

Chapter 3

Pursuing Low-Dimensional Distributions via Denoising

“Information is the resolution of uncertainty.”

—Claude Shannon, 1948

In Chapter 2, we showed how to learn simple classes of distributions whose supports are assumed to be either a single low-dimensional subspace, a mixture of such subspaces, or low-rank Gaussians. For simplicity, the different (hidden) linear or Gaussian modes are taken to be orthogonal or independent,¹ as illustrated in Figure 2.1. For these special distributions, simple and effective learning algorithms with correctness and efficiency guarantees can be derived, and the geometric and statistical interpretations of the associated algorithmic operations are clear.

In practice, both linearity and independence are idealistic assumptions that distributions of real-world high-dimensional data *rarely* satisfy. The only assumption we can make is that the intrinsic dimension of the distribution is very low compared with the dimension of the ambient space in which the data are embedded. Hence, in this chapter we show how to learn a more general class of low-dimensional distributions in a high-dimensional space that is *not* necessarily (piecewise) linear.

Real-world distributions typically contain multiple components or modes, for example, corresponding to different object classes in images. These modes may not be statistically independent and can even have different intrinsic dimensions. Therefore, we generally assume that the data are distributed on a mixture of low-dimensional nonlinear submanifolds embedded in a high-dimensional space. Figure 3.1 illustrates such a distribution. Another constraint is that in real-world situations we typically have access only to a finite number of samples

¹Or can be reduced to such idealistic cases under benign conditions.

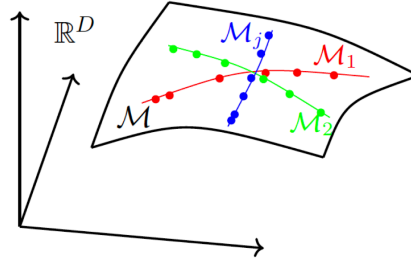


Figure 3.1: Data \mathbf{X} are distributed on a mixture of low-dimensional submanifolds $\cup_j \mathcal{M}_j$ in a very high-dimensional ambient space, say \mathbb{R}^D . This is a much more general case than the special cases studied in Chapter 2, where the submanifolds were assumed to be piecewise linear and thus admitted parametric forms and analytical solutions.

from the data distribution, rather than prior knowledge of the exact distribution itself.

To learn such a distribution under these conditions, we must address three fundamental questions:

- What is a general approach to learning and representing a low-dimensional distribution in a high-dimensional space?
- How do we measure the complexity of the distribution so that we can exploit its low dimensionality effectively for learning?
- How do we make the learning process computationally tractable and scalable, especially when the ambient dimension is high and the sample size is large?

As we will see, the fundamental idea of *compression*, or *dimension reduction*, which we have already shown to be effective in the linear/independent case, remains a guiding principle for developing computational models and methods for general nonlinear low-dimensional distributions.

Due to its theoretical and practical importance, we examine in greater depth how this compression-based framework materializes when the distribution of interest can be well modeled or approximated by a mixture of low-dimensional subspaces or low-rank Gaussians. As later chapters show, these seemingly restrictive model classes serve as basic building blocks for algorithms that perform compression on general distributions.

3.1 Entropy Minimization and Compression

3.1.1 Entropy and Coding Rate

In Chapter 1 we mentioned that the goal of learning is to find the simplest way to generate a given set of data. Conceptually, Kolmogorov complexity was intended to provide such a measure of complexity, but it is not computable and is not associated with any implementable scheme that can actually reproduce the data. Hence, we need an alternative, computable, and realizable measure of complexity. This leads us to the notion of *entropy*, introduced by Shannon in 1948 [Sha48].

To illustrate the constructive nature of entropy, let us start with the simplest case. Suppose we have a discrete random variable that takes N distinct values, or *tokens*, $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ with equal probability $1/N$. We can encode each token \mathbf{x}_i using the $\log_2 N$ -bit binary representation of i . This coding scheme can be generalized to encoding arbitrary discrete distributions [CT91]: given a distribution p such that $\sum_{i=1}^N p(\mathbf{x}_i) = 1$, one can assign each token \mathbf{x}_i with probability $p(\mathbf{x}_i)$ a binary code of length $\log_2[1/p(\mathbf{x}_i)] = -\log_2 p(\mathbf{x}_i)$ bits. Hence the average number of bits, or *the coding rate*, needed to encode any sample from the distribution $p(\cdot)$ is given by the expression²:

$$H(p) \doteq \mathbb{E}_{\mathbf{x} \sim p}[-\log p(\mathbf{x})] = -\sum_{i=1}^N p(\mathbf{x}_i) \log p(\mathbf{x}_i). \quad (3.1.1)$$

This is known as the *entropy* of the (discrete) distribution $p(\cdot)$. Note that this entropy is always nonnegative and it is zero if and only if $p(\mathbf{x}_i) = 1$ for some \mathbf{x}_i with $i \in [N]$.³ For a random variable \mathbf{x} , we will often write $H(\mathbf{x})$ to mean the entropy of the (marginal) distribution of \mathbf{x} .

3.1.2 Differential Entropy

When the random variable $\mathbf{x} \in \mathbb{R}^D$ is continuous and has a probability density p , one may view the limit of the above sum (3.1.1) as related to an integral:

$$h(p) \doteq \mathbb{E}_{\mathbf{x} \sim p}[-\log p(\mathbf{x})] = -\int_{\mathbb{R}^D} p(\boldsymbol{\xi}) \log p(\boldsymbol{\xi}) d\boldsymbol{\xi}. \quad (3.1.2)$$

More precisely, given a continuous probability distribution p on \mathbb{R}^D , we may discretize \mathbb{R}^D into a union of disjoint cubes C_i^ϵ of side length $\epsilon^{1/D}$, each (say) centered at \mathbf{x}_i , and define the discretized distribution p^ϵ by $p^\epsilon(\mathbf{x}_i) = \mathbb{P}_{\mathbf{x} \sim p}[\mathbf{x} \in C_i^\epsilon]$. Then one can show that $H(p^\epsilon) + \log \epsilon \approx h(p)$. By taking ϵ to be small, we observe that the differential entropy $h(p)$ can be negative. Interested readers may refer to [CT91] for a more detailed explanation. Similar to the discrete entropy,

²By the convention of information theory [CT91], the log here has base 2. Hence, entropy is measured in (binary) bits.

³Here we use the fact $\lim_{p \rightarrow 0} p \log p = 0$.

we will often write $h(\mathbf{x})$ to mean the entropy of the (marginal) distribution of \mathbf{x} .

Example 3.1 (Entropy of Gaussian distributions). Through direct calculation, one can show that the entropy of a Gaussian random variable $x \sim \mathcal{N}(\mu, \sigma^2)$ is

$$h(x) = \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2}. \quad (3.1.3)$$

The entropy of a multivariate Gaussian random variable $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ in \mathbb{R}^D is

$$h(\mathbf{x}) = \frac{D}{2} (1 + \log(2\pi)) + \frac{1}{2} \log \det(\boldsymbol{\Sigma}). \quad (3.1.4)$$

■

Similar to the entropy of a discrete distribution, we would like the differential entropy to be associated with the coding rate of some realizable coding scheme. For example, as mentioned above, if we discretize the domain of the distribution with a grid of size $\epsilon > 0$ with centers \mathbf{x}_i , one possible coding scheme is to map each \mathbf{x} to its nearest center \mathbf{x}_i , and the entropy of the resulting discrete distribution approximates the differential entropy [CT91].

There are some caveats associated with the definition of differential entropy. For a distribution in a high-dimensional space, when its support becomes degenerate (low-dimensional), its differential entropy diverges to $-\infty$. This fact is proved in Theorem B.1, but even in the simple explicit case of Gaussian distributions (3.1.4), when the covariance $\boldsymbol{\Sigma}$ is singular, we can see that $\log \det(\boldsymbol{\Sigma}) = -\infty$ so we have $h(\mathbf{x}) = -\infty$ (note that by Theorem B.1, Gaussian distributions have the largest entropy among all distributions with the same mean and covariance). In such a situation, it is not obvious how to properly quantize or encode such a distribution. Nevertheless, degenerate (Gaussian) distributions are precisely the simplest possible, and arguably the most important, instances of low-dimensional distributions in a high-dimensional space. In the next chapter, we will discuss a complete resolution to this seeming difficulty with degeneracy.

3.1.3 Minimizing Coding Rate

The learning problem entails recovering a (potentially continuous) distribution $p(\mathbf{x})$ from samples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ drawn from it. For ease of exposition, we write $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$. Because the distributions of interest are (nearly) low-dimensional, we expect their (differential) entropy to be very small. However, unlike the situations studied in the previous chapter, we do not know the family of (analytical) low-dimensional models to which $p(\mathbf{x})$ belongs. Thus, checking whether the entropy is small seems to be the only guideline we can rely on to identify and model the distribution.

Now, given the samples alone without knowing what $p(\mathbf{x})$ is, in theory they could be interpreted as samples from any generic distribution. In particular, they could be interpreted as any of the following:

1. as samples from the empirical distribution $p^{\mathbf{X}}$ itself, which assigns probability $1/N$ to each of the N samples $\mathbf{x}_i, i \in [N]$;
2. as samples from a standard normal distribution $\mathbf{x}^n \sim p^n \doteq \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ with variance σ^2 large enough (say larger than the squared sample norms);
3. as samples from a normal distribution $\mathbf{x}^e \sim p^e \doteq \mathcal{N}(\mathbf{0}, \hat{\Sigma})$ with covariance $\hat{\Sigma} = \frac{1}{N} \mathbf{X} \mathbf{X}^\top$ being the empirical covariance of the samples;
4. as samples from a distribution $\hat{\mathbf{x}} \sim \hat{q}(\mathbf{x})$ that closely approximates the ground-truth distribution p .

The question is: which interpretation is better, and in what sense? Suppose that you believe these data \mathbf{X} are drawn from a particular distribution $q(\mathbf{x})$, which may be one of the above. Then we could encode the data points with the optimal codebook for the distribution $q(\mathbf{x})$. Then, as the number of samples N becomes large, the required average coding length (or coding rate) is given by the so-called *cross-entropy*:

$$\text{CE}(p \parallel q) \doteq \mathbb{E}_{\mathbf{x} \sim p}[-\log q(\mathbf{x})]. \quad (3.1.5)$$

If we have identified the correct distribution $p(\mathbf{x})$, the coding rate is given by the entropy $h(p) = \mathbb{E}_{\mathbf{x} \sim p}[-\log p(\mathbf{x})]$. It turns out that the above coding length $\text{CE}(p \parallel q)$ is always greater than or equal to the entropy $h(\mathbf{x})$ unless $p = q$.

Hence, given a set of sampled data \mathbf{X} , to determine which distribution best characterizes the data among p^n, p^e , and \hat{q} , we may compare their coding rates for \mathbf{X} and see which yields the lowest rate. We know from the above that the (theoretically achievable) coding rate for a distribution is closely related to its entropy. In general,

$$h(p^n) > h(p^e) > h(\hat{q}). \quad (3.1.6)$$

If the data \mathbf{X} were encoded by the codebook associated with each of these distributions, the coding rate for \mathbf{X} would therefore decrease in the same order:

$$\text{CE}(p \parallel p^n) > \text{CE}(p \parallel p^e) > \text{CE}(p \parallel \hat{q}). \quad (3.1.7)$$

This observation gives a general guideline for pursuing a distribution $p(\mathbf{x})$ that has low-dimensional structure. It suggests two possible approaches:

- A1. Starting with a general distribution (say a normal distribution) with high entropy, gradually transform the distribution toward the (empirical) distribution of the data by reducing entropy;
- A2. Among a large family of (parametric or non-parametric) distributions with explicit coding schemes that encode the given data, progressively search for better coding schemes that give lower coding rates.

Conceptually, both approaches A1 and A2 try to do the same thing. For A1, we need to ensure such a path of transformation exists and is computable. For A2, it is necessary that the chosen family is rich enough and can closely

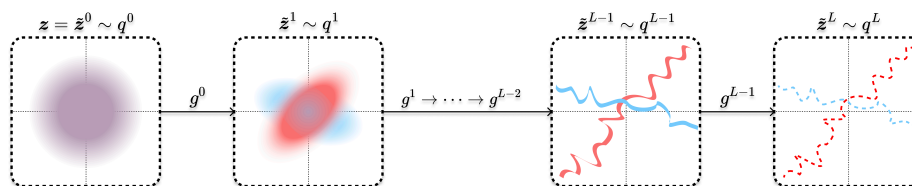


Figure 3.2: Illustration of an iterative denoising process that, starting from an isotropic Gaussian distribution, converges to an arbitrary data distribution.

approximate (or contain) the ground-truth distribution. For either approach, we need to ensure that solutions with lower entropy or better coding rates can be efficiently computed and converge to the desired distribution quickly.⁴ We will explore both approaches in the remaining sections of this chapter.

3.2 Compression via Denoising

In this section we describe a *natural* and *computationally tractable* way to learn a distribution $p_{\mathbf{x}}$ by learning a parametric encoding that minimizes the entropy or coding rate of the representation, then using this encoding to transform high-entropy samples from a standard Gaussian into low-entropy samples from the target distribution, as illustrated in Figure 3.2. This methodology combines approaches A1 and A2 above to learn and sample from the distribution.

3.2.1 Diffusion and Denoising Processes

We first want a procedure that reduces the entropy of a very noisy random variable, producing a lower-entropy sample from the data distribution. Here we describe one potential approach, perhaps the most natural. First, we devise a way to *gradually increase* the entropy of samples already drawn from the data distribution. Next, we seek an *approximate inverse* of this process. In general, increasing entropy has no inverse, because information from the original distribution may be destroyed. We therefore restrict ourselves to a special case in which (1) the entropy-increasing operation is simple, computable, and reversible, and (2) we have a (parametric) encoding of the data distribution, as alluded to above. As we will see, these two conditions make our approach feasible.

Diffusion Processes

We will increase the entropy in arguably the simplest possible way, i.e., *adding isotropic Gaussian noise*. More precisely, given the random variable \mathbf{x} , we can consider the *stochastic process* $(\mathbf{x}_t)_{t \in [0, T]}$ that adds gradual noise to it, i.e.,

$$\mathbf{x}_t \doteq \mathbf{x} + t\mathbf{g}, \quad \forall t \in [0, T], \quad (3.2.1)$$

⁴Say the distribution of real-world data such as images and texts.

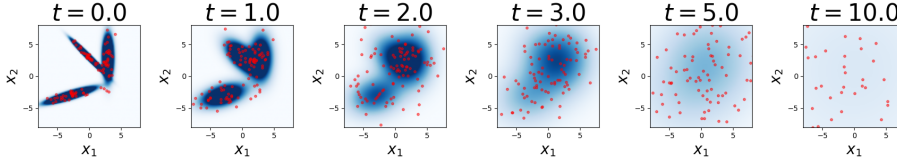


Figure 3.3: **Diffusing a mixture of Gaussians \mathbf{x} .** From left to right, we observe the evolution of the density p_t of \mathbf{x}_t as t grows from 0 to 10, along with some representative samples (red). The plane is colored by the probability density p_t ; high-density regions are colored darker blue. We observe that the probability mass becomes less concentrated as t increases, signaling that entropy increases.

where $T \in [0, \infty)$ is a time horizon and $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is drawn independently of \mathbf{x} . This process is an example of a *diffusion process*, so named because it spreads (i.e., *diffuses*) the probability mass over all of \mathbb{R}^D as time goes on, increasing the entropy. This intuition is confirmed graphically by Figure 3.3, and rigorously via the following theorem.

Theorem 3.1 (Simplified version of Theorem B.2). *Suppose that $(\mathbf{x}_t)_{t \in [0, T]}$ follows the model (3.2.1). For any $t \in (0, T]$, the random variable \mathbf{x}_t has a probability density p_t and differential entropy $h(\mathbf{x}_t) > -\infty$. Moreover, under certain technical conditions on \mathbf{x} ,*

$$\frac{d}{dt} h(\mathbf{x}_t) > 0, \quad \forall t \in (0, T], \quad (3.2.2)$$

i.e., the entropy of the distribution of \mathbf{x}_t increases over time t .

The proof is elementary, but it is rather long, so we postpone it to Section B.2.1. The main as-yet unstated implication of this result is that $h(\mathbf{x}_t) > h(\mathbf{x}_s)$ for every $t > s \geq 0$. To see this, note that if $s = 0$ and $h(\mathbf{x}) = h(\mathbf{x}_0) = -\infty$ then $h(\mathbf{x}_t) > -\infty$ for all $t > 0$, so the conclusion is true; otherwise $h(\mathbf{x}_t) = h(\mathbf{x}_s) + \int_s^t \left[\frac{d}{du} h(\mathbf{x}_u) \right] du > h(\mathbf{x}_s)$ by the fundamental theorem of calculus, so in both cases $h(\mathbf{x}_t) > h(\mathbf{x}_s)$ for every $t > s \geq 0$.

Example 3.2 (Adding noise to a mixture of Gaussians). In this example, we consider adding noise to a mixture of K Gaussians. Choose K means $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K \in \mathbb{R}^D$, K covariance matrices $\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K \in \text{PSD}(D)$, and probabilities $\pi_1, \dots, \pi_K \in [0, 1]$ such that $\sum_{k=1}^K \pi_k = 1$. Then \mathbf{x} is sampled as follows.

- First, an index (or *label*) $y \in [K]$ is sampled s.t. $y = k$ with probability π_k .
- Second, \mathbf{x} is sampled from the Gaussian $\mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$; i.e., if $y = k$ then $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$.

Thus \mathbf{x} has the Gaussian mixture distribution

$$\mathbf{x} \sim \sum_{k=1}^K \pi_k \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k). \quad (3.2.3)$$

Gaussian mixture models are relevant because they can approximate many distributions: given any distribution for \mathbf{x} , there exists a Gaussian mixture that approximates it arbitrarily well. For this reason, Gaussian mixture models are a popular, analytically tractable surrogate for multi-modal and complex data distributions. In this book we use them extensively as a “model problem” to understand and develop almost every concept. For now, we represent the data \mathbf{x} by a Gaussian mixture. After adding noise, the distribution of \mathbf{x}_t is again a Gaussian mixture:

$$\mathbf{x}_t = \mathbf{x} + t\mathbf{g} \sim \sum_{k=1}^K \pi_k \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k + t^2 \mathbf{I}). \quad (3.2.4)$$

To see by example how the entropy evolves, suppose $K = 1$, so $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. By (3.1.4), the entropy of \mathbf{x}_t is

$$h(\mathbf{x}_t) = \frac{D}{2}(1 + \log(2\pi)) + \frac{1}{2} \log \det(\boldsymbol{\Sigma} + t^2 \mathbf{I}). \quad (3.2.5)$$

Differentiating with respect to time (see Chapter A) gives

$$\frac{d}{dt} h(\mathbf{x}_t) = t \operatorname{tr}((\boldsymbol{\Sigma} + t^2 \mathbf{I})^{-1}) = \sum_{i=1}^D \frac{t}{\lambda_i(\boldsymbol{\Sigma}) + t^2} > 0, \quad (3.2.6)$$

where $\lambda_i(\boldsymbol{\Sigma})$ are the (non-negative) eigenvalues of $\boldsymbol{\Sigma}$. Because this derivative is positive, the entropy $h(\mathbf{x}_t)$ of the noised distribution increases over time. ■

Denoising Processes

The inverse operation to adding noise is known as *denoising*. It is a classical and well-studied topic in signal processing and system theory, instantiated (for example) as the Wiener filter and the Kalman filter (see Section 1.3.1). Several problems discussed in Chapter 2, such as PCA, ICA, and Dictionary Learning, are specific instances of the denoising problem. For a fixed t and the additive Gaussian noise model (3.2.1), the denoising problem can be formulated as attempting to learn a function $\bar{\mathbf{x}}^*(t, \cdot)$ that forms the best possible approximation (in expectation) of the true random variable \mathbf{x} , given both t and \mathbf{x}_t :

$$\bar{\mathbf{x}}^*(t, \cdot) \in \arg \min_{\bar{\mathbf{x}}(t, \cdot)} \mathbb{E} \|\mathbf{x} - \bar{\mathbf{x}}(t, \mathbf{x}_t)\|_2^2. \quad (3.2.7)$$

The solution to this problem, when optimizing $\bar{\mathbf{x}}(t, \cdot)$ over all possible (square-integrable) functions, is the so-called *Bayes optimal denoiser*, and turns out to be the conditional expectation:

$$\bar{\mathbf{x}}^*(t, \boldsymbol{\xi}) \doteq \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}]. \quad (3.2.8)$$

This expression justifies the notation $\bar{\mathbf{x}}$, which computes a conditional expectation (i.e., conditional mean or conditional average). In short, it attempts

to remove the noise from the noisy input, outputting the best possible guess (in expectation and with respect to the ℓ^2 distance) of the denoised original random variable.

Example 3.3 (Denoising Gaussian Noise from a Mixture of Gaussians). In this example we compute the Bayes optimal denoiser for an important class of distributions, the Gaussian mixture model in Example 3.2. Let us define $\varphi(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ as the probability density of $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ evaluated at \mathbf{x} . In this notation, the density of \mathbf{x}_t is

$$p_t(\mathbf{x}_t) = \sum_{k=1}^K \pi_k \varphi(\mathbf{x}_t; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k + t^2 \mathbf{I}). \quad (3.2.9)$$

Conditioned on y , the variables are jointly Gaussian: if we write $\mathbf{x} = \boldsymbol{\mu}_y + \boldsymbol{\Sigma}_y^{1/2} \mathbf{u}$ where $(\cdot)^{1/2}$ is the matrix square root and $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ independently of y (and \mathbf{g}), then we have

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{x}_t \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu}_y \\ \boldsymbol{\mu}_y \end{bmatrix} + \begin{bmatrix} \boldsymbol{\Sigma}_y^{1/2} & \mathbf{0} \\ \boldsymbol{\Sigma}_y^{1/2} & t\mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{g} \end{bmatrix}. \quad (3.2.10)$$

This shows that \mathbf{x} and \mathbf{x}_t are jointly Gaussian (conditioned on y) as claimed. Thus we can write

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{x}_t \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_y \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_y & \boldsymbol{\Sigma}_y \\ \boldsymbol{\Sigma}_y & \boldsymbol{\Sigma}_y + t^2 \mathbf{I} \end{bmatrix}\right). \quad (3.2.11)$$

Thus the conditional expectation of \mathbf{x} given \mathbf{x}_t (i.e., the Bayes optimal denoiser conditioned on y) is famously (see also Exercise 3.2)

$$\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t, y] = \boldsymbol{\mu}_y + \boldsymbol{\Sigma}_y(\boldsymbol{\Sigma}_y + t^2 \mathbf{I})^{-1}(\mathbf{x}_t - \boldsymbol{\mu}_y). \quad (3.2.12)$$

To find the overall Bayes optimal denoiser, we use the law of iterated expectation, obtaining

$$\bar{\mathbf{x}}^*(t, \mathbf{x}_t) = \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t] \quad (3.2.13)$$

$$= \mathbb{E}[\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t, y] \mid \mathbf{x}_t] \quad (3.2.14)$$

$$= \sum_{k=1}^K \mathbb{P}[y = k \mid \mathbf{x}_t] \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t, y = k]. \quad (3.2.15)$$

The probability can be dealt with as follows. Let $p_{t|y}$ be the probability density of \mathbf{x}_t conditioned on the value of y . Then

$$\mathbb{P}[y = k \mid \mathbf{x}_t] = \frac{p_{t|y}(\mathbf{x}_t \mid k) \pi_k}{p_t(\mathbf{x}_t)} \quad (3.2.16)$$

$$= \frac{\pi_k \varphi(\mathbf{x}_t; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k + t^2 \mathbf{I})}{\sum_{i=1}^K \pi_i \varphi(\mathbf{x}_t; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i + t^2 \mathbf{I})}. \quad (3.2.17)$$

On the other hand, the conditional expectation is as described before:

$$\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t, y = k] = \boldsymbol{\mu}_k + \boldsymbol{\Sigma}_k(\boldsymbol{\Sigma}_k + t^2\mathbf{I})^{-1}(\mathbf{x}_t - \boldsymbol{\mu}_k). \quad (3.2.18)$$

So putting this all together, the true Bayes optimal denoiser is

$$\bar{\mathbf{x}}^*(t, \mathbf{x}_t) = \sum_{k=1}^K \left\{ \frac{\pi_k \varphi(\mathbf{x}_t; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k + t^2\mathbf{I})}{\sum_{i=1}^K \pi_i \varphi(\mathbf{x}_t; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i + t^2\mathbf{I})} \right. \quad (3.2.19)$$

$$\left. \cdot (\boldsymbol{\mu}_k + \boldsymbol{\Sigma}_k(\boldsymbol{\Sigma}_k + t^2\mathbf{I})^{-1}(\mathbf{x}_t - \boldsymbol{\mu}_k)) \right\}. \quad (3.2.20)$$

This example is particularly important, and several special cases will give us great conceptual insight later. For now, let us attempt to extract some geometric intuition from the functional form of the optimal denoiser (3.2.19).

To understand (3.2.19) intuitively, let us first set $K = 1$ (i.e., one Gaussian) such that $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Let us then diagonalize $\boldsymbol{\Sigma} = \mathbf{V}\boldsymbol{\Lambda}\mathbf{V}^\top$. Then the Bayes optimal denoiser is

$$\bar{\mathbf{x}}^*(t, \mathbf{x}_t) = \boldsymbol{\mu} + \boldsymbol{\Sigma}(\boldsymbol{\Sigma} + t^2\mathbf{I})^{-1}(\mathbf{x}_t - \boldsymbol{\mu}) \quad (3.2.21)$$

$$= \boldsymbol{\mu} + \mathbf{V} \begin{bmatrix} \lambda_1/(\lambda_1 + t^2) & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \lambda_D/(\lambda_D + t^2) \end{bmatrix} \mathbf{V}^\top(\mathbf{x}_t - \boldsymbol{\mu}), \quad (3.2.22)$$

where $\lambda_1, \dots, \lambda_D$ are the eigenvalues of $\boldsymbol{\Sigma}$. We can observe that this denoiser has three steps:

- Translate the input \mathbf{x}_t by $\boldsymbol{\mu}$.
- Contract the (translated) input $\mathbf{x}_t - \boldsymbol{\mu}$ in each eigenvector direction by a quantity $\lambda_i/(\lambda_i + t^2)$. If the translated input is low-rank and some eigenvalues of $\boldsymbol{\Sigma}$ are zero, these directions get immediately contracted to 0 by the denoiser, ensuring that the output of the contraction is similarly low-rank.
- Translate the output back by $\boldsymbol{\mu}$.

This is the geometric interpretation of the denoiser of a *single* Gaussian. The overall denoiser of the Gaussian mixture model (3.2.19) uses K such denoisers, weighting their output by the posterior probabilities $\mathbb{P}[y = k \mid \mathbf{x}_t]$. If the means of the Gaussians are well separated, these posterior probabilities are very close to 0 or 1 near each mean or cluster. In this regime, the overall denoiser (3.2.19) has the same geometric interpretation as the above single Gaussian denoiser. ■

Intuitively, and as we can see from Example 3.3, the Bayes-optimal denoiser $\bar{\mathbf{x}}^*(t, \cdot)$ should move its input \mathbf{x}_t toward the modes of the distribution of \mathbf{x} . It turns out that we can quantify this by showing that the Bayes-optimal denoiser *takes a gradient-ascent step* on the (log-)density of \mathbf{x}_t , which (recall) we denoted

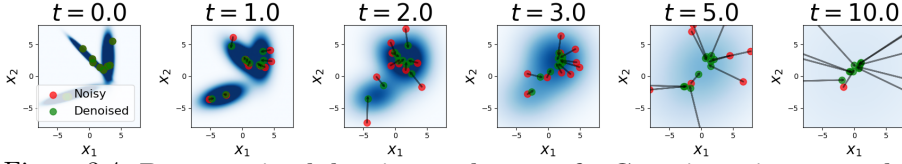


Figure 3.4: **Bayes-optimal denoiser and score of a Gaussian mixture model.** In the same setting as Figure 3.3, we demonstrate the effect of the Bayes-optimal denoiser $\bar{\mathbf{x}}^*$ by plotting \mathbf{x}_t (red) and $\bar{\mathbf{x}}^*(t, \mathbf{x}_t)$ (green) for some choice of t and \mathbf{x}_t . By Tweedie’s formula Theorem 3.2, the residual between them is proportional to the so-called (Hyvärinen) score $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$. We can see that the score points toward the modes of the distribution of \mathbf{x}_t .

p_t . That is, following the denoiser means moving from the input iterate to a region of higher probability within this (perturbed) distribution. For small t , the perturbation is small, so our initial intuition is therefore almost exactly right. The picture is visualized in Figure 3.4 and rigorously formulated as Tweedie’s formula [Rob56].

Theorem 3.2 (Tweedie’s Formula). *Suppose that $(\mathbf{x}_t)_{t \in [0, T]}$ obeys (3.2.1). Let p_t be the density of \mathbf{x}_t (as previously declared). Then*

$$\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t] = \mathbf{x}_t + t^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t). \quad (3.2.23)$$

Proof. For the proof, suppose that \mathbf{x} has a density,⁵ and call this density p . Let $p_{0|t}$ and $p_{t|0}$ be the conditional densities of $\mathbf{x} = \mathbf{x}_0$ given \mathbf{x}_t and \mathbf{x}_t given \mathbf{x} respectively. Let $\varphi(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ be the density of $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ evaluated at \mathbf{x} , so that $p_{t|0}(\mathbf{x}_t \mid \mathbf{x}) = \varphi(\mathbf{x}_t; \mathbf{x}, t^2 \mathbf{I})$. Then a simple calculation gives

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) = \frac{\nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t)}{p_t(\mathbf{x}_t)} \quad (3.2.24)$$

$$= \frac{1}{p_t(\mathbf{x}_t)} \nabla_{\mathbf{x}_t} \int_{\mathbb{R}^D} p(\mathbf{x}) p_{t|0}(\mathbf{x}_t \mid \mathbf{x}) d\mathbf{x} \quad (3.2.25)$$

$$= \frac{1}{p_t(\mathbf{x}_t)} \nabla_{\mathbf{x}_t} \int_{\mathbb{R}^D} p(\mathbf{x}) \varphi(\mathbf{x}_t; \mathbf{x}, t^2 \mathbf{I}) d\mathbf{x} \quad (3.2.26)$$

$$= \frac{1}{p_t(\mathbf{x}_t)} \int_{\mathbb{R}^D} p(\mathbf{x}) [\nabla_{\mathbf{x}_t} \varphi(\mathbf{x}_t; \mathbf{x}, t^2 \mathbf{I})] d\mathbf{x} \quad (3.2.27)$$

$$= \frac{1}{p_t(\mathbf{x}_t)} \int_{\mathbb{R}^D} p(\mathbf{x}) \varphi(\mathbf{x}_t; \mathbf{x}, t^2 \mathbf{I}) \left[-\frac{\mathbf{x}_t - \mathbf{x}}{t^2} \right] d\mathbf{x} \quad (3.2.28)$$

$$= \frac{1}{t^2 p_t(\mathbf{x}_t)} \int_{\mathbb{R}^D} p(\mathbf{x}) \varphi(\mathbf{x}_t; \mathbf{x}, t^2 \mathbf{I}) [\mathbf{x} - \mathbf{x}_t] d\mathbf{x} \quad (3.2.29)$$

$$= \frac{1}{t^2 p_t(\mathbf{x}_t)} \int_{\mathbb{R}^D} p(\mathbf{x}) \varphi(\mathbf{x}_t; \mathbf{x}, t^2 \mathbf{I}) \mathbf{x} d\mathbf{x} \quad (3.2.30)$$

⁵Note that the theorem is true without this assumption, and it can be proven through a minimal modification of this proof to use probability measures instead of densities.

$$\begin{aligned}
& - \frac{\mathbf{x}_t}{t^2 p_t(\mathbf{x}_t)} \int_{\mathbb{R}^D} p(\mathbf{x}) \varphi(\mathbf{x}_t; \mathbf{x}, t^2 \mathbf{I}) d\mathbf{x} \\
&= \frac{1}{t^2 p_t(\mathbf{x}_t)} \int_{\mathbb{R}^D} p(\mathbf{x}) p_{t|0}(\mathbf{x}_t | \mathbf{x}) \mathbf{x} d\mathbf{x} - \frac{\mathbf{x}_t}{t^2 p_t(\mathbf{x}_t)} p_t(\mathbf{x}_t) \quad (3.2.31)
\end{aligned}$$

$$= \frac{1}{t^2 p_t(\mathbf{x}_t)} \int_{\mathbb{R}^D} p_t(\mathbf{x}_t) p_{0|t}(\mathbf{x} | \mathbf{x}_t) \mathbf{x} d\mathbf{x} - \frac{\mathbf{x}_t}{t^2 p_t(\mathbf{x}_t)} p_t(\mathbf{x}_t) \quad (3.2.32)$$

$$= \frac{1}{t^2} \int_{\mathbb{R}^D} p_{0|t}(\mathbf{x} | \mathbf{x}_t) \mathbf{x} d\mathbf{x} - \frac{\mathbf{x}_t}{t^2} \quad (3.2.33)$$

$$= \frac{1}{t^2} \mathbb{E}[\mathbf{x} | \mathbf{x}_t] - \frac{\mathbf{x}_t}{t^2} \quad (3.2.34)$$

$$= \frac{\mathbb{E}[\mathbf{x} | \mathbf{x}_t] - \mathbf{x}_t}{t^2}. \quad (3.2.35)$$

Simple rearrangement of the above equality proves the theorem. \square

This result develops a connection between denoising and optimization: the Bayes-optimal denoiser takes a single step of gradient ascent on the perturbed data density p_t , and the step size adaptively becomes smaller (i.e., takes more precise steps) as the perturbation to the data distribution grows smaller. The quantity $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ is called the (*Hyvärinen*) *score* and frequently appears in discussions about denoising, etc.; it first appeared in a paper of Aapo Hyvärinen in the context of ICA [Hyv05].

Similar to how one step of gradient descent is almost never sufficient to minimize an objective in practice when initializing far from the optimum, the output of the Bayes-optimal denoiser $\bar{\mathbf{x}}^*(t, \cdot)$ is almost never contained in a high-probability region of the data distribution when t is large, *especially* when the data have low-dimensional structure. We illustrate this point explicitly in the following example.

Example 3.4 (Denoising a two-point mixture). Let x be uniform on the two-point set $\{-1, +1\}$ and let $(x_t)_{t \in [0, T]}$ follow (3.2.1). This is precisely a degenerate Gaussian mixture model with priors equal to $\frac{1}{2}$, means $\{-1, +1\}$, and covariances both equal to 0. For a fixed $t > 0$ we can use the calculation of the Bayes-optimal denoiser in (3.2.19) to obtain (proof as exercise)

$$\bar{x}^*(t, x_t) = \frac{\varphi(x_t; +1, t^2) - \varphi(x_t; -1, t^2)}{\varphi(x_t; +1, t^2) + \varphi(x_t; -1, t^2)} = \tanh\left(\frac{x_t}{t^2}\right). \quad (3.2.36)$$

For t near 0, this quantity is near $\{-1, +1\}$ for almost all inputs $\bar{x}^*(t, x_t)$. However, for t large, this quantity is not necessarily even approximately in the original support of x , which, remember, is $\{-1, +1\}$. In particular, for $x_t \approx 0$ it holds $\bar{x}^*(t, x_t) \approx 0$, which lies completely in between the two possible points. Thus \bar{x}^* will not output “realistic” x . More rigorously, the distribution of $\bar{x}^*(t, x_t)$ is very different from the distribution of x . \blacksquare

Incremental Denoising

Therefore, to denoise the very noisy sample \mathbf{x}_T (where, recall, T is the maximum time), we cannot use the denoiser just once. Instead, we must apply it many times, analogously to gradient descent with decaying step sizes, to converge to a stationary point $\hat{\mathbf{x}}$. Namely, we use the denoiser to go from \mathbf{x}_T to $\hat{\mathbf{x}}_{T-\delta}$, which approximates $\mathbf{x}_{T-\delta}$, then from $\hat{\mathbf{x}}_{T-\delta}$ to $\hat{\mathbf{x}}_{T-2\delta}$, and so on, all the way from $\hat{\mathbf{x}}_\delta$ to $\hat{\mathbf{x}} = \hat{\mathbf{x}}_0$. Each denoising step makes the action of the denoiser more like a gradient step on the original (log-)density.

More formally, we uniformly discretize $[0, T]$ into $L + 1$ timesteps $0 = t_0 < t_1 < \dots < t_L = T$, i.e.,

$$t_\ell = \frac{\ell}{L}T, \quad \ell \in \{0, 1, \dots, L\}. \quad (3.2.37)$$

Then for each $\ell \in [L] = \{1, 2, \dots, L\}$, going from $\ell = L$ to $\ell = 1$, we run the iteration

$$\hat{\mathbf{x}}_{t_{\ell-1}} = \mathbb{E}[\mathbf{x}_{t_{\ell-1}} \mid \mathbf{x}_{t_\ell} = \hat{\mathbf{x}}_{t_\ell}] \quad (3.2.38)$$

$$= \mathbb{E}[\mathbf{x} + t_{\ell-1}\mathbf{g} \mid \mathbf{x}_{t_\ell} = \hat{\mathbf{x}}_{t_\ell}] \quad (3.2.39)$$

$$= \mathbb{E}\left[\mathbf{x} + t_{\ell-1} \cdot \frac{\mathbf{x}_{t_\ell} - \mathbf{x}}{t_\ell} \mid \mathbf{x}_{t_\ell} = \hat{\mathbf{x}}_{t_\ell}\right] \quad (3.2.40)$$

$$= \frac{t_{\ell-1}}{t_\ell} \hat{\mathbf{x}}_{t_\ell} + \left(1 - \frac{t_{\ell-1}}{t_\ell}\right) \mathbb{E}[\mathbf{x} \mid \mathbf{x}_{t_\ell} = \hat{\mathbf{x}}_{t_\ell}] \quad (3.2.41)$$

$$= \left(1 - \frac{1}{\ell}\right) \hat{\mathbf{x}}_{t_\ell} + \frac{1}{\ell} \bar{\mathbf{x}}^*(t_\ell, \hat{\mathbf{x}}_{t_\ell}). \quad (3.2.42)$$

At the beginning of the iteration, when ℓ is large, we barely trust the denoiser's output and mostly keep the current iterate. This makes sense, as the denoiser can have huge variance (cf. Example 3.4). When ℓ is small, the denoiser “locks on” to the modes of the data distribution, since a denoising step is essentially a gradient step on the true distribution's log-density, and we can trust it not to produce unreasonable samples; thus the denoising step is dominated by the denoiser's output. At $\ell = 1$ we even discard the current iterate and keep only the denoiser's output.

We visualize the convergence process in \mathbb{R}^3 in Figure 3.5. While we will develop some rigorous results about convergence later, we now give a closed-form example of the iterative denoising process applied to denoising a Gaussian data distribution.

Example 3.5 (Iterative denoising for a single Gaussian). To understand the above iterative denoising iteration intuitively, we apply it to denoise \mathbf{x}_T defined in (3.2.1), where $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Using (3.2.42), we compute

$$\hat{\mathbf{x}}_{t_{\ell-1}} - \boldsymbol{\mu} = \left(1 - \frac{1}{\ell}\right) \hat{\mathbf{x}}_{t_\ell} + \frac{1}{\ell} \bar{\mathbf{x}}^*(t_\ell, \hat{\mathbf{x}}_{t_\ell}) - \boldsymbol{\mu} \quad (3.2.43)$$

$$= \left(1 - \frac{1}{\ell}\right) (\hat{\mathbf{x}}_{t_\ell} - \boldsymbol{\mu}) + \frac{1}{\ell} (\bar{\mathbf{x}}^*(t_\ell, \hat{\mathbf{x}}_{t_\ell}) - \boldsymbol{\mu}). \quad (3.2.44)$$

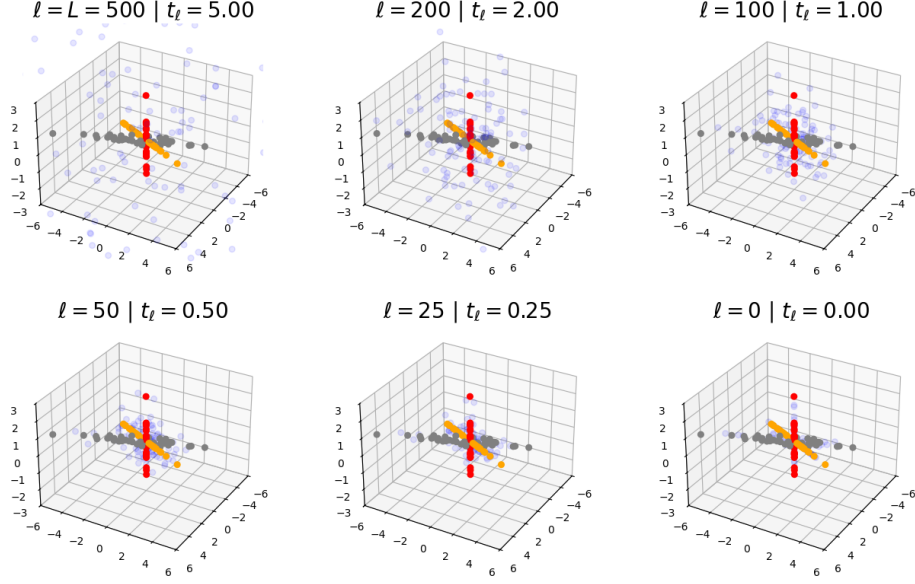


Figure 3.5: **Denoising a low-rank mixture of Gaussians.** Each figure shows samples from the true data distribution (gray, orange, red) and samples undergoing the denoising process (3.2.82) (light blue). At top left, the process has just started and the noise is very large. As the process continues, the noise is pushed further toward the support of the low-rank data distribution. Finally, in the bottom right, the generated samples are perfectly aligned with the support of the data and closely resemble samples drawn from the low-rank Gaussian mixture model.

Then, plugging in the form for $\bar{\mathbf{x}}^*$ that we obtain from (3.2.22) and performing some tedious matrix algebra, we obtain

$$\mathbf{x}_{t_{\ell-1}} - \boldsymbol{\mu} = \mathbf{V} \begin{bmatrix} 1 - \frac{\ell T^2}{\lambda_1 L^2 + \ell^2 T^2} & & \\ & \ddots & \\ & & 1 - \frac{\ell T^2}{\lambda_D L^2 + \ell^2 T^2} \end{bmatrix} \mathbf{V}^\top (\mathbf{x}_{t_\ell} - \boldsymbol{\mu}). \quad (3.2.45)$$

Notice that all entries in this diagonal matrix are in $(0, 1)$. Taking the ℓ^2 -norm and using a cheap yet tight upper bound, it holds

$$\|\mathbf{x}_{t_{\ell-1}} - \boldsymbol{\mu}\|_2 \leq \underbrace{\max_{i \in [D]} \left\{ 1 - \frac{\ell T^2}{\lambda_i L^2 + \ell^2 T^2} \right\}}_{c_{\ell, L}} \cdot \|\mathbf{x}_{t_\ell} - \boldsymbol{\mu}\|_2. \quad (3.2.46)$$

At first glance, this inequality seems to imply that the denoising iteration is *very close*, conceptually, to a *contraction mapping* such as power iteration (2.1.17) sending the iterate \mathbf{x}_t to the mean $\boldsymbol{\mu}$, since $c_{\ell, L}$ are all in $(0, 1)$ for any ℓ and L .

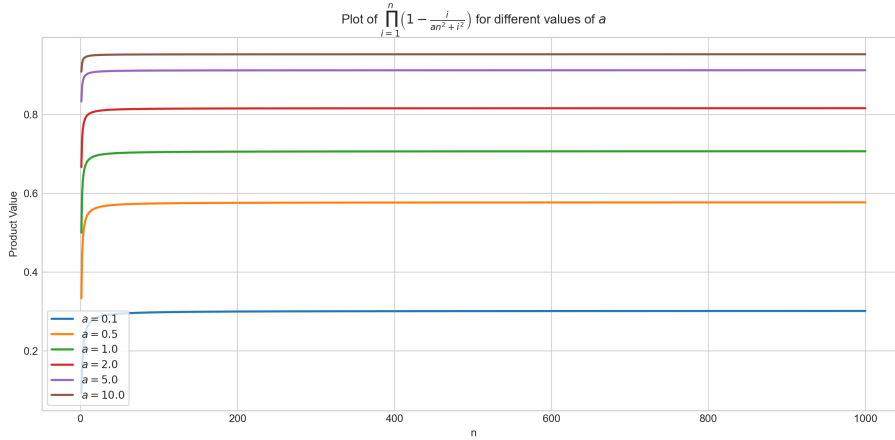


Figure 3.6: **The denoising iteration from Example 3.5 is not a contraction mapping.** Plot of the contraction coefficient products $\prod_{\ell=1}^L c_{\ell,L}$ in the case where the data are in $D = 1$, the time horizon $T = 1$, and the data variance $\lambda_1 = a$, for varying values of a . If the denoising iteration were (conceptually similar to) a contraction mapping, these products would go to 0, but instead they converge to non-zero values.

However, this is false. To see why, observe that by iterating (3.2.46) it holds

$$\|\mathbf{x}_0 - \boldsymbol{\mu}\|_2 \leq \prod_{\ell=1}^L c_{\ell,L} \cdot \|\mathbf{x}_T - \boldsymbol{\mu}\|_2. \quad (3.2.47)$$

If each denoising update had similar properties to a contraction mapping, then $\prod_{\ell=1}^L c_{\ell,L}$ would go to 0 as $L \rightarrow \infty$, and the iterate would converge to $\boldsymbol{\mu}$ extremely fast, rendering the sampler completely useless or vacuous. But in reality the denoising iteration is *actually not a contraction mapping* since $c_{\ell,L}$ changes with both the indices ℓ and L , so that $\prod_{\ell=1}^L c_{\ell,L}$ does not go to 0 as $L \rightarrow \infty$, as demonstrated in Figure 3.6 (it is also possible to compute the limit analytically). In reality, the sampler ensures that the terminal iterates spread out amongst the support of the data distribution, enabling actual sampling from the data distribution, which is quite unlike the objective of power iteration (for example). ■

Recall that we wanted to build the iterative denoising process to reduce the entropy of the iterates. While we did this in a roundabout way by inverting a process that adds entropy, it is now time to pay the piper and confirm that our iterative denoising process reduces the entropy.

Theorem 3.3 (Simplified Version of Theorem B.3). *Suppose that $(\mathbf{x}_t)_{t \in [0, T]}$ obeys (3.2.1). Then, under certain technical conditions on \mathbf{x} , for every $s < t$ with $s, t \in (0, T]$,*

$$h(\mathbb{E}[\mathbf{x}_s \mid \mathbf{x}_t]) < h(\mathbf{x}_t). \quad (3.2.48)$$

The full statement of the theorem and the proof itself require some technicality, so they are postponed to Section [B.2.2](#).

Computing the Denoiser

The last thing we discuss here is that many times, we will *not be able to compute* $\bar{\mathbf{x}}^*(t, \cdot)$ for any t , since we do not have the distribution p_t . But we can try to *learn one from data*. Recall that the denoiser $\bar{\mathbf{x}}^*$ is defined in [\(3.2.7\)](#) as minimizing the mean-squared error $\mathbb{E} \|\bar{\mathbf{x}}(t, \mathbf{x}_t) - \mathbf{x}\|_2^2$. We can use this mean-squared error as a loss or objective function to learn the denoiser. For example, we can parameterize $\bar{\mathbf{x}}(t, \cdot)$ by a neural network, writing it as $\bar{\mathbf{x}}_\theta(t, \cdot)$, and optimize the loss over the parameter space Θ :

$$\min_{\theta \in \Theta} \mathbb{E} \|\bar{\mathbf{x}}_\theta(t, \mathbf{x}_t) - \mathbf{x}\|_2^2. \quad (3.2.49)$$

The solution to this optimization problem, implemented via gradient descent or a similar algorithm, will give us a $\bar{\mathbf{x}}_{\theta^*}(t, \cdot)$ which is a good approximation to $\bar{\mathbf{x}}^*(t, \cdot)$ (at least if the training works) and which we will use as our denoiser.

What is a good architecture for this neural network $\bar{\mathbf{x}}_{\theta^*}(t, \cdot)$? To answer this question, we will first attempt to understand the case of data \mathbf{x} that are *low-rank*, i.e., having low-rank support. The most canonical way to represent low-rank data is through a *low-rank Gaussian*, and so we start our exploration by assuming that \mathbf{x} is drawn (approximately) from a low-rank Gaussian:

$$\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{U}\mathbf{U}^\top), \quad (3.2.50)$$

where $\mathbf{U} \in \mathcal{O}(D, P)$ is a tall orthogonal matrix with $P \leq D$ (using zero-mean for essentially notational convenience). The optimal denoiser under [\(3.2.50\)](#) is (from [Example 3.3](#)):

$$\bar{\mathbf{x}}^*(t, \mathbf{x}_t) = \mathbf{U}\mathbf{U}^\top (\mathbf{U}\mathbf{U}^\top + t^2\mathbf{I})^{-1} \mathbf{x}_t. \quad (3.2.51)$$

The matrix $\mathbf{U}\mathbf{U}^\top + t^2\mathbf{I}$ is a low-rank perturbation of the full-rank matrix $t^2\mathbf{I}$, and thus ripe for simplification via the *Sherman-Morrison-Woodbury identity*, i.e., for matrices $\mathbf{A}, \mathbf{C}, \mathbf{U}, \mathbf{V}$ such that \mathbf{A} and \mathbf{C} are invertible,

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}\mathbf{VA}^{-1}. \quad (3.2.52)$$

We prove this identity in [Exercise 3.3](#). For now we apply this identity with $\mathbf{A} = t^2\mathbf{I}$, $\mathbf{U} = \mathbf{U}$, $\mathbf{V} = \mathbf{U}^\top$, and $\mathbf{C} = \mathbf{I}$, obtaining

$$(\mathbf{U}\mathbf{U}^\top + t^2\mathbf{I})^{-1} = \frac{1}{t^2}\mathbf{I} - \frac{1}{t^4}\mathbf{U} \left(\mathbf{I} + \frac{1}{t^2}\mathbf{U}^\top\mathbf{U} \right)^{-1} \mathbf{U}^\top \quad (3.2.53)$$

$$= \frac{1}{t^2}\mathbf{I} - \frac{1}{t^4 \left(1 + \frac{1}{t^2}\right)} \mathbf{U}\mathbf{U}^\top \quad (3.2.54)$$

$$= \frac{1}{t^2} \left(\mathbf{I} - \frac{1}{1 + t^2} \mathbf{U}\mathbf{U}^\top \right). \quad (3.2.55)$$

This implies

$$\mathbf{U}\mathbf{U}^\top(\mathbf{U}\mathbf{U}^\top + t^2\mathbf{I})^{-1} = \frac{1}{t^2}\mathbf{U}\mathbf{U}^\top \left(\mathbf{I} - \frac{1}{1+t^2}\mathbf{U}\mathbf{U}^\top \right) \quad (3.2.56)$$

$$= \frac{1}{t^2} \left(1 - \frac{1}{1+t^2} \right) \mathbf{U}\mathbf{U}^\top \quad (3.2.57)$$

$$= \frac{1}{1+t^2} \mathbf{U}\mathbf{U}^\top. \quad (3.2.58)$$

Our optimal denoiser is a (rescaled) linear projection onto the low-rank subspace supporting our data distribution. This is essentially the same thing as probabilistic PCA from Chapter 2 (cf. Section 2.1.4). We can formalize this intuition by considering a class of subspace-parameterized denoisers:

$$\bar{\mathbf{x}}_{\mathbf{V}}(t, \mathbf{x}_t) = \frac{1}{1+t^2} \mathbf{V}\mathbf{V}^\top \mathbf{x}_t, \quad (3.2.59)$$

and showing that the denoising problem $\min_{\mathbf{V} \in \mathcal{O}(D,P)} \mathbb{E} \|\bar{\mathbf{x}}_{\mathbf{V}}(t, \mathbf{x}_t) - \mathbf{x}\|_2^2$ obtains the same solution as probabilistic PCA (Section 2.1.4). Indeed, plugging our subspace denoiser into the training loss (3.2.7), we obtain

$$\min_{\mathbf{V} \in \mathcal{O}(D,P)} \mathbb{E} \|\bar{\mathbf{x}}_{\mathbf{V}}(t, \mathbf{x}_t) - \mathbf{x}\|_2^2 = \mathbb{E} \left\| \frac{1}{1+t^2} \mathbf{V}\mathbf{V}^\top (\mathbf{x} + t\mathbf{g}) - \mathbf{x} \right\|_2^2, \quad (3.2.60)$$

where the equality is due to (3.2.1). Conditioned on \mathbf{x} (thus integrating over \mathbf{g}), we compute

$$\mathbb{E}_{\mathbf{g}} \left\| \frac{1}{1+t^2} \mathbf{V}\mathbf{V}^\top (\mathbf{x} + t\mathbf{g}) - \mathbf{x} \right\|_2^2 \quad (3.2.61)$$

$$= \left\| \frac{1}{1+t^2} \mathbf{V}\mathbf{V}^\top \mathbf{x} - \mathbf{x} \right\|_2^2 - \frac{t}{1+t^2} \mathbb{E}_{\mathbf{g}} \left\langle \mathbf{x} - \frac{1}{1+t^2} \mathbf{V}\mathbf{V}^\top \mathbf{x}, \mathbf{V}\mathbf{V}^\top \mathbf{g} \right\rangle \quad (3.2.62)$$

$$+ \frac{t^2}{(1+t^2)^2} \mathbb{E}_{\mathbf{g}} \|\mathbf{V}\mathbf{V}^\top \mathbf{g}\|_2^2 \quad (3.2.63)$$

$$= \left\| \frac{1}{1+t^2} \mathbf{V}\mathbf{V}^\top \mathbf{x} - \mathbf{x} \right\|_2^2 + \frac{t^2 P}{(1+t^2)^2} \quad (3.2.64)$$

where the second equality follows from $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and the fact that $\mathbf{V} \in \mathcal{O}(D, P)$ so that $\mathbb{E}_{\mathbf{g}} \|\mathbf{V}\mathbf{V}^\top \mathbf{g}\|_2^2 = \mathbb{E}_{\mathbf{g}} \|\mathbf{V}^\top \mathbf{g}\|_2^2 = P$. Therefore, the denoising problem is equivalent to

$$\min_{\mathbf{V} \in \mathcal{O}(D,P)} \left\{ \mathbb{E} \left\| \frac{1}{1+t^2} \mathbf{V}\mathbf{V}^\top \mathbf{x} - \mathbf{x} \right\|_2^2 \right. \quad (3.2.65)$$

$$\left. = \mathbb{E} \|\mathbf{x}\|_2^2 + \left(\frac{1}{(1+t^2)^2} - \frac{2}{1+t^2} \right) \mathbb{E} \|\mathbf{V}^\top \mathbf{x}\|_2^2 \right\}, \quad (3.2.66)$$

which in turn is equivalent to the problem

$$\max_{\mathbf{V} \in \mathcal{O}(D, P)} \mathbb{E} \|\mathbf{V}^\top \mathbf{x}\|_2^2 \quad (3.2.67)$$

which was shown in Section 2.1.1 to be equivalent to the problem

$$\min_{\mathbf{V} \in \mathcal{O}(D, P)} \mathbb{E} \|\mathbf{x} - \mathbf{V}\mathbf{V}^\top \mathbf{x}\|_2^2, \quad (3.2.68)$$

i.e., the probabilistic PCA formulation (2.1.30).

Thus, if we were to assume our data were drawn from a low-rank Gaussian, our denoising algorithm would essentially consist of many PCA iterations in a row applied to progressively less-noisy data. While this method sounds very simple, it can actually produce complex and relatively powerful denoising behaviors when applied to practical data, as demonstrated by PCANet [CJG+15]. Nevertheless, to obtain a more general denoising operator that may work on complex and multi-modal data distributions, we will consider a broader class of data distributions: the ubiquitous class of *low-rank Gaussian mixture models*, whose denoiser we computed in Example 3.2 and Example 3.3. In the context of denoising, since Gaussian mixture models can approximate a very broad class of distributions arbitrarily well, in practice optimizing among the class of denoisers for Gaussian mixture models can potentially give us something close to the optimal denoiser for the real data distribution.

Formally, we now assume that \mathbf{x} is drawn from a low-rank Gaussian mixture model $\frac{1}{K} \sum_{k=1}^K \mathcal{N}(\mathbf{0}, \mathbf{U}_k \mathbf{U}_k^\top)$, where $\mathbf{U}_k \in \mathcal{O}(D, P)$ are tall orthogonal matrices whose columns span the P -dimensional subspaces that support the K low-rank Gaussian components. We write

$$\mathbf{x} \sim \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\mathbf{0}, \mathbf{U}_k \mathbf{U}_k^\top). \quad (3.2.69)$$

Then the optimal denoiser under (3.2.1) is (from Example 3.3)

$$\bar{\mathbf{x}}^*(t, \mathbf{x}_t) = \sum_{k=1}^K \left\{ \frac{\varphi(\mathbf{x}_t; \mathbf{0}, \mathbf{U}_k \mathbf{U}_k^\top + t^2 \mathbf{I})}{\sum_{i=1}^K \varphi(\mathbf{x}_t; \mathbf{0}, \mathbf{U}_i \mathbf{U}_i^\top + t^2 \mathbf{I})} \right. \quad (3.2.70)$$

$$\left. \cdot (\mathbf{U}_k \mathbf{U}_k^\top (\mathbf{U}_k \mathbf{U}_k^\top + t^2 \mathbf{I})^{-1} \mathbf{x}_t) \right\}. \quad (3.2.71)$$

Applying the same Sherman-Morrison-Woodbury identity as in (3.2.55), we obtain

$$(\mathbf{U}_k \mathbf{U}_k^\top + t^2 \mathbf{I})^{-1} = \frac{1}{t^2} \left(\mathbf{I} - \frac{1}{1+t^2} \mathbf{U}_k \mathbf{U}_k^\top \right). \quad (3.2.72)$$

We can then compute the posterior probabilities as follows. Since the \mathbf{U}_k are orthogonal, $\det(\mathbf{U}_k \mathbf{U}_k^\top + t^2 \mathbf{I})$ is the same for every k . Thus

$$\frac{\varphi(\mathbf{x}_t; \mathbf{0}, \mathbf{U}_k \mathbf{U}_k^\top + t^2 \mathbf{I})}{\sum_{i=1}^K \varphi(\mathbf{x}_t; \mathbf{0}, \mathbf{U}_i \mathbf{U}_i^\top + t^2 \mathbf{I})} \quad (3.2.73)$$

$$= \frac{\exp\left(-\frac{1}{2}\mathbf{x}_t^\top(\mathbf{U}_k\mathbf{U}_k^\top + t^2\mathbf{I})^{-1}\mathbf{x}_t\right)}{\sum_{i=1}^K \exp\left(-\frac{1}{2}\mathbf{x}_t^\top(\mathbf{U}_i\mathbf{U}_i^\top + t^2\mathbf{I})^{-1}\mathbf{x}_t\right)} \quad (3.2.74)$$

$$= \frac{\exp\left(-\frac{1}{2t^2}\mathbf{x}_t^\top\left(\mathbf{I} - \frac{1}{1+t^2}\mathbf{U}_k\mathbf{U}_k^\top\right)\mathbf{x}_t\right)}{\sum_{i=1}^K \exp\left(-\frac{1}{2t^2}\mathbf{x}_t^\top\left(\mathbf{I} - \frac{1}{1+t^2}\mathbf{U}_i\mathbf{U}_i^\top\right)\mathbf{x}_t\right)} \quad (3.2.75)$$

$$= \frac{\exp\left(-\frac{1}{2t^2}\|\mathbf{x}_t\|_2^2 + \frac{1}{2t^2(1+t^2)}\|\mathbf{U}_k^\top\mathbf{x}_t\|_2^2\right)}{\sum_{i=1}^K \exp\left(-\frac{1}{2t^2}\|\mathbf{x}_t\|_2^2 + \frac{1}{2t^2(1+t^2)}\|\mathbf{U}_i^\top\mathbf{x}_t\|_2^2\right)} \quad (3.2.76)$$

$$= \frac{\exp\left(-\frac{1}{2t^2}\|\mathbf{x}_t\|_2^2\right) \exp\left(\frac{1}{2t^2(1+t^2)}\|\mathbf{U}_k^\top\mathbf{x}_t\|_2^2\right)}{\exp\left(-\frac{1}{2t^2}\|\mathbf{x}_t\|_2^2\right) \sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)}\|\mathbf{U}_i^\top\mathbf{x}_t\|_2^2\right)} \quad (3.2.77)$$

$$= \frac{\exp\left(\frac{1}{2t^2(1+t^2)}\|\mathbf{U}_k^\top\mathbf{x}_t\|_2^2\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)}\|\mathbf{U}_i^\top\mathbf{x}_t\|_2^2\right)}. \quad (3.2.78)$$

This is a softmax operation weighted by the projection of \mathbf{x}_t onto each subspace measured by $\|\mathbf{U}_i^\top\mathbf{x}_t\|_2$ (tempered by the temperature $2t^2(1+t^2)$). Meanwhile, the component denoisers can be written in the same way as (3.2.58) to obtain

$$\mathbf{U}_k\mathbf{U}_k^\top(\mathbf{U}_k\mathbf{U}_k^\top + t^2\mathbf{I})^{-1}\mathbf{x}_t = \frac{1}{1+t^2}\mathbf{U}_k\mathbf{U}_k^\top\mathbf{x}_t. \quad (3.2.79)$$

Putting these together, we have

$$\bar{\mathbf{x}}^*(t, \mathbf{x}_t) = \frac{1}{1+t^2} \sum_{k=1}^K \frac{\exp\left(\frac{1}{2t^2(1+t^2)}\|\mathbf{U}_k^\top\mathbf{x}_t\|_2^2\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)}\|\mathbf{U}_i^\top\mathbf{x}_t\|_2^2\right)} \mathbf{U}_k\mathbf{U}_k^\top\mathbf{x}_t, \quad (3.2.80)$$

i.e., a projection of \mathbf{x}_t onto each of the K subspaces, weighted by a softmax of a quadratic function of \mathbf{x}_t . This functional form is similar to an *attention mechanism* in a transformer architecture! As we will see in Chapter 5, this is no coincidence at all; the deep link between denoising and lossy compression (to be covered in Section 4.1) makes transformer denoisers so effective in practice. Thus, our Gaussian mixture model theory motivates the use of transformer-like neural networks for denoising.

Overall, the learned denoiser architecture corresponds to an (implicit parametric) encoding scheme of the given data, since it can be used to denoise/project onto the data distribution. Training a denoiser is equivalent to finding a better coding scheme, and this partially implements approach A2 (i.e., pursuing a distribution by searching for schemes with better coding rates) at the end of Section 3.1.3. In the sequel, we discuss how to implement the other approach.

3.2.2 Sampling a Distribution via Iterative Denoising

Remember that at the end of Section 3.1.3, we discussed two approaches for pursuing a distribution with low-dimensional structure. The first approach (A1)

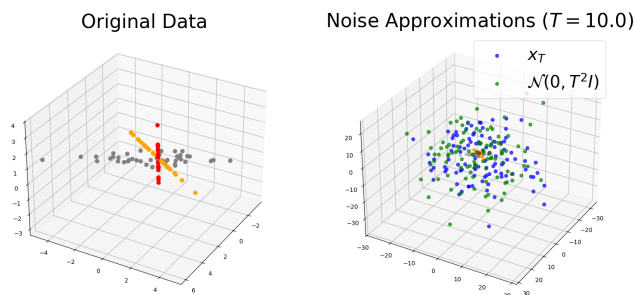


Figure 3.7: **Visualizing \mathbf{x}_T versus $\mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$.** *Left:* A plot of Gaussian mixture model data \mathbf{x} . *Right:* A plot of \mathbf{x} as well as \mathbf{x}_T and an independent sample of $\mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$, for $T = 10$. On the right plot, \mathbf{x} is plotted in the same colors as the left; however, samples from \mathbf{x}_T and $\mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$ are both much larger, on average, than samples from \mathbf{x} , and so it appears much smaller because of the scaling. Despite this large blow-up, we clearly observe the similarities in the distributions of \mathbf{x}_T and $\mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$.

starts with a normal distribution of high entropy and gradually reduces its entropy until it reaches the data distribution. We call this procedure *sampling*, since we generate new samples. It is now time to discuss how to implement this with the toolkit we have built.

We know how to denoise very noisy samples \mathbf{x}_T to obtain approximations $\hat{\mathbf{x}}$ whose distributions resemble that of the target random variable \mathbf{x} . Yet the sampling procedure requires us to start from a template distribution that has *no* influence from the distribution of \mathbf{x} and to use the denoiser to guide the iterates toward the distribution of \mathbf{x} . How can we do this? One way is motivated as follows:

$$\frac{\mathbf{x}_T}{T} = \frac{\mathbf{x} + T\mathbf{g}}{T} = \frac{\mathbf{x}}{T} + \mathbf{g} \xrightarrow{T \rightarrow +\infty} \mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (3.2.81)$$

Thus, $\mathbf{x}_T \approx \mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$. This approximation is quite good for almost all practical distributions, and is visualized in Figure 3.7.

So, discretizing $[0, T]$ into $0 = t_0 < t_1 < \dots < t_L = T$ uniformly using $t_\ell = T\ell/L$ (as in the previous section), one possible way to sample from pure noise is:

- Sample $\hat{\mathbf{x}}_T \sim \mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$ (independently of everything else).
- Run the denoising iteration as in Section 3.2.1, i.e.,

$$\hat{\mathbf{x}}_{t_{\ell-1}} = \left(1 - \frac{1}{\ell}\right) \hat{\mathbf{x}}_{t_\ell} + \frac{1}{\ell} \bar{\mathbf{x}}^*(t_\ell, \hat{\mathbf{x}}_{t_\ell}). \quad (3.2.82)$$

- Output $\hat{\mathbf{x}} = \hat{\mathbf{x}}_0$.

This is conceptually all there is behind *diffusion models*, which transform noise into data samples in accordance with the first desideratum. However, a few

steps remain before we obtain models that can actually sample from real data distributions like images under practical resource constraints. In the sequel, we introduce and motivate several such steps.

Step 1: Different discretizations. The first step is motivated by the observation that *we do not need to spend so many denoising iterations* at large t . Figure 3.5 shows that the first 200 or 300 iterations out of the 500 iterations of the sampling process are spent contracting the noise toward the data distribution as a whole, before the remaining iterations push the samples toward a subspace. Given a fixed iteration count L , this suggests that we should allocate more timesteps t_ℓ near $t = 0$ than near $t = T$. During sampling (and training), we can therefore use an alternative discretization of $[0, T]$ into $0 \leq t_0 < t_1 < \dots < t_L \leq T$, such as an *exponential discretization*:

$$t_\ell = C_1(e^{C_2\ell} - 1), \quad \forall \ell \in \{0, 1, \dots, L\} \quad (3.2.83)$$

where $C_1, C_2 > 0$ are constants that can be tuned for optimal performance in practice; theoretical analyses often specify such optimal constants as well. The denoising/sampling iteration then becomes

$$\hat{\mathbf{x}}_{t_{\ell-1}} \doteq \frac{t_{\ell-1}}{t_\ell} \hat{\mathbf{x}}_{t_\ell} + \left(1 - \frac{t_{\ell-1}}{t_\ell}\right) \bar{\mathbf{x}}^*(t_\ell, \hat{\mathbf{x}}_{t_\ell}), \quad (3.2.84)$$

with $\hat{\mathbf{x}}_{t_L} \sim \mathcal{N}(\mathbf{0}, t_L^2 \mathbf{I})$.

Step 2: Different noise models. The second step is to consider slightly different models compared to (3.2.1). The basic motivation for this is as follows. In practice, the noise distribution $\mathcal{N}(\mathbf{0}, t_L^2 \mathbf{I})$ becomes an increasingly poor estimate of the true covariance in high dimensions, i.e., (3.2.81) becomes an increasingly poor approximation, especially with anisotropic high-dimensional data. The increased distance between $\mathcal{N}(\mathbf{0}, t_L^2 \mathbf{I})$ and the true distribution of \mathbf{x}_{t_L} may cause the denoiser to perform worse in such circumstances. Theoretically, \mathbf{x}_{t_L} never converges to any distribution as t_L increases, so this setup is difficult to analyze end-to-end. In this case, our remedy is to *simultaneously add noise and shrink the contribution of \mathbf{x} , such that \mathbf{x}_T converges as $T \rightarrow \infty$* . The rate of added noise is denoted $\sigma: [0, T] \rightarrow \mathbb{R}_{\geq 0}$, and the rate of shrinkage is denoted $\alpha: [0, T] \rightarrow \mathbb{R}_{\geq 0}$, such that σ is *increasing* and α is (not strictly) *decreasing*, and

$$\mathbf{x}_t \doteq \alpha_t \mathbf{x} + \sigma_t \mathbf{g}, \quad \forall t \in [0, T]. \quad (3.2.85)$$

The previous setup has $\alpha_t = 1$ and $\sigma_t = t$, and this is called the *variance-exploding (VE) process*. A popular choice that decreases the contribution of \mathbf{x} , as we described originally, has $T = 1$ (so that $t \in [0, 1]$), $\alpha_t = \sqrt{1 - t^2}$ and $\sigma_t = t$; this is the *variance-preserving (VP) process*. Note that under the VP process, $\mathbf{x}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ exactly, so we can just sample from this standard

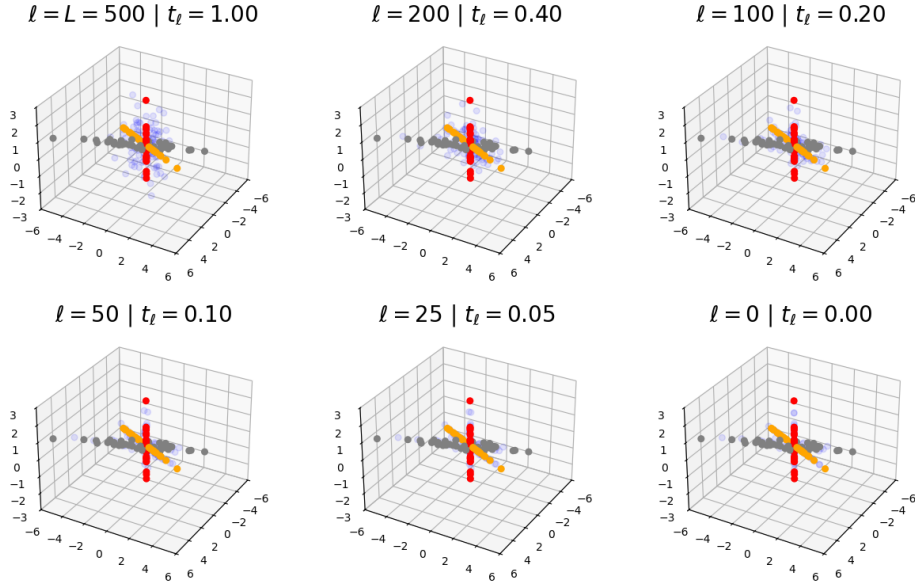


Figure 3.8: **Denoising a mixture of Gaussians using the VP diffusion process.** We use the same figure setup and data distribution as Figure 3.5. Note that compared to Figure 3.5, the noise distribution is much more concentrated around the origin.

distribution and iteratively denoise. As a result, the VP process is much easier to analyze theoretically and more stable empirically.⁶

With this more general setup, Tweedie’s formula (3.2.23) becomes

$$\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t] = \frac{1}{\alpha_t} (\mathbf{x}_t + \sigma_t^2 \nabla \log p_t(\mathbf{x}_t)). \quad (3.2.86)$$

The denoising iteration (3.2.84) becomes

$$\hat{\mathbf{x}}_{t_{\ell-1}} = \frac{\sigma_{t_{\ell-1}}}{\sigma_{t_{\ell}}} \hat{\mathbf{x}}_{t_{\ell}} + \left(\alpha_{t_{\ell-1}} - \frac{\sigma_{t_{\ell-1}}}{\sigma_{t_{\ell}}} \alpha_{t_{\ell}} \right) \bar{\mathbf{x}}^*(t_{\ell}, \hat{\mathbf{x}}_{t_{\ell}}). \quad (3.2.87)$$

Example 3.6 (Gaussian mixture model). For a Gaussian mixture model (see Example 3.2), its denoiser (3.2.19) becomes

$$\bar{\mathbf{x}}^*(t, \mathbf{x}_t) = \sum_{k=1}^K \left\{ \frac{\pi_k \varphi(\mathbf{x}_t; \alpha_t \boldsymbol{\mu}_k, \alpha_t^2 \boldsymbol{\Sigma}_k + \sigma_t^2 \mathbf{I})}{\sum_{i=1}^K \pi_i \varphi(\mathbf{x}_t; \alpha_t \boldsymbol{\mu}_i, \alpha_t^2 \boldsymbol{\Sigma}_i + \sigma_t^2 \mathbf{I})} \cdot (\boldsymbol{\mu}_k + \alpha_t \boldsymbol{\Sigma}_k (\alpha_t^2 \boldsymbol{\Sigma}_k + \sigma_t^2 \mathbf{I})^{-1} (\mathbf{x}_t - \alpha_t \boldsymbol{\mu}_k)) \right\}. \quad (3.2.88)$$

⁶Why use the whole α, σ setup? As we will see in Exercise 3.5, it encapsulates and unifies many proposed processes, including the recently popular so-called *flow matching* process. Despite this, the VE and VP processes are still the most popular empirically and theoretically (so far), and so we will consider them in this Section.

Figure 3.8 demonstrates iterations of the sampling procedure. ■

Note that the denoising iteration (3.2.87) gives a sampling algorithm called the DDIM (“Denoising Diffusion Implicit Model”) sampler [SME20], and is one of the most popular sampling algorithms used in diffusion models today. We summarize it in Algorithm 3.1.

Algorithm 3.1 Sampling using a denoiser.

Input: An ordered list of timesteps $0 \leq t_0 < \dots < t_L \leq T$ to use for sampling.

Input: A denoiser $\bar{\mathbf{x}}: \{t_\ell\}_{\ell=1}^L \times \mathbb{R}^D \rightarrow \mathbb{R}^D$.

Input: Scale and noise level functions $\alpha, \sigma: \{t_\ell\}_{\ell=0}^L \rightarrow \mathbb{R}_{\geq 0}$.

Output: A sample $\hat{\mathbf{x}}$, approximately from the distribution of \mathbf{x} .

```

1: function DDIMSAMPLER( $\bar{\mathbf{x}}, (t_\ell)_{\ell=0}^L$ )
2:   Initialize  $\hat{\mathbf{x}}_{t_L} \sim$  approximate distribution of  $\mathbf{x}_{t_L}$ 
   # VP  $\implies \mathcal{N}(\mathbf{0}, \mathbf{I})$ , VE  $\implies \mathcal{N}(\mathbf{0}, t_L^2 \mathbf{I})$ .

3:   for  $\ell = L, L - 1, \dots, 1$  do
4:     Compute

```

$$\hat{\mathbf{x}}_{t_{\ell-1}} \doteq \frac{\sigma_{t_{\ell-1}}}{\sigma_{t_\ell}} \hat{\mathbf{x}}_{t_\ell} + \left(\alpha_{t_{\ell-1}} - \frac{\sigma_{t_{\ell-1}}}{\sigma_{t_\ell}} \alpha_{t_\ell} \right) \bar{\mathbf{x}}(t_\ell, \hat{\mathbf{x}}_{t_\ell})$$

```

5:   end for
6:   return  $\hat{\mathbf{x}}_{t_0}$ 
7: end function

```

Step 3: Optimizing training pipelines. If we use the procedure dictated by Section 3.2.1 to learn a separate denoiser $\bar{\mathbf{x}}(t, \cdot)$ for each time t to be used in the sampling algorithm, *we would have to learn L separate denoisers!* This is highly inefficient—the usual case is that we have to train L separate neural networks, taking up L times the training time and storage memory, and then be locked into using these timesteps for sampling forever. Instead, we can *train a single neural network* to denoise across all times t , taking as input the continuous variables \mathbf{x}_t and t (instead of just \mathbf{x}_t before). Mechanically, our training loss averages over t , i.e., solves the following problem:

$$\min_{\theta} \mathbb{E} \|\bar{\mathbf{x}}_{\theta}(\tau, \mathbf{x}_{\tau}) - \mathbf{x}\|_2^2, \quad (3.2.89)$$

where $\tau \sim \mathcal{U}([0, T])$, i.e., is drawn uniformly from $[0, T]$. Similar to Step 1, where we used more timesteps closer to $t = 0$ to ensure a better sampling process, we may want to ensure that the denoiser is higher quality closer to $t = 0$, and thereby *weight the loss* so that t near 0 has higher weight.⁷ Letting w_t be the

⁷In practice, the distribution of sampled t may also be non-uniform, in order to ensure (in conjunction with the weighting scheme) that the Monte Carlo estimate for the loss has small variance, greatly improving training stability.

weight at time t , the weighted loss would look like

$$\min_{\theta} \mathbb{E}[w_{\tau} \|\bar{\mathbf{x}}_{\theta}(\tau, \mathbf{x}_{\tau}) - \mathbf{x}\|_2^2]. \quad (3.2.90)$$

One reasonable choice of weight in practice is $w_t = \alpha_t^2 / \sigma_t^2$. The precise reason will be covered in the next paragraph, but generally it serves to up-weight the losses corresponding to t near 0 while still remaining reasonably numerically stable. Also, of course, we cannot compute the expectation in practice, so we use the most straightforward Monte Carlo average to estimate it. The series of changes made here have several conceptual and computational benefits: we do not need to train multiple denoisers, we can train on one set of timesteps and sample using a subset (or others entirely), etc. The full pipeline is discussed in Algorithm 3.2. Certainly, this pipeline works in practice; we discuss some results in Sections 8.6 to 8.10.

Step 4: Changing the estimation target. It is common to instead reorient the whole denoising pipeline around *noise predictors*, i.e., estimates of $\mathbb{E}[\mathbf{g} \mid \mathbf{x}_t]$. In practice, noise predictors are slightly easier to train because their output is (almost) always of comparable size to a Gaussian random variable, so training is more numerically stable, but they may require larger models since they predict a full-dimensional target instead of potentially low-dimensional data [LH25]. By (3.2.85) we have

$$\mathbf{x}_t = \alpha_t \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t] + \sigma_t \mathbb{E}[\mathbf{g} \mid \mathbf{x}_t] \implies \mathbb{E}[\mathbf{g} \mid \mathbf{x}_t] = \frac{1}{\sigma_t} (\mathbf{x}_t - \alpha_t \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t]), \quad (3.2.91)$$

so any predictor for \mathbf{x} can be turned into a predictor for \mathbf{g} using the above relation, i.e.,

$$\bar{\mathbf{g}}(t, \mathbf{x}_t) = \frac{1}{\sigma_t} \mathbf{x}_t - \frac{\alpha_t}{\sigma_t} \bar{\mathbf{x}}(t, \mathbf{x}_t), \quad (3.2.92)$$

and vice-versa. Thus a good network for estimating $\bar{\mathbf{g}}$ is the same as a good network for estimating $\bar{\mathbf{x}}$ *plus a residual connection* (as seen in, e.g., transformers). Their losses are also the same as the denoiser, up to the factor of α_t / σ_t , i.e.,

$$\mathbb{E}[w_{\tau} \|\mathbf{g} - \bar{\mathbf{g}}(\tau, \mathbf{x}_{\tau})\|_2^2] = \mathbb{E}\left[w_{\tau} \frac{\alpha_{\tau}^2}{\sigma_{\tau}^2} \|\mathbf{x} - \bar{\mathbf{x}}(\tau, \mathbf{x}_{\tau})\|_2^2\right], \quad (3.2.93)$$

recalling that $\tau \sim \mathcal{U}([0, T])$. Other targets have been proposed for different tasks, e.g., $\mathbb{E}[\mathbf{v}_t \mid \mathbf{x}_t]$ where $\mathbf{v}_t = \frac{d}{dt} \mathbf{x}_t$ (called *v-prediction* or *velocity prediction*), etc. This turns out to be equivalent to denoising and noise prediction, but also fundamentally useful for generalizations of diffusion such as flow matching, which we will discuss in Remark 3.3.

Sampling Error Rate of the Denoising Process

We have made many changes to our original platonic noising/denoising process. To assure ourselves that the new process still works in practice, we can compute

Algorithm 3.2 Learning a denoiser from data.

Input: Dataset $\mathcal{D} \subseteq \mathbb{R}^D$.

Input: An ordered list of timesteps $0 \leq t_0 < \dots < t_L \leq T$ to use for sampling.

Input: A weighting function $w: \{t_\ell\}_{\ell=1}^L \rightarrow \mathbb{R}_{\geq 0}$.

Input: Scale and noise level functions $\alpha, \sigma: \{t_\ell\}_{\ell=0}^L \rightarrow \mathbb{R}_{\geq 0}$.

Input: A parameter space Θ and a denoiser architecture $\bar{\mathbf{x}}_\theta$.

Input: An optimization algorithm for the parameters.

Input: The number of optimization iterations M .

Input: The number of Monte Carlo draws N per training step.

Output: A trained denoiser $\bar{\mathbf{x}}_{\theta^*}$.

```

1: function TRAINDENOISER( $\mathcal{D}, \Theta$ )
2:   Initialize  $\theta^{(1)} \in \Theta$ 
3:   for  $m \in [M]$  do
4:     for  $i \in [N]$  do
5:       # Draw a sample from the dataset.
6:        $\mathbf{x}^{(m,i)} \sim \mathcal{D}$ 
7:
8:       # Sample a timestep.
9:        $t^{(m,i)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}(\{t_\ell\}_{\ell=1}^L)$ 
10:
11:      # Compute weights and scales.
12:       $(\alpha^{(m,i)}, \sigma^{(m,i)}, w^{(m,i)}) = (\alpha_{t^{(m,i)}}, \sigma_{t^{(m,i)}}, w_{t^{(m,i)}})$ 
13:
14:      # Sample a noise vector.
15:       $\mathbf{g}^{(m,i)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
16:
17:      # Compute the noised sample.
18:       $\mathbf{x}_t^{(m,i)} \doteq \alpha^{(m,i)} \mathbf{x}^{(m,i)} + \sigma^{(m,i)} \mathbf{g}^{(m,i)}$ 
19:
20:     end for
21:     # Compute the training loss.
22:      $\hat{\mathcal{L}}^{(m)} \doteq \frac{1}{N} \sum_{i=1}^N w^{(m,i)} \|\mathbf{x}^{(m,i)} - \bar{\mathbf{x}}_{\theta^{(m)}}(t^{(m,i)}, \mathbf{x}_t^{(m,i)})\|_2^2$ 
23:
24:     # Update the model parameters.
25:      $\theta^{(m+1)} \doteq \text{OptimUpdate}^{(m)}(\theta^{(m)}, \nabla_{\theta^{(m)}} \hat{\mathcal{L}}^{(m)})$ 
26:
27:   end for
28:   # Return the trained denoiser.
29:   return  $\bar{\mathbf{x}}_{\theta^{(M+1)}}$ 
30: end function

```

numerical examples (such as Figure 3.8). To assure ourselves that it is theoretically sound, we can prove a *bound on the error rate* for the sampling algorithm, which shows that the error rate is small. We now furnish such a rate from the literature, which shows that the output distribution of the sampler converges in the so-called *total variation (TV) distance* to the true distribution. The TV distance between two random variables \mathbf{x} and \mathbf{y} is defined as:⁸

$$\mathrm{TV}(\mathbf{x}, \mathbf{y}) \doteq \sup_{A \subseteq \mathbb{R}^d} |\mathbb{P}[\mathbf{x} \in A] - \mathbb{P}[\mathbf{y} \in A]|. \quad (3.2.94)$$

If \mathbf{x} and \mathbf{y} are very close (uniformly), then the supremum will be small. So the TV distance measures the closeness of random variables. (It is indeed a metric, as the name suggests; the proof is an exercise.)

Theorem 3.4 ([LY24] Theorem 1, Simplified). *Suppose that $\mathbb{E} \|\mathbf{x}\|_2 < \infty$. If \mathbf{x} is denoised according to the VP process with an exponential discretization⁹ as in (3.2.83), the output $\hat{\mathbf{x}}$ of Algorithm 3.1 satisfies the total variation bound*

$$\mathrm{TV}(\mathbf{x}, \hat{\mathbf{x}}) = \tilde{\mathcal{O}} \left(\underbrace{\frac{D}{L}}_{\text{discretization error}} + \underbrace{\sqrt{\frac{1}{L} \sum_{\ell=1}^L \frac{\alpha_{t_\ell}^2}{\sigma_{t_\ell}^4} \mathbb{E} \|\bar{\mathbf{x}}^*(t_\ell, \mathbf{x}_{t_\ell}) - \bar{\mathbf{x}}(t_\ell, \mathbf{x}_{t_\ell})\|_2^2}}_{\text{average excess error of the denoiser}} \right) \quad (3.2.95)$$

where $\bar{\mathbf{x}}^*$ is the Bayes optimal denoiser for \mathbf{x} , and $\tilde{\mathcal{O}}$ is a version of the big- \mathcal{O} notation that ignores logarithmic factors in L .

The very high-level proof technique is, as discussed earlier, to bound the error at each step, distinguish the error sources (between discretization and denoiser error), and carefully ensure that the errors do not accumulate too much (or even cancel out).

Note that if $L \rightarrow \infty$ and we correctly learn the Bayes-optimal denoiser $\bar{\mathbf{x}} = \bar{\mathbf{x}}^*$ (so the excess error is 0), then the sampling process in Algorithm 3.1 yields a *perfect (in distribution) inverse* of the noising process, since the error rate in Theorem 3.4 goes to 0,¹⁰ as argued heuristically above.

Remark 3.1. What if the data is low-dimensional, say supported on a *low-rank subspace* of the high-dimensional space \mathbb{R}^D ? If the data distribution is compactly supported—say if the data is normalized to the unit hypercube, which is often ensured as a preprocessing step for real data such as images—it is possible to

⁸The supremum here is only over *measurable* sets A , not absolutely all sets. As with all other measure-theoretic caveats in this book, feel free to ignore it for practical purposes.

⁹The precise definition is rather lengthy in our notation and only defined up to various absolute constants, so we omit it here for brevity. It is in the original paper [LY24].

¹⁰Similar results hold for VE processes; the theoretical toolkits used for the proofs differ because of the structure of the underlying stochastic processes.

do better. Namely, the authors of [LY24] also define a measure of *approximate intrinsic dimension* κ using the asymptotics of the so-called covering number, which is extremely similar in intuition (if not in implementation) to the rate-distortion function presented in the next Chapter. Then they show that using a particular small modification of the DDIM sampler in Algorithm 3.1 (i.e., slightly perturbing the update coefficients), the discretization error becomes $\tilde{\mathcal{O}}(\kappa/L)$ instead of $\tilde{\mathcal{O}}(D/L)$ as in Theorem 3.4. Therefore, using this modified algorithm, L does not have to be too large even as D reaches the thousands or millions, since real data have low-dimensional structure. However, in practice we use the DDIM sampler instead, so L should have a mild dependence on D to achieve consistent error rates. The exact choice of L trades off between the computational complexity (e.g., runtime or memory consumption) of sampling and the statistical complexity of learning a denoiser for low-dimensional structures. The value of L is often different at training time (where a larger L allows better coverage of the interval $[0, T]$, which helps the network learn a relationship that generalizes over t) and sampling time (where L being smaller means more efficient sampling). One can even pick the timesteps adaptively at sampling time to optimize this tradeoff [BLZ+22].

Remark 3.2. Various other works define the reverse process as moving backward in the time index t using an explicit difference equation, or differential equation in the limit $L \rightarrow \infty$, or forward in time using the transformation $\mathbf{y}_t = \mathbf{x}_{T-t}$, such that if t increases then \mathbf{y}_t becomes closer to \mathbf{x}_0 . Each paper uses its own notation, and it can sometimes be confusing to understand the relationships between different papers’ results. In this work we strive to keep consistency: we move forward in time to noise, and backward in time to denoise. If you are reading another work that is not clear on the time index, or trying to implement an algorithm that is similarly unclear, there is one way to disambiguate every time: the sampling process should always have a *positive* coefficient on the denoiser term (or equivalently the score function) when taking a step.

Remark 3.3. The diffusion/denoising processes and associated models provide a way to transport the data distribution to and from a high-entropy (Gaussian or Gaussian-like) distribution. On the other hand, it may be useful to transport the data distribution to and from an arbitrary target distribution, which could be Gaussian, a Gaussian mixture, or even another type of data distribution. If we have this desideratum, we can use the *flow matching* framework to achieve this. Suppose the target random variable is \mathbf{y} , and our flow is

$$\mathbf{x}_t = \alpha_t \mathbf{x} + \sigma_t \mathbf{y}, \quad \forall t \in [0, 1]. \quad (3.2.96)$$

such that $\alpha_0 = \sigma_1 = 1$ and $\alpha_1 = \sigma_0 = 0$, so that $\mathbf{x}_0 = \mathbf{x}$ and $\mathbf{x}_1 = \mathbf{y}$. Nothing about our strategy changes: we can train a denoiser $\bar{\mathbf{x}}_\theta$ to denoise \mathbf{x}_t back to \mathbf{x} (or equivalently a “noise” predictor to predict \mathbf{y}), and then apply our iterative denoising algorithms to start from a realization of \mathbf{y} to obtain \mathbf{x} . The only thing missing is Tweedie’s formula which connects the score $\nabla \log p_t$ of the denoiser to the denoiser; such a formula does not exist in general.

In practice, inspired by the concept of learning a so-called *transport map* between data pairs \mathbf{x} and \mathbf{y} , practitioners often wish to train a *velocity predictor* $\bar{\mathbf{v}}_\theta$ to predict $\mathbf{v}_t = \frac{d}{dt}\mathbf{x}_t = \alpha'_t\mathbf{x} + \sigma'_t\mathbf{y}$ instead of a denoiser $\bar{\mathbf{x}}_\theta$, or rather train $\bar{\mathbf{v}}_\theta$ to predict $\mathbb{E}[\mathbf{v}_t | \mathbf{x}_t]$. In this case, it is simplest to take the choice of coefficients $\alpha_t = 1-t$ and $\sigma_t = t$ (for $t \in [0, 1]$), so that the prediction target is $\mathbf{v}_t = \mathbf{y} - \mathbf{x}$ for each time t . Given the estimator $\bar{\mathbf{v}}_\theta$, the standard sampling algorithm follows the field backwards in time, i.e., starting from a realization of \mathbf{y} , we can iteratively update the iterate as follows:

$$\hat{\mathbf{x}}_{t_{\ell-1}} = \hat{\mathbf{x}}_{t_\ell} - \bar{\mathbf{v}}_\theta(t_\ell, \hat{\mathbf{x}}_{t_\ell})\Delta t, \quad (3.2.97)$$

where $\Delta t = t_\ell - t_{\ell-1}$ is the timestep. Note that velocity prediction is algebraically equivalent to denoising (see Exercise 3.6). See Sections 8.6 and 8.9 for practical applications of flow matching.

Remark 3.4. In practice, diffusion usually takes many steps and many function evaluations of the denoiser. It is therefore natural to ask: “can we speed up the sampling process by using a smaller number of steps and function evaluations?” The answer is yes, and such models are often called *consistency models* or, more descriptively, *few-step models*. The core idea is to use a denoiser that takes in two times $\bar{\mathbf{x}}_\theta(s, t, \mathbf{x}_t)$ and estimates $\mathbb{E}[\mathbf{x}_s | \mathbf{x}_t]$, which in some sense is an easier task, or alternatively to use a velocity predictor $\bar{\mathbf{v}}_\theta(s, t, \mathbf{x}_t)$ and estimate the average velocity $\mathbb{E}[(\mathbf{x}_t - \mathbf{x}_s)/(t - s) | \mathbf{x}_t]$. Such denoisers have additional consistency conditions such as $\bar{\mathbf{x}}_\theta(u, t, \bar{\mathbf{x}}_\theta(s, u, \boldsymbol{\xi})) = \bar{\mathbf{x}}_\theta(s, t, \boldsymbol{\xi})$, which motivate auxiliary training losses. Once such a predictor is learned, the sampling process becomes very simple and efficient, as it only requires a few steps of the form

$$\hat{\mathbf{x}}_s = \bar{\mathbf{x}}_\theta(s, t, \hat{\mathbf{x}}_t) = \hat{\mathbf{x}}_t - \bar{\mathbf{v}}_\theta(s, t, \hat{\mathbf{x}}_t) \cdot (t - s). \quad (3.2.98)$$

Relevant formulations include flow maps [BAV25] and mean flows [GDB+25].

3.3 Memorization, Generalization, and Coding Rates

The calculation presented at the end of Section 3.2.1 seems to suggest (loosely speaking) that in practice, using a transformer-like network is a good choice for learning or approximating a denoiser. This is reasonable, but what is the problem with using any old neural network (such as a multi-layer perceptron (MLP)) and simply scaling it up to infinity? To answer this question, let us consider a constraint we have ignored or glossed over so far in this Chapter: we only have finite data with which to learn the denoiser. Indeed, let us consider the training problem (3.2.90), except this time we assume that our data \mathbf{x} are drawn from the empirical distribution $p^{\mathbf{X}}$ on N training samples $\mathbf{X} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}]$. Because $p^{\mathbf{X}}$ is supported on a 0-dimensional subset of the ambient space \mathbb{R}^D , we have that $h(p^{\mathbf{X}}) = -\infty$, so in some sense $p^{\mathbf{X}}$ is the simplest (i.e., lowest-entropy) distribution that could generate our observations \mathbf{X} . In short, this

assumption arguably *has the least inductive bias* about the data distribution! Many popular commentaries on deep learning and generative models would have you believe that, due to this property, the best way to solve the generative modeling problem is to scale a generic denoiser as large as possible and reap the corresponding benefits. However, this advice is (provably!) wrong in this setting. To understand why, let us note that the empirical distribution $p^{\mathbf{X}}$ is actually a mixture of Gaussians with zero covariances:

$$p^{[\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}]} = \frac{1}{N} \sum_{i=1}^N \mathcal{N}(\mathbf{x}^{(i)}, \mathbf{0}). \quad (3.3.1)$$

As such, we know that the (globally) optimal solution to the denoising problem (3.2.90) with \mathbf{x} drawn from the empirical distribution, written out as follows:

$$\min_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \mathbb{E}[w_\tau \|\mathbf{x}^{(i)} - \bar{\mathbf{x}}_\theta(\tau, \mathbf{x}_\tau^{(i)})\|_2^2] \quad (3.3.2)$$

where $\mathbf{x}_\tau^{(i)}$ is the noised version of the training sample $\mathbf{x}^{(i)}$. The optimizer of this loss over all (square-integrable) denoisers $\bar{\mathbf{x}}$ is the associated optimal Gaussian denoiser from Example 3.3, i.e.,

$$\bar{\mathbf{x}}^*(t, \mathbf{x}_t) = \sum_{i=1}^N \frac{\exp(-\|\mathbf{x}_t - \alpha_t \mathbf{x}^{(i)}\|_2^2 / (2\sigma_t^2))}{\sum_{j=1}^N \exp(-\|\mathbf{x}_t - \alpha_t \mathbf{x}^{(j)}\|_2^2 / (2\sigma_t^2))} \mathbf{x}^{(i)}. \quad (3.3.3)$$

This is a convex combination of the data $\mathbf{x}^{(i)}$, and the coefficients get “sharper” (i.e., closer to 0 or 1) as $t \rightarrow 0$. In particular, for small t , the denoiser is almost exactly equal to the closest sample in the training set, and the update moves the iterate closer to the closest point, such that at the end of sampling, the iterate $\hat{\mathbf{x}}$ is (nearly) equal to a training sample $\mathbf{x}^{(i)}$. Let us re-emphasize this point: *if we use the denoiser (3.3.3) and a good sampler, our samples will re-generate the training data.* This is confirmed by our theoretical guarantee on sampling Theorem 3.4, which, when instantiated in this context, says that if we take the number of sampling steps to ∞ and use the prescribed timestep schedule, the diffusion model associated with the denoiser (3.3.3) will always sample from the empirical distribution and thus always reproduce a training sample. This phenomenon is a particularly strict instantiation of *memorization* in generative modeling.¹¹

Since in practice we only have a fixed finite dataset on which we can train our denoiser, the learning problem we aim to solve in this case is the finite-sample denoising loss (3.3.2) (compared to the infinite-sample denoising loss (3.2.90), for example). As such, *if we have an arbitrarily large model and optimize it as well as possible, it will learn to memorize and re-generate the training set.* In other words, this very large-scale model is not much more interesting than a simple database of the training data.

¹¹Memorization is a much broader phenomenon whose exact characterization depends on the domain the model is trained on (e.g., text, images, videos, audio, etc.). Later in this Section we will discuss more heuristic measures of memorization.

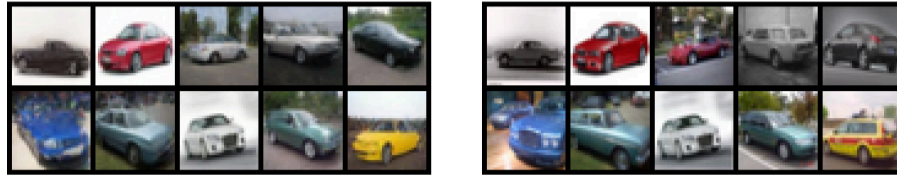


Figure 3.9: **Generated samples (left) from a diffusion model trained on 6,000 samples and (right) their closest points in the training dataset (in the Euclidean norm).** We observe that the generated samples are essentially pixel-perfect equivalents to the most similar training data. Both sets of images are from [YCK+23].

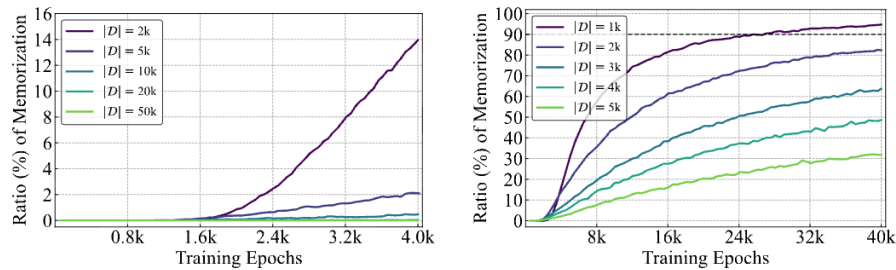


Figure 3.10: **Memorization versus dataset size for a fixed model, on (left) 4,000 training epochs and (right) 40,000 training epochs with smaller datasets.** We observe from the left plot that large enough datasets will never be memorized (e.g., 10,000 to 50,000 samples), whereas small datasets (1,000 to 5,000 samples) will be memorized after enough training. This means that as the dataset becomes larger relative to the model, the memorization ratio at the end of training goes from near-1 to near-0 relatively rapidly, indicating the possible existence of a phase transition from memorization to generalization as the dataset size increases. Figure taken from [GDP+23].

The above argument holds for arbitrarily powerful models that can fit the (true) memorizing denoiser. It is therefore relevant to ask whether it actually bears out in practice: do very large models (relative to the dataset) learn to memorize the training data? It turns out that this is empirically true. One of the first studies on this was [YCK+23]. To understand their results, we define the memorization ratio as the probability of memorization:

$$\text{MemRatio}(\bar{\mathbf{x}}) = \mathbb{P}_{\hat{\mathbf{x}} \sim \text{DM}(\bar{\mathbf{x}})}[\hat{\mathbf{x}} \text{ is memorized}], \quad (3.3.4)$$

where $\text{DM}(\bar{\mathbf{x}})$ is the distribution induced by the diffusion model that uses denoiser $\bar{\mathbf{x}}$, and the sample $\hat{\mathbf{x}}$ is memorized using the following criterion:

$$\hat{\mathbf{x}} \text{ is memorized} \iff \frac{\|\hat{\mathbf{x}} - \mathbf{x}_{(1)}(\hat{\mathbf{x}})\|_2}{\|\hat{\mathbf{x}} - \mathbf{x}_{(2)}(\hat{\mathbf{x}})\|_2} < \frac{1}{3}, \quad (3.3.5)$$

where $\mathbf{x}_{(i)}(\hat{\mathbf{x}})$ is the i^{th} closest point to $\hat{\mathbf{x}}$ in the training dataset (in the Euclidean norm), and $1/3$ is a constant that is empirically shown to agree well with human perception of memorization. Notice that this definition exactly corresponds to our earlier intuitive definition of memorization as sampled points each being much closer to a particular training sample than all other training samples. They showed that this kind of memorization occurs in practice when the model is much larger (in terms of the number of parameters) relative to the data; Figure 3.9 demonstrates that in this case, generated samples are essentially pixel-perfect copies of training samples. Slightly later, [GDP+23] performed a careful study to characterize how various factors such as architectures and training time impact memorization, showing among other things that as the dataset size increases relative to a fixed model, the memorization ratio rapidly changes from near-1 to near-0, i.e., the model quickly stops memorizing. This indicates the possibility of a phase transition between memorization and generalization as the model becomes smaller relative to the data. [BPM+25] investigates the existence and theoretical mechanism of this phase transition in synthetic data drawn from a Gaussian mixture model, demonstrating that it exists and providing a partial characterization of the model size at which models will start or stop memorizing. Other theoretical work such as [GVM25; KG24; NZM+24; ZLL+25] also attempts to provide a theoretical characterization of the phase transition in terms of model and dataset size in different settings.

In order to empirically prove the existence of a phase transition between memorization and generalization in real data, [ZZL+24] rethought the definition of memorization using the following reasoning. In the very high-dimensional space of natural images, Euclidean distances between points become much more strict [WM22], in the sense that if two images are not pixel-perfect matches to each other, then their Euclidean distances will concentrate sharply around a fixed value (scaling with the images' dimension). This does not catch the case where certain *concepts* (such as copyrighted cartoon characters) are generated in novel positions, making it very difficult to form effective and practical memorization criteria around the Euclidean distance. To fix this, [ZZL+24] uses a particular function f_{SSCD} (a pre-trained neural network called a Self-Supervised

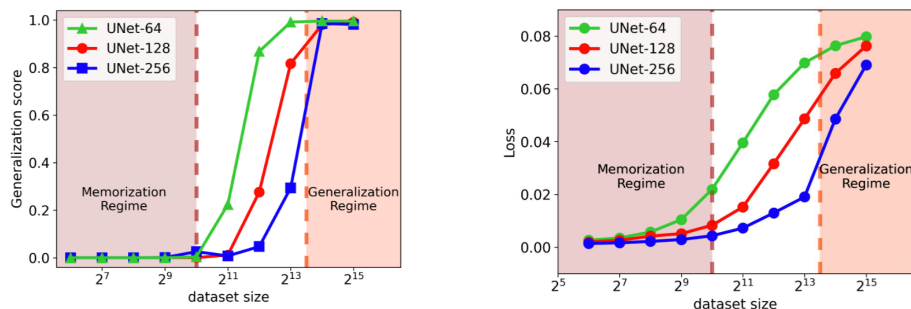


Figure 3.11: **The generalization score (left) and training loss (right) for a fixed model and increasing data**, from [ZZL+24]. The generalization score is 1– the fraction of samples that are memorized, and the different lines correspond to differently-sized denoiser models (UNet-64 is smaller than UNet-128, which in turn is smaller than UNet-256). The figures show that, as the dataset size becomes larger relative to the model size, the loss increases significantly and memorization stops occurring, showing that our theory and intuition apply to real models.

Copy Detector [PRR+22]) to map images into a Euclidean space, then takes the normalized inner product to obtain their (cosine) similarity, such that images that are similar due to this cosine similarity are claimed to be conceptually similar in the sense needed for memorization. Then, [ZZL+24] says that a sample $\hat{\mathbf{x}}$ is memorized if its similarity with any training point is $> 3/5$, i.e.,

$$\hat{\mathbf{x}} \text{ is memorized} \iff \max_{i \in [N]} \frac{f_{\text{SSCD}}(\hat{\mathbf{x}})^\top f_{\text{SSCD}}(\mathbf{x}^{(i)})}{\|f_{\text{SSCD}}(\hat{\mathbf{x}})\|_2 \|f_{\text{SSCD}}(\mathbf{x}^{(i)})\|_2} > \frac{3}{5} \quad (3.3.6)$$

Figure 3.11 shows the outcome of taking many (around 10,000) samples from diffusion models trained on differently-sized training datasets in terms of estimating memorization and computing the training loss. The outcome is in line with the above discussion; when the model is too large relative to the training data and the training loss is too small, memorization occurs and samples reproduce the training data. Thus, our intuition applies in practice: *for generative diffusion models, scaling the model is not all you need!*

So now the question is: what kind of denoiser should we use? The key is to choose a network architecture that can approximate the true denoiser (say, the one corresponding to a low-rank distribution as in (3.2.80)) without approximating the memorizing denoiser (3.3.3). Some work has explored this fine line and why modern diffusion models, which use Transformer- and convolution-based architectures, can memorize and generalize in different regimes [KG24; NZM+24; ZLL+25].

Let us at last return to the initial premise of this Chapter, described in Section 3.1.3: finding the minimal-entropy distribution that models the data is the only way to learn the distribution. The above discussion has hopefully shown that this claim is false. However, a related claim, also argued in Section 3.1.3

faint orange balls.¹³ For increasingly larger well-trained diffusion models, the perturbation radius ϵ shrinks to 0, and the model learns the empirical distribution. On the other hand, ϵ large enough that balls corresponding to different points overlap or are close may imply that the diffusion model generalizes.

Thus, our successful diffusion model can be viewed as minimizing a *lossy coding rate*, i.e., coding rate up to a precision ϵ (so as to be able to disregard the ϵ -balls around each point). This notion of the lossy coding rate turns out to be the right thing to minimize, morally speaking, and will be *extremely* useful. We will tease out its uses in Chapter 4 and throughout the rest of this book.

3.4 Summary and Notes

Key messages. In this Chapter, we have shown that, to pursue a general low-dimensional distribution, a universal computational approach is to start with a generic distribution and compress it toward the distribution through (iterative) noise or entropy reduction. As we discussed in the preceding Chapter, this is even the case for arguably the simplest distributions of a single subspace, such as computing PCA via the power-iteration algorithm. As we will see in future chapters, the fundamental roles of modern deep networks are not to fit arbitrary (noisy) data distributions or functions; instead, they are to realize or approximate those (lossy) compressing or denoising operators. We will derive precise forms of those operators and the associated deep architectures in future chapters.

More extensions and references. The use of denoising and diffusion for sampling has a rich history. The first work clearly about a diffusion model is probably [SWM+15], but before this there are many works on denoising as a computational and statistical problem. The most relevant of these is probably [Hyv05], which explicitly uses the score function to denoise and to perform independent component analysis. The most popular follow-ups are essentially co-occurring: [HJA20; SE19]. Since then, thousands of papers have built on diffusion models; we revisit this topic in Chapter 6.

Many of these works use a different stochastic process than the simple linear combination (3.2.85). In fact, all works listed above emphasize the need to add *independent* Gaussian noise at the beginning of each step of the forward process. Theoretically-minded work uses Brownian motion or stochastic differential equations to formulate the forward process [SSK+21]. However, since linear combinations of Gaussians still result in Gaussians, the *marginal distributions* of such processes still take the form of (3.2.85). Most of our discussion requires only that the marginal distributions are what they are, so our overly simplistic model is actually enough for almost everything. The only time marginal distributions are not enough is when we derive an expression for $\mathbb{E}[\mathbf{x}_s | \mathbf{x}_t]$ in terms

¹³In high dimensions, it is hypothesized that a better model for these perturbations is not being contained by balls but rather by much more anisotropic geometry-adaptive shapes [FPH+25], but the intuition remains similar.

of $\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t]$. Different (noising) processes give different such expressions, which can be used for sampling (and of course there are other ways to derive efficient samplers, such as the ever-popular DDPM sampler). The process in (3.2.85) is a bona fide stochastic process whose “natural” denoising iteration takes the form of the popular DDIM algorithm [SME20]. (Even this equivalence is not trivial; we cite [DGG+25] as justification.)

On top of the theoretical work [LY24] covered in Section 3.2.2, and the lineage of work that it builds on, which studies the *sampling* efficiency of diffusion models when the data has low-dimensional structure, there is a large body of work that studies the *training* efficiency of diffusion models when the data has low-dimensional structure. Specifically, Chen et al. [CHZ+23] and Wang et al. [WZZ+24] characterized the approximation and estimation error of denoisers when the data belongs to a mixture of low-rank Gaussians, showing that the number of training samples required to accurately learn the distribution scales with the intrinsic dimension of the data rather than the ambient dimension. There is considerable *methodological* work that attempts to use the low-dimensional structure of the data to perform various tasks with diffusion models. We highlight three: image editing [CZG+24], watermarking [LZQ24], and unlearning [CZL+25], though this list is inexhaustive.

3.5 Exercises and Extensions

Exercise 3.1. Show that (3.2.8) is the optimal solution of Problem (3.2.7).

Exercise 3.2. Consider random vectors $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{y} \in \mathbb{R}^d$, such that the pair $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{D+d}$ is jointly Gaussian. This means that

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_x & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{xy}^\top & \boldsymbol{\Sigma}_y \end{bmatrix} \right),$$

where the mean and covariance parameters are given by

$$\boldsymbol{\mu}_x = \mathbb{E}[\mathbf{x}], \quad \boldsymbol{\mu}_y = \mathbb{E}[\mathbf{y}], \quad \begin{bmatrix} \boldsymbol{\Sigma}_x & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{xy}^\top & \boldsymbol{\Sigma}_y \end{bmatrix} = \mathbb{E} \left[\begin{bmatrix} \mathbf{x} - \mathbb{E}[\mathbf{x}] \\ \mathbf{y} - \mathbb{E}[\mathbf{y}] \end{bmatrix} \begin{bmatrix} \mathbf{x} - \mathbb{E}[\mathbf{x}] \\ \mathbf{y} - \mathbb{E}[\mathbf{y}] \end{bmatrix}^\top \right]$$

Assume that $\boldsymbol{\Sigma}_y$ is positive definite (hence invertible); then positive semidefiniteness of the covariance matrix is equivalent to the Schur complement condition $\boldsymbol{\Sigma}_x - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_y^{-1} \boldsymbol{\Sigma}_{xy}^\top \succeq \mathbf{0}$.

In this exercise, we will prove that the conditional distribution $p_{\mathbf{x}|\mathbf{y}}$ is Gaussian: namely,

$$p_{\mathbf{x}|\mathbf{y}} \sim \mathcal{N} \left(\boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_y^{-1} (\mathbf{y} - \boldsymbol{\mu}_y), \boldsymbol{\Sigma}_x - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_y^{-1} \boldsymbol{\Sigma}_{xy}^\top \right). \quad (3.5.1)$$

A direct path to prove this result manipulates the defining ratio of densities $p_{\mathbf{x},\mathbf{y}}/p_{\mathbf{y}}$. We sketch an algebraically concise argument of this form below.

1. Verify the following matrix identity for the covariance:

$$\begin{bmatrix} \Sigma_x & \Sigma_{xy} \\ \Sigma_{xy}^\top & \Sigma_y \end{bmatrix} = \begin{bmatrix} I_D & \Sigma_{xy}\Sigma_y^{-1} \\ \mathbf{0} & I_d \end{bmatrix} \begin{bmatrix} \Sigma_x - \Sigma_{xy}\Sigma_y^{-1}\Sigma_{xy}^\top & \mathbf{0} \\ \mathbf{0} & \Sigma_y \end{bmatrix} \begin{bmatrix} I_D & \mathbf{0} \\ \Sigma_y^{-1}\Sigma_{xy}^\top & I_d \end{bmatrix}. \quad (3.5.2)$$

One arrives at this identity by performing two rounds of (block) Gaussian elimination on the covariance matrix.

2. Based on the previous identity, show that

$$\begin{bmatrix} \Sigma_x & \Sigma_{xy} \\ \Sigma_{xy}^\top & \Sigma_y \end{bmatrix}^{-1} = P \begin{bmatrix} (\Sigma_x - \Sigma_{xy}\Sigma_y^{-1}\Sigma_{xy}^\top)^{-1} & \mathbf{0} \\ \mathbf{0} & \Sigma_y^{-1} \end{bmatrix} P^\top, \quad (3.5.3)$$

where

$$P = \begin{bmatrix} I_D & \mathbf{0} \\ -\Sigma_y^{-1}\Sigma_{xy}^\top & I_d \end{bmatrix}, \quad (3.5.4)$$

and whenever the relevant inverses are defined.¹⁴ Conclude that

$$\begin{bmatrix} \mathbf{x} - \boldsymbol{\mu}_x \\ \mathbf{y} - \boldsymbol{\mu}_y \end{bmatrix}^\top \begin{bmatrix} \Sigma_x & \Sigma_{xy} \\ \Sigma_{xy}^\top & \Sigma_y \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{x} - \boldsymbol{\mu}_x \\ \mathbf{y} - \boldsymbol{\mu}_y \end{bmatrix} \quad (3.5.5)$$

$$= \begin{bmatrix} \mathbf{x} - \mathbf{m} \\ \mathbf{y} - \boldsymbol{\mu}_y \end{bmatrix}^\top \begin{bmatrix} (\Sigma_x - \Sigma_{xy}\Sigma_y^{-1}\Sigma_{xy}^\top)^{-1} & \mathbf{0} \\ \mathbf{0} & \Sigma_y^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{x} - \mathbf{m} \\ \mathbf{y} - \boldsymbol{\mu}_y \end{bmatrix}, \quad (3.5.6)$$

where $\mathbf{m} = \boldsymbol{\mu}_x + \Sigma_{xy}\Sigma_y^{-1}(\mathbf{y} - \boldsymbol{\mu}_y)$.

3. By dividing $p_{x,y}/p_y$, prove Equation (3.5.1). (*Hint: Using the previous identities, only minimal algebra should be necessary. For the normalizing constant, use Equation (3.5.3) to factor the determinant similarly.*)

Exercise 3.3. Show the Sherman–Morrison–Woodbury identity, i.e., for matrices \mathbf{A} , \mathbf{C} , \mathbf{U} , \mathbf{V} such that \mathbf{A} , \mathbf{C} , and $\mathbf{A} + \mathbf{UCV}$ are invertible,

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}\mathbf{VA}^{-1}. \quad (3.5.7)$$

Exercise 3.4. Rederive the following, assuming \mathbf{x}_t follows the generalized noise model (3.2.85).

- Tweedie's formula: (3.2.86).
- The DDIM iteration: (3.2.87).
- The Bayes-optimal denoiser for a Gaussian mixture model: (3.2.88).

Exercise 3.5.

¹⁴In cases where the Schur complement term is not invertible, the same result holds with its inverse replaced by the Moore–Penrose pseudoinverse. In particular, the conditional distribution (3.5.1) becomes a degenerate Gaussian distribution.

1. Implement the formulae derived in Exercise 3.4, building a sampler for Gaussian mixtures.
2. Reproduce Figure 3.4 and Figure 3.8.
3. We now introduce *Flow Matching (FM)* (see Remark 3.3), as follows:

$$\alpha_t = 1 - t, \quad \sigma_t = t. \quad (3.5.8)$$

Implement this process using the same framework, and test it for sampling in high dimensions. Which process gives better or more stable results?

Exercise 3.6. Consider the somewhat familiar interpolation

$$\mathbf{x}_t = \alpha_t \mathbf{x} + \sigma_t \mathbf{y}, \quad \forall t \in [0, T]. \quad (3.5.9)$$

where \mathbf{y} is a (not necessarily Gaussian) random variable. We define the *velocity* as

$$\mathbf{v}_t = \frac{d}{dt} \mathbf{x}_t = \alpha'_t \mathbf{x} + \sigma'_t \mathbf{y}, \quad \forall t \in [0, T]. \quad (3.5.10)$$

Compute an algebraic relationship writing $\mathbb{E}[\mathbf{v}_t | \mathbf{x}_t]$ in terms of $\mathbb{E}[\mathbf{x} | \mathbf{x}_t]$ and vice versa, showing that velocity prediction is algebraically equivalent to denoising and thus also to noise prediction (see Remark 3.3).

