be observed clearly in Figure 3, displaying the exact asymptotics of the risks, obtained in the proofs in Section B.2. It demonstrates that interpolation and generalization are mutually incompatible. We also study the dependence on t in Figure 3, showing that the rate of explosion of the excess risk as $t \rightarrow 0^+$ depends on the model complexity.

Comparison with regression. To understand how these results differ from traditional regression, we consider a multi-output regression problem with the same data and architecture. Given a target $B_\star \in \mathbb{R}^{d \times d}$ and noiseless labels $Y_i := B_\star Z_i$, we minimize the empirical risk $A \mapsto n^{-1} \sum_{i=1}^n \|s_{W,A}(Z_i) - Y_i\|^2 + \lambda \|AW\|_{\text{F}}^2$. We focus on the overparameterized regime $p > d > n$ and $\lambda = 0$, where minimizers A satisfy $(AW - B_\star) \widehat{\Sigma} = 0$. Since $\widehat{\Sigma}$ has rank $n < d$ a.s., this gives a $d(d-n)$ -dimensional affine subspace of solutions for AW , contrasting sharply with Lemma 4.1, where score matching produces a unique value of $AW = -\widehat{\Sigma}_t^{-1}$ even without regularization.

The excess risk for the regression problem is given directly by $\|(AW - B_\star) \Sigma^{1/2}\|_{\text{F}}^2$. The expression has a similar form

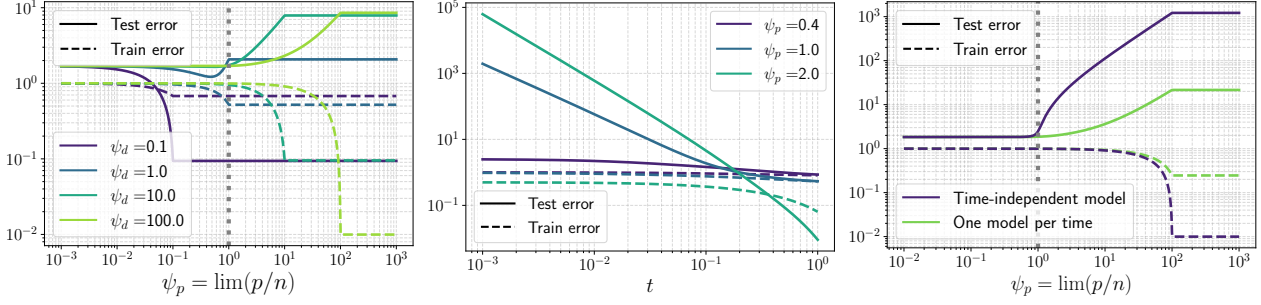


Figure 3. Linear random features under Assumption 4.1. (Left) time is fixed at $t = 10^{-1}$ and ψ_p varies. (Center) ψ_d is fixed at 10^2 and t varies. (Right) time integrated losses (as in Section 5.1) with π the uniform distribution on $[0, T]$. Train error denotes the empirical DSM loss with the weighting used in training in practice. Test error is the population ESM loss with ELBO weighting.

as the score matching ESM loss in Lemma 4.1 once replacing B_\star with $-\Sigma_t^{-1}$, and Σ with Σ_t in the weighting. To resolve the under-determination, we consider the ridgeless limit $\lambda \rightarrow 0^+$, recovering the minimum-norm interpolator. This projects AW on the orthogonal of $\ker(\widehat{\Sigma})$, leading to the excess risk $\|B_\star \Pi_{\ker(\widehat{\Sigma})} \Sigma^{1/2}\|_F^2$, where $\Pi_{\ker(\widehat{\Sigma})}$ is the column-wise projection onto the kernel of $\widehat{\Sigma}$. Thus, generalization can occur in the overparameterized regime due to two factors: (i) the minimum-norm bias projects AW on the orthogonal of $\ker(\widehat{\Sigma})$; (ii) the residual is small when B_\star aligns with $\widehat{\Sigma}^{1/2}$ within the kernel of $\widehat{\Sigma}$. This is made precise in the benign overfitting literature that identify a form of *self-induced regularization* resulting from alignment (Bartlett et al., 2020; Hastie et al., 2022).

The picture is fundamentally different for score matching. First, since the empirical problem exactly determines the solution $AW = -\widehat{\Sigma}_t^{-1}$, there can be no minimum-norm bias that projects AW orthogonally to $\ker(\widehat{\Sigma})$. Moreover, whereas B_\star could reasonably align with $\Sigma^{1/2}$ on the null space in such a way that makes their product small, it is less likely that $-\Sigma_t^{-1}$ aligns in the same way with $\Sigma_t^{1/2}$. Indeed, if Σ is not full rank, we have $\|\Sigma_t^{-1} \Pi_{\ker(\widehat{\Sigma})} \Sigma_t^{1/2}\|_F^2 \geq (d - n)\sigma_t^{-2}$ which explodes as $t \rightarrow 0$. This is made precise in the proof of Proposition 4.1, leading to the first term of Equation (6), corresponding to the contribution of the null space of $\widehat{\Sigma}$ to the ESM loss.

5. Preventing overfitting in diffusion models

The absence of benign overfitting in diffusion models might be seen as a discouraging result. However, in this section, we extend the analysis of Section 4 and show that the score matching framework contains key components that act as implicit regularization mechanisms and help prevent overfitting.

5.1. Implicit regularization due to time smoothness

The analysis of Section 4 is done at a fixed noise scale $t > 0$. Here, we integrate in this analysis the time integral defining the ESM losses. Critically, we study the impact of the time-smoothness of the score estimator, considering the extreme case where it is time-independent to simplify the derivations.

Setup. We parameterize the score network as in Section 4 with $p > d$ and we consider a time weighting ϖ which is a probability distribution on $[0, T]$ such that we can define the probability measure $\pi(dt) \propto \sigma_t^{-2} \varpi(dt)$. For any distribution $\omega \in \{\pi, \varpi\}$ on $[0, T]$, we write

$$\alpha_\omega^2 := \int \alpha_t^2 d\omega(t), \quad \sigma_\omega^2 := \int \sigma_t^2 d\omega(t),$$

$$\Sigma_\omega := \int \Sigma_t d\omega(t), \quad \widehat{\Sigma}_\omega := \int \widehat{\Sigma}_t d\omega(t).$$

In the next lemma, we characterize the empirical risk minimization in this setting.

Lemma 5.1. *Assume that $p > d$ and W is full-rank. Then, any $\widehat{A} \in \arg \min_A \mathcal{L}_{\text{DSM}}(s_{A,W}, \varpi)$ satisfies $\widehat{A}W = -\widehat{\Sigma}_\varpi^{-1}$ and we have*

$$\mathcal{L}_{\text{ESM}}(s_{\widehat{A},W}, \pi) = \|(\widehat{\Sigma}_\varpi^{-1} - \Sigma_\pi^{-1}) \Sigma_\pi^{1/2}\|_F^2 + C_\pi,$$

$$\widehat{\mathcal{L}}_{\text{DSM}}(s_{\widehat{A},W}, \varpi) = d\sigma_\pi^{-2} - \text{Tr}(\widehat{\Sigma}_\varpi^{-1}),$$

where $C_\pi \geq 0$ is given in the proof in Equation (11). C_π is 0 only when π is a Dirac distribution.

Remark 5.1 (Equivalence with ridge). By Lemma 5.1, we have $\widehat{A}W = -(\alpha_\varpi^2 \widehat{\Sigma} + \sigma_\varpi \text{Id})^{-1}$. As ϖ is a probability distribution, there exists $t_\varpi > 0$ such that $\alpha_{t_\varpi}^2 = \alpha_\varpi^2$. Thus, for any $t < t_\varpi$, we have

$$\widehat{\Sigma}_\varpi = \eta_{t,\varpi}^{-1} (\widehat{\Sigma}_t + 2\lambda_{t,\varpi} \text{Id}),$$

$$2\lambda_{t,\varpi} := \eta_{t,\varpi} (\sigma_\varpi^2 - \sigma_t^2) > 0, \quad \eta_{t,\varpi} := \alpha_t^2 / \alpha_\varpi^2.$$

Therefore, up to a time reparameterization by changing $\kappa \geq 0$ in Equation (1) to $\tilde{\kappa} := \kappa \sigma_\varpi^2 \mu_t^2 / (\sigma_t^2 \mu_\varpi^2)$, we see that $\widehat{A} \in$

$\arg \min_A (\widehat{\mathcal{L}}_{\text{DSM}}(s_{A,W}, t) + \lambda_{t,\varpi} \|AW\|^2)$. Therefore, in this case, considering the time-integrated loss is formally similar to adding ridge regularization for small times $t < t_\varpi$, which by Proposition 4.1 are responsible for most of the (not benign) overfitting.

Next, we make this observation more precise by characterizing the asymptotic score matching losses.

Proposition 5.1. *Under Assumption 4.1 and in the limit $n, d, p \rightarrow \infty$ with $\lim(n/d) = \psi_d$ and $\lim(p/n) = \psi_p$. Let $\widehat{A} \in \arg \min_A \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,A}, \varpi)$. When $\psi_p > \psi_d > 1$, we have, almost-surely,*

$$\begin{aligned} \frac{1}{d} \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,\widehat{A}}, \varpi) &\rightarrow \frac{1}{\sigma_\pi^2} - \frac{1}{\sigma_\varpi^2} + \mathcal{O}(\psi_d^{-1}), \\ \frac{1}{d} \mathcal{L}_{\text{ESM}}(s_{W,\widehat{A}}, \pi) &\rightarrow L_\pi + \mathcal{O}(\psi_d^{-1}), \end{aligned}$$

where $L_\pi \geq 0$ is a constant whose expression is given in Equation (14), and $\sigma_\pi^{-2} - \sigma_\varpi^{-2} \geq 0$.

Proposition 5.1 shows that the time-smoothness and the time-integrated loss effectively prevent the explosion of the population loss observed in Proposition 4.1. It also shows that the limiting empirical risk cannot be arbitrarily small as soon as ϖ is not a Dirac distribution (by similar calculations, we reach the same conclusion for $\widehat{\mathcal{L}}_{\text{ESM}}(s_{W,\widehat{A}}, \varpi)$). This contrasts with the fixed noise scale case of Proposition 4.1, even when ψ_d is large, suggesting that the time-smoothness prevents overfitting. Moreover, while in Proposition 4.1 the asymptotic ESM loss could explode as $t \rightarrow 0^+$. In above proposition, this explosion is prevented, and the ESM loss stays bounded, showing that the time-smoothness acts as a form of implicit regularization, as hinted in Remark 5.1.

In the proof of Proposition 5.1 (see Section B.3), we derive exact asymptotics for the population and empirical losses, which we report in the right figure in Figure 3 (green curves). Using the asymptotics of the previous section, we compare it against an overfit linear score network that has independent weights at each time, corresponding to the completely non-time-smooth case (purple curves). We observe that the time smoothness increases the empirical risk while decreasing the population loss, showing that it is indeed acting as an implicit regularization mechanism.

We verify this in a practical image generation setting similar to Figure 2, but considering a time-independent DDPM U-Net model trained across a range $[t, t+r]$ for varying r . We evaluate train and test error at t to see how changing r affects the error. We present the results in Figure 4 where we find that increasing r increases the train error at t and therefore is, indeed, preventing overfitting. Interestingly, the test error first decreases and then increases for large r , suggesting a tradeoff on r .

5.2. Implicit regularization via early stopping

In this part, we study the impact of early stopping, *i.e.*, stopping the training of the score network before reaching overfitting. To formalize this, we fix a time $t \in [0, T]$ and study the dynamics of the parameter A under gradient flow,

$$\frac{d}{d\tau} A(\tau) = -p^{-1} \nabla_A (\sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,A(\tau)}, t)), \quad (7)$$

where $\tau \geq 0$, $A \in \mathbb{R}^{d \times p}$ and $W \in \mathbb{R}^{p \times d}$ as in Section 4, and the p^{-1} factor is ensuring a finite limit as $p \rightarrow \infty$ (Bonnaire et al., 2025). Without loss of generality, we assume that $A(0) = 0$. We show in Section B.4 that in the limit $p \rightarrow \infty$ with (n, d) fixed, the matrix $A(\tau)W$ converges to a finite limit following the gradient flow of a linear model (as shown in Section B.4). Moreover, when $\tau \rightarrow \infty$, $A(\tau)W \rightarrow -\widehat{\Sigma}_t^{-1}$, which leads to the same score matching bounds as in Proposition 4.1 for $p > d$. In the next proposition, we show that small enough values of τ avoid the explosion of the test loss.

Proposition 5.2. *Suppose that Assumption 4.1 holds and define $\mathcal{E}_t(\tau, n, d, p) := \mathcal{L}_{\text{ESM}}(s_{W,A(\tau)}, t)$. Then, $\mathcal{E}_t(\tau, n, d, p)$ has an almost-sure limit $\bar{\mathcal{E}}_t(\tau, n, d) := \lim_{p \rightarrow \infty} \mathcal{L}_{\text{ESM}}(s_{W,A(\tau)}, t)$. Moreover, in the limit $n, d \rightarrow \infty$ with $\lim(d/n) = \psi_d > 1$, we have, almost-surely*

$$\begin{aligned} \frac{1}{d} \bar{\mathcal{E}}_t(\tau, n, d) &\rightarrow \\ &\left(1 - \frac{1}{\psi_d}\right) \frac{\left((\alpha_t^2 \beta + \sigma_t^2) \sigma_t^{-2} (1 - e^{-2\tau \sigma_t^4}) - 1\right)^2}{\alpha_t^2 \beta + \sigma_t^2} + \mathcal{O}\left(\frac{1}{\psi_d}\right), \end{aligned}$$

In particular, if $\tau = \Omega(\sigma_t^{-2})$, then the limit above is bounded as $t \rightarrow 0^+$.

This result shows that early stopping prevents the explosion of the test loss as $t \rightarrow 0$ for a fixed τ . As a corollary of Section 3, it also prevents overfitting. More precisely, we observe that if the early stopping τ is of order σ_t^{-2} , then the population loss does not diverge as $t \rightarrow 0$. This suggests that, despite not being benign, overfitting and memorization appear later in training, allowing early stopping to effectively prevent it, which aligns with recent results (Bonnaire et al., 2025; Merger & Goldt, 2025). The advantages of early stopping can also be observed in high-dimensional experiments: in Figure 5, we report the evolution of the test and train errors for a U-Net across iterations, observing that there is a unique optimal stopping time where the ESM loss is controlled as the number of parameters grow

6. Related works & Conclusion

Random features analysis & linear models. While we focus on quantifying overfitting, we also derive exact asymptotics of the losses in the proof of Proposition 4.1, complementing prior analyses (George & Macris, 2026; George

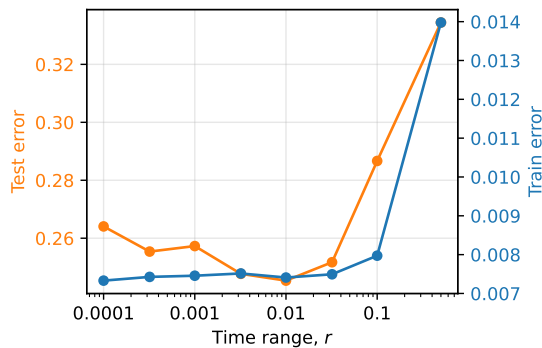


Figure 4. Time-independent DDPM model evaluated at $t = 0.1$. We train it on the range $[t, t + r]$ and observe how increasing r can benefit generalization at time t . See Appendix C.2 for details.

et al., 2025; Merger & Goldt, 2025). As discussed in Section 4, the two-layer linear RFNNs analysis recovers linear diffusion models as a particular case, extending the regularized linear model analysis of (Merger & Goldt, 2025). In particular, we recover the conclusion that overfitting is driven by the null space. Recent work also studied the learning dynamics of diffusion models (Biroli et al., 2024; Merger & Goldt, 2025), including RFNN settings (Bonnaire et al., 2025; Li et al., 2023). Section 5.2 supplements these approaches by quantifying the role of early stopping in mitigating overfitting in linear networks. Overall, we provide a fine-grained characterization of overfitting in linear RFNNs and, crucially, incorporate the **time-smoothness** of the models with time-integrated losses (Section 5.1), which is beyond the scope of prior work.

Diffusion model generalization. In addition to the aforementioned works, there has a lot of recent interest in the generalization properties of diffusion models (Kadkhodaie et al., 2024; Li et al., 2023; Chakraborty et al., 2026). In particular, minimax statistical rates have been obtained (Oko et al., 2023; Zhang et al., 2024; Azangulov et al., 2024), the relationship with data geometry has been explored (Farghly et al., 2026; Pidstrigach, 2022; He et al., 2026) and, more recently, algorithm-dependent generalization bounds were proposed for diffusion models (Dupuis et al., 2025; Farghly et al., 2025). Our work offers a new but complementary perspective, by suggesting that the mechanisms causing generalization in diffusion models are fundamentally different from those in classical settings.

Statistical lower bounds. Recent work analyzed memorization (Buchanan et al., 2025) from a lower bounds perspective. In particular Lemma 3.1 can be compared with (Ye et al., 2026), who have shown that $\mathcal{J}(\hat{p}_t||p_t) = \widehat{\mathcal{L}}_{\text{DSM}}(\nabla \log p_t, t) - \widehat{\mathcal{L}}_{\text{DSM}}(\nabla \log \hat{p}_t, t)$. This is of a different nature than the left-hand side of Equation (4), which is the sum of the empirical and population ESM losses for a generic score. That being said, the lower bounds

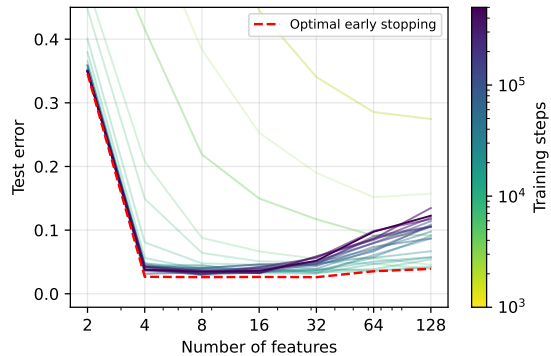


Figure 5. The same setting as Figure 2 but we observe how the test error curve changes with training steps. We find that early stopping produces a test error that does not grow with NF.

on $\mathcal{J}(\hat{p}_t||p_t)$ derived in (Ye et al., 2026) for mixture data could easily be combined with Lemma 3.1. Our results in Section 3.2 differ in that they consider a very generic setting and exploit the intrinsic dimension of the data in statistical guarantees. Note that (4) does not require any restricting assumption on ν , such as log-Sobolev inequality as in (Koehler et al., 2023).

Conclusion. In this work, we showed that benign overfitting does not occur in diffusion models and precisely quantified how harmful it is through generic theoretical and empirical results. We provided a fine-grained analysis of this behavior in a 2-layer linear random feature model, showing its relation with model complexity, dimension, and sample size. Finally, by incorporating the time-integrated losses and gradient dynamics into our experiments and linear network analysis, we demonstrated that these key components of score matching can help prevent overfitting.

Limitations & future works. Our analysis of the mechanisms preventing overfitting in Section 5 is limited to linear networks. Extending this perspective to non-linear networks is an important direction for future research. The results of Propositions 4.1, 5.1 and 5.2 are also limited to simple data distributions, and extending these results to more complex settings is an interesting direction. More generally, while our work shows that benign overfitting cannot happen in diffusion models, the question remains partially open to understand how the data distribution can be well-approximated by practically relevant architectures, which we leave for future work.

References

- Anderson, B. D. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982. ISSN 0304-4149. doi: [https://doi.org/10.1016/0304-4149\(82\)90051-5](https://doi.org/10.1016/0304-4149(82)90051-5). URL <https://www.sciencedirect.com/>

- 440 [science/article/pii/0304414982900515](https://arxiv.org/abs/2409.18804).
- 441
- 442 Azangulov, I., Deligiannidis, G., and Rousseau, J. Conver-
- 443 gence of diffusion models under the manifold hypothe-
- 444 sis in high-dimensions, 2024. URL <https://arxiv.org/abs/2409.18804>.
- 445
- 446 Bakry, D., Gentil, I., and Ledoux, M. *Analysis and Geome-*
- 447 *try of Markov Diffusion Operators*. Springer, 2014.
- 448
- 449 Bartlett, P. L., Long, P. M., Lugosi, G., and Tsigler, A.
- 450 Benign Overfitting in Linear Regression. *Proceedings*
- 451 *of the National Academy of Sciences*, 117(48):30063–
- 452 30070, December 2020. ISSN 0027-8424, 1091-6490.
- 453 doi: 10.1073/pnas.1907378117.
- 454
- 455 Belkin, M., Hsu, D., Ma, S., and Mandal, S. Reconciling
- 456 modern machine learning practice and the bias-variance
- 457 trade-off. *Proceedings of the National Academy of Sci-*
- 458 *ences*, 116(32):15849–15854, August 2019. ISSN 0027-
- 459 8424, 1091-6490. doi: 10.1073/pnas.1903070116.
- 460
- 461 Benton, J., Bortoli, V. D., Doucet, A., and Deligianni-
- 462 dis, G. Nearly d-linear convergence bounds for diffu-
- 463 sion models via stochastic localization. In *The Twelfth*
- 464 *International Conference on Learning Representations*,
- 465 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=r5njv3BsuD)
- 466 [id=r5njv3BsuD](https://openreview.net/forum?id=r5njv3BsuD).
- 467
- 468 Biroli, G., Bonnaire, T., de Bortoli, V., and Mézard, M.
- 469 Dynamical regimes of diffusion models. *Nature Commu-*
- 470 *nications*, 15, 2024.
- 471
- 472 Bonnaire, T., Urfin, R., Biroli, G., and Mezard, M. Why Dif-
- 473 fusion Models Don’t Memorize: The Role of Implicit Dy-
- 474 namical Regularization in Training. In *The Thirty-ninth*
- 475 *Annual Conference on Neural Information Processing*
- 476 *Systems*, 2025.
- 477
- 478 Bortoli, V. D. Convergence of denoising diffusion
- 479 models under the manifold hypothesis. *Transactions*
- 480 *on Machine Learning Research*, 2022. ISSN 2835-
- 481 8856. URL [https://openreview.net/forum?](https://openreview.net/forum?id=MhK5aXo3gB)
- 482 [id=MhK5aXo3gB](https://openreview.net/forum?id=MhK5aXo3gB). Expert Certification.
- 483
- 484 Buchanan, S., Pai, D., Ma, Y., and Bortoli, V. D. On the
- 485 edge of memorization in diffusion models, 2025. URL
- 486 <https://arxiv.org/abs/2508.17689>.
- 487
- 488 Carlini, N., Hayes, J., Nasr, M., Jagielski, M., Sehwa-
- 489 g, V., Tramèr, F., Balle, B., Ippolito, D., and Wallace, E.
- 490 Extracting training data from diffusion models, 2023.
- 491
- 492 Chafai, D. Entropies, convexity, and functional inequali-
- 493 ties. *Kyoto Journal of Mathematics*, 44(2), January 2004.
- 494 ISSN 2156-2261. doi: 10.1215/kjm/1250283556.
- Chakraborty, S., Berthet, Q., and Bartlett, P. L. Gener-
alization properties of score-matching diffusion mod-
els for intrinsically low-dimensional data, 2026. URL
<https://arxiv.org/abs/2603.03700>.
- Conforti, G., Durmus, A., and Silveri, M. G. Kl convergence
guarantees for score diffusion models under minimal data
assumptions. *SIAM Journal on Mathematics of Data
Science*, 7(1):86–109, 2025.
- Couillet, R. and Liao, Z. *Random Matrix Methods for Ma-*
chine Learning. Cambridge University Press, 1 edition,
July 2022. doi: 10.1017/9781009128490.
- Dhariwal, P. and Nichol, A. Diffusion models beat gans
on image synthesis. In *Advances in Neural Information
Processing Systems*, volume 34, pp. 8780–8794. Curran
Associates, Inc., 2021.
- Dupuis, B., Shariatian, D., Haddouche, M., Durmus, A., and
Simsekli, U. Algorithm- and data-dependent generaliza-
tion bounds for score-based generative models. In *39th
Conference on Neural Information Processing Systems
(NeurIPS 2025)*, 2025.
- Esser, P., Kulal, S., Blattmann, A., Entezari, R., Müller,
J., Saini, H., Levi, Y., Lorenz, D., Sauer, A., Boesel, F.,
Podell, D., Dockhorn, T., English, Z., and Rombach, R.
Scaling rectified flow transformers for high-resolution
image synthesis. In *Forty-first International Conference
on Machine Learning*, 2024.
- Falconer, K. *Fractal Geometry - Mathematical Foundations
and Applications*. Wiley, 2014.
- Farghly, T., Rebeschini, P., Deligiannidis, G., and Doucet, A.
Implicit regularisation in diffusion models: An algorithm-
dependent generalisation analysis, 2025. URL <https://arxiv.org/abs/2507.03756>.
- Farghly, T., Potapchik, P., Howard, S., Deligiannidis, G.,
and Pidstrigach, J. Diffusion models and the manifold
hypothesis: Log-domain smoothing is geometry adap-
tive. In *The Thirty-ninth Annual Conference on Neural
Information Processing Systems*, 2026.
- George, A. J. and Macris, N. Asymptotic Learning Curves
for Diffusion Models with Random Features Score and
Manifold Data, March 2026.
- George, A. J., Veiga, R., and Macris, N. Denoising Score
Matching with Random Features: Insights on Diffusion
Models from Precise Learning Curves. *arXiv preprint at
arXiv:2502.00336*, 2025.
- Hastie, T., Montanari, A., Rosset, S., and Tibshirani, R. J.
Surprises in high-dimensional ridgeless least squares in-
terpolation. *The Annals of Statistics*, 50(2):949–986,
2022.

- 495 Haussmann, U. G. and Pardoux, E. Time reversal of diffu-
 496 sions. *The Annals of Probability*, 14(4):1188–1205, 1986.
 497 doi: 10.1214/aop/1176992362.
 498
- 499 He, Y., Qiu, Y., and Tao, M. Diffusion model’s generaliza-
 500 tion can be characterized by inductive biases toward a
 501 data-dependent ridge manifold. *arXiv [stat.ML]*, Febru-
 502 ary 2026.
 503
- 504 Ho, J., Jain, A., and Abbeel, P. Denoising diffusion proba-
 505 bilistic models. In Larochelle, H., Ranzato, M., Hadsell,
 506 R., Balcan, M., and Lin, H. (eds.), *Advances in Neural
 507 Information Processing Systems*, volume 33, pp. 6840–
 508 6851. Curran Associates, Inc., 2020.
 509
- 510 Hyvärinen, A. Estimation of non-normalized statistical
 511 models by score matching. *Journal of Machine Learning
 512 Research*, 6(24):695–709, 2005.
 513
- 514 Kadkhodaie, Z., Guth, F., Simoncelli, E. P., and Mallat, S.
 515 Generalization in diffusion models arises from geometry-
 516 adaptive harmonic representations. In *The Twelfth Inter-
 517 national Conference on Learning Representations*, 2024.
 518
- 519 Kingma, D. P. and Ba, J. Adam: A method for stochastic
 520 optimization. In *ICLR (Poster)*, 2015. URL [http://
 521 arxiv.org/abs/1412.6980](http://arxiv.org/abs/1412.6980).
 522
- 523 Koehler, F., Heckett, A., and Risteski, A. Statistical effi-
 524 ciency of score matching: The view from isoperimetry.
 525 In *The Eleventh International Conference on Learning
 526 Representations*, 2023.
 527
- 528 Krizhevsky, A. Learning multiple layers of features from
 529 tiny images. *Master’s thesis, University of Tront*, 2009.
 530
- 531 Li, P., Li, Z., Zhang, H., and Bian, J. On the generalization
 532 properties of diffusion models. In *Advances in Neural
 533 Information Processing Systems*, volume 36, 2023.
 534
- 535 Marchenko, V. A. and Pastur, L. A. Distribution of eigen-
 536 values for some sets of random matrices. *Mathematics of
 537 the USSR-Sbornik*, 1:457–483, 1967.
 538
- 539 Mattila, P., Moran, M., and Rey, J.-m. Dimension of a
 540 measure. *Studia mathematica* 142 (3), 2000.
 541
- 542 Mei, S. and Montanari, A. The Generalization Error of
 543 Random Features Regression: Precise Asymptotics and
 544 the Double Descent Curve. *Communications on Pure and
 545 Applied Mathematics*, 75(4):667–766, April 2022. ISSN
 546 0010-3640, 1097-0312. doi: 10.1002/cpa.22008.
 547
- 548 Merger, C. and Goldt, S. Generalization dynamics of linear
 549 diffusion models. In *EurIPS 2025 Workshop on Princi-
 550 ples of Generative Modeling (PriGM)*, 2025.
- 551 Nakkiran, P., Kaplun, G., Bansal, Y., Yang, T., Barak, B.,
 552 and Sutskever, I. Deep Double Descent: Where Bigger
 553 Models and More Data Hurt. *ICLR 2020*, December
 554 2019.
- 555 Oko, K., Akiyama, S., and Suzuki, T. Diffusion models are
 556 minimax optimal distribution estimators. In *Proceedings
 557 of the 40th International Conference on Machine Learn-
 558 ing*, volume 202, pp. 26517–26582. PMLR, 23–29 Jul
 559 2023.
- 560 Øksendal, B. *Stochastic Differential Equations*. Universi-
 561 text. Springer, Berlin, Heidelberg, 2003. doi: 10.1007/
 562 978-3-642-14394-6.
- 563 Otto, F. and Villani, C. Comment on: “Hypercontractivity
 564 of Hamilton–Jacobi equations”, by S. Bobkov, I. Gentil
 565 and M. Ledoux. *Journal de Mathématiques Pures et
 566 Appliquées*, 80(7):697–700, 2001.
- 567 Pesin, Y. B. *Dimension Theory in Dynamical Systems -
 568 Contemporary Views and Applications*. Chicago Lectures
 569 in Mathematics. The University of Chicago Press, 1997.
- 570 Pesin, Ya. B. On rigorous mathematical definitions of cor-
 571 relation dimension and generalized spectrum for dimen-
 572 sions. *Journal of Statistical Physics*, 71(3):529–547, May
 573 1993. ISSN 1572-9613. doi: 10.1007/BF01058436.
- 574 Pidstrigach, J. Score-based generative models detect mani-
 575 folds. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave,
 576 D., Cho, K., and Oh, A. (eds.), *Advances in Neural Infor-
 577 mation Processing Systems*, volume 35, pp. 35852–35865.
 578 Curran Associates, Inc., 2022.
- 579 Procaccia, I., Grassberger, P., and Hentschel, H. G. E. On the
 580 characterization of chaotic motions. In *Dynamical Sys-
 581 tems and Chaos*, volume 179 of *Lecture Notes in Physics*,
 582 pp. 212–222, Berlin, 1983. Springer. Proceedings of the
 583 conference held in Sitges/Barcelona, 1982.
- 584 Rahimi, A. and Recht, B. Random features for large-scale
 585 kernel machines. In *Advances in Neural Information
 586 Processing Systems*, volume 20, 2007.
- 587 Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and
 588 Ommer, B. High-resolution image synthesis with la-
 589 tent diffusion models. In *Proceedings of the IEEE/CVF
 590 Conference on Computer Vision and Pattern Recognition
 591 (CVPR)*, pp. 10684–10695, June 2022.
- 592 Ronneberger, O., Fischer, P., and Brox, T. U-net: Convolu-
 593 tional networks for biomedical image segmentation. In
 594 *International Conference on Medical image computing
 595 and computer-assisted intervention*, pp. 234–241, 2015.
- 596 Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton,
 597 E. L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan,

- 550 B., Salimans, T., Ho, J., Fleet, D., and Norouzi, M. Photo-
551 realistic text-to-image diffusion models with deep lan-
552 guage understanding. In *Advances in Neural Information*
553 *Processing Systems*, volume 35, 2022.
- 554 Shalev-Shwartz, S. and Ben-David, S. *Understanding Ma-*
555 *chine Learning: From Theory to Algorithms*. Cambridge
556 University Press, 2014.
- 558 Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and
559 Ganguli, S. Deep unsupervised learning using nonequi-
560 librium thermodynamics. In *Proceedings of the 32nd*
561 *International Conference on Machine Learning*, pp. 2256–
562 2265, 2015.
- 564 Song, Y. and Ermon, S. Generative modeling by estimating
565 gradients of the data distribution. In *Advances in Neural*
566 *Information Processing Systems*, volume 32, 2019.
- 567 Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Er-
568 mon, S., and Poole, B. Score-based generative modeling
569 through stochastic differential equations. In *International*
570 *Conference on Learning Representations*, 2021.
- 572 Vincent, P. A Connection Between Score Matching and
573 Denoising Autoencoders. *Neural Computation*, 23(7):
574 1661–1674, July 2011. ISSN 0899-7667, 1530-888X.
575 doi: 10.1162/NECO_a.00142.
- 577 Watson, J. L., Juergens, D., Bennett, N. R., Trippe, B. L.,
578 Yim, J., Eisenach, H. E., Ahern, W., Borst, A. J., Ragotte,
579 R. J., Milles, L. F., Wicky, B. I. M., Hanikel, N., Pellock,
580 S. J., Courbet, A., Sheffler, W., Wang, J., Venkatesh,
581 P., Sappington, I., Torres, S. V., Lauko, A., De Bortoli,
582 V., Mathieu, E., Ovchinnikov, S., Barzilay, R., Jaakkola,
583 T. S., DiMaio, F., Baek, M., and Baker, D. De novo
584 design of protein structure and function with RFDiffusion.
585 *Nature*, 620(7976):1089–1100, August 2023. ISSN 1476-
586 4687. doi: 10.1038/s41586-023-06415-8.
- 587 Yang, L., Zhang, Z., Song, Y., Hong, S., Xu, R., Zhao, Y.,
588 Zhang, W., Cui, B., and Yang, M.-H. Diffusion models:
589 A comprehensive survey of methods and applications.
590 *ACM Comput. Surv.*, 56(4), 2023.
- 592 Ye, Z., Zhu, Q., Tao, M., and Chen, M. Provable separations
593 between memorization and generalization in diffusion
594 models. In *The Fourteenth International Conference on*
595 *Learning Representations*, 2026.
- 597 Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals,
598 O. Understanding deep learning requires rethinking gen-
599 eralization. In *International Conference on Learning*
600 *Representations (ICLR)*, 2017.
- 601 Zhang, K., Yin, H., Liang, F., and Liu, J. Minimax op-
602 timality of score-based diffusion models: Beyond the
603 density lower bound assumptions. In Salakhutdinov, R.,
604
- Kolter, Z., Heller, K., Weller, A., Oliver, N., Scarlett,
J., and Berkenkamp, F. (eds.), *Proceedings of the 41st*
International Conference on Machine Learning, volume
235 of *Proceedings of Machine Learning Research*, pp.
60134–60178. PMLR, 21–27 Jul 2024.

A. Additional Technical Background

A.1. Intrinsic dimensions

Our main results in Section 3 are expressed in terms of the Rényi dimension of a Borel measure μ on \mathbb{R}^d . This notion was first developed by (Procaccia et al., 1983) and then rigorously treated in (Pesin, 1997; 1993). It is a standard notion in the theory of chaotic dynamical systems. We provide below the definition.

Definition A.1 (Rényi dimension of a measure). Let μ be a Borel probability measure on \mathbb{R}^d . The upper Rényi dimension of μ is defined by

$$\bar{\gamma}(\mu) := \limsup_{\delta \rightarrow 0} \left(\frac{1}{\log \delta} \log \left(\int \mu(B(x, \delta)) d\mu(x) \right) \right).$$

Similarly, the lower Rényi dimension of μ is defined by

$$\underline{\gamma}(\mu) := \liminf_{\delta \rightarrow 0} \left(\frac{1}{\log \delta} \log \left(\int \mu(B(x, \delta)) d\mu(x) \right) \right).$$

Remark A.1. The definition above is actually a particular case of the larger family of Rényi dimension, where it actually corresponds to the Rényi dimension of order 2. Note that Definition A.1 is also sometimes called the *correlation dimension* in the literature (Pesin, 1993).

Remark A.2. In particular, any measure that is supported on a smooth manifold and is Ahlfors regular (see (Mattila et al., 2000; Falconer, 2014)) on this manifold, then the Rényi dimension of Definition A.1 is equal to the manifold dimension.

A.2. Random covariance matrices

We recall here some basic random matrix theory results related to random covariance matrices and the Marchenko-Pastur theorem. We invite the reader to consult (Couillet & Liao, 2022) for additional details.

Given a symmetric matrix $\Sigma \in \mathbb{R}^{d \times d}$, we denote its spectrum by $\text{Spec}(\Sigma)$ (where the eigenvalues are counted with their multiplicity). The spectral measure of Σ is defined by

$$\hat{\mu}_{\Sigma} := \frac{1}{d} \sum_{\lambda \in \text{Spec}(\Sigma)} \delta_{\lambda}. \quad (8)$$

We recall below the celebrated Marchenko-Pastur theorem (Marchenko & Pastur, 1967). To this end, we first define the Marchenko-Pastur distribution below.

Definition A.2 (Marchenko-Pastur distribution). Let $\psi > 0$. The Marchenko-Pastur distribution with shape parameter ψ , denoted here $\mu_{\text{MP}}^{(\psi)}$, is the probability distribution on \mathbb{R}_+ defined by

$$\mu_{\text{MP}}^{(\psi)}(dx) := \left(1 - \frac{1}{\psi}\right)_+ \delta_0 + \frac{\sqrt{(\psi^+ - x)(x - \psi^-)}}{2\pi\psi x} \mathbb{1}_{[\psi^-, \psi^+]}(x) dx,$$

where $y_+ := \max(0, y)$, $\psi^- = |1 - \sqrt{\psi}|^2$, and $\psi^+ := (1 + \sqrt{\psi})^2$.

The Marchenko-Pastur's theorem is given below.

Theorem A.1 (Marchenko-Pastur's theorem). Let X be random vector \mathbb{R}^d with i.i.d. components with zero mean, unit variance, and uniformly bounded moments¹ of order $4 + \epsilon$ for some $\epsilon > 0$. Define the empirical covariance matrix,

$$\hat{\Sigma} := \frac{1}{n} \sum_{i=1}^n X_i X_i^T,$$

with X_1, \dots, X_n i.i.d. copies of X . Consider that $d = d_n$ with $\lim_{n \rightarrow \infty} (d_n/n) = \psi > 0$. Then, with probability one, the empirical measure $\hat{\mu}_{\hat{\Sigma}}$ converges weakly to the Marchenko-Pastur distribution with shape parameter ψ .

¹We refer to Theorem 2.4 of (Couillet & Liao, 2022) for a discussion on the different conditions required by the Marchenko-Pastur theorem.

B. Omitted Proofs

In this section we present all the omitted proofs of our main results of Sections 3 to 5.

B.1. Omitted proofs of Section 3

B.1.1. PROOF OF LEMMA 3.1.

We present below the proof of Lemma 3.1.

Proof. Let $S := (Z_1, \dots, Z_n) \sim \nu^{\otimes n}$ be a dataset sampled from the data distribution. Let $(\widehat{X}_t)_{t \geq 0}$ the Ornstein-Uhlenbeck process $d\widehat{X}_t = -\kappa \widehat{X}_t dt + \sqrt{2} dB_t$ initialized at the empirical data distribution $\widehat{X}_0 \sim \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$. We denote by \hat{p}_t the probability density of \widehat{X}_t , conditioned on $S := (Z_1, \dots, Z_n)$. Let $s : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ be an arbitrary measurable score function. By definition of the Fisher information and the triangle inequality (note that the expectation is over both the noise and the dataset), we have, for any $t > 0$,

$$\begin{aligned} \mathbb{E} [\mathcal{J}(\hat{p}_t | p_t)] &= \mathbb{E} \left[\left\| \nabla \log \hat{p}_t(\widehat{X}_t) - \nabla \log p_t(\widehat{X}_t) \right\|^2 \right] \\ &\leq 2\mathbb{E} \left[\left\| \nabla \log \hat{p}_t(\widehat{X}_t) - s(t, \widehat{X}_t) \right\|^2 \right] + 2\mathbb{E} \left[\left\| s(t, \widehat{X}_t) - \nabla \log p_t(\widehat{X}_t) \right\|^2 \right] \\ &= 2\mathbb{E} \left[\left\| \nabla \log \hat{p}_t(\widehat{X}_t) - s(t, \widehat{X}_t) \right\|^2 \right] + 2\mathbb{E} \left[\left\| s(t, X_t) - \nabla \log p_t(X_t) \right\|^2 \right], \end{aligned}$$

where the last equality follows as $\text{law}(X_t)$ is the marginal distribution of \widehat{X}_t under $S \sim \nu^{\otimes n}$. We immediately deduce that for any $t > 0$, we have

$$\mathbb{E} [\mathcal{J}(\hat{p}_t | p_t)] \leq 2\mathbb{E} \left[\mathcal{L}_{\text{ESM}}(s, t) + \widehat{\mathcal{L}}_{\text{ESM}}(s, t) \right].$$

For the second part of the statement, by the triangle inequality, we have that

$$\mathbb{E} [\mathcal{J}(\hat{p}_t | p_t)] = \int_{\mathbb{R}^d} \mathbb{E} \left[\left\| \nabla \log \frac{\hat{p}_t}{p_t}(x) \right\|^2 \hat{p}_t(x) \right] dx \geq \frac{1}{2} \mathbb{E} [\mathcal{J}(\hat{p}_t)] - \mathcal{J}(p_t).$$

Note that $\mathcal{J}(p_t) \leq \mathcal{J}(\nu) < +\infty$ under our assumptions. Then, the fact that $\mathbb{E} [\mathcal{J}(\hat{p}_t)] \rightarrow +\infty$ as $t \rightarrow 0^+$ follows from standard arguments. For the sake of completeness, let us provide a short proof in the case $\kappa = 0$ (without loss of generality, the other cases follow by the same computation by using the invariant distribution of the forward process instead of the Lebesgue measure as a reference measure). Let $h(f) := -\int f(x) \log f(x) dx$ denote the entropy functional, as soon as it is well-defined. As a consequence of the De-Bruijn's identity and the fact that the Fisher information is non-increasing along the semigroup (Bakry et al., 2014), we have, for $t > s > 0$, that

$$h(\hat{p}_t) - h(\hat{p}_s) = \int_s^t \mathcal{J}(\hat{p}_u) du \leq (t - s) \mathcal{J}(\hat{p}_s).$$

Given $(Z_1, \dots, Z_n) \sim \nu^{\otimes n}$, we easily see that

$$-h(\hat{p}_s) \geq -\frac{1}{n} \sum_{i=1}^n h(f_{Z_i, \sigma_t}) - \log(n) = -\frac{d}{2} \log(2\pi e) + \frac{d}{2} \log(\sigma_t^{-2}) - \log(n),$$

where $f_{\mu, \sigma}$ is the density of $N(\mu, \sigma I_d)$. Thus $\mathbb{E} [\mathcal{J}(\hat{p}_s)] \rightarrow +\infty$ as $s \rightarrow 0^+$.

This concludes the proof. \square

Remark B.1. As an immediate consequence, for any positive Borel measure ϖ on $[0, T]$, we have

$$\mathbb{E} \left[\mathcal{L}_{\text{ESM}}(s, \varpi) + \widehat{\mathcal{L}}_{\text{ESM}}(s, \varpi) \right] \geq \frac{1}{2} \int_{[0, T]} \mathbb{E} [\mathcal{J}(\hat{p}_t | p_t)] d\varpi(t).$$

B.1.2. LOWER BOUND UNDER LOG-SOBOLEV INEQUALITIES

In this subsection, we present some additional results to complement Lemma 3.1. Our goal is to show that, under the assumption that the data distribution satisfies a logarithmic Sobolev inequality, we can obtain a similar lower bound than Lemma 3.1, where the relative Fisher information term can be replaced by other divergences, such as the total variation distance.

Conventions for probabilistic inequalities. We define the total variation distance as

$$\text{TV}(\mu, \nu) := \sup_A |\mu(A) - \nu(A)|.$$

We say that a Borel probability distribution ν satisfies the log-Sobolev inequality with constant ρ if for all differentiable $f \in L^1(\nu)$, we have

$$\text{Ent}_\nu(f) \leq \frac{\rho}{2} \int \frac{\|\nabla f\|^2}{f} d\nu.$$

Equipped with these definitions, we have the following proposition.

Proposition B.1. *Assume that the data distribution satisfies the log-Sobolev inequality with constant ρ_0 . Let ϖ be a positive Borel measure on $[0, T]$, we have*

$$\mathbb{E} \left[\mathcal{L}_{\text{ESM}}(\theta, \varpi) + \widehat{\mathcal{L}}_{\text{ESM}}(\theta, \varpi) \right] \geq \int_{[0, T]} \frac{2}{\mu_t^2 \rho_0 + \sigma_t^2} \text{TV}(p_t, \hat{p}_t)^2 d\varpi(t).$$

Proof. By stability of the log-Sobolev inequality under Lipschitz mappings and convolution (Chafai, 2004), we have that p_t satisfies the log-Sobolev inequality with constant $\rho_t := \mu_t^2 \rho_0 + \sigma_t^2$. Therefore, for any $t > 0$, we have

$$\begin{aligned} \frac{1}{2} \mathcal{J}(\hat{p}_t \| p_t) &\geq \frac{1}{\mu_t^2 \rho_0 + \sigma_t^2} \text{KL}(\hat{p}_t \| p_t) \quad (\text{log-Sobolev inequality}) \\ &\geq \frac{2}{\mu_t^2 \rho_0 + \sigma_t^2} \text{TV}(\hat{p}_t, p_t)^2 \quad (\text{Pinsker's inequality}). \end{aligned}$$

The result immediately follows by integrating over ϖ and applying Lemma 3.1. \square

Remark B.2. We immediately see that, if $\alpha > 0$, we have,

$$\mathbb{E} \left[\mathcal{L}_{\text{ESM}}(\theta, \varpi) + \widehat{\mathcal{L}}_{\text{ESM}}(\theta, \varpi) \right] \geq \frac{2}{\max(\alpha^{-1}, \rho_0)} \int_{[0, T]} \text{TV}(p_t, \hat{p}_t)^2 d\varpi(t).$$

B.1.3. PROOF OF THEOREM 3.1

Before giving the proof of Theorem 3.1, we provide some technical lemmas below. Recall that we denote by $p_{t|0}(\cdot|x)$ the conditional density of X_t given $X_0 = x$, defined by Equation (2).

Lemma B.1. *Given any $x \in \mathbb{R}^d$ and $Y \sim p_{\epsilon|0}(\cdot|x)$, we have that for any $\Delta \geq 0$,*

$$\begin{aligned} \mathbb{P}(\|Y - \mu_\epsilon x\| \geq \sigma_\epsilon m + \Delta) &\leq \exp\left(-\frac{\Delta^2}{2\sigma_\epsilon^2}\right), \\ \mathbb{P}(\|Y - \mu_\epsilon x\| \leq \sigma_\epsilon m - \Delta) &\leq \exp\left(-\frac{\Delta^2}{2\sigma_\epsilon^2}\right), \end{aligned}$$

where $m = \mathbb{E}_{Z \sim \mathcal{N}(0, I_d)}[\|Z\|]$.

Proof. By definition, we have that the random variable $Z = \sigma_\epsilon^{-1}(Y - \mu_\epsilon X)$ is a standard multivariate Gaussian. Thus, by Gaussian concentration of measure (see e.g., Section 5.4.2 of (Bakry et al., 2014)), we have that for any 1-Lipschitz function f ,

$$\mathbb{P}(f(Z) - \mathbb{E}[f(Z)] \geq t) \leq \exp(-t^2/2).$$

Choosing $f(z) = \|z\|$ and $f(z) = -\|z\|$ along with $t = \Delta/\sigma_\epsilon$. \square

The proof of Theorem 3.1 is based on the following proposition, providing a lower bound on the sample size, provided that two KL divergences associated to the population and empirical forward processes are both smaller than an absolute constant. Below, we use the following convention and notation for the total variation distance:

$$\|\mu - \nu\|_{\text{TV}} := \sup_A |\mu(A) - \nu(A)|, \quad (9)$$

for probability measure μ and ν . We also recall that we denote by p_t the probability density of X_t following Equation (1) and by \hat{p}_t the probability density of the process $d\hat{X}_t = -\kappa\hat{X}_t dt + \sqrt{2}dB_t$ initialized at $\hat{X}_0 \sim \frac{1}{n} \sum_{i=1}^n \delta_{Z_i}$, for $t > 0$, where $S := (Z_1, \dots, Z_n) \sim \nu^{\otimes n}$ is the dataset.

Proposition B.2. *Let $\epsilon > 0$. For any distribution $\hat{\nu}$ on \mathbb{R}^d , which might depend on the dataset S , if $\mathbb{E}[\text{KL}(\hat{p}_\epsilon|\hat{\nu})] \leq 1/16$ and $\mathbb{E}[\text{KL}(p_\epsilon|\hat{\nu})] \leq 1/16$, then we have $n \geq n_{\min}(\epsilon)$, with, for ϵ small enough*

$$\log n_{\min}(\epsilon) \geq -\frac{1}{2} \log \left(8\sqrt{d} \vee \log(64)\sigma_\epsilon \right) \underline{\gamma}(\nu),$$

where $\underline{\gamma}(\nu)$ is the lower correlation dimension defined in Definition A.1.

Proof of Proposition B.2. We begin with Pinsker's inequality and the reverse triangle inequality for the total variation distance, which are used to obtain

$$\begin{aligned} \mathbb{E}[\text{KL}(p_\epsilon|\hat{\nu})] &\geq 2\mathbb{E}[\|p_\epsilon - \hat{\nu}\|_{\text{TV}}^2] \\ &\geq 2\mathbb{E}\left[\left(\|p_\epsilon - \hat{p}_\epsilon\|_{\text{TV}} - \|\hat{p}_\epsilon - \hat{\nu}\|_{\text{TV}}\right)^2\right] \\ &\geq 2\left(\sqrt{\mathbb{E}[\|p_\epsilon - \hat{p}_\epsilon\|_{\text{TV}}^2]} - \sqrt{\mathbb{E}[\|\hat{p}_\epsilon - \hat{\nu}\|_{\text{TV}}^2]}\right)^2, \end{aligned}$$

where the final inequality follows from the triangle inequality in weighted L^2 -norm. Since, we have $2\mathbb{E}[\|\hat{p}_\epsilon - \hat{\nu}\|_{\text{TV}}^2] \leq \mathbb{E}[\text{KL}(\hat{p}_\epsilon|\hat{\nu})] \leq 1/16$, we conclude from this that if $\mathbb{E}[\text{KL}(p_\epsilon|\hat{\nu})] \leq 1/16$ then $\mathbb{E}[\|p_\epsilon - \hat{p}_\epsilon\|_{\text{TV}}^2] \leq \frac{1}{8}$.

We now identify a test set to lower bound the total variation using that for any Borel set A ,

$$\|p_\epsilon - \hat{p}_\epsilon\|_{\text{TV}} \geq |p_\epsilon(A) - \hat{p}_\epsilon(A)|.$$

As in Lemma B.1, we define $m = \mathbb{E}_{Z \sim N(0, I_d)}[\|Z\|]$. We define the set $A_R := \{y \in \mathbb{R}^d : y \notin B_{\sigma_\epsilon m + R}(\mu_\epsilon Z_i), \forall i = 1, \dots, n\}$ so that the empirical density is upper bounded by

$$\hat{p}_\epsilon(A_R) \leq \frac{1}{N} \sum_{i=1}^N p_{\epsilon|0}(B_{m\sigma_\epsilon + R}(\mu_\epsilon Z_i)^c | Z_i) = 1 - p_{\epsilon, R},$$

where $p_{\epsilon, R} := p_{\epsilon|0}(B_{m\sigma_\epsilon + R}(\mathbf{0})|0)$. Define the set

$$\tilde{A}_R := \{y \in \mathbb{R}^d : \|y - Z_i\| \geq 2(R + m\sigma_\epsilon)/\mu_\epsilon, \forall i = 1, \dots, n\}.$$

The population density is lower bounded using that for any $x \in \tilde{A}_R$, we have $B_{m\sigma_\epsilon + R}(\mu_\epsilon x) \subset A_R$ and hence,

$$\begin{aligned} p_\epsilon(A_R) &\geq \int_{\tilde{A}_R} p_{\epsilon|0}(A_R|x)\nu(dx) \\ &\geq \int_{\tilde{A}_R} p_{\epsilon|0}(B_{m\sigma_\epsilon + R}(\mu_\epsilon x)|x)\nu(dx) \\ &= \nu(\tilde{A}_R)p_{\epsilon, R}. \end{aligned}$$

We further bound this using the union bound,

$$\begin{aligned} \mathbb{E}_S[\nu(\tilde{A}_R)] &= 1 - \mathbb{E}_S \left[\mathbb{P}_{X \sim \nu} \left(\bigcup_{i=1}^n \{\|X - x_i\| < 2(R + m\sigma_\epsilon)/\mu_\epsilon\} \right) \right] \\ &\geq \left(1 - \sum_{i=1}^n \mathbb{E}[\mathbb{P}_{X \sim \nu}(\|X - x_i\| < 2(R + m\sigma_\epsilon)/\mu_\epsilon)] \right)_+ \\ &= (1 - n\delta(2(R + m\sigma_\epsilon)/\mu_\epsilon))_+, \end{aligned}$$

where we define $\delta(r) = \mathbb{E}_{X \sim \nu}[\nu(B_r(X))]$. The total variation is then lower bounded using A_R as a test set: Correct equation:

$$\mathbb{E}[\|p_\epsilon - \hat{p}_\epsilon\|_{\text{TV}}] \geq \mathbb{E}[p_\epsilon(A_R) - \hat{p}_\epsilon(A_R)] \geq \left(\left(1 - n\delta \left(\frac{2(R + m\sigma_\epsilon)}{\mu_\epsilon} \right) \right)_+ p_{\epsilon,R} - (1 - p_{\epsilon,R}) \right)_+.$$

To guarantee that this bound is non-trivial, we now choose R using Lemma B.1. Setting $R^2 = 2 \log(8)\sigma_\epsilon^2$, leads to the bound,

$$p_{\epsilon,R} \geq 1 - \exp \left(- \frac{R^2}{2\sigma_\epsilon^2} \right) \geq \frac{3}{4}.$$

We also choose ϵ such that $\mu_\epsilon \geq 1/2$ so that $2(R + m\sigma_\epsilon)/\mu_\epsilon \leq 8\sqrt{d \vee \log(64)}\sigma_\epsilon$ and $8\sqrt{d \vee \log(64)}\sigma_\epsilon < 1$.

Thus, by Jensen's inequality, we have the lower bound,

$$\mathbb{E}[\|p_\epsilon - \hat{p}_\epsilon\|_{\text{TV}}^2]^{1/2} \geq \frac{1}{4} \left(3(1 - n\delta(8\sqrt{d \vee \log(64)}\sigma_\epsilon))_+ - 1 \right)_+.$$

From this bound we we deduce that,

$$\begin{aligned} \mathbb{E}[\|p_\epsilon - \hat{p}_\epsilon\|_{\text{TV}}^2] &\leq \frac{1}{8} \implies (3(1 - n\delta(8\sqrt{d \vee \log(64)}\sigma_\epsilon))_+ \leq 1 + \sqrt{2} \\ &\implies n \geq \frac{2 - \sqrt{2}}{3\delta(8\sqrt{d \vee \log(64)}\sigma_\epsilon)} =: n_{\min}(\epsilon). \end{aligned}$$

By the definition of the lower Rényi dimension in Definition A.1, we have that

$$\underline{\gamma}(\nu) = \liminf_{\epsilon \rightarrow 0^+} \frac{\log \delta(8\sqrt{d \vee \log(64)}\sigma_\epsilon)}{\log(8\sqrt{d \vee \log(64)}\sigma_\epsilon)}.$$

Let $\eta > 0$ be fixed, by definition, we have that there exists ϵ_0 such that, for all $\epsilon < \epsilon_0$, we have

$$\log n_{\min}(\epsilon) \geq -\log \left(\frac{2 - \sqrt{2}}{3} \right) - \log \left(8\sqrt{d \vee \log(64)}\sigma_\epsilon \right) (\underline{\gamma}(\nu) - \eta).$$

In particular, if $\underline{\gamma}(\nu) > 0$, for ϵ small enough, we have that

$$\log n_{\min}(\epsilon) \geq -\frac{1}{2} \log \left(8\sqrt{d \vee \log(64)}\sigma_\epsilon \right) \underline{\gamma}(\nu).$$

Moreover, if $\underline{\gamma}(\nu) = 0$, then the above inequality is always verified. This concludes the proof. \square

We can now present the proof of Theorem 3.1.

Proof of Theorem 3.1. Suppose that $\mathbb{E}_{X \sim \nu}[\|X\|^2] = \sigma^2 < \infty$, $T \geq 1 + \frac{1}{2\kappa} \log(32(\sigma^2 + \kappa^{-1}d))$, and $d \geq 5$. We begin by expressing the weighted ESM loss in terms of the standard ESM loss by using the fact that the density w of ϖ is non-decreasing on its support. We have, for $\delta \geq \epsilon$,

$$\widehat{\mathcal{L}}_{\text{ESM}}(\hat{s}, \varpi) = \int_\epsilon^T \widehat{\mathcal{L}}_{\text{ESM}}(\hat{s}, t) w(t) dt \geq w(\delta) \widehat{\mathcal{L}}_{\text{ESM}}(\hat{s}, \varpi_{\text{ELBO}}^\delta),$$

where we recall $\delta = w^{-1}(32\epsilon)$. By the assumption that $\epsilon \leq \sup(w)/32$, we have,

$$\mathbb{E}[\widehat{\mathcal{L}}_{\text{ESM}}(\hat{s}, \varpi_{\text{ELBO}}^\delta)] \leq \frac{\epsilon}{w(\delta)} \leq \frac{1}{32}.$$

Via the same argument, we also obtain that $\mathbb{E}[\mathcal{L}_{\text{ESM}}(\hat{s}, \varpi)] \leq \epsilon$ implies that $\mathbb{E}[\mathcal{L}_{\text{ESM}}(\hat{s}, \varpi_{\text{ELBO}}^\delta)] \leq 1/32$. Therefore it is sufficient to show that $\mathbb{E}[\widehat{\mathcal{L}}_{\text{ESM}}(\hat{s}, \varpi_{\text{ELBO}}^\delta)] \leq 1/32$, $\mathbb{E}[\mathcal{L}_{\text{ESM}}(\hat{s}, \varpi_{\text{ELBO}}^\delta)] \leq 1/32$ leads to the lower bound on n .

It follows from standard bounds based on Girsanov's theorem that

$$\begin{aligned} \text{KL}(p_\delta || q_{T-\delta}) &\leq \mathcal{L}_{\text{ESM}}(\hat{s}, \varpi_{\text{ELBO}}^\delta) + \text{KL}(p_T || q_0) , \\ \text{KL}(\hat{p}_\delta || q_{T-\delta}) &\leq \widehat{\mathcal{L}}_{\text{ESM}}(\hat{s}, \varpi_{\text{ELBO}}^\delta) + \text{KL}(\hat{p}_T || q_0) . \end{aligned}$$

Using that q_0 satisfies a logarithmic Sobolev inequality with constant 2κ , in combination with Corollary 2 of (Otto & Villani, 2001), we obtain the upper bound,

$$\begin{aligned} \text{KL}(p_T || q_0) &\leq \exp(-2\kappa(T-s)) \text{KL}(p_s || q_0) \\ &\leq \frac{\exp(-2\kappa(T-s))}{4s} W_2(p_0, q_0)^2 , \end{aligned}$$

any $s \in (0, T]$. Setting $s = 1$, we obtain,

$$\text{KL}(p_T || q_0) \leq \frac{\exp(-2\kappa(T-1))}{2} (\kappa^{-1}\sigma^2 + d) .$$

By the lower bound assumption on T , we have that $\text{KL}(p_T || q_0) \leq 1/32$ and, via a similar argument, $\mathbb{E}[\text{KL}(\hat{p}_T || q_0)] \leq 1/32$ also. Thus, it follows that

$$\mathbb{E}[\text{KL}(p_\delta || q_{T-\delta})], \mathbb{E}[\text{KL}(\hat{p}_\delta || q_{T-\delta})] \leq 1/16 .$$

The result then follows immediately from Proposition B.2. \square

B.2. Omitted proofs of Section 4

In this section, we present the proofs of Lemma 4.1 and Proposition 4.1.

B.2.1. PROOF OF LEMMA 4.1

We present below the proof of Lemma 4.1.

Proof. Let $(Z_1, \dots, Z_n) \in (\mathbb{R}^d)^N$ be a dataset and $W \in \mathbb{R}^{p \times d}$ be fixed. By assumption, W has full rank. Let $\Xi \sim \text{N}(0, I_d)$. We can write the empirical risk as

$$A \mapsto \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\|\sigma_t A W (\alpha_t Z_i + \sigma_t \Xi) + \Xi\|^2 \right] .$$

As $\mathbb{E}[\Xi] = 0$, we have

$$\begin{aligned} \sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,A}, t) &= d + \sigma_t^2 \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\text{Tr} (W^T A^T A W (\alpha_t Z_i + \sigma_t \Xi) (\alpha_t Z_i + \sigma_t \Xi)^T)] \\ &\quad + 2\sigma_t^2 \mathbb{E} [\text{Tr} (A W \Xi \Xi^T)] \\ &= d + \sigma_t^2 \text{Tr} (A^T A W (\alpha_t^2 \widehat{\Sigma} + \sigma_t^2 I_d) W^T) + 2\sigma_t^2 \text{Tr}(A W) \\ &= d + \sigma_t^2 \text{Tr} (A^T A W \widehat{\Sigma}_t W^T) + 2\sigma_t^2 \text{Tr}(A W) . \end{aligned}$$

Let $A \in \arg \min \widehat{\mathcal{R}}_{W,t}(\widehat{A})$, by taking the gradient above, we have

$$A W \widehat{\Sigma}_t W^T + W^T = 0 . \tag{10}$$

Case 1 ($p < d$). By our assumptions, the matrix W has full rank, so its rank is p in this case. The matrix $\widehat{\Sigma}_t$ is always invertible, as soon as $t > 0$; we also know that $\text{rank}(W) = \text{rank}(W \widehat{\Sigma}_t^{1/2})$, therefore, we conclude that $\text{rank}(W \widehat{\Sigma}_t W^T) = p$, so the matrix $W \widehat{\Sigma}_t W^T \in \mathbb{R}^{p \times p}$ is invertible.

Therefore, the empirical risk minimizer is given by

$$\hat{A} := -W^T(W\hat{\Sigma}_tW^T)^{-1}.$$

Using Equation (10), we have that

$$\sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,\hat{A}}, t) = d - \sigma_t^2 \text{Tr} \left(W^T(W\hat{\Sigma}_tW^T)^{-1}W \right).$$

To compute the score matching loss, we note that when $\vec{X}_0 \sim \nu = \text{N}(0, \Sigma)$, then $X_t \sim \text{N}(0, \alpha_t^2 \Sigma + \sigma_t^2 \text{I}_d) = \text{N}(0, \Sigma_t)$. Therefore, the score can be expressed as

$$\nabla \log p_t(x) = -\Sigma_t^{-1}x.$$

The score matching loss is then equal to

$$\mathcal{L}_{\text{ESM}}(s_{W,\hat{A}}, t) = \mathbb{E} \left[\left\| \hat{A}W(\alpha_t Z + \sigma_t \Xi) + \Sigma_t^{-1}(\alpha_t Z + \sigma_t \Xi) \right\|^2 \right],$$

with $(Z, \Xi) \sim \nu \otimes \text{N}(0, \text{I}_d)$. By symmetry of $\hat{A}W$, we obtain that

$$\begin{aligned} \mathcal{L}_{\text{ESM}}(s_{W,\hat{A}}, t) &= \mathbb{E} \left[\left\| \hat{A}W(\alpha_t Z + \sigma_t \Xi) + \Sigma_t^{-1}(\alpha_t Z + \sigma_t \Xi) \right\|^2 \right] \\ &= \text{Tr} \left(W^T \hat{A}^T \hat{A}W \Sigma_t + 2\hat{A}W + \Sigma_t^{-1} \right) \\ &= \left\| \left(W^T(W\hat{\Sigma}_tW^T)^{-1}W - \Sigma_t^{-1} \right) \Sigma_t^{1/2} \right\|_{\text{F}}^2. \end{aligned}$$

Case 2. ($p \geq d$) In this case, by similar argument as in the first case, we note that the matrix $W \in \mathbb{R}^{p \times d}$ has full rank by assumption. Consider the linear map $\Phi_W : \mathbb{R}^{d \times p} \rightarrow \mathbb{R}^{d \times d}$ defined by $\Phi_W(A) = AW$. Then, we easily show that this map is almost-surely surjective, because the matrix W admits a left inverse W^+ such that $W^+W = \text{I}_d$ (which can be obtained from the singular value decomposition of W).

As both the empirical risk and the score matching loss depend only on A through the product AW , we have that all the empirical risk minimizers \hat{A} satisfy $\hat{A}W = \hat{H} \in \mathbb{R}^{d \times d}$, with \hat{H} the empirical minimizer of the following linear model:

$$\hat{H} \in \arg \min_H \left\{ \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\left\| \sigma_t H(\alpha_t Z_i + \sigma_t \Xi) + \Xi \right\|^2 \right] \right\}, \quad (Z, \Xi) \sim \nu \otimes \text{N}(0, \text{I}_d).$$

By similar computations as before, we easily show that $\hat{H} = -\hat{\Sigma}_t^{-1}$. Therefore, the minimum empirical risk is

$$\sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,\hat{A}}, t) = d - \sigma_t^2 \text{Tr} \left(\hat{\Sigma}_t^{-1} \right).$$

Finally, the score matching loss can be expressed as

$$\mathcal{L}_{\text{ESM}}(s_{W,\hat{A}}, t) = \text{Tr} \left((\hat{A}W + \Sigma_t^{-1})^2 \Sigma_t \right) = \left\| (\hat{\Sigma}_t^{-1} - \Sigma_t^{-1}) \Sigma_t^{1/2} \right\|_{\text{F}}^2.$$

This concludes the proof. \square

B.2.2. PROOF OF PROPOSITION 4.1

We prove Proposition 4.1 as a corollary of the following more general result, which gives the exact asymptotic behavior of the score matching loss and the empirical risk.

Proposition B.3. Consider the same setup as in Proposition 4.1. Let $\hat{A} \in \arg \min_A \widehat{\mathcal{R}}_{W,t}(A)$. If $\psi_p > \psi_d$, we have

$$\lim_{n \rightarrow \infty} \frac{1}{d} \inf_A \left(\sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,A}, t) \right) = 1 - \sigma_t^2 \int \frac{d\mu_{\text{MP}}^{(\psi_d)}(\lambda)}{\alpha_t^2 \beta \lambda + \sigma_t^2},$$

and

$$\lim_{n \rightarrow \infty} \mathcal{L}_{\text{ESM}}(s_{W, \hat{A}}, t) = (\alpha_t^2 \beta + \sigma_t^2) \int \left(\frac{1}{\alpha_t^2 \beta \lambda + \sigma_t^2} - \frac{1}{\alpha_t^2 \beta + \sigma_t^2} \right)^2 d\mu_{\text{MP}}^{(\psi_d)}(\lambda).$$

If $\psi_p < \psi_d$, we have

$$\frac{1}{d} \sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W, \hat{A}}, t) \xrightarrow{n \rightarrow \infty} 1 - \sigma_t^2 \frac{\psi_p}{\psi_d} \int \frac{1}{\alpha_t^2 \beta \lambda + \sigma_t^2} d\mu_{\text{MP}}^{(\psi_p)}(\lambda),$$

and

$$\frac{1}{d} \mathcal{L}_{\text{ESM}}(s_{W, \hat{A}}, t) \xrightarrow{n \rightarrow \infty} (\alpha_t^2 \beta + \sigma_t^2) \frac{\psi_p}{\psi_d} \int \left(\frac{1}{\alpha_t^2 \beta \lambda + \sigma_t^2} - \frac{1}{\alpha_t^2 \beta + \sigma_t^2} \right)^2 d\mu_{\text{MP}}^{(\psi_p)}(\lambda) + \frac{1 - \frac{\psi_p}{\psi_d}}{\alpha_t^2 \beta + \sigma_t^2}.$$

Proof. First, we note that in that case we have $\Sigma_t = (\alpha_t^2 \beta + \sigma_t^2) \mathbf{I}_d$. Moreover, we observe that, as the matrix $W \in \mathbb{R}^{p \times d}$ has i.i.d. standard Gaussian entries, it has almost-surely full rank.

We can therefore apply Lemma 4.1, we distinguish two cases.

Case 1. ($\psi_p > \psi_d$) For n large enough, we have $p > d$. In this case, by Lemma 4.1, the normalized empirical risk can be written as

$$\frac{1}{d} \inf_A \left(\sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W, A}, t) \right) = 1 - \sigma_t^2 \int \frac{d\mu_{\beta^{-1/2} \Sigma}(\lambda)}{\alpha_t^2 \beta \lambda + \sigma_t^2},$$

where, for a positive definite matrix M , $\widehat{\mu}_M$ denotes the empirical spectral measure of M , see Section A.2.

By the Marchenko-Pastur theorem (e.g. Theorem 2.4 of (Couillet & Liao, 2022)), we know that, almost-surely, the measure $\mu_{\beta^{-1/2} \Sigma}$ converges weakly to the Marchenko-Pastur distribution with shape parameter ψ_d , denoted $\mu_{\text{MP}}^{(\psi_d)}$, see Section A.2. The map $\lambda \mapsto (\alpha_t^2 \beta \lambda + \sigma_t^2)$ is bounded and continuous, therefore, we have

$$\lim_{n \rightarrow \infty} \frac{1}{d} \inf_A \left(\sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W, A}, t) \right) = 1 - \sigma_t^2 \int \frac{d\mu_{\text{MP}}^{(\psi_d)}(\lambda)}{\alpha_t^2 \beta \lambda + \sigma_t^2}.$$

Similarly, we have, with $\hat{A} \in \arg \min_A \widehat{\mathcal{R}}_{W, t}(A)$,

$$\begin{aligned} \mathcal{L}_{\text{ESM}}(s_{W, \hat{A}}, t) &= \frac{1}{d} (\alpha_t^2 \beta + \sigma_t^2) \left\| \left(\widehat{\Sigma}_t^{-1} - \Sigma_t^{-1} \right) \Sigma_t^{1/2} \right\|_{\text{F}}^2 \\ &\xrightarrow{n \rightarrow \infty} (\alpha_t^2 \beta + \sigma_t^2) \int \left(\frac{1}{\alpha_t^2 \beta \lambda + \sigma_t^2} - \frac{1}{\alpha_t^2 \beta + \sigma_t^2} \right)^2 d\mu_{\text{MP}}^{(\psi_d)}(\lambda). \end{aligned}$$

Case 2. ($\psi_p < \psi_d$) For n large enough, we have $d > p$. In this case, we note (see the proof of Lemma 4.1) that the matrix W is almost surely of rank p . Therefore, we can write the singular value decomposition of W in the form $W = U \Lambda V^T$, where $U \in \mathbb{R}^{p \times p}$ is orthonormal, $\Lambda \in \mathbb{R}^{p \times p}$ is a diagonal matrix with positive diagonal coefficients, and $V \in \mathbb{R}^{d \times p}$ is column orthonormal ($V^T V = \mathbf{I}_p$).

With these notations and noting that $V^T V = \mathbf{I}_p$, we have, almost-surely, that

$$\begin{aligned} \frac{1}{d} \sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W, \hat{A}}, t) &= 1 - \sigma_t^2 \text{Tr} \left(V \Lambda U^T (U \Lambda V^T \widehat{\Sigma}_t V \Lambda U^T)^{-1} U \Lambda V^T \right) \\ &= 1 - \sigma_t^2 \text{Tr} \left(V (V^T \widehat{\Sigma}_t V)^{-1} V^T \right) \\ &= 1 - \sigma_t^2 \text{Tr} \left((V^T \widehat{\Sigma}_t V)^{-1} \right). \end{aligned}$$

Now we observe that

$$V^T \widehat{\Sigma}_t V = \alpha_t^2 V^T \widehat{\Sigma} V + \sigma_t^2 \mathbf{I}_p = \frac{\alpha_t^2}{n} \sum_{i=1}^n V^T Z_i (V^T Z_i)^T + \sigma_t^2 \mathbf{I}_p.$$

Recall that $(Z_1, \dots, Z_n) \sim N(0, \beta \mathbf{I}_d)^{\otimes n}$ and that V is column orthonormal (actually, it is even uniformly distributed) and independent of (Z_1, \dots, Z_n) . By independence and invariance by rotation of $N(0, \beta \mathbf{I}_d)$, we have that $V^T Z_i \sim N(0, \beta \mathbf{I}_p)$. Moreover, we have that $(V^T Z_1, \dots, V^T Z_n) \sim N(0, \beta \mathbf{I}_p)^{\otimes n}$. Thus, we can apply the p -dimensional version of the Marchenko-Pastur theorem, which gives that

$$\frac{1}{d} \sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W, \hat{A}}, t) \xrightarrow{n \rightarrow \infty} 1 - \sigma_t^2 \frac{\psi_p}{\psi_d} \int \frac{1}{\alpha_t^2 \beta \lambda + \sigma_t^2} d\mu_{\text{MP}}^{(\psi_p)}(\lambda).$$

By Lemma 4.1, we have

$$\begin{aligned} \frac{1}{d} \mathcal{L}_{\text{ESM}}(s_{W, \hat{A}}, t) &= \frac{1}{d} \left\| \left(W^T (W \widehat{\Sigma}_t W^T)^{-1} W - \Sigma_t^{-1} \right) \Sigma_t^{1/2} \right\|_{\text{F}}^2 \\ &= \frac{\alpha_t^2 \beta + \sigma_t^2}{d} \left\| \left(V (V^T \widehat{\Sigma}_t V)^{-1} V^T - \Sigma_t^{-1} \right) \right\|_{\text{F}}^2. \end{aligned}$$

Recall that $\Sigma_t = (\alpha_t^2 \beta + \sigma_t^2) \mathbf{I}_d$. By the Pythagorean theorem, we have

$$\begin{aligned} \frac{1}{d} \mathcal{L}_{\text{ESM}}(s_{W, \hat{A}}, t) &= \frac{\alpha_t^2 \beta + \sigma_t^2}{d} \left(\left\| V (V^T \widehat{\Sigma}_t V)^{-1} V^T - \frac{V V^T}{\alpha_t^2 \beta + \sigma_t^2} \right\|_{\text{F}}^2 + \frac{\|\mathbf{I}_d - V V^T\|_{\text{F}}^2}{(\alpha_t^2 \beta + \sigma_t^2)^2} \right) \\ &= \frac{\alpha_t^2 \beta + \sigma_t^2}{d} \left(\left\| (V^T \widehat{\Sigma}_t V)^{-1} - \frac{\mathbf{I}_p}{\alpha_t^2 \beta + \sigma_t^2} \right\|_{\text{F}}^2 + \frac{\text{Tr}(\mathbf{I}_d - V V^T)}{(\alpha_t^2 \beta + \sigma_t^2)^2} \right) \\ &= \frac{\alpha_t^2 \beta + \sigma_t^2}{d} \left(\left\| (V^T \widehat{\Sigma}_t V)^{-1} - \frac{\mathbf{I}_p}{\alpha_t^2 \beta + \sigma_t^2} \right\|_{\text{F}}^2 + \frac{d - p}{(\alpha_t^2 \beta + \sigma_t^2)^2} \right). \end{aligned}$$

By the arguments above (for the empirical risk), we observe that we can apply the p dimensional version of the Marchenko-Pastur theorem to the first term, this gives, almost-surely,

$$\frac{1}{d} \mathcal{L}_{\text{ESM}}(s_{W, \hat{A}}, t) \xrightarrow{n \rightarrow \infty} (\alpha_t^2 \beta + \sigma_t^2) \frac{\psi_p}{\psi_d} \int \left(\frac{1}{\alpha_t^2 \beta \lambda + \sigma_t^2} - \frac{1}{\alpha_t^2 \beta + \sigma_t^2} \right)^2 d\mu_{\text{MP}}^{(\psi_p)}(\lambda) + \frac{1 - \frac{\psi_p}{\psi_d}}{\alpha_t^2 \beta + \sigma_t^2}.$$

This concludes the proof. \square

We present below the proof of Proposition 4.1.

Proof. (of Proposition 4.1) We apply Proposition B.3 and focus on the case where $\psi_p \wedge \psi_d > 1$. Then, in both cases below, the Marchenko-Pastur distributions appearing in Proposition B.3 have a mass at zero, which we exploit to obtain our estimates. Let $f_{\text{MP}}^{(\psi)}$ denote the bulk density of the Marchenko-Pastur distribution with shape parameter $\psi > 0$. Note that this is only a probability density when $\psi \leq 1$, see Section A.2.

Case 1. ($\psi_p > \psi_d$). Then, we have, almost-surely,

$$\lim_{n \rightarrow \infty} \frac{1}{d} \inf_A \left(\sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W, A}, t) \right) = 1 - \sigma_t^2 \left(1 - \frac{1}{\psi_d} \right) \frac{1}{\sigma_t^2} - \sigma_t^2 \int \frac{f_{\text{MP}}^{(\psi_d)}(\lambda)}{\alpha_t^2 \beta \lambda + \sigma_t^2} d\lambda.$$

The bulk density $f_{\text{MP}}^{(\psi_d)}$ is supported on $[(\sqrt{\psi_d} - 1)^2, (\sqrt{\psi_d} + 1)^2]$ and we have $\int f_{\text{MP}}^{(\psi_d)}(\lambda) d\lambda = \psi_d^{-1}$. Therefore, we have, almost-surely

$$\lim_{n \rightarrow \infty} \frac{1}{d} \inf_A \left(\sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W, A}, t) \right) = \frac{1}{\psi_d} + \mathcal{O} \left(\frac{1}{\psi_d^2} \right).$$

For the score matching loss, we have similarly that, almost-surely,

$$\begin{aligned} \lim_{n \rightarrow \infty} \left(\frac{1}{d} \mathcal{L}_{\text{ESM}}(s_{W, \hat{A}}, t) \right) &= \left(1 - \frac{1}{\psi_d} \right) \frac{\alpha_t^4 \beta^2}{\sigma_t^4 (\alpha_t^2 \beta + \sigma_t^2)} \\ &\quad + (\alpha_t^2 \beta + \sigma_t^2) \int \left(\frac{1}{\alpha_t^2 \beta \lambda + \sigma_t^2} - \frac{1}{\alpha_t^2 \beta + \sigma_t^2} \right)^2 f_{\text{MP}}^{(\psi_d)}(\lambda) d\lambda \\ &= \left(1 - \frac{1}{\psi_d} \right) \frac{\alpha_t^4 \beta^2}{\sigma_t^4 (\alpha_t^2 \beta + \sigma_t^2)} + \mathcal{O} \left(\frac{1}{\psi_d} \right). \end{aligned}$$

Case 2. ($\psi_p < \psi_d$) In that case, by very similar computations, we obtain

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{d} \inf_A \left(\sigma_t^2 \widehat{\mathcal{L}}_{\text{DSM}}(s_{W,A}, t) \right) &= 1 - \frac{\psi_p}{\psi_d} \left(1 - \frac{1}{\psi_p} \right) - \sigma_t^2 \frac{\psi_p}{\psi_d} \int \frac{f_{\text{MP}}^{(\psi_p)}(\lambda)}{\alpha_t^2 \beta \lambda + \sigma_t^2} d\lambda \\ &= 1 + \frac{1}{\psi_d} - \frac{\psi_p}{\psi_d} + \mathcal{O}_t \left(\frac{1}{\psi_p \psi_d} \right). \end{aligned}$$

Finally, for the score matching loss, we obtain

$$\begin{aligned} \lim_{n \rightarrow \infty} \left(\frac{1}{d} \mathcal{L}_{\text{ESM}}(s_{W,\hat{A}}, t) \right) &= \frac{\psi_p}{\psi_d} \left(1 - \frac{1}{\psi_p} \right) \frac{\alpha_t^4 \beta^2}{\sigma_t^4 (\alpha_t^2 \beta + \sigma_t^2)} + \frac{1 - \frac{\psi_p}{\psi_d}}{\alpha_t^2 \beta + \sigma_t^2} \\ &\quad + (\alpha_t^2 \beta + \sigma_t^2) \frac{\psi_p}{\psi_d} \int \left(\frac{1}{\alpha_t^2 \beta \lambda + \sigma_t^2} - \frac{1}{\alpha_t^2 \beta + \sigma_t^2} \right)^2 f_{\text{MP}}^{(\psi_p)}(\lambda) d\lambda \\ &= \frac{\psi_p}{\psi_d} \left(1 - \frac{1}{\psi_p} \right) \frac{\alpha_t^4 \beta^2}{\sigma_t^4 (\alpha_t^2 \beta + \sigma_t^2)} + \frac{1 - \frac{\psi_p}{\psi_d}}{\alpha_t^2 \beta + \sigma_t^2} + \mathcal{O} \left(\frac{1}{\psi_d} \right). \end{aligned}$$

The result follows by rearranging the terms. \square

B.3. Omitted proofs of Section 5.1

In this section, we present the proofs of Section 5.1, related to the impact of the time integration in the score matching loss.

B.3.1. PROOF OF LEMMA 5.1

Proof. (of Lemma 5.1) By calculations similar to the proof of Lemma 4.1, we can write the denoising score matching loss as

$$\begin{aligned} \widehat{\mathcal{L}}_{\text{DSM}}(s_{A,W}, \varpi) &= \sigma_\pi^{-2} \int_{[0,T]} \left(d + \sigma_t^2 \text{Tr}(A^T A W \widehat{\Sigma}_t W^T) + 2\sigma_t^2 \text{Tr}(A W) \right) d\pi(t) \\ &= \frac{d}{\sigma_\pi^2} + \text{Tr}(W^T A^T A W \widehat{\Sigma}_\varpi) + 2\text{Tr}(A W). \end{aligned}$$

By reasoning as in the proof Lemma 4.1, we observe that (because $p > d$) the empirical risk minimization is equivalent to a linear model with variable all the empirical risk minimizers \hat{A} satisfy

$$\hat{A}W = -\widehat{\Sigma}_\varpi^{-1}.$$

Therefore, the minimum empirical risk is equal to

$$\inf_A \widehat{\mathcal{L}}_{\text{DSM}}(s_{A,W}, \varpi) = \frac{d}{\sigma_\pi^2} - \text{Tr} \left(\widehat{\Sigma}_\varpi^{-1} \right).$$

Again, by similar calculations as in the proof of Lemma 4.1, we can express the population score matching loss as

$$\mathcal{L}_{\text{ESM}}(s_{\hat{A},W}, \pi) = \int_{[0,T]} \mathbb{E} \left[\left\| \widehat{A}W(\alpha_t Z + \sigma_t \Xi) + \Sigma_t^{-1}(\alpha_t Z + \sigma_t \Xi) \right\|^2 \right] d\pi(t).$$

with $(Z, \Xi) \sim \nu \otimes \mathcal{N}(0, \Sigma)$. Therefore, we have

$$\begin{aligned} \mathcal{L}_{\text{ESM}}(s_{\hat{A},W}, \pi) &= \int_{[0,T]} \text{Tr} \left(\widehat{\Sigma}_\varpi^{-2} \Sigma_t - 2\widehat{\Sigma}_\varpi^{-1} + \Sigma_t^{-1} \right) d\pi(t) \\ &= \text{Tr} \left(\widehat{\Sigma}_\varpi^{-2} \Sigma_\pi - 2\widehat{\Sigma}_\varpi^{-1} + \int_{[0,T]} \Sigma_t^{-1} d\pi(t) \right) \\ &= \left\| \left(\widehat{\Sigma}_\varpi^{-1} - \Sigma_\pi^{-1} \right) \Sigma_\pi^{1/2} \right\|_{\text{F}}^2 + \text{C}_\pi, \end{aligned}$$

with

$$C_\pi := \text{Tr} \left(\int_{[0,T]} \Sigma_t^{-1} d\pi(t) - \Sigma_\pi^{-1} \right). \quad (11)$$

By diagonalizing Σ_t and using Jensen's inequality, we see that $C_\pi \geq 0$. \square

Remark B.3 (Case where $p \leq d$). By proceeding like in the proof of Lemma 4.1, we have that

$$\inf_A \widehat{\mathcal{L}}_{\text{DSM}}(s_{A,W}, \varpi) = \frac{d}{\sigma_\pi^2} - \text{Tr} \left(W^T (W \widehat{\Sigma}_\varpi W^T)^{-1} W \right).$$

Similarly, for the explicit score matching loss, we have

$$\mathcal{L}_{\text{ESM}}(s_{\widehat{A},W}, \pi) \left\| \left(W^T (W \widehat{\Sigma}_\varpi W^T)^{-1} W - \Sigma_\pi^{-1} \right) \Sigma_\pi^{1/2} \right\|_{\text{F}}^2 + C_\pi,$$

where $\widehat{A} = -(W \widehat{\Sigma}_\varpi W^T)^{-1}$ is the empirical risk minimizer. These formulas are used to obtain the exact asymptotics presented in Figure 3, using very similar computations as in the proofs of Section 4.

Remark B.4 (Case where π is the uniform distribution and $\kappa = 1$). In that case, we have

$$\frac{1}{d} C_\pi = \frac{1}{T} \int_0^T \frac{dt}{\alpha_t^2 \beta + \sigma_t^2} - \frac{1}{\alpha_\pi^2 \beta + \sigma_\pi^2}.$$

We easily see that

$$\alpha_\pi^2 = \frac{1}{T} \int_0^T e^{-2t} dt = \frac{1 - e^{-2T}}{2T}, \quad \sigma_\pi^2 = 1 - \frac{1 - e^{-2T}}{2T}.$$

Therefore, we have

$$\begin{aligned} \frac{1}{d} C_\pi &= \frac{1}{T} \int_0^T \frac{dt}{e^{-2t} \beta + 1 - e^{-2t}} - \frac{2T}{\beta(1 - e^{-2T}) + 2T - 1 + e^{-2T}} \\ &= \frac{1}{2T} \log \left(\frac{e^{2T} + \beta - 1}{\beta} \right) - \frac{2T}{\beta(1 - e^{-2T}) + 2T - 1 + e^{-2T}}. \end{aligned}$$

Therefore, we have

$$\frac{1}{d} C_\pi = \frac{1}{2T} \log \left(\frac{1 + e^{-2T}(\beta - 1)}{\beta} \right) + \frac{\beta(1 - e^{-2T}) - 1 + e^{-2T}}{\beta(1 - e^{-2T}) + 2T - 1 + e^{-2T}} = \mathcal{O} \left(\frac{1}{T} \right). \quad (12)$$

B.3.2. PROOF OF PROPOSITION 5.1

We provide below the proof of Proposition 5.1. We first provide, in the next proposition, the exact asymptotics of the score matching losses in the case of time-integrated losses.

Proposition B.4. *Consider the same setup as in Proposition 5.1, we have*

$$\lim_{n \rightarrow \infty} \left(\frac{1}{d} \inf_A \widehat{\mathcal{L}}_{\text{DSM}}(s_{A,W}, \varpi) \right) = \frac{1}{\sigma_\pi^2} - \int \frac{d\mu_{\text{MP}}^{(\psi_d)}(\lambda)}{\alpha_\varpi^2 \beta \lambda + \sigma_\varpi^2},$$

and

$$\lim_{n \rightarrow \infty} \left(\frac{1}{d} \mathcal{L}_{\text{ESM}}(s_{\widehat{A},W}, \pi) \right) = (\alpha_\pi^2 \beta + \sigma_\pi^2) \int \left(\frac{1}{\alpha_\varpi^2 \beta \lambda + \sigma_\varpi^2} - \frac{1}{\alpha_\pi^2 \beta + \sigma_\pi^2} \right)^2 d\mu_{\text{MP}}^{(\psi_d)}(\lambda) + \overline{C}_\pi,$$

where, as before, $\mu_{\text{MP}}^{(\psi_d)}$ is the Marchenko-Pastur distribution with shape parameter ψ_d , and

$$\overline{C}_\pi := \int_{[0,T]} \frac{\pi(dt)}{\alpha_t^2 \beta + \sigma_t^2} - \frac{1}{\alpha_\pi^2 \beta + \sigma_\pi^2} \geq 0. \quad (13)$$

Proof. Assume that $\psi_p > \psi_d$. Then, by Lemma 4.1 (using the same notations), we have

$$\frac{1}{d} \inf_A \widehat{\mathcal{L}}_{\text{DSM}}(s_{A,W}, \varpi) = \frac{1}{\sigma_\pi^2} - \int \frac{d\mu_{\beta^{-1/2}\widehat{\Sigma}}(\lambda)}{\alpha_\varpi^2 \beta \lambda + \sigma_\varpi^2}.$$

By the Marchenko-Pastur's theorem, we have

$$\lim_{n \rightarrow \infty} \left(\frac{1}{d} \inf_A \widehat{\mathcal{L}}_{\text{DSM}}(s_{A,W}, \varpi) \right) = \frac{1}{\sigma_\pi^2} - \int \frac{d\mu_{\text{MP}}^{(\psi_d)}(\lambda)}{\alpha_\varpi^2 \beta \lambda + \sigma_\varpi^2}.$$

Similarly, the score matching loss can be expressed as

$$\frac{1}{d} \mathcal{L}_{\text{ESM}}(s_{\widehat{A},W}, \pi) = \frac{\alpha_\pi^2 \beta + \sigma_\pi^2}{d} \left\| \widehat{\Sigma}_\varpi^{-1} - \frac{\mathbf{I}_d}{\alpha_\pi^2 \beta + \sigma_\pi^2} \right\|_{\text{F}}^2 + \frac{\text{C}_\pi}{d},$$

where C_π is defined in the proof of Lemma 5.1. Under Assumption 4.1, we observe that C_π/d is independent of d , and therefore we can define

$$\overline{\text{C}}_\pi := \frac{\text{C}_\pi}{d} = \int_{[0,T]} \frac{\pi(dt)}{\alpha_t^2 \beta + \sigma_t^2} - \frac{1}{\alpha_\pi^2 \beta + \sigma_\pi^2} \geq 0,$$

where the inequality follows from Jensen's inequality. By applying the Marchenko-Pastur theorem, we obtain

$$\lim_{n \rightarrow \infty} \left(\frac{1}{d} \mathcal{L}_{\text{ESM}}(s_{\widehat{A},W}, \pi) \right) = (\alpha_\pi^2 \beta + \sigma_\pi^2) \int \left(\frac{1}{\alpha_\varpi^2 \beta \lambda + \sigma_\varpi^2} - \frac{1}{\alpha_\pi^2 \beta + \sigma_\pi^2} \right)^2 d\mu_{\text{MP}}^{(\psi_d)}(\lambda) + \overline{\text{C}}_\pi.$$

This concludes the proof. \square

Remark B.5. Based on Remark B.3, we can easily extend these derivations to the case $\psi_p < \psi_d$.

We present below the proof of Proposition 5.1.

Proof. (of Proposition 5.1) Assume that $\psi_p > \psi_d > 1$, then the Marchenko-Pastur has a mass at zero in all the cases below. By Proposition B.4, we immediately have

$$\begin{aligned} \lim_{n \rightarrow \infty} \left(\frac{1}{d} \inf_A \widehat{\mathcal{L}}_{\text{DSM}}(s_{A,W}, \varpi) \right) &= \frac{1}{\sigma_\pi^2} - \frac{1}{\sigma_\varpi^2} \left(1 - \frac{1}{\psi_d} \right) + \int \frac{f_{\text{MP}}^{(\psi_d)}(\lambda)}{\alpha_\varpi^2 \beta \lambda + \sigma_\varpi^2} d\lambda \\ &= \frac{1}{\sigma_\pi^2} - \frac{1}{\sigma_\varpi^2} + \frac{1}{\psi_d \sigma_\varpi^2} + \mathcal{O} \left(\frac{1}{\psi_d^2} \right). \end{aligned}$$

Note that, by the Cauchy-Schwarz inequality, we always have $\sigma_\pi^{-2} \geq \sigma_\varpi^{-2}$. Similarly, the asymptotic score matching loss satisfies

$$\lim_{n \rightarrow \infty} \left(\frac{1}{d} \mathcal{L}_{\text{ESM}}(s_{\widehat{A},W}, \pi) \right) = \left(1 - \frac{1}{\psi_d} \right) \frac{(\alpha_\pi^2 \beta + \sigma_\pi^2 - \sigma_\varpi^2)^2}{\sigma_\varpi^4 (\alpha_\pi^2 \beta + \sigma_\pi^2)} + \overline{\text{C}}_\pi + \mathcal{O} \left(\frac{1}{\psi_d} \right).$$

The proof follows by defining

$$\text{L}_\pi := \frac{(\alpha_\pi^2 \beta + \sigma_\pi^2 - \sigma_\varpi^2)^2}{\sigma_\varpi^4 (\alpha_\pi^2 \beta + \sigma_\pi^2)} + \overline{\text{C}}_\pi \tag{14}$$

This concludes the proof. \square

B.4. Proof of Section 5.2

We present below the proof of Proposition 5.2.

Proof. (of Proposition 5.2) **Step 1.** We start by solving the gradient flow dynamics of Equation (7).

We can expand the gradient flow dynamics as

$$\frac{d}{d\tau} A(\tau) = -2p^{-1}\sigma_t^2 A(\tau) W \widehat{\Sigma}_t W^T - 2p^{-1}\sigma_t^2 W^T.$$

We can assume, without loss of generality, that $p > d$. Then, by reasoning as in the proof of Proposition 4.1, we observe that the matrix $W^T W \in \mathbb{R}^{d \times d}$ is almost-surely invertible. Therefore, we can define $\tilde{A} := (W^T W)^{-1} A W$ and we have almost-surely that

$$\frac{d}{d\tau} \tilde{A}(\tau) = -2p^{-1}\sigma_t^2 \tilde{A}(\tau) \widehat{\Sigma}_t W^T W - 2p^{-1}\sigma_t^2 \mathbf{I}_d.$$

Taking into account that the initialization satisfies $A(0) = 0$, we deduce that

$$\tilde{A}(\tau) = -(\widehat{\Sigma}_t W^T W)^{-1} \left(\mathbf{I}_d - e^{-2\tau p^{-1}\sigma_t^2 \widehat{\Sigma}_t W^T W} \right).$$

Therefore, we have

$$A(\tau) W = -\widehat{\Sigma}_t^{-1} \left(\mathbf{I}_d - e^{-2\tau p^{-1}\sigma_t^2 \widehat{\Sigma}_t W^T W} \right) \in \mathbb{R}^{d \times d}.$$

By the strong law of large numbers, we have, almost-surely (for fixed n, d , and τ), that $p^{-1} W^T W \rightarrow \mathbf{I}_d$ as $p \rightarrow \infty$. Thus, almost-surely, we have

$$\lim_{p \rightarrow \infty} (A(\tau) W) = -\widehat{\Sigma}_t^{-1} \left(\mathbf{I}_d - e^{-2\tau \sigma_t^2 \widehat{\Sigma}_t} \right) \in \mathbb{R}^{d \times d}. \quad (15)$$

Step 2. By computations similar to the proof of Lemma 4.1 and Proposition 4.1, we have (for fixed (p, d, n))

$$\mathcal{E}_t(\tau, n, d, p) := \mathcal{L}_{\text{ESM}}(s_{W, A(\tau)}, t) = \left\| \left(\widehat{\Sigma}_t^{-1} \left(\mathbf{I}_d - e^{-2\tau p^{-1}\sigma_t^2 \widehat{\Sigma}_t W^T W} \right) - \Sigma_t^{-1} \right) \Sigma_t^{1/2} \right\|_{\mathbb{F}}^2.$$

Therefore, by the law of large numbers, we have almost surely that

$$\begin{aligned} \lim_{p \rightarrow \infty} \mathcal{E}_t(\tau, n, d, p) &= \left\| \left(\widehat{\Sigma}_t^{-1} \left(\mathbf{I}_d - e^{-2\tau \sigma_t^2 \widehat{\Sigma}_t} \right) - \Sigma_t^{-1} \right) \Sigma_t^{1/2} \right\|_{\mathbb{F}}^2 \\ &= (\alpha_t^2 \beta + \sigma_t^2) \left\| \widehat{\Sigma}_t^{-1} \left(\mathbf{I}_d - e^{-2\tau \sigma_t^2 \widehat{\Sigma}_t} \right) - \Sigma_t^{-1} \right\|_{\mathbb{F}}^2 =: \bar{\mathcal{E}}_t(\tau, n, d). \end{aligned}$$

Recall that we take $n, d \rightarrow \infty$ with $\psi_d := \lim(d/n) > 1$. Therefore, we can apply the Marchenko-Pastur theorem and obtain, through similar computations as in the proof of Proposition 4.1, that

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{d} \bar{\mathcal{E}}_t(\tau, n, d) &= (\alpha_t^2 \beta + \sigma_t^2) \int \left(\frac{1 - e^{-2\tau \sigma_t^2 (\alpha_t^2 \beta \lambda + \sigma_t^2)}}{\alpha_t^2 \beta \lambda + \sigma_t^2} - \frac{1}{\alpha_t^2 \beta + \sigma_t^2} \right)^2 d\mu_{\text{MP}}^{(\psi_d)}(\lambda) \\ &= \left(1 - \frac{1}{\psi_d} \right) \frac{\left((\alpha_t^2 \beta + \sigma_t^2) \left(1 - e^{-2\tau \sigma_t^4} \right) - \sigma_t^2 \right)^2}{\sigma_t^4 (\alpha_t^2 \beta + \sigma_t^2)} + \mathcal{O} \left(\frac{1}{\psi_d} \right). \end{aligned}$$

This concludes the proof. \square

C. Experimental details

In this section, we provide some additional details on the experiments presented in the main text.

C.1. Random features experiments

In this short subsection, we quickly describe the empirical setup for the diffusion random feature experiment presented in Figure 1.

Regression experiments. In the regression experiment presented in Figure 1, we use a classical random feature regression setting. We consider a data distribution of the form

$$y_i = \langle \beta_\star, x_i \rangle + \sigma \epsilon_i, \quad x_i \sim \mathcal{N}(0, \mathbf{I}_d), \quad \epsilon \sim \mathcal{N}(0, 1),$$

where ϵ_i and x_i are independent and we sample n i.i.d. pairs (x_i, y_i) for $i \in \{1, \dots, n\}$. We consider a score network of the form

$$f_{A,W}(x) = \frac{A}{\sqrt{p}} \varrho \left(\frac{W}{\sqrt{d}} x \right),$$

where $A \in \mathbb{R}^{d \times p}$, $W \in \mathbb{R}^{p \times d}$ is a random matrix with i.i.d. standard Gaussian entries, and $\varrho(x) := \max(x, 0)$ is the ReLU activation.

We consider the minimization of the empirical risk

$$\hat{A} \in \arg \min_A \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - f_{A,W}(x_i))^2 + \frac{\lambda}{\sqrt{p}} \|A\|_F^2 \right\},$$

with $\Xi \sim \mathcal{N}(0, \mathbf{I}_d)$ and λ a ridge regularization parameter.

Then, we evaluate the test (using a validation set) and train risk evaluated at the empirical risk minimizer \hat{A} and report it in Figure 1.

Hyperparameters details. In Figure 1, we use $n = 10^2$, $d = 20$, $N_g = 10$, $\lambda = 10^{-4}$, and $p \in [5, 4 \cdot 10^2]$.

Diffusion experiments. We consider the following family of random features score networks, which is similar to existing works (Bonnaire et al., 2025; George & Macris, 2026), for $x \in \mathbb{R}^d$

$$f_{A,W}(x) = \frac{A}{\sqrt{p}} \varrho \left(\frac{W}{\sqrt{d}} x \right),$$

where $A \in \mathbb{R}^{d \times p}$, $W \in \mathbb{R}^{p \times d}$ is a random matrix with i.i.d. standard Gaussian entries, and $\varrho(x) := \max(x, 0)$ is the ReLU activation. For the data distribution, we take the data distribution to be $\nu = \mathcal{N}(0, d^{-1} \mathbf{I}_d)$.

Given a fixed $t > 0$, a random W and a dataset $S^{(n)} := (Z_1, \dots, Z_n) \sim \nu^{\otimes n}$, we consider the minimization of the regularized empirical risk, defined by

$$A \mapsto \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\Xi} \left[\|\sigma_t f_{A,W}(\alpha_t Z_i + \sigma_t \Xi) + \Xi\|^2 \right] + \frac{\lambda}{\sqrt{p}} \|A\|_F^2,$$

with $\Xi \sim \mathcal{N}(0, \mathbf{I}_d)$ and λ a ridge regularization parameter. An important difference with the regression case above is that, in order to make the empirical risk minimization tractable, we use a Monte Carlo approximation of the expectation over $\Xi \sim \mathcal{N}(0, \mathbf{I}_d)$. Given an integer $N_g \in \mathbb{N}^*$ and $(\xi_{ij})_{1 \leq i \leq n, 1 \leq j \leq N_g} \sim \mathcal{N}(0, \mathbf{I}_d)^{\otimes (n N_g)}$, we therefore compute

$$\hat{A} \in \arg \min_A \left\{ \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{N_g} \|\sigma_t f_{A,W}(\alpha_t Z_i + \sigma_t \xi_{ij}) + \xi_{ij}\|^2 + \frac{\lambda}{\sqrt{p}} \|A\|_F^2 \right\}$$

by noting that it can be written as the solution of a linear regression problem. This procedure is similar to (George et al., 2025).

After estimating A , we compute the empirical and population denoising score matching loss, where the latter is estimated using a validation set. The hyperparameters details are provided below.

Hyperparameters details. In Figure 1, we use $\kappa = 1$, $n = 2 \cdot 10^4$, $t = 10^{-1}$, $d = 3$, $N_g = 10$, $\lambda = 10^{-4}$, and $p \in [2, 10^4]$.

C.2. High-dimensional experiments

Our experiments are performed using an implementation of the DDPM model from (Song et al., 2021). We use the configuration titled `vp.ddpm.cifar10_continuous` which implements the DDPM model of (Ho et al., 2020) but for the continuous-time variance preserving setting. The architecture is a U-Net (Ronneberger et al., 2015) with the encoder and decoder each consisting of four resolution levels ($32 \times 32, 16 \times 16, 8 \times 8, 4 \times 4$), with two residual blocks per level and utilizes self-attention at a resolution of 16×16 . The model is conditioned on time, with the timestep encoded via a sinusoidal embedding and injected into each residual block.

For figures 2 and 5 we modify the number of features (NF) parameter (denoted `nf` in the configuration file and taking a default value of 128). This parameter defines the base channel width used throughout model which is multiplied by $m = 1, 2, 2, 2$ to produce the actual channel widths for each of the four resolution levels. The number of parameters in the model scales quadratically with channel width. We consider $NF \in \{2, 4, 8, 16, 32, 64, 128\}$ to see how the number of parameters affects overparameterization. Some technical modifications to the code were required to support smaller values of NF.

For all high-dimensional experiments we use CIFAR-10 but restrict to the first 1,000 examples to make overfitting easier to achieve. Evaluation is performed on the full 10,000 held-out examples. We train the model with the default setup in the codebase (ADAM with 10% dropout and EMA decay 0.9999) and we train for 1,000,000 iterations, which is larger than typical to better guarantee overfitting. In the experiments of Figure 2 we choose the checkpoint with smallest train loss. The samples generated on the right-hand side of Figure 2 uses the DDIM implementation in the codebase to generate samples with a fixed seed across NF values. We perform log-likelihood calculations using the `bpd` implementation in the codebase.

For the experiment in Figure 4, we consider the DDPM network but we restrict the access of the neural network to the time input. However, we do still scale the network properly according to the time input. We then train the model as usual but restrict the time variable to the range $[t, t + r]$ for $t = 0.1$ and $r \in [0.0001, 0.00032, 0.001, 0.0032, 0.01, 0.032, 0.1, 0.5]$. We then evaluate the train and test error at t to obtain the plot.

C.3. Compute resources

Our experiments were ran on an Amazon EC2 G6e (g6e.xlarge) instance which has a single NVIDIA L40S Tensor Core GPU.

C.4. Licenses

Codebases:

- Score-Based Generative Modeling through Stochastic Differential Equations (Song et al., 2021): Apache License 2.0.

Datasets:

- CIFAR-10 (Krizhevsky, 2009): MIT license.