

LOSING DIMENSIONS: GEOMETRIC MEMORIZATION IN GENERATIVE DIFFUSION

Anonymous authors

Paper under double-blind review

ABSTRACT

Generative diffusion processes are state-of-the-art machine learning models deeply connected with fundamental concepts in statistical physics. Depending on the dataset size and the capacity of the network, their behavior is known to transition from an associative memory regime to a generalization phase in a phenomenon that has been described as a glassy phase transition. Here, using statistical physics techniques, we extend the theory of memorization in generative diffusion to manifold-supported data. Our theoretical and experimental findings indicate that different tangent subspaces are lost due to memorization effects at different critical times and dataset sizes, which depend on the local variance of the data along their directions. Perhaps counterintuitively, we find that, under some conditions, subspaces of higher variance are lost first due to memorization effects. This leads to a selective loss of dimensionality where some prominent features of the data are memorized without a full collapse on any individual training point. We validate our theory with a comprehensive set of experiments on networks trained both in image datasets and on linear manifolds, which result in a remarkable qualitative agreement with the theoretical predictions.

1 INTRODUCTION

Generative diffusion models (Sohl-Dickstein et al., 2015) have achieved spectacular performance in image (Ho et al., 2020; Song and Ermon, 2019; S. et al., 2021) and video (Ho et al., 2022; Singer et al., 2022; Blattmann et al., 2023; B. et al., 2024b) generation and currently form the backbone of most state-of-the-art image generation software (Betker et al., 2023; Esser et al., 2024). The defining feature of these models is their remarkable ability to generalize on complex high-dimensional data distributions. However, diffusion models are known to be capable of fully memorizing the training set in the low-data regime (Somepalli et al., 2023a;b; Kadkhodaie et al., 2024), and in this regime have been shown (Ambrogioni, 2024; Hoover et al., 2023) to be mathematically equivalent to Dense Associative Memory networks, which are modern variants of the celebrated Hopfield model admitting very large memory storage capacity (Krotov and Hopfield, 2016). The ability of the models to memorize has widespread societal implications, as in case of memorization their adoption would likely violate existing copyright laws. Therefore, understanding the factors that lead to generalization and memorization has great practical and theoretical value and drives future developments. It is often conjectured that natural images and other high-dimensional natural data have most of their variability confined on a relatively small sub-space of the ambient space of all possible pixel-values (Peyré, 2009; Fefferman et al., 2016). This *latent manifold* charts the space of possible images, which is embedded in a significantly larger space of meaningless configurations of pixels. The dimensionality of the latent manifold provides an estimate of the richness of generalizations since it quantifies the number of (linearly independent) ways in which an image can be altered while remaining in the space of possible generations. From this geometric perspective, a network that fully memorized the training set corresponds to a zero-dimensional latent manifold (i.e., a collection of points), because the individual “memories” cannot be altered without leaving the space of possible generated images. Generative diffusion models have a rich mathematical structure that closely mirrors systems that have been heavily studied in statistical physics (Montanari,

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

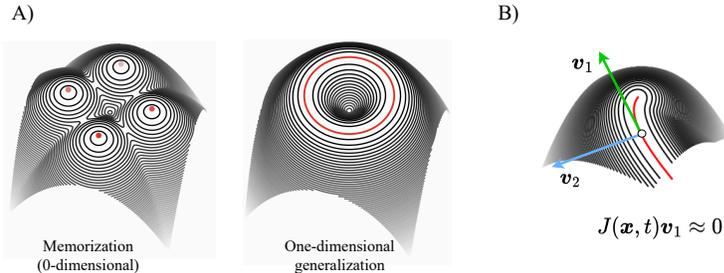


Figure 1: Visualization of the latent manifold of a diffusion model. The contour lines denote the log-probability (i.e. the (negative) ‘energy’). The manifold of fixed points is drawn as a red line. A) Manifolds corresponding to memorization and one-dimensional generalization. B) Tangent and orthogonal singular vectors of the score.

2023; Raya and Ambrogioni, 2024; B. et al., 2024a; Biroli and Mézard, 2023; Ambrogioni, 2023). Recent works have shown that the transition from generalization to memorization in diffusion models is the result of a glass phase transition in the underlying energy function (B. et al., 2024a; Ambrogioni, 2023), which turns generative diffusion in a form of Dense Associative Memory network (Ambrogioni, 2024; Hoover et al., 2023). When the data is supported on a latent manifold, it is then natural to ask if all dimensions are lost at once due to memorization, or if instead they are lost gradually in separate transition events. In this paper, we refer to this intermediate memorization phenomenon as *geometric memorization* (Ross et al., 2024).

2 RELATED WORK

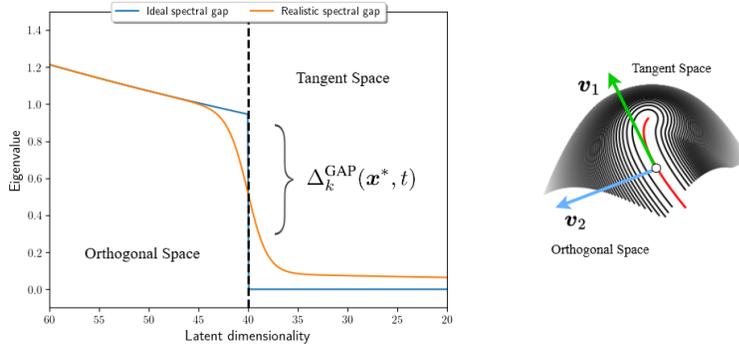
Diffusion Models (DMs) are the current state-of-the-art generative models in several domains, and works such as (De Bortoli, 2022; Pidstrigach, 2022) show that they are capable of learning the manifold structure of the data. There are several methods that, given a trained DM, can then estimate the Local Intrinsic Dimensionality (LID) at individual datapoints, defined as the dimension of the manifold in the component corresponding to the datapoint (Stanczuk et al., 2022; Horvat and Pfister, 2024; Kamkari et al., 2024b). Our analysis makes use of the method from (Stanczuk et al., 2022) to extend the analysis of manifold learning by DMs. The idea that memories can be detected as datapoints with lower LID is formalized in (Kamkari et al., 2024a). While our work shares similar conclusions, our focus is on the interplay between memorization and generalization for diffusion models, and provide a theoretical explanation of such phenomenon. Furthermore, (Kadkhodaie et al., 2024) show, both theoretically and empirically, how generalization arises in diffusion models as a function of the number of datapoints as well as the way the manifold dimension varies along the diffusion trajectory. While their findings are related to our results, our analysis differs as we analyze the nonparametric empirical score and we provide for a statistical physics perspective similar to the work of (Biroli and Mézard, 2023) on mixture of Gaussian data. This differs from (W. et al., 2024) which considers parameterized low-rank score functions. The advantage of the non-parametric empirical formulation is that it provides insight on the large capacity regime, which can cast insight on the behavior of modern large scale architectures. Regarding the usage of statistical mechanics, our work takes inspiration from the analysis conducted by (B. et al., 2024a; Raya and Ambrogioni, 2024; Ambrogioni, 2023) to describe the backward process of DMs as a series of phase transitions similar to those ones occurring in disordered systems.

3 GENERATIVE DIFFUSION MODELS

We will consider a Brownian process where an ensemble of “particles” \mathbf{x}_0 starts from an initial distribution $p_0(\mathbf{x})$ and then evolves through the stochastic equation

$$d\mathbf{x}_t = d\mathbf{Z}_t \tag{1}$$

108
109
110
111
112
113
114
115
116
117
118
119



120
121
122
123
124
125
126
127
128
129
130
131
132
133

Figure 2: Illustration of the gaps in the singular values of the Jacobian of the score function in the presence of a latent manifold. The small values correspond to the tangent manifold (possible generations) while the high values correspond to the orthogonal manifold (forbidden generations). The singular values determine the steepness of the potential well along each eigen-direction.

where $d\mathbf{Z}_t$ is a standard Wiener process. In generative modeling applications, \mathbf{x}_0 is usually chosen to be the distribution of data such as natural images. Eq. (1) is “solved” by the following formal expression: $p(\mathbf{x}_t, t) = \mathbb{E}_{\mathbf{x}_0 \sim p_0} \left[\frac{1}{(2\pi t)^{d/2}} e^{-\frac{\|\mathbf{x}_t - \mathbf{x}_0\|_2^2}{2t}} \right]$, where d is the dimensionality of the ambient space. The *target distribution* $p_0(\mathbf{x})$ can be recovered by using a reversed diffusion process (Anderson, 1982). At large time t_f , we start from a sample $\mathbf{x}_{t_f} \sim \mathcal{N}(0, t_f I_d)$, and let it evolve through the backward process defined by

$$d\mathbf{x}_t = -\nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t) dt + d\mathbf{Z}_t \tag{2}$$

134
135
136
137
138
139
140
141
142
143
144
145
146
147

backward from t_f to $t_0 = 0$. In the machine learning literature, the term $s(\mathbf{x}, t) = \nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ is commonly referred to as the score function. These formulas are given according to what is known as the variance-exploding framework in the generative diffusion literature. However, all results can be ported directly to the variance-preserving (i.e. Ornstein–Uhlenbeck) case.

Given a training set $\{\mathbf{y}^1, \dots, \mathbf{y}^N\}$ sampled from $p_0(\mathbf{x})$, we can obtain a neural approximation of the score function by training a noise-prediction network $\hat{\epsilon}_{\theta}(\mathbf{x}, t)$ (parameterized by θ) using the empirical denoising score matching objective (Hyvärinen and Dayan, 2005; Vincent, 2011; Ho et al., 2020). The network $\hat{\epsilon}_{\theta}(\mathbf{x}_t, t)$ is trained to predict the added noise ϵ from a noisy state $\mathbf{x}_t = \mathbf{x}_0 + \sqrt{t}\epsilon$. The estimated score is then given by $\hat{s}_{\theta}(\mathbf{x}, t) = -\frac{\hat{\epsilon}_{\theta}(\mathbf{x}, t)}{\sqrt{t}}$. Learning $\hat{\epsilon}_{\theta}$ instead of \hat{s}_{θ} has the advantage of maintaining finite the output of the network for small $t \rightarrow 0$, where we know that the magnitude of the score becomes infinite outside of the support of the data.

4 LOCAL GEOMETRY AND LATENT MANIFOLDS

148
149
150
151
152
153

Here, we assume that the data distribution $p_0(\mathbf{x})$ is supported on a d -dimensional manifold \mathcal{M}_0 embedded in \mathbb{R}^n . We study the family of latent manifolds \mathcal{M}_t for different values of the diffusion time t . When the distribution inside the manifold is uniform, these manifolds can be defined as sets of stable fixed-points of the score function:

$$\mathcal{M}_t = \{\mathbf{x}^* \mid s(\mathbf{x}^*, t) = \mathbf{0}, \text{ with } s(\mathbf{x}^*, t) \text{ n.s.d.}\} \tag{3}$$

154
155
156
157
158
159
160
161

where the negative semi-definiteness (n.s.d.) condition refers to the Jacobian of the score function at \mathbf{x}^* . Due to the Gaussian annulus theorem, during high-dimensional generative diffusion the samples “orbit” a shell around the manifold. However, while the manifold itself is not usually visited by the particles, its shape reflects the shape of the potential around it, which guides the generated trajectories. In order to understand the generative process, it is therefore important to study how the manifold emerges and changes during generation, and

162 how its shape depends on the distribution of the data and on the number of training samples.
 163 Its local geometry can be studied by analyzing the linearization of the score function around
 164 each point (Stanczuk et al., 2023). More precisely, for a given point \mathbf{x}^* in \mathcal{M}_t , we can
 165 analyze the effect of adding a small perturbation vector \mathbf{p} with magnitude in the order of \sqrt{t} :

$$166 \quad s(\mathbf{x}^* + \mathbf{p}, t) \approx J(\mathbf{x}^*, t) \mathbf{p} , \quad (4)$$

167 where $J(\mathbf{x}^*, t)$ is the *smoothed Jacobian matrix* $J(\mathbf{x}, t)$, whose columns are defined as

$$168 \quad \mathbf{J}_j(\mathbf{x}, t) = \left[s(\mathbf{x} + \sqrt{t}\mathbf{e}_j, t) - s(\mathbf{x}, t) \right] / \sqrt{t} , \quad (5)$$

169 where \mathbf{e}_j is a vector in an orthonormal basis set, which converges to the exact Jacobian of the
 170 score for $t \rightarrow 0$. This discretization in the definition has the advantage of only considering
 171 perturbations on the scale given by \sqrt{t} , since smaller perturbations cannot be resolved at
 172 the given noise level (Stanczuk et al., 2023). The effect of the linear perturbations can be
 173 expressed in terms of the singular values of $J(\mathbf{x}^*, t)$:

$$174 \quad J(\mathbf{x}^*, t) \mathbf{p} = \sum_j \lambda_j(\mathbf{x}^*, t) \mathbf{v}_j (\mathbf{w}_j \cdot \mathbf{p}) , \quad (6)$$

175 where the \mathbf{w}_j (\mathbf{v}_j) is the j -th right (left) singular vector and $\lambda_j(\mathbf{x}^*, t)$ is its associated (non-
 176 negative) singular value. The right singular vectors \mathbf{w}_j define a set of linearly independent
 177 perturbations of \mathbf{x}^* , while the scaled left singular vectors $-\lambda_j(\mathbf{x}^*, t)\mathbf{v}_j$ give a linearization of
 178 the score at the perturbed point. Generally, \mathbf{v}_j is roughly aligned to \mathbf{w}_j , meaning that the
 179 tend to push the diffusing particle back towards the point on the manifold with a strength
 180 determined by $\lambda_j(\mathbf{x}^*, t)$. For $t \rightarrow 0$, the matrix $J(\mathbf{x}, t)$ becomes symmetric and the singular
 181 values correspond to the negative of the eigenvalues. We refer the set of singular values of
 182 $J(\mathbf{x}^*, t)$ as its *spectrum*. By analyzing these spectra, we can extract important information
 183 concerning the local geometry of \mathcal{M}_t . For a given \mathbf{x}^* , consider the set of left singular vectors
 184 that are orthogonal to $T_{\mathcal{M}_t}(\mathbf{x}^*)$. These vectors correspond to perturbations that move
 185 the particle orthogonally outside of the manifold, which are therefore associated with large
 186 singular values. On the other hand, perturbations that lie inside the tangent space correspond
 187 to a set of vanishing singular values since along these directions the particles can diffuse
 188 almost freely. This is visualized geometrically in Fig. 2 B. From these considerations, it
 189 follows the dimensionality of $T_{\mathcal{M}_t}$ can be estimated as the dimensionality of the right singular
 190 space of $J(\mathbf{x}^*, t)$ with 0 singular value. The drop between non-zero and zero singular values
 191 can be detected as a discontinuity (i.e. a gap) in its spectrum (see Fig. 2) (Stanczuk et al.,
 192 2023). In real datasets, the target distribution is generally not uniform when restricted to the
 193 manifold. In this case, the latent manifold cannot be defined as a set of fixed-points, we can
 194 still define \mathcal{M}_t by annealing the target distribution into a uniform distribution defined over
 195 its support. Aside from definitional issues, the main consequence of having a non-uniform
 196 distribution is that the Jacobian will generally not have zero singular values. Instead, the
 197 singular value corresponding to tangent right singular vectors $\mathbf{v} \in \mathcal{M}_0$ will have a magnitude
 198 that depends inversely on the local variance of the data along that direction. While we
 199 cannot simply count the space corresponding to zero singular values, we can still quantify
 200 the dimensionality of the latent manifold by analyzing the time-dependency of gaps (i.e.
 201 sharp jumps) in the spectrum:

$$202 \quad \Delta_{\text{GAP}}^{(k)}(\mathbf{x}^*, t) = \lambda_{k+1}(\mathbf{x}^*, t) - \lambda_k(\mathbf{x}^*, t) . \quad (7)$$

203 In this case, the spectra of $J(\mathbf{x}^*, t)$ contain more information than just the local dimensionality
 204 of the manifold \mathcal{M}_t . Vectors in the tangent space $T_{\mathcal{M}}(\mathbf{x}^*, t)$ define locally linear sub-spaces
 205 with different variance. The variances of these local linear sub-spaces characterize the
 206 variability of the data along the corresponding directions.

207 5 THEORETICAL ANALYSIS

208 Here, we will provide a theoretical explanation of geometric memorization (dimensionality
 209 loss) using concepts from random-matrix theory and the statistical physics of disordered
 210 systems. To study the dynamics of the latent manifold evolution analytically, we will consider

data distributed according to the linear model $\mathbf{x}_0 = F\mathbf{z}$ where \mathbf{z} is a m -dimensional standard Gaussian vector and F is a $d \times m$ projection matrix. Equivalently, $\mathbf{x}_0 \sim \mathcal{N}(0, FF^\top)$. The choice of the linear model simplifies the statistical analysis while qualitatively capturing important features of the local phenomenology of the tangent spaces. In fact, at time t the curvature of \mathcal{M}_t is suppressed by the smoothing induced by the forward process, which, roughly speaking, linearizes the geometry at the same scale as our local Jacobian analysis (see Eq. (5)).

5.1 THE EXACT SCORE

In the linear Gaussian case, the Jacobian of the exact score function is independent of \mathbf{x} and is given by the formula

$$J_t/t = \frac{1}{t}F \left[I_m + \frac{1}{t}F^\top F \right]^{-1} F^\top - I_d. \quad (8)$$

where we considered J_t/t to control the divergence of the score at $t \rightarrow 0$. If the columns of F are mutually orthogonal with $\|F_k\|_2^2 = \sigma^2$ for all indices, the singular spectrum of J_t has a single gap at $d - m$, which encodes the dimensionality of the manifold. In a parallel paper (Anonymous, 2024), it is shown that the size of the gap at time t is given by

$$\Delta_{GAP}(t; \sigma) = \frac{\sigma^2(1 + \alpha_m^{-1/2})^2}{t + \sigma^2(1 + \alpha_m^{-1/2})^2}. \quad (9)$$

where $\alpha_m = m/d$ denotes the ratio between the dimensionality of the manifold and the ambient dimensionality. While the gap is always present, in practice it only becomes distinguishable from the background noise at a finite time. Therefore, this formula indicates that subspaces with large variance emerge earlier during the generative reverse process. In the presence of m different subspaces with dimensions $\{d_1, \dots, d_m\}$ and variance $\{\sigma_1^2, \dots, \sigma_m^2\}$, we will have one total manifold gap at $d - m$ and $m - 1$ intermediate gaps whose size approaches $(\sigma_k^{-2} - \sigma_{k+1}^{-2})/t$ for $t \rightarrow 0$. Only the total manifold gap remains for $t \rightarrow 0$ due to the $1/t$ normalization of the score. This reflects the fact that the component of the score orthogonal to the manifold diverges while the parallel component reaches a constant value, leaving us with: $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \simeq \frac{1}{t} [\Pi - I_d] \mathbf{x}$, where $\Pi = F(F^\top F)^{-1}F^\top$ is the projector on the linear space \mathcal{M} .

5.2 THE EMPIRICAL SCORE

Computing the exact score function involves an average with respect to the true target distribution p_0 . In real applications, we do not have direct access on p_0 , whose behavior can only be inferred through a finite training set comprised of N samples $\{\mathbf{y}^1, \dots, \mathbf{y}^N\}$, with $\mathbf{y}^\mu \stackrel{iid}{\sim} p_0$. When training, we sample from the empirical distribution which we use as a proxy of the true distribution. The empirical distribution at time t in the variance exploding framework is

$$p_t^N(\mathbf{x}) = \frac{1}{N\sqrt{(2\pi t)^d}} \sum_{\mu=1}^N e^{-\frac{\|\mathbf{x} - \mathbf{y}^\mu\|^2}{2t}}. \quad (10)$$

From the empirical distribution, we can write down the empirical score:

$$\nabla \log p_t^N(\mathbf{x}) = \sum_{\mu=1}^N w_\mu(\mathbf{x}, t) (\mathbf{y}^\mu - \mathbf{x}) / t, \quad (11)$$

where the weight $w_\mu(\mathbf{x}, t) = p(\mathbf{y}^\mu | \mathbf{x}) / \sum_{\nu=1}^N p(\mathbf{y}^\nu | \mathbf{x})$ is the posterior probability of the pattern \mathbf{y}^μ given the noisy state \mathbf{x} , were the possible states are restricted to the empirical set. This estimator is consistent, meaning that its bias approached the true score for $N \rightarrow \infty$. The random sampling of the dataset introduces statistical fluctuations that we can quantify by considering the estimator variance, which for large N can be approximated as

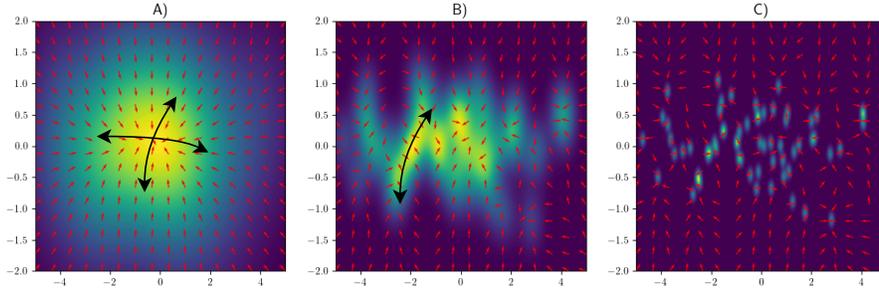


Figure 3: Visualization of the dimensionality loss phenomenon. Manifold sub-spaces with higher variance are lost due to ‘condensation’ (i.e. memorization). Panels A,B and C show the score estimated from a bivariate distribution with unequal variances for $\beta = 1$, $\beta = 10$ and $\beta = 100$ respectively. The red arrows show the empirical score while the heat-map visualizes the density.

$\text{var}[\nabla \log p_t^N(\mathbf{x})] \approx \text{var}(\mathbf{x}_0 | \mathbf{x}) / \mathbb{E}[\tilde{N}_t(\mathbf{x})]$, where $\text{var}(\mathbf{x}_0 | \mathbf{x})$ is the true posterior variance and $\tilde{N}_t(\mathbf{x}) = \left(\sum_{\mu=1}^N w_{\mu}^2(\mathbf{x}, t)\right)^{-1} \leq N$ is the effective number of data points used to estimate the score. When $t \rightarrow 0$, we always have that $\tilde{N}_t(\mathbf{x}) \rightarrow 1$, because the empirical score always fully memorizes in this limit. However, the empirical score exhibits generalization when the expected value is larger than the standard deviation induced by $\tilde{N}_t(\mathbf{x})$. Note that $\tilde{N}_t(\mathbf{x})$ is a function of the state \mathbf{x} and that, consequently, the fluctuations in the empirical score depend on the ‘‘location’’ \mathbf{x} . This property is fundamental in our analysis of geometric memorization.

5.3 MEMORIZATION AS GLASSY PHASE TRANSITION

The statistical behavior of the empirical score can be analyzed in the large N limit by interpreting Eq. (10) as proportional to the partition function of a Random Energy Model (REM) (B. et al., 2024a; Lucibello and Mézard, 2024), which offers a simple model of disordered thermodynamic systems. The thermodynamic analysis of generative diffusion models is outlined in (Ambrogioni, 2023). In summary, each energy level E_{μ} is associated with a data point \mathbf{y}^{μ} in the training set and its energy depends on its Euclidean distance with the current state \mathbf{x}_t (Ambrogioni, 2023), with the energy given by

$$E_{\mu}(\mathbf{x}) = -\frac{1}{2}\|\mathbf{y}^{\mu}\|^2 + \mathbf{x} \cdot \mathbf{y}^{\mu} \quad (12)$$

which leads to the partition function

$$Z_N(\mathbf{x}, t) = \sum_{\mu=1}^N e^{-\frac{1}{t}E_{\mu}(\mathbf{x})} \quad (13)$$

where the time parameter t is analogous to the temperature of the system, which can be used to express the weights as a Boltzmann distribution: $w_{\mu}(\mathbf{x}, t) = \frac{1}{Z_N(\mathbf{x}, t)} e^{-\frac{1}{t}E_{\mu}(\mathbf{x})}$. Since the empirical score is a Boltzmann average according to Eq. (13), studying its fluctuation under the random sampling of the data allows us to quantify the deviations from the exact score due to memorization effects. In our case, the energy levels are distributed according to

$$p(E; \mathbf{x}) = \int_{\mathbb{R}^n} \delta\left(E + \frac{1}{2}\|\mathbf{y}\|^2 - \mathbf{x} \cdot \mathbf{y}\right) dP_0(\mathbf{y}) \quad (14)$$

For small values of t and large dataset sizes, the empirical score can be shown to be self-averaging, meaning that it is insensitive to the specific sampling of the training points, resulting to generalization of the underlying distribution. More formally, from the physical theory of REMs (Derrida, 1981), we know that, at the asymptotic limit of $d \rightarrow \infty$, the

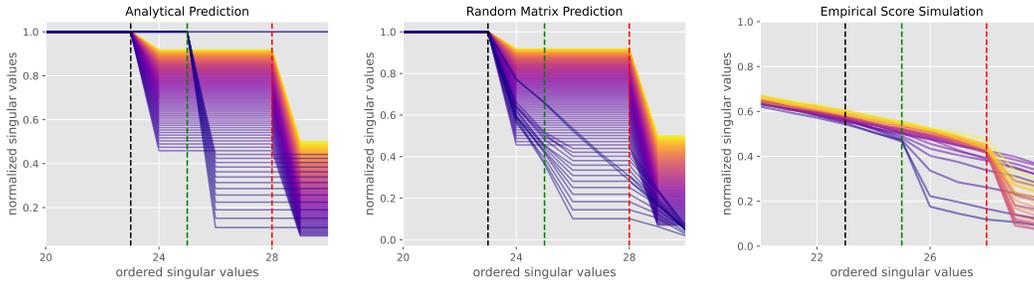


Figure 4: The ordered singular values of the Jacobian of the empirical score function of a linear manifold model as a function of the diffusion time t . Lighter colours are associated to larger times in the colour map. The parameters for the model are $d = 30$, $m = 7$, $\log(N)/d = 0.23$ with subspaces associated to variances $\sigma_1^2 = 1$ and $\sigma_2^2 = 0.3$ with dimensions $m_1 = 2$ and $m_2 = 5$ respectively. Left: approximated theoretical prediction in the memorization phase according to Eq. (19). Center: prediction from the approximated Jacobian in Eq. (18). Right: singular values obtained by the numerical measure of the Jacobian of the empirical score function (as described in Supp. B), evaluated from a synthetic data set of $N = 10^3$ points.

statistical system specified by Eq. (10) undergoes a random phase transition that separates a self-averaging high-temperature regime to a *condensation* regimes where Boltzmann averages depend on a small (i.e. sub-exponential) fraction of energy levels (Montanari and Mézard, 2009). In Supp. C, we show that, for d much smaller than N , the condensation time for linear manifold data is, to a leading exponential order, equal to

$$t_c(\mathbf{x}) = \sqrt{\frac{d}{2\log(N)} \left(\frac{r_{4,\sigma}}{2} + \omega^2(\mathbf{x}) \right)}, \quad (15)$$

which demarcates the diffusion time when the empirical score becomes susceptible to fluctuations introduced by the random sampling of the dataset. In the formula, the term $r_{4,\sigma} = d^{-1} \sum_{i=1}^d \sigma_i^4$ captures the fluctuations in the norm of the data, while the directional quantity $\omega^2(\mathbf{x}) = d^{-1} \sum_{i=1}^d x_i^2 \sigma_i^2$ is the *variance density* along the direction \mathbf{x} . As we shall see, the balance between these two quantities plays a crucial role in geometric memorization effects. From standard REM calculations (see Supp. D), we can express the effective number of data points used to estimate the score at \mathbf{x} at time t as

$$\tilde{N}_t(\mathbf{x}) = \min\left(N, \frac{t}{1 - t_c^{-1}(\mathbf{x})}\right). \quad (16)$$

where we introduced the minimum operator heuristically to account for the finite size of the system. The exact asymptotic theory is recovered for $N \rightarrow \infty$. Note that, since these quantities scale to the leading exponential order, they are therefore neglected quantities that scale sub-exponentially in N .

5.4 A THEORY OF GEOMETRIC MEMORIZATION

The large N analysis outlined in the previous sections give us a description of the fluctuations in the empirical score as a function of the state \mathbf{x} . The “spatial” dependency of these fluctuations ultimately depend on the data distribution $p_0(\mathbf{x})$, which outlines a rich geometric landscape that interacts in a complex way with the “spatial” variations in the exact score $\nabla \log p_t(\mathbf{x})$. To study the effect of this “spatially” non-homogeneous random fluctuations on the spectrum, we start from an approximate formula for the empirical score obtained by restricting the average to only $\tilde{N}_t(\mathbf{x})$ “active samples”:

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \approx \frac{1}{\tilde{N}_t(\mathbf{x})} \sum_{\mu=1}^{\tilde{N}_t(\mathbf{x})} (\mathbf{y}^\mu - \mathbf{x}) / t \quad (17)$$

where the “active samples” \mathbf{y}^μ are sampled from the posterior distribution $p(\mathbf{x}_0 | \mathbf{x}; t)$. In the linear Gaussian case, Eq. (17) follows a Normal distribution, since the posterior is itself

Normal. In the following, for the sake of simplicity, we will assume that F is a diagonal $d \times d$ with diagonal entries σ_k , with $\sigma_k = 0$ for $d - m$ indices. This corresponds to a rotation to the basis of eigenvectors of $F^\top F$. If we assume that the fluctuations in the score are uncorrelated for a separation of the order of \sqrt{t} , we can quantify the statistical variability of the (smoothed) Jacobian (Eq. (5)) through the formula

$$J_{ij}(t) \sim \mathcal{N} \left(-\delta_{ij} (t + \sigma_i^2)^{-1}, \frac{\sigma_i^2}{t(t + \sigma_i^2)} \left[\phi(t, \mathbf{0}) + \phi(t, \mathbf{e}_j \cdot \sqrt{t}) \right] \right), \quad (18)$$

where we defined the function $\phi(t, \mathbf{x}) = \max(1/N, t^{-1} - t_c^{-1}(\mathbf{x}))$. The expected value of this expression is just the Jacobian of the exact score, which determines the opening of the spectral gaps as explained in Sec. 5.1. On the other hand, in this model gaps can close due to the variance term. We can see this phenomenon qualitatively by considering the singular values spectrum of the expected value of Eq. (18):

$$\bar{s}_i = \sqrt{\frac{1}{(t + \sigma_i^2)^2} + \sum_{k=1}^d \frac{\sigma_k^2}{t^2 (t + \sigma_k^2)^2} \left[\phi(t, \mathbf{0}) + \phi(t, \mathbf{e}_i \cdot \sqrt{t}) \right]^2}. \quad (19)$$

Remember that we see a gap in the sorted spectrum if there is a large difference between two consecutive sorted singular values s_k and s_{k+1} . This gap can disappear if I) $\phi(t, \mathbf{e}_k \cdot \sqrt{t})$ is larger than $\phi(t, \mathbf{e}_{k+1} \cdot \sqrt{t})$, or II) if the contribution of these variance terms make the contribution of the expected value negligible. Case I) is directional, as it depends on the direction of perturbations \mathbf{e}_k and \mathbf{e}_{k+1} and it leads to the selective suppression of a particular subspace. On the other hand, case II) is non-directional: it induces a synchronous suppression of all gaps and leads to complete memorization. The phenomenon of selective memorization is visualized in Fig. 3 for a two-dimensional distribution. For linear Gaussian, the closing times are determined by the critical time $t_c^{-1}(\mathbf{x})$, which itself depends on the constant term $r_{4,\sigma} = d^{-1} \sum_{i=1}^d \sigma_i^4$ and on the ‘‘directional’’ term $\omega^2(\mathbf{x}) = d^{-1} \sum_{i=1}^d x_i^2 \sigma_i^2$. This latter term is proportional to the variance along the subspace spanned by \mathbf{x} and plays a crucial role in determining the differential disappearance of different subspaces at different times. Perhaps counter-intuitively, the subspace spanned by \mathbf{e}_k is more vulnerable to memorization when $\omega^2(\mathbf{e}_k)$ is large. Therefore, subspaces that are more prominent in the distribution of the data and that emerge earlier during the diffusion process are also more vulnerable to memorization in the later stages of diffusion. This correspond to the form of feature memorization suggested in (Ross et al., 2024).

6 EXPERIMENTS

6.1 TESTING NUMERICALLY THE THEORY OF GEOMETRIC MEMORIZATION

Fig. 4 shows the evolution in time of the spectrum of the singular values obtained from Eq. (19) (left panel) Eq. (18) (central panel) and the direct computation of the empirical score introduced in section 5.2 (right panel). The experimental curves obtained from the empirical score look consistent with the theory, both in signaling the dimension of the subspaces and the opening times for the gaps, displaying the hierarchical structure associated to the ordering of the variances. Additional experiments are reported in Supp. D.1.

6.2 DIFFUSION NETWORKS TRAINED ON LINEAR MANIFOLD DATA

While our theory analyzes the empirical score, our experiments show that our results cast fundamental insight on geometric memorization in trained networks. Fig. 5 shows the spectra estimated from network trained on a linear manifold where the matrix F is built such that datapoints live in two sub-spaces with variances equal to 1 (high variance) and 0.3 (low variance) respectively. The details of the experiment are given in Supp. E. We can compare these results with the spectra obtained in Fig. 4. The behavior of the trained network has several of the qualitative features predicted by our theoretical analysis. When N is large, the spectra show the total manifold gap, as predicted by the exact theory. This gap remains present in the network even for $t = 1e^{-5}$, which shows that that trained network have a

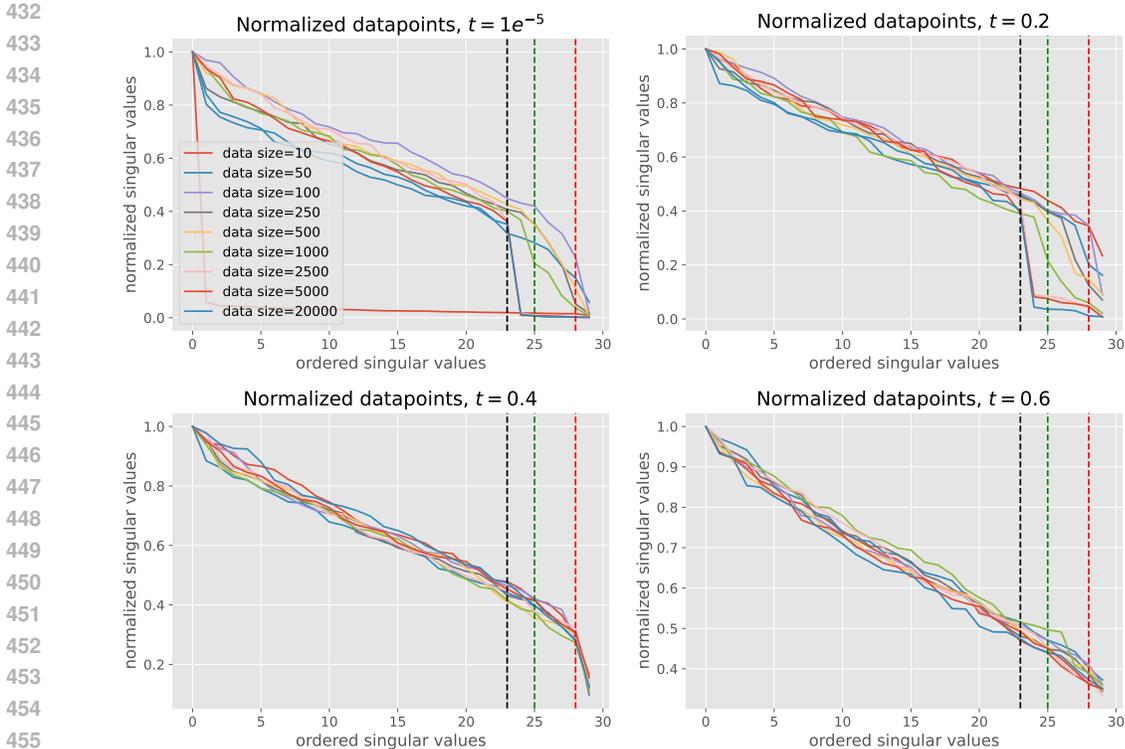


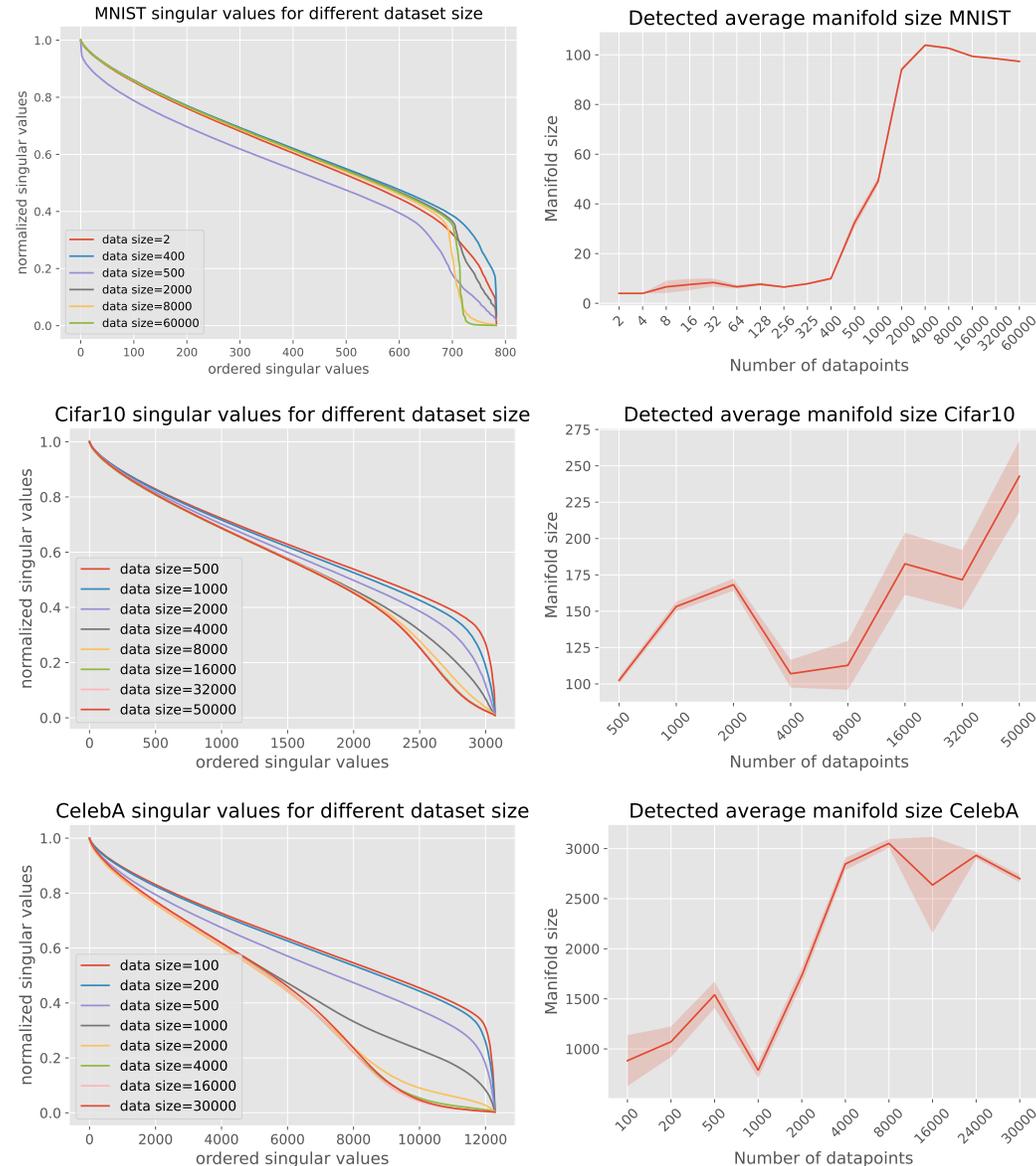
Figure 5: Spectra for different t estimated from diffusion networks trained on linear Normal data with two subspaces with different variance using a range of dataset sizes. The red (green) dashed line correspond to the location of the theoretical spectral gap for the high (low) variance subspace. The black dashed line corresponds to the total manifold gap.

tendency to generalize, likely due to their finite capacity and by the implicit regularization induced by the parameterization. For intermediate values of N (e.g. $N = 1000$), the theory predicts that, for small t , only the low variance gap should be observable, as the high variance sub-space is lost due to geometric memorization. Interestingly, this counterintuitive behavior seems to be present in the trained networks, where for $N = 250$ and $N = 500$, the drop in the spectrum is roughly aligned with the low-variance gap (green dashed line). Furthermore, the final generalization profile at $t = 1e^{-5}$ still behaves according to this prediction, suggesting that the ultimate generalization of the network can be predicted by the temporal dynamic of the empirical score. Finally, for small dataset sizes the behavior of the network disaligns from the theoretical prediction as the network returns to exhibit the high variance gap even for $t = 1e^{-5}$. This is not surprising since in this regime our large N analysis is outside its domain of applicability, and the behavior of the network seems to reflect the global fit of a parameterized linear model. This is also visible in an additional experiment given in Supp. E.

6.3 GEOMETRIC MEMORIZATION IN NATURAL IMAGE DATASETS

We will now report the results of a series of experimental analysis where we trained on a series of increasingly large sub-datasets extracted from MNIST, Cifar10 and Celeb10. For each dataset size and time point, we estimate the latent dimensionality by locating the largest spectral gap and we study how dimensionality changes with the dataset size. The full details of dimensionality estimation and training are given in Supp. A and Supp. B. Fig. 6 shows the average spectra (left) and the average detected dimensionality (right) for the whole dataset (both train and test set), and see how that changes as a function of dataset size and diffusion time. Sharp spectral gaps are visible in the network trained on the MNIST dataset, where the estimated dimensionality increases sharply starting from 400 data points and reaches its peak at around 4000. The other datasets show less clear gaps in their spectra.

486 However, the location of their spectral inflection point decreases predictably as function of
 487 the dataset size, revealing an increasing trend in the estimated dimensionality. Interestingly,
 488 the dimensionality in Cifar10 does not seem to saturate, suggesting that the total dataset
 489 size still results in partial (geometric) memorization.
 490



529 Figure 6: Spectra at $t = 1e^{-5}$ estimated on deep networks trained on natural image datasets
 530 with different dataset sizes. The estimated dimensionality tend to increase with the dataset
 531 size, suggesting a phenomenon of geometric memorization.
 532

534 **7 CONCLUSIONS**

535

536 Our work opens the door for further analysis of generative diffusion using the tools of
 537 statistical physics, differential geometry and random matrix theory, which may cast light on
 538 generalization in these fascinating generative methods. Our theoretical analysis was focused
 539 on the empirical score function, further research may analyze theoretical Jacobian spectra in
 simple trained models and elucidate their deviations from the empirical theory.

REFERENCES

- 540
541
542 Ambrogioni, L. (2023). The statistical thermodynamics of generative diffusion models. *arXiv*
543 *preprint arXiv:2310.17467*.
- 544 Ambrogioni, L. (2024). In search of dispersed memories: Generative diffusion models are
545 associative memory networks. *Entropy*, 26(5):381.
- 546
547 Anderson, B. D. (1982). Reverse-time diffusion equation models. *Stochastic Processes and*
548 *their Applications*, 12(3):313–326.
- 549 Anonymous (2024). Manifolds, random matrices and spectral gaps: The geometric phases of
550 generative diffusion. *Supplementary*.
- 551
552 B., G., B., T., B., V., and M., M. (2024a). Dynamical regimes of diffusion models. *arXiv*
553 *preprint arXiv:2402.18491*.
- 554 B., T. et al. (2024b). Video generation models as world simulators.
- 555
556 Betker, J. et al. (2023). Improving image generation with better captions. *Computer Science.*,
557 2(3):8.
- 558
559 Biroli, G. and Mézard, . (2023). Generative diffusion in very large dimensions. *Journal of*
560 *Statistical Mechanics: Theory and Experiment*, 2023(9):093402.
- 561 Blattmann, A. et al. (2023). Align your latents: High-resolution video synthesis with latent
562 diffusion models. In *Conference on Computer Vision and Pattern Recognition*.
- 563
564 De Bortoli, V. (2022). Convergence of denoising diffusion models under the manifold
565 hypothesis. *Transactions on Machine Learning Research*.
- 566
567 Deng, L. (2012). The mnist database of handwritten digit images for machine learning
568 research [best of the web]. *IEEE Signal Processing Magazine*, 29(6):141–142.
- 569 Derrida, B. (1981). Random-energy model: An exactly solvable model of disordered systems.
570 *Physical Review B*, 24:2613–2626.
- 571 Esser, P., Kulal, S., Blattmann, A., Entezari, R., Müller, J., Saini, H., Levi, Y., Lorenz, D.,
572 Sauer, A., Boesel, F., et al. (2024). Scaling rectified flow transformers for high-resolution
573 image synthesis. In *International Conference on Machine Learning*.
- 574
575 Fefferman, C., Mitter, S., and Narayanan, H. (2016). Testing the manifold hypothesis.
576 *Journal of the American Mathematical Society*, 29(4):983–1049.
- 577
578 Ho, J., Jain, A., and Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances*
579 *in Neural Information Processing Systems*.
- 580
581 Ho, J., Salimans, T., Gritsenko, A., Chan, W., Norouzi, M., and Fleet, D. J. (2022). Video
582 diffusion models. *arXiv preprint arXiv:2204.03458*.
- 583 Hoover, B., Strobelt, H., Krotov, D., Hoffman, J., Kira, Z., and Chau, D. H. (2023). Memory
584 in plain sight: A survey of the uncanny resemblances between diffusion models and
585 associative memories. *arXiv preprint arXiv:2309.16750*.
- 586
587 Horvat, C. and Pfister, J. (2024). On gauge freedom, conservativity and intrinsic dimension-
588 ality estimation in diffusion models. In *The Twelfth International Conference on Learning*
Representations.
- 589
590 Hyvärinen, A. and Dayan, P. (2005). Estimation of non-normalized statistical models by
591 score matching. *Journal of Machine Learning Research*, 6(4).
- 592
593 Kadkhodaie, Z., Guth, F., Simoncelli, E. P., and Mallat, S. (2024). Generalization in
diffusion models arises from geometry-adaptive harmonic representations. In *International*
Conference on Learning Representations.

- 594 Kamkari, H., Ross, B. L., Cresswell, J. C. and Caterini, A. L., Krishnan, R., and Loaiza-
595 Ganem, G. (2024a). A geometric explanation of the likelihood ood detection paradox. In
596 *International Conference on Machine Learning*.
597
- 598 Kamkari, H., Ross, B. L., Hosseinzadeh, R., Cresswell, J. C., and Loaiza-Ganem, G. (2024b).
599 A geometric view of data complexity: Efficient local intrinsic dimension estimation with
600 diffusion models. *arXiv preprint arXiv:2406.03537*.
- 601 Krizhevsky, A., Nair, V., and Hinton, G. (2014). The cifar-10 dataset. 55(5):2.
602
- 603 Krotov, D. and Hopfield, J. J. (2016). Dense associative memory for pattern recognition.
604 *Advances in neural information processing systems*.
- 605 Liu, Z., Luo, P., Wang, X., and Tang, X. (2015). Deep learning face attributes in the wild.
606 In *Proceedings of International Conference on Computer Vision*.
607
- 608 Lucibello, C. and Mézard, M. (2024). Exponential capacity of dense associative memories.
609 *Physical Review Letters*, 132(7):077301.
- 610 Mezard, M., Parisi, G., and Virasoro, M. (1986). *Spin Glass Theory and Beyond*. World
611 Scientific.
- 612 Montanari, A. (2023). Sampling, diffusions, and stochastic localization. *arXiv preprint*
613 *arXiv:2305.10690*.
614
- 615 Montanari, A. and Mézard, M. (2009). *Information, physics, and computation*. Oxford Univ.
616 Press.
- 617 Peyré, G. (2009). Manifold models for signals and images. *Computer Vision and Image*
618 *Understanding*, 113(2):249–260.
619
- 620 Pidstrigach, J. (2022). Score-based generative models detect manifolds. *Advances in Neural*
621 *Information Processing Systems*, 35:35852–35865.
- 622 Raya, G. and Ambrogioni, L. (2024). Spontaneous symmetry breaking in generative diffusion
623 models. *Advances in Neural Information Processing Systems*, 36.
624
- 625 Ross, B. L., Kamkari, H., Liu, Z., Wu, T., Stein, G., Loaiza-Ganem, G., and Cresswell, J. C.
626 (2024). A geometric framework for understanding memorization in generative models. In
627 *ICML 2024 Next Generation of AI Safety Workshop*.
- 628 S., Y. et al. (2021). Score-based generative modeling through stochastic differential equations.
629 In *International Conference on Learning Representations*.
630
- 631 Salimans, T., Karpathy, A., Chen, X., and Kingma, D. P. (2017). Pixelcnn++: Improving
632 the pixelcnn with discretized logistic mixture likelihood and other modifications. *arXiv*
633 *preprint arXiv:1701.05517*.
- 634 Singer, U. et al. (2022). Make-a-video: Text-to-video generation without text-video data.
635 *arXiv preprint arXiv:2209.14792*.
636
- 637 Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and Ganguli, S. (2015). Deep unsu-
638 pervised learning using nonequilibrium thermodynamics. In *International Conference on*
639 *Machine Learning*.
- 640 Somepalli, G., Singla, V., Goldblum, M., Geiping, J., and Goldstein, T. (2023a). Diffusion
641 art or digital forgery? investigating data replication in diffusion models. In *IEEE/CVF*
642 *Conference on Computer Vision and Pattern Recognition*.
- 643 Somepalli, G., Singla, V., Goldblum, M., Geiping, J., and Goldstein, T. (2023b). Under-
644 standing and mitigating copying in diffusion models. *Advances in Neural Information*
645 *Processing Systems*.
646
- 647 Song, Y. and Ermon, S. (2019). Generative modeling by estimating gradients of the data
distribution. *Advances in Neural Information Processing Systems*.

648 Stanczuk, J., Batzolis, G., Deveney, T., and Schönlieb, C. (2022). Your diffusion model
649 secretly knows the dimension of the data manifold. *arXiv preprint arXiv:2212.12611*.
650

651 Stanczuk, J. P., Batzolis, G., Deveney, T., and Schönlieb, C. (2023). Diffusion models
652 encode the intrinsic dimension of data manifolds. In *International Conference on Machine*
653 *Learning*.

654 Vincent, P. (2011). A connection between score matching and denoising autoencoders. *Neural*
655 *Computation*, 23(7):1661–1674.
656

657 W., P., Z., H., Z., Z., C., S., M., Y., and Q., Q. (2024). Diffusion models learn low-dimensional
658 distributions via subspace clustering. *arXiv preprint arXiv:2409.02426*.

659 Y., T., C., J. Y., K., S., and R., E. K. (2023). Diffusion probabilistic models generalize
660 when they fail to memorize. In *ICML Workshop on Structured Probabilistic Inference &*
661 *Generative Modeling*.
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A EXPERIMENTAL METHODOLOGY: TRAINING AND MODEL ARCHITECTURE DETAILS

Dataset	Image Size	Latent Dim.	Channel Mult.	Param. Count	Batch size	Iterations
Cifar10	32	128	(1, 2, 2, 2)	35.7M	128	500,000
Mnist	28	128	(1, 2, 2)	24.5M	128	400,000
CelebA-HQ	64	64	(1, 1, 2, 2, 4, 4)	27.4M	64	800,000

Table 1: Table displaying both model and training configurations for each dataset.

For the toy examples, we train a Variance Exploding continuous score model with 2M training steps with batch size 128. We use a Residual Multi Layer Perceptron with hidden size of 128, with two residual blocks. Each block is composed by two linear layers with SiLu activation.

For the image models, we follow the diffusion setting in (Ho et al., 2020). We kept the variance scheduler, where $\beta_{\min} = 10^{-4}$ and $\beta_{\max} = 2 \times 10^{-2}$, the time steps $T = 1000$, and the score model backbone (PixelCNN++ (Salimans et al., 2017)) the same. In addition, for each of the datasets, we varied the model’s channel multipliers, latent dimension, batch size, and training iterations to account for the complexity of the dataset and our available computing resources; see Table 1. For context, we primarily utilized NVIDIA Tesla V100 GPUs with 32 GB of memory for the training of our models.

For consideration of the data sizes we chosen for our experiment, we closely followed the setup of (Y. et al., 2023). Specifically, we first trained multiple diffusion models using different data split sizes from $\{0.5k, 1k, 2k, 4k, \dots, |S|\}$, where $|S|$ is the full size of a given dataset. We trained our models without random flipping and utilized the exponential moving average version of the trained models, where we set the decay value to 0.9999 during training. For Cifar10 (Krizhevsky et al., 2014) and Mnist (Deng, 2012), we did not center-crop or resize the images. However, for CelebA-HQ (Liu et al., 2015), we center-cropped and downsized the images to 64×64 resolution. Moreover, we only use dropout for Cifar10 with the value of 0.1. To get a more comprehensive view of the reduction in the manifold size, in Fig. 6, we provide additional points for the low data size region. Specifically, for Mnist, we have $\{2, 4, 8, 16, 32, 64, 128, 256, 325, 400, 500\}$, and $\{100, 200\}$ for CelebA-HQ.

B EXPERIMENTAL METHODOLOGY: MEASURING THE INTRINSIC DIMENSION OF THE DATA MANIFOLD

The method used to **geometrically visualize** the intrinsic dimension of the data manifold is the same as in (Stanczuk et al., 2022): the score function is measured across several independent positions in the vicinity of the manifold and ordered as the columns of a rectangular matrix S ; the singular values of the matrix S are computed and collected; the intrinsic dimension of the manifold is given by the $d - \ker(S)$, with the kernel is estimated directly from the spectrum of the singular values of S .

On the other hand, we propose a new method to **geometrically estimate** the intrinsic dimension of the data manifold. The procedure is based on empirically computing the absolute value of the second derivative of the singular values, selecting the first bigger value with respect to the median multiplied by a threshold factor. We further discard the initial singular value as it tends to be large, resulting in instabilities. We found this method to be more robust than the one proposed in (Stanczuk et al., 2022), especially for high dimensional datasets where there is no sharp drop in the spectrum of the singular values.

B.1 COMPUTING THE SINGULAR VALUES OF THE JACOBIAN OF THE SCORE

For computing the singular values, we use the procedure described in (Stanczuk et al., 2022) reported in algorithm 1. For the linear models and MNIST models we used a symmetrized version which we empirically found to be more stable, reported in algorithm 2.

Algorithm 1 Estimate singular values at x_0

Require: s_θ - trained diffusion model (score), t_0 - sampling time, K - number of score vectors.

- 1: Sample $x_0 \sim p_0(x)$ from the data set
- 2: $d \leftarrow \dim(x_0)$
- 3: $S \leftarrow$ empty matrix
- 4: **for** $i = 1, \dots, K$ **do**
- 5: Sample $x_{t_0}^{(i)} \sim \mathcal{N}(x_{t_0} | x_0, \sigma_{t_0}^2 I)$
- 6: Append $s_\theta(x_{t_0}^{(i)}, t_0)$ as a new column to S
- 7: **end for**
- 8: $(s_i)_{i=1}^d, (v_i)_{i=1}^d, (w_i)_{i=1}^d \leftarrow \text{SVD}(S)$

Algorithm 2 Estimate singular values at x_0 with central difference

Require: s_θ - trained diffusion model (score), t_0 - sampling time, K - number of score vectors.

- 1: Sample $x_0 \sim p_0(x)$ from the data set
- 2: $d \leftarrow \dim(x_0)$
- 3: $S \leftarrow$ empty matrix
- 4: **for** $i = 1, \dots, K$ **do**
- 5: Sample $x_{t_0}^{+(i)} \sim \mathcal{N}(x_{t_0} | x_0, \sigma_{t_0}^2 I)$
- 6: Sample $x_{t_0}^{-(i)} \sim \mathcal{N}(x_{t_0} | x_0, -\sigma_{t_0}^2 I)$
- 7: Append $\frac{s_\theta(x_{t_0}^{+(i)}, t_0) - s_\theta(x_{t_0}^{-(i)}, t_0)}{2}$ as a new column to S
- 8: **end for**
- 9: $(s_i)_{i=1}^d, (v_i)_{i=1}^d, (w_i)_{i=1}^d \leftarrow \text{SVD}(S)$

B.2 COMPUTING THE INTRINSIC DIMENSION AT x_0

We report in 3 the algorithm used to find the intrinsic manifold dimension given the singular values. For MNIST we use $\bar{d} = 100$, $c = 15$ and $t = 0$, for Cifra10 $\bar{d} = 1000$, $c = 10$ and $t = 15$, and for CelebA $\bar{d} = 1000$, $c = 10$ and $t = 0$. Here t correspond to the diffusion index in DDPM. We report an example of second derivative in Fig. 7.

Algorithm 3 Estimate intrinsic manifold dimension at x_0

Require: $(s_i)_{i=1}^d$ from algorithms 1 or 2 - diffusion time, t - datapoint size, d - threshold, c - discard values, \bar{d} .

- 1: $d_{svd}^2 \leftarrow | \frac{d^2}{ds} s_t[\bar{d} :] |$
- 2: $m \leftarrow \text{median}(d_{svd}^2)$
- 3: $n \leftarrow \arg \text{where}(d_{svd}^2 > c * m)$
- 4: $k \leftarrow d - n + \bar{d}$
- 5: **return** manifold dimension k

C CONDENSATION TIME FOR POSITIONAL REM

For simplicity, we will perform the analysis for coordinate-aligned linear manifolds. Consider d -dimensional normally distributed vector-valued data \mathbf{y}^μ where each component y_k^μ follows a centered normal distribution with variance σ_k^2 . In the linear manifold case, number $d - m$ of these variances is equal to zero, meaning that the distribution spans a m -dimensional linear manifold. The number of datapoints are taken to be exponential in the size of the ambient space, i.e. $N = \exp(\alpha d)$, with $\alpha > 0$. Notice that σ_k^2 correspond to the eigenvalues of $F^\top F$ in the random projection model and we assume we have changed the coordinate system. Let us take a fixed \mathbf{x} . Hence, in the variance exploding framework we have

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

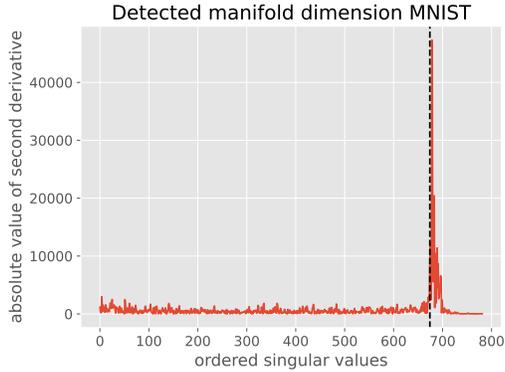


Figure 7: The figure shows the absolute values of the second derivatives computed with algorithm 3. The vertical line is the detected manifold dimension.

$$p_t(\mathbf{x}) = \frac{1}{N\sqrt{2\pi t}^d} \sum_{\mu=0}^N e^{-\frac{1}{2t}\|\mathbf{x}-\mathbf{y}^\mu\|^2} \quad (20)$$

$$= \frac{1}{N\sqrt{2\pi t}^d} e^{-\frac{\|\mathbf{x}\|^2}{2t}} \sum_{\mu=0}^N \exp\left(-\frac{1}{2t}\|\mathbf{y}^\mu\|^2 + \frac{1}{t}\mathbf{x}\mathbf{y}^\mu\right). \quad (21)$$

It is useful, at this point, to introduce the Random Energy Model (REM), firstly proposed by physicists (Derrida, 1981; Montanari and Mézard, 2009), now imported to computer science to characterize diffusion models (B. et al., 2024a). The REM consists in a collection of energy levels $\{E_\mu\}_{\mu \leq N}$ that interact with an external heat-bath at an inverse temperature β . The energy levels are random variables generated from a probability density function $p(E|\boldsymbol{\theta})$ where $\boldsymbol{\theta}$ can be some parameters of the model and source of disorder for the system. The thermodynamics of the model shows a condensation phase at a critical temperature β_c that shares similarities with glassy transitions in spin-glass models (Mezard et al., 1986). Condensation, in turn, is analogous to memorization in diffusion models. The main thermodynamic quantities, such as the condensation temperature, can be fully recovered starting from the *partition function* of the system, given by

$$Z_N(\beta) = \sum_{\mu=1}^N e^{-\beta E_\mu}. \quad (22)$$

We can now map our diffusion model into a REM by redefining

$$\beta(t) = 1/t, \quad (23)$$

and

$$E_\mu(\mathbf{x}) = \frac{1}{2}\|\mathbf{x} - \mathbf{y}^\mu\|^2. \quad (24)$$

We call this model, *positional* REM, because the occurrence of condensation will depend on a position in the d -dimensional Euclidean space. Standard REM calculations are now performed to compute the free energy of the model and then the condensation time. The moment generating function of the energies is

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

$$\zeta(\lambda) = \lim_{d \rightarrow \infty} \frac{1}{d} \log \mathbb{E}_{\mathbf{y}} e^{-\frac{\lambda}{2t} \|\mathbf{y}\|^2 + \frac{\lambda}{t} \mathbf{x} \mathbf{y}} \quad (25)$$

$$= \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=1}^d \log \int \frac{dy_i}{\sqrt{2\pi\sigma_i^2}} \exp -\frac{y_i^2}{2} \left(\frac{1}{\sigma_i^2} + \frac{\lambda}{t} \right) + \frac{\lambda}{t} x_i y_i \quad (26)$$

$$= \lim_{d \rightarrow \infty} \frac{1}{d} \left[-\frac{1}{2} \sum_{i=1}^d \log \left(1 + \lambda \frac{\sigma_i^2}{t} \right) + \frac{\lambda^2}{2t^2} \sum_{i=1}^d \frac{x_i^2 \sigma_i^2}{1 + \lambda \frac{\sigma_i^2}{t}} \right] \quad (27)$$

The derivative of the zeta function is

$$\zeta'(\lambda) = \lim_{d \rightarrow \infty} \frac{1}{d} \left[-\frac{1}{2t} \sum_i \frac{\sigma_i^2}{1 + \lambda \frac{\sigma_i^2}{t}} + \frac{\lambda}{t^2} \sum_i \frac{x_i^2 \sigma_i^2}{1 + \lambda \frac{\sigma_i^2}{t}} - \frac{\lambda^2}{2t^3} \sum_i \frac{x_i^2 \sigma_i^4}{(1 + \lambda \frac{\sigma_i^2}{t})^2} \right]. \quad (28)$$

At large times, $\zeta(\lambda)$ and $\zeta'(\lambda)$ become respectively

$$\zeta(\lambda) = -\frac{\lambda}{2t} r_{2,\sigma} + \frac{\lambda^2}{4t^2} r_{4,\sigma} + \frac{\lambda^2}{2t^2} \omega^2(\mathbf{x}), \quad (29)$$

$$\zeta'(\lambda) = -\frac{\lambda}{2t} r_{2,\sigma} + \frac{\lambda^2}{2t^2} r_{4,\sigma} + \frac{\lambda^2}{t^2} \omega^2(\mathbf{x}). \quad (30)$$

Where

$$r_{2,\sigma} = \lim_{d \rightarrow \infty} \frac{1}{d} \sum_i \sigma_i^2 \quad (31)$$

$$r_{4,\sigma} = \lim_{d \rightarrow \infty} \frac{1}{d} \sum_i (\sigma_i^2)^2 \quad (32)$$

$$\omega^2(\mathbf{x}) = \lim_{d \rightarrow \infty} \frac{1}{d} \sum_i (x_i)^2 \sigma_i^2. \quad (33)$$

The condition for the condensation time is $\alpha + \zeta(1) - \zeta'(1) = 0$, from which we obtain

$$t_c(\mathbf{x}) = \sqrt{\frac{r_{4,\sigma} + \omega^2(\mathbf{x})}{2\alpha}}. \quad (34)$$

As clear from the formula, this time depends on the variance $\omega^2(\mathbf{x})$ along the direction of \mathbf{x} . This implies that, when \mathbf{x} is aligned to a linear sub-manifold with higher variance, condensation around this state will happen earlier, leading to a decrease in the estimated commonality of the latent manifold. Fig. 8 shows a comparison between the exact approach for computing $t_c(x)$ (i.e. using Eqs. (25), (28)) and the small α expansion (i.e. Eq. (15)), showing a good qualitative agreement between the two quantities at all values of α . The right panel of the same figure also displays a strong dependence of the exactly computed condensation time.

If each dimension has equal variance σ^2 , the directional variance density is just σ^2 , which implies that the critical condensation time depends linearly on the dimensionality but only logarithmically on the number of data points. This implies that in isotropic case, in order to avoid condensation an exponential number of data points is needed. However, if only α_m dimensions have non-zero variance, it is straightforward to see that the exponential dependency will scale with α_m instead of m . More generally, the exponential scaling depends on the total variance $m \omega(\mathbf{x})$, which implies that it is realistic to learn high-dimensional spaces as far as most of these dimensions have vanishing variance.

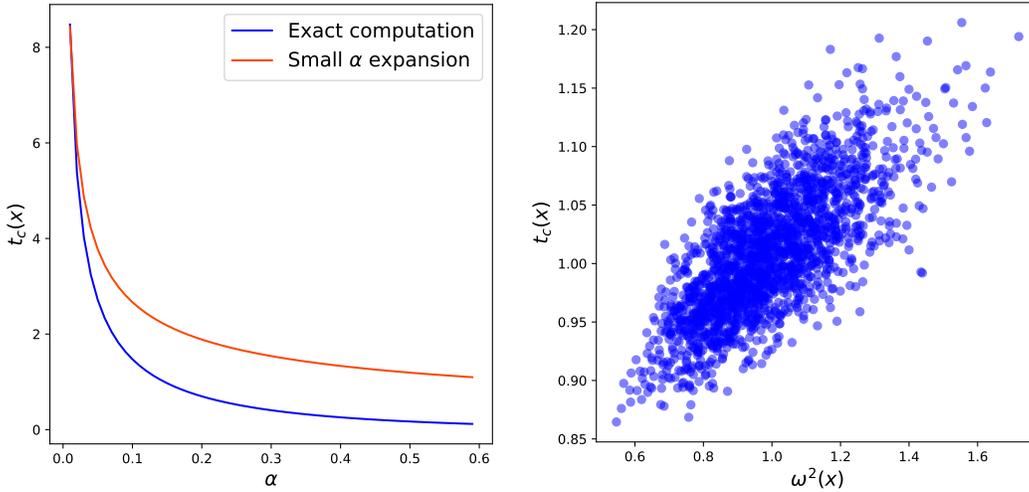


Figure 8: Condensation time as a function of position x computed according to the REM calculations. Left: we have generated one single position \mathbf{x} in a ambient space of dimension $d = 100$ and one single matrix F of dimensions 100×50 (with $m = 50$ dimension of the latent space). Both \mathbf{x} and F are generated according to a Gaussian process with zero mean and unitary variance; we show the comparison between the exact calculation of the positional condensation time and the approximated version that is fully explicit in the directional variance $\omega^2(\mathbf{x})$. Right: we generate 2000 random positions \mathbf{x} around the origin of the ambient space of dimension $d = 100$; the latent space dimension is $m = 50$ and $\alpha = 0.15$; we show the dependence of the exact positional condensation time as a function of $\omega^2(\mathbf{x})$, showing a qualitatively similar behaviour with respect to the approximated expression of t_c .

D ANALYSIS OF THE EMPIRICAL JACOBIAN

We can relate this random energy analysis to the spectra of Jacobian eigenvalues using a heuristic argument. In the linear manifold example, the Jacobian of the true score function at $t = 0$ is diagonal with eigenvalues equal to $-1/\sigma_k^2$. This results in spectral gaps when different sub-spaces have different variances. For a finite value of the inverse temperature $\beta(t) = 1/t$, the eigenvalues are $-1/\sigma_k^2 - \beta$. After the critical condensation time, the empirical score gives a good approximation of the true score. On the other hand, in the condensation phase the empirical score is dominated by the (quenched) fluctuations in the data distribution. First, we can introduce the participation ratio

$$Y(\beta, \mathbf{x}) = \frac{Z(2\beta, \mathbf{x})}{Z(\beta, \mathbf{x})^2}. \quad (35)$$

This thermodynamic quantity can be roughly interpreted as the inverse of the number of energy levels with non-vanishing weights. In the condensation phase, this will be a finite number while it becomes infinite in the high temperature phase.

In the thermodynamic limit and for $\beta(t) \geq \beta_c(t)$, the participation ratio of our REM model is given by

$$\mathbb{E}[Y(\beta, \mathbf{x})] = 1 - \frac{\beta_c(t, \mathbf{x})}{\beta(t)}, \quad (36)$$

which implies that the number of datapoints that contribute to the score function at \mathbf{x} is

$$\tilde{N} = e^{\alpha \tilde{d}(\beta, \mathbf{x})} = 1/Y(\beta, \mathbf{x}) = \frac{\beta(t)}{\beta(t) - \beta_c(t, \mathbf{x})}. \quad (37)$$

Note that this number tends to one for $\beta(t) \rightarrow \infty$, meaning that in the low-time limit the score depends on a single pattern.

In this phase, the score is dominated by approximately $e^{\alpha \bar{d}(\beta, \mathbf{x})} = \frac{\beta(t)}{\beta(t) - \beta_c(t, \mathbf{x})}$, leading to the expression

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \approx \frac{\beta}{e^{\alpha \bar{d}(\beta, \mathbf{x})}} \sum_{\mu=1}^{e^{\alpha \bar{d}(\beta, \mathbf{x})}} (\mathbf{y}^\mu - \mathbf{x}) \quad (38)$$

where $\mathbf{y}^\mu \sim p(\mathbf{y}^\mu | \mathbf{x}, \beta) \propto e^{-\mathbf{y}^T (\Lambda^{-1} + \beta I_d) \mathbf{y} / 2 + \beta \mathbf{x} \cdot \mathbf{y}}$. Therefore, the empirical score approximately follows the distribution

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \sim \mathcal{N}(-M(\beta)\mathbf{x}, \beta(\Lambda^{-1} + \beta I_d)^{-1} \max(0, \beta - \beta_c(\mathbf{x}))) . \quad (39)$$

where $M(\beta) = \beta(\Lambda^{-1} + \beta I_d)^{-1} \Lambda^{-1}$ and Λ being the diagonal matrix collecting the variances σ_k^2 and we used the fact that $e^{\alpha \bar{d}(\beta, \mathbf{x})} = \beta / (\beta - \beta_c(\mathbf{x}))$. The minimum in the formula is due to the fact that, for $\beta < \beta_c$, an exponentially large number of patterns participate in the estimation of the score, which leads to a complete suppression of the variance. On the other hand, the variance of the empirical score estimator diverges for $\beta \rightarrow \infty$. In fact, during condensation, the fluctuations in the random sampling of the datapoints are not suppressed due to the small number of non-vanishing weights.

We can finally estimate the distribution of the eigenvalues estimated from the empirical Jacobian matrix. Let us set ourselves on $\mathbf{x} = \mathbf{0}$ and perturb along the directions of the eigenvectors of $F^T F$. We estimate the elements of the Jacobian of the score function with respect to the orthogonal direction \mathbf{e}_j using a perturbative approach, i.e.

$$J_{ij}(\beta) \approx \sqrt{\beta} \left(\partial_{x_i} \log p_t(\mathbf{e}_j / \sqrt{\beta}) - \partial_{x_i} \log p_t(\mathbf{0}) \right) . \quad (40)$$

Using Eq. (39), we can then write an approximate distribution for the elements of the Jacobian as

$$J_{ij}(\beta) \sim \mathcal{N} \left(-\beta \delta_{ij} (1 + \beta \sigma_i^2)^{-1}, \beta^2 (\sigma_i^{-2} + \beta)^{-1} \left[\phi(\beta, \mathbf{0}) + \phi(\beta, \mathbf{e}_j / \sqrt{\beta}) \right] \right) . \quad (41)$$

where we assumed that the fluctuations in $\nabla_{\mathbf{x}} \log p_t(\mathbf{e}_j / \sqrt{\beta})$ are independent from the fluctuations in $\nabla_{\mathbf{x}} \log p_t(\mathbf{0})$ and $\nabla_{\mathbf{x}} \log p_t(\mathbf{e}_k / \sqrt{\beta})$ for all k s. In this expression, we introduced the function

$$\phi(\beta, \mathbf{x}) = \max(0, \beta - \beta_c(\mathbf{x})) . \quad (42)$$

We can now recover the singular values of $J(\beta)$ as minus the square roots of the eigenvalues of $J(\beta)^T J(\beta)$. In general, the matrix $J(\beta)^T J(\beta)$ can have a complex spectral distribution. An approximate formula for the singular values of $J(\beta)$ is

$$s_i \approx -\sqrt{\beta^2 (1 + \beta \sigma_i^2)^{-2} + \beta^4 \sum_{k=1}^d (\sigma_k^{-2} + \beta)^{-2} \left[\phi(\beta, \mathbf{0}) + \phi(\beta, \mathbf{e}_i / \sqrt{\beta}) \right]^2} . \quad (43)$$

To obtain this formula, we write J as

$$J = A + B \quad (44)$$

where A is a diagonal matrix corresponding to the mean of Eq. (41), while B corresponds to the variance. Therefore, $J^T J$ becomes

$$J^T J = A^T A + A^T B + B^T A + B^T B . \quad (45)$$

This expression is dominated by the two symmetric terms, so we can write

$$J^T J \approx A^T A + B^T B . \quad (46)$$

Then, the term $A^T A = A^2$ is, of course, still diagonal, while the term $B^T B$ is diagonally dominant. Calling $C = \sum_{ik} B_{ik} B_{ik}$, we can approximate the singular values as $\sqrt{A^2 + C^2}$, obtaining Eq. (43). Note however that the distribution of the spectrum does not concentrate exactly to Eq. 43 in the large N . Nevertheless, Eq. 43 gives an accurate picture of the qualitative behavior, as shown in Sec. D.1.

These results also show that in some regimes Eq. 43 is more in agreement with the numerical empirical score than the correct spectrum of Eq. 41, which is likely due to the fact that Eq. 41 overestimates the fluctuations by ignoring the correlations of the score at different points.

D.1 NUMERICAL ANALYSIS

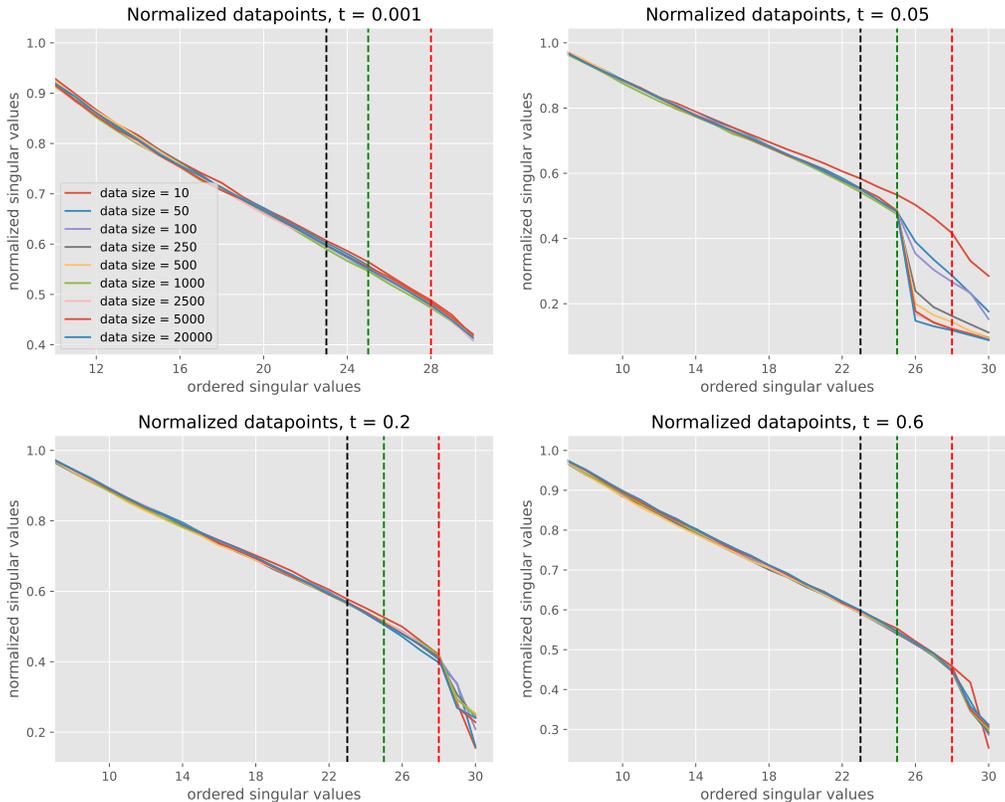


Figure 9: Ordered singular values of the Jacobian of the empirical score in the case of the linear data model. The parameters for the model are $d = 30$, $m = 7$, $\alpha = \log(N)/d = 0.23$ with a subspace associated to a variance $\sigma_1^2 = 1$ of dimension $m_1 = 2$ and another subspace with variance $\sigma_2^2 = 0.3$ and dimension $m_2 = 5$. Different lines are associated to different sizes of the training set. Measures have been averaged over 30 realizations of the experiment.

In order to test our theory of the empirical score, we plot the singular values of the Jacobian in an ordered fashion, as done experimentally in the previous sections: this allows to visualize the drops forming due to the described memorization phase transition. As a first test, a set of N data have been generated according to the linear manifold model with two variances, and the empirical score has been computed out of these points as in Eq. (10). Therefore, we measured the Jacobian according to same method used for the trained models, described in section B. The condensation time appearing in the formulas has been computed according to the method explained in Supp. C. Figs. 9 and 10 report the profiles of the ordered spectra at different times when data-sets of different sizes are employed to assemble the empirical score. We notice the same phenomenology of the gaps predicted by the theory: the gaps indicating the dimensions of both the two subspaces progressively open starting from the largest variance and ending to the smallest one.

As a second experiment we confront the evolution of the ordered singular values in time, as obtained from three methods: the functional form in Eq. (19); by extraction from the random matrix expressed in Eq. (18); by computing the empirical score function from a synthetic set of N datapoints (i.e. as performed in the previous experiment). The time evolution of the gaps is reported in figures 11 and 12 for two choices of the parameters of the model. We conclude that both the random matrix prediction and the simulation capture the phenomenology predicted by the analytical expression of the singular values obtained in Eq. (19).

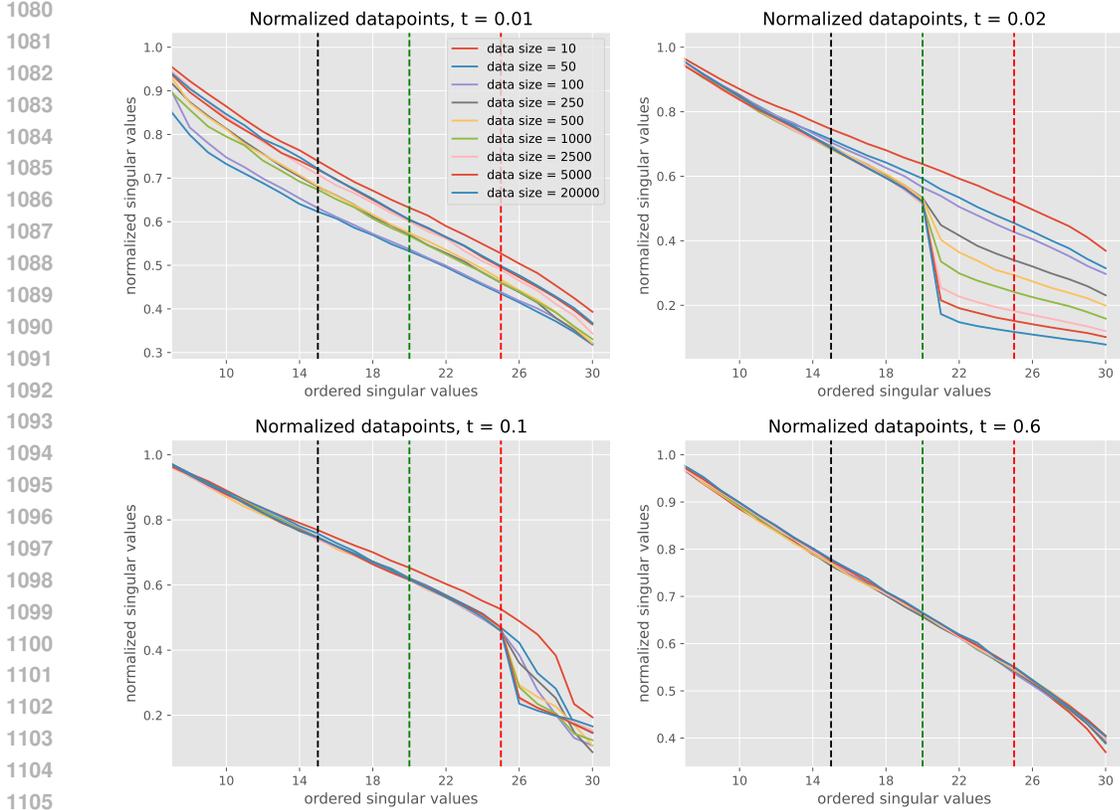


Figure 10: Ordered singular values of the Jacobian of the empirical score in the case of the linear data model. $d = 30$, $m = 15$, $\alpha = \log(N)/d = 0.23$ with a subspace associated to a variance $\sigma_1^2 = 1$ of dimension $m_1 = 5$ and another subspace with variance $\sigma_2^2 = 0.3$ and dimension $m_2 = 10$. Different lines are associated to different sizes of the training set. Measures have been averaged over 30 realizations of the experiment.

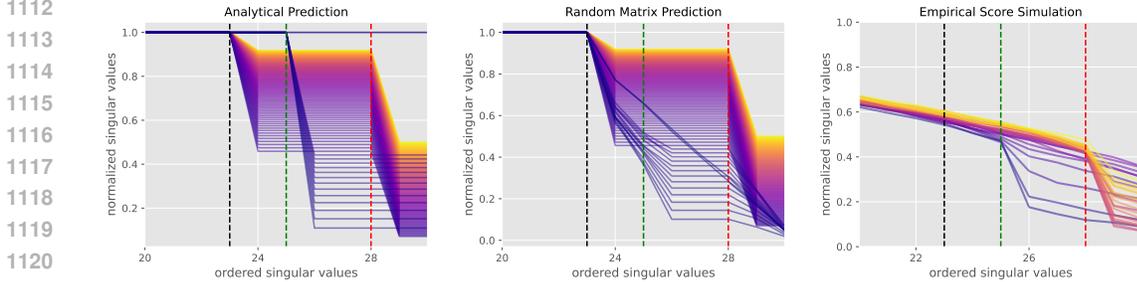
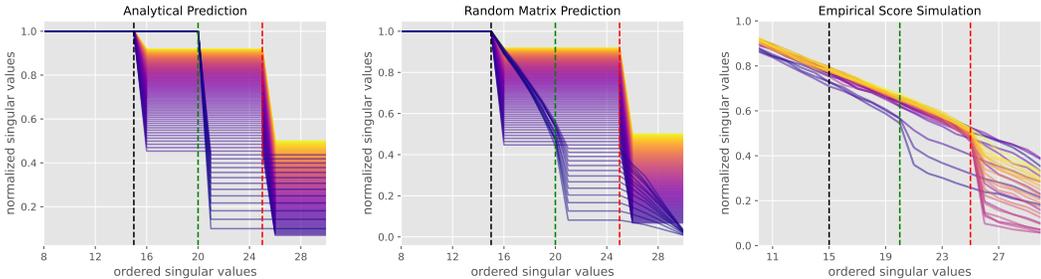


Figure 11: The ordered singular values of the Jacobian of the empirical score function of a linear manifold model as a function of the diffusion time t . Lighter colours are associated to larger times in the colour map. The parameters for the model are $d = 30$, $m = 7$, $\alpha = \log(N)/d = 0.23$ with a subspace associated to a variance $\sigma_1^2 = 1$ of dimension $m_1 = 2$ and another subspace with variance $\sigma_2^2 = 0.3$ and dimension $m_2 = 5$. Left: approximated theoretical prediction in the memorization phase according to Eq. (19). Center: prediction from the approximated Jacobian in Eq. (18). Right: singular values obtained by the numerical measure of the Jacobian of the empirical score function (as described in section B), evaluated from a synthetic data set of $N = 10^3$ points.

Finally, figures 13 and 14 report the evolution of the gaps according to, respectively, the closed formula for the singular values in Eq. (19) and the approximated random matrix in Eq. (10), when we change the size of the data-set.

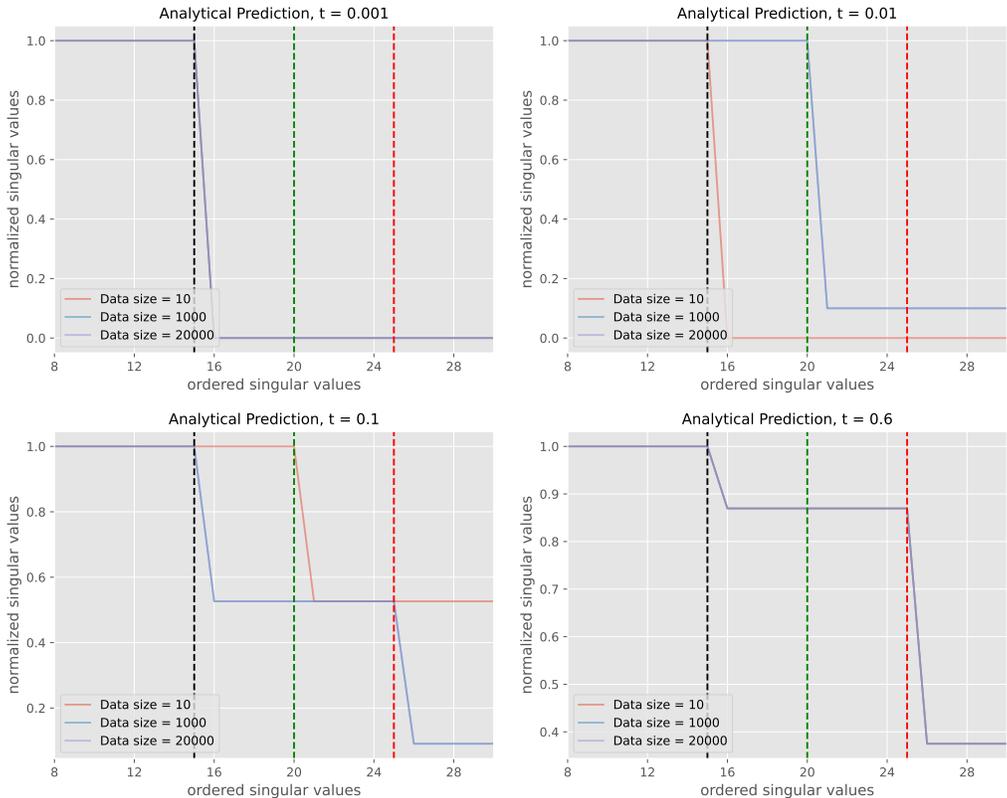
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144



1145
1146
1147
1148
1149
1150
1151
1152
1153

Figure 12: The ordered singular values of the Jacobian of the empirical score function of a linear manifold model as a function of the diffusion time t . Lighter colours are associated to larger times in the colour map. The parameters for the model are $d = 30$, $m = 15$, $\alpha = \log(N)/d = 0.23$ with a subspace associated to a variance $\sigma_1^2 = 1$ of dimension $m_1 = 5$ and another subspace with variance $\sigma_2^2 = 0.3$ and dimension $m_2 = 10$. Left: approximated theoretical prediction in the memorization phase according to Eq. (19). Center: prediction from the approximated Jacobian in Eq. (18). Right: singular values obtained by the numerical measure of the Jacobian of the empirical score function (as described in section B), evaluated from a synthetic data set of $N = 10^3$ points.

1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181



1182
1183
1184
1185
1186
1187

Figure 13: Ordered singular values of the Jacobian of the empirical score in the case of the linear data model estimated from Eq. (19). The parameters for the model are $d = 30$, $m = 7$, $\alpha = \log(N)/d = 0.23$ with a subspace associated to a variance $\sigma_1^2 = 1$ of dimension $m_1 = 2$ and another subspace with variance $\sigma_2^2 = 0.3$ and dimension $m_2 = 5$. Different lines are associated to different sizes of the training set.

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

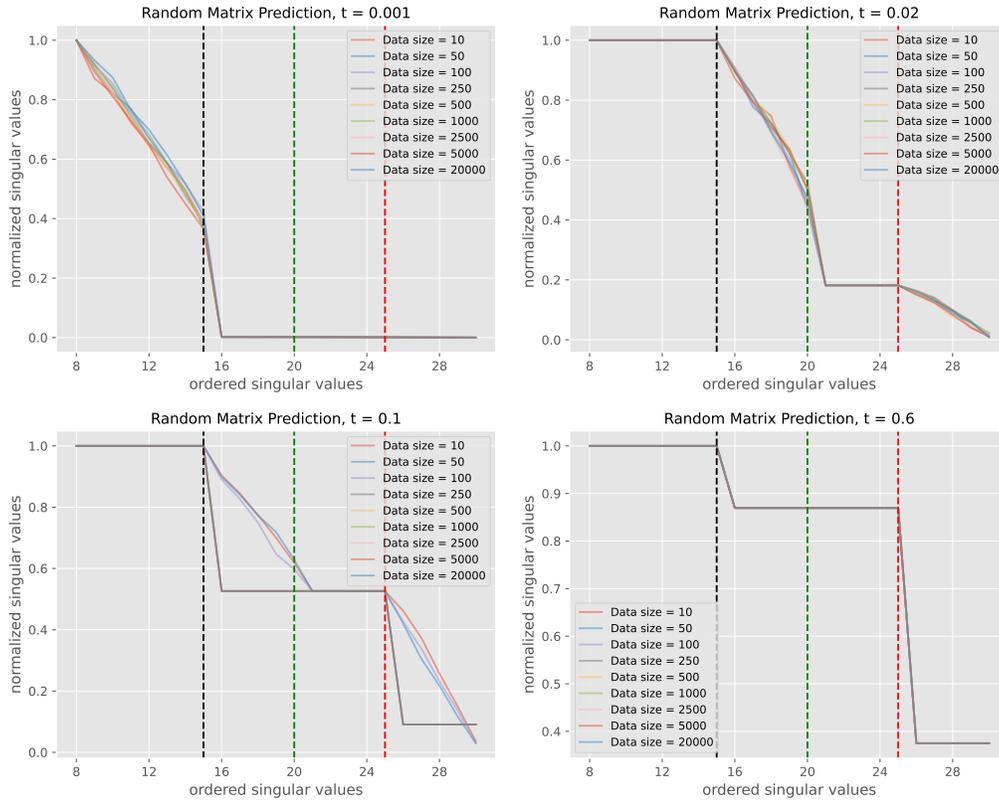


Figure 14: Ordered singular values of the Jacobian of the empirical score in the case of the linear data model, estimated from the random matrix in Eq. (10). $d = 30$, $m = 15$, $\alpha = \log(N)/d = 0.23$ with a subspace associated to a variance $\sigma_1^2 = 1$ of dimension $m_1 = 5$ and another subspace with variance $\sigma_2^2 = 0.3$ and dimension $m_2 = 10$. Different lines are associated to different sizes of the training set.

E ADDITIONAL EXPERIMENTAL RESULTS ON TRAINED NETWORKS

In addition to the experiments described in section 6 we report here some tests on synthetic data generated on the linear manifold introduced in section 5.

First, we report numerical results from experiments in an ambient space of dimension $d = 100$, while the manifold lives in a space of dimension $m = 40$. On the same manifold, through the choice of a diagonal F matrix, we define subspaces of different variances, which will result in the opening of gaps at different times in the spectrum of the singular values of the matrix obtained by sampling the score functions. In the specific, we choose the case of two subspaces associated to two variances of the data and the particular scenario of m different variances sampled uniformly at random. In Fig. 15 we plot the spectra at the smallest diffusion time ϵ , for models trained on datasets with different amount of training samples. With few training samples, we cannot see any gap opening. However, as the training samples increase, the model starts generalizing to the sub-spaces with higher variance, indicating both a smooth transition between generalization and memorization, and that the subspaces with higher variance are learned first by the model. As we shall see, as predicted by the theory the network will instead generalize to subspaces of low variance for parameter settings where σ_2^2 is lower than σ_1^2 but not negligible. Furthermore, in Fig. 16 report a similar experiment with a smaller ambient dimension, i.e. $d = 30$. Now the geometric memorization phenomenon is more evident for medium data-set sizes at small times. Moreover, the phenomenology emerging from the trained model is fully consistent with the one resulting from the empirical score, obtained through the same choice of the parameters, as showed in figures 12, 13 and 14. This conclusion suggests two powerful insights about diffusion models: the network behaviour is consistent with our theory of memorization derived from the physics of Random Energy Models; the trained score function behaves consistently with the empirical score for a certain choice of the parameters.

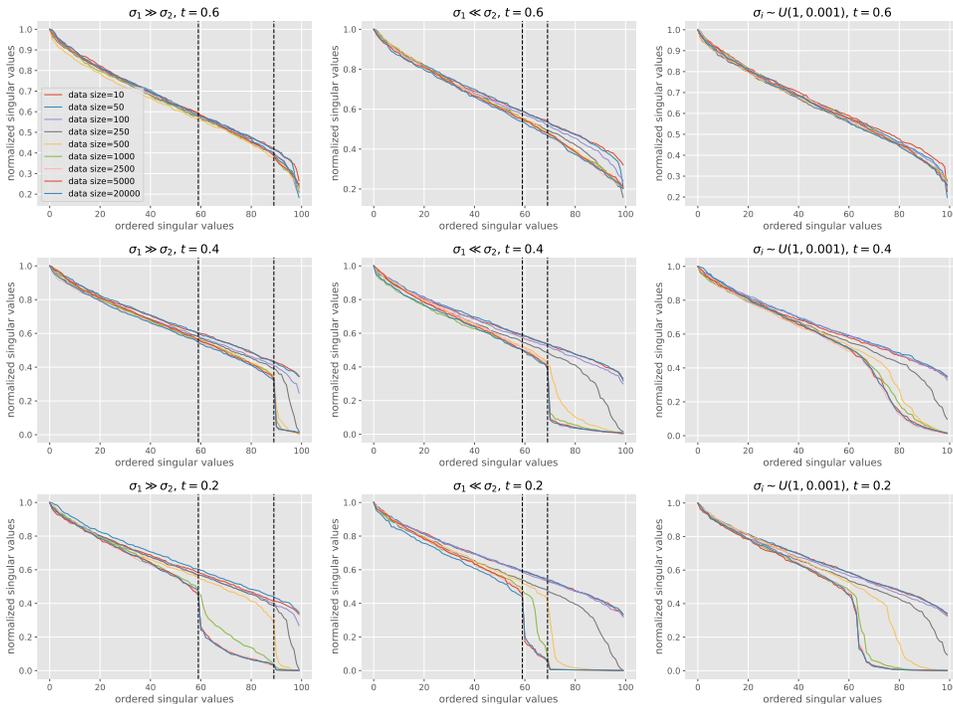


Figure 15: Singular values for models trained on different number of samples from the dataset (in the legend) at $t = 0.6$, $t = 0.4$ and $t = 0.2$, from top to bottom respectively. From left to right: model with $\sigma_1^2 = 1$, $\sigma_2^2 = 0.01$; model with $\sigma_1^2 = 0.01$, $\sigma_2^2 = 1$; model with uniformly sampled variances.

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

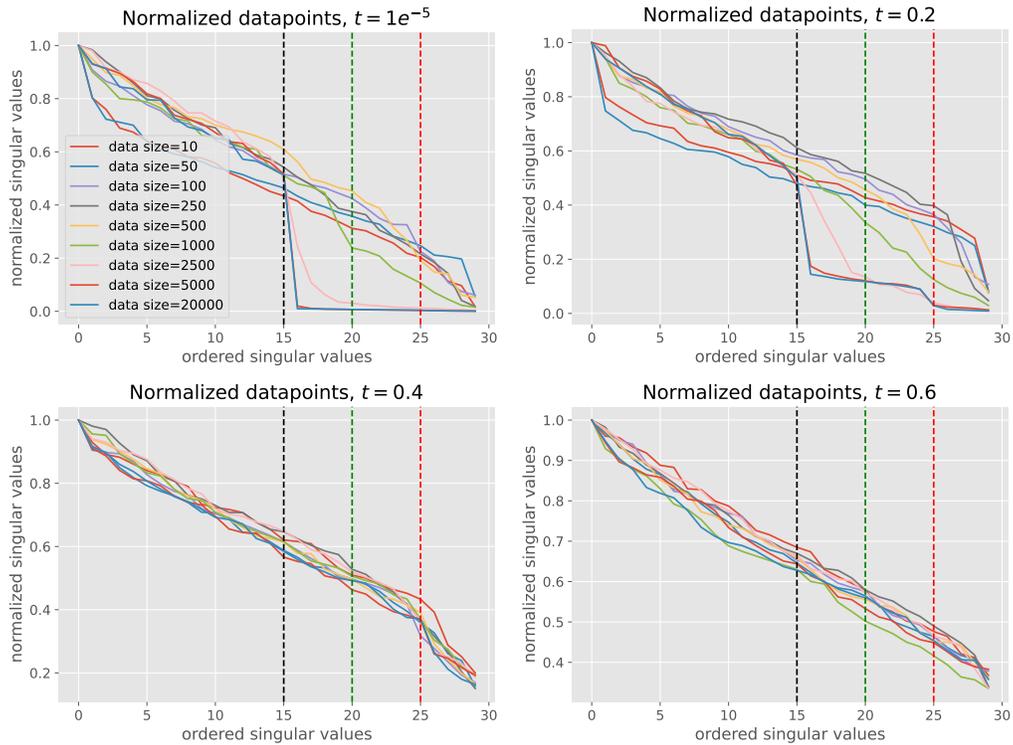


Figure 16: Ordered singular values of the Jacobian of the trained score function in the case of the linear data model. $d = 30$, $m = 15$, $\alpha = \log(N)/d = 0.23$ with a subspace associated to a variance $\sigma_1^2 = 1$ of dimension $m_1 = 5$ and another subspace with variance $\sigma_2^2 = 0.3$ and dimension $m_2 = 10$. Different lines are associated to different sizes of the training set.