

Chapter 7

Inference with Low-Dimensional Distributions

“Mathematics is the art of giving the same name to different things.”

— Henri Poincaré

In the previous Chapters of this book, we have studied how to effectively and efficiently learn a representation for a variable \boldsymbol{x} in the world with a distribution $p(\boldsymbol{x})$ that has a low-dimensional support in a high-dimensional space. So far, we have mainly developed the methodology for learning representation and autoencoding in a general, distribution or task-agnostic fashion. With such a learned representation, one can already use it to perform some generic and basic tasks such as classification (if the encoding is supervised with the class) and generation of random samples that have the same distribution as the given data (say natural images or natural languages).

More generally, however, the universality and scalability of the theoretical and computational framework presented in this book has enabled us to learn the distribution of a variety of important real-world high-dimensional data such as natural languages, human poses, natural images, videos, and even 3D scenes. Once the intrinsically rich and low-dimensional structures of these real data can be learned and represented correctly, they start to enable a broad family of powerful, often seemingly miraculous, tasks. Hence, from here onwards, we will start to show how to connect and tailor the general methods presented in previous Chapters to learn useful representations for specific structured data distributions and for many popular tasks in modern practice of machine intelligence.

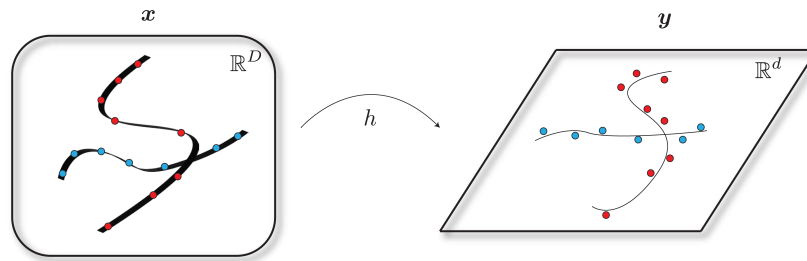


Figure 7.1: **Inference with low-dimensional distributions.** This is the generic picture for this chapter: we have a low-dimensional distribution for $\mathbf{x} \in \mathbb{R}^D$ (here depicted as a union of two 2-dimensional manifolds in \mathbb{R}^3) and a measurement model $\mathbf{y} = h(\mathbf{x}) + \mathbf{w} \in \mathbb{R}^d$. We want to infer various things about this model, including the conditional distribution of \mathbf{x} given \mathbf{y} , or the conditional expectation $\mathbb{E}[\mathbf{x} | \mathbf{y}]$, given various information about the model and (potentially finite) samples of either \mathbf{x} or \mathbf{y} .

7.1 Bayesian Inference and Constrained Optimization

7.1.1 Bayesian Inference with Low-Dimensional Distributions

Leveraging low-dimensionality for stable and robust inference. Generally speaking, a good representation or autoencoding should enable us to utilize the learned low-dimensional distribution of the data \mathbf{x} and its representation \mathbf{z} for various subsequent classification, estimation, and generation tasks under different conditions. As we have alluded to earlier in Chapter 1 Section 1.2.2, the importance of the *low-dimensionality* of the distribution is the key for us to conduct stable and robust inference related to the data \mathbf{x} , as illustrated by the few simple examples in Figure 1.11, from incomplete, noisy, and even corrupted observations. As it turns out, the very same concept carries over to real-world high-dimensional data whose distributions have a low-dimensional support, such as natural images and languages.

Despite a dazzling variety of applications in the practice of machine learning with data such as languages, images, videos and many other modalities, almost all practical applications can be viewed as a special case of the following inference problem: given an observation \mathbf{y} that depends on \mathbf{x} , say

$$\mathbf{y} = h(\mathbf{x}) + \mathbf{w}, \quad (7.1.1)$$

where $h(\cdot)$ represents measurements of a part of \mathbf{x} or certain observed attributes and \mathbf{w} represents some measurement noise and even (sparse) corruptions, solve the “inverse problem” of obtaining a most likely estimate $\hat{\mathbf{x}}(\mathbf{y})$ of \mathbf{x} or generating a sample $\hat{\mathbf{x}}$ that is at least consistent with the observation $\mathbf{y} \approx h(\hat{\mathbf{x}})$. Figure 7.1 illustrates the general relationship between \mathbf{x} and \mathbf{y} .

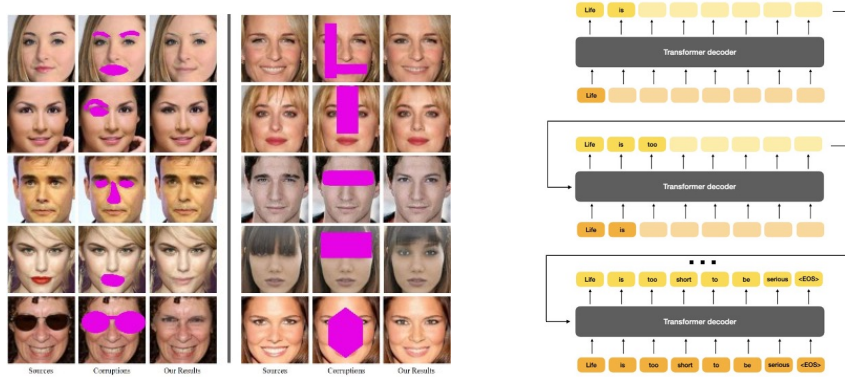


Figure 7.2: Left: image completion. Right: text prediction. In particular, text prediction is the inspiration for the popular Generative Pre-trained Transformer (GPT).

Example 7.1 (Image Completion and Text Prediction). The popular natural image completion and natural language prediction are two typical tasks that require us to recover a full data \mathbf{x} from its partial observations \mathbf{y} , with parts of \mathbf{x} masked out and to be completed based on the rest. Figure 7.2 shows some examples of such tasks. In fact, it is precisely these tasks which have inspired how to train modern large models for text generation (such as GPT) and image completion (such as the masked autoencoder) that we will study in greater detail later. ■

Statistical interpretation via Bayes' rule. Generally speaking, to accomplish such tasks well, we need to get ahold of the conditional distribution $p(\mathbf{x} | \mathbf{y})$. If we had this, then we would be able to find the maximal likelihood estimate (prediction):

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} p(\mathbf{x} | \mathbf{y}); \quad (7.1.2)$$

or compute the conditional expectation estimate:

$$\hat{\mathbf{x}} = \mathbb{E}[\mathbf{x} | \mathbf{y}] = \int \mathbf{x} p(\mathbf{x} | \mathbf{y}) d\mathbf{x}; \quad (7.1.3)$$

or sample from the conditional distribution:

$$\hat{\mathbf{x}} \sim p(\mathbf{x} | \mathbf{y}). \quad (7.1.4)$$

Notice that if the conditional distribution $p(\mathbf{x} | \mathbf{y})$ has a low-dimensional support that is nonlinear, these three different estimates can be rather different, as illustrated in Figure 7.3. Conceptually, the maximum a posteriori (MAP) estimate is the most desired one—it is the sample of the highest probability

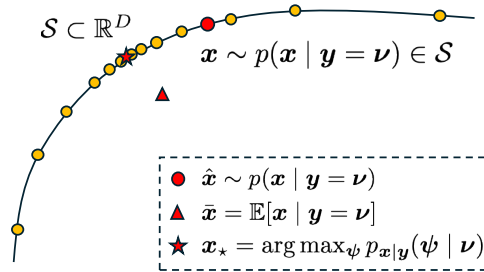


Figure 7.3: Comparison of three different surrogates for an estimate of \mathbf{x} given \mathbf{y} .

to have produced the observation \mathbf{y} , but it is typically the most expensive to compute.

Notice that from Bayes' rule, we have

$$p(\mathbf{x} | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{x})p(\mathbf{x})}{p(\mathbf{y})}. \quad (7.1.5)$$

For instance, the maximal likelihood estimate can be computed by solving the following (maximal log likelihood) program:

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} [\log p(\mathbf{y} | \mathbf{x}) + \log p(\mathbf{x})], \quad (7.1.6)$$

say via gradient ascent:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \cdot (\nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{x})). \quad (7.1.7)$$

Efficiently computing the conditional distribution $p(\mathbf{x} | \mathbf{y})$ naturally depends on how we learn and exploit the low-dimensional distribution $p(\mathbf{x})$ of the data \mathbf{x} and the observation model $\mathbf{y} = h(\mathbf{x}) + \mathbf{w}$ that determines the conditional distribution $p(\mathbf{y} | \mathbf{x})$.

Remark 7.1 (End-to-End versus Bayesian Inference). In the modern practice of data-driven machine learning, for certain popular tasks people often directly learn the conditional distribution $p(\mathbf{x} | \mathbf{y})$ or a (probabilistic) mapping or a regressor. Such a mapping is often modeled by some deep networks and trained end-to-end with sufficient paired samples (\mathbf{x}, \mathbf{y}) . Such an approach is very different from the above Bayesian approach in which both the distribution of $\mathbf{x} \sim p(\mathbf{x})$ and the (observation) mapping are needed. The benefit of the Bayesian approach is that the learned distribution $p(\mathbf{x})$ can facilitate many different tasks with varied observation models and conditions.

7.1.2 Constrained Optimization with Submanifolds

Geometric interpretation as constrained optimization. As the support $\mathcal{S}_{\mathbf{x}}$ of the distribution of \mathbf{x} is low-dimensional, we may assume that there exists

a function F such that

$$F(\mathbf{x}) = 0 \iff \mathbf{x} \in \mathcal{S}_{\mathbf{x}} \quad (7.1.8)$$

such that $\mathcal{S}_{\mathbf{x}} = F^{-1}(\{0\})$ is the low-dimensional support of the distribution $p(\mathbf{x})$. Geometrically, one natural choice of $F(\mathbf{x})$ is the “distance function” to the support $\mathcal{S}_{\mathbf{x}}$:

$$F(\mathbf{x}) = \min_{\mathbf{x}_p \in \mathcal{S}_{\mathbf{x}}} \|\mathbf{x} - \mathbf{x}_p\|_2. \quad (7.1.9)$$

Notice that, in reality, we only have discrete samples on the support of the distribution. In the same spirit of continuation, through diffusion or lossy coding studied in Chapters 3 and 4, we may approximate the distance function as $F(\mathbf{x}) \approx \min_{\mathbf{x}_p \in \mathcal{C}_{\mathbf{x}}^\epsilon} \|\mathbf{x} - \mathbf{x}_p\|_2$ where $\mathcal{S}_{\mathbf{x}}$ is replaced by a covering $\mathcal{C}_{\mathbf{x}}^\epsilon$ of the samples with ϵ -balls. But if the analytical form of a simple distribution is given, sometimes $F(\mathbf{x})$ can be computed explicitly.

Example 7.2 (Distance to a Line in \mathbb{R}^3). Let us assume that the support of a low-dimensional distribution in \mathbb{R}^3 is the x_3 -axis. Then the distance function $F(\mathbf{x})$ is given by:

$$F(\mathbf{x}) = \sqrt{x_1^2 + x_2^2}. \quad (7.1.10)$$

■

To see how such a function plays an important role in exploiting the distribution, for simplicity, we will assume that, for the rest of the subsection, the distance function is already given.

Now given $\mathbf{y} = h(\mathbf{x}) + \mathbf{w}$, to solve for \mathbf{x} , we can solve the following constrained optimization problem:

$$\max_{\mathbf{x}} -\frac{1}{2} \|h(\mathbf{x}) - \mathbf{y}\|_2^2 \quad \text{s.t.} \quad F(\mathbf{x}) = 0. \quad (7.1.11)$$

Using the method of augmented Lagrange multipliers, we can solve the following unconstrained program:

$$\max_{\mathbf{x}} \left[-\frac{1}{2} \|h(\mathbf{x}) - \mathbf{y}\|_2^2 + \lambda F(\mathbf{x}) - \frac{\mu}{2} F(\mathbf{x})^2 \right] \quad (7.1.12)$$

for some constant Lagrange multiplier λ . This is equivalent to the following program:

$$\max_{\mathbf{x}} \left[\log \exp \left(-\frac{1}{2} \|h(\mathbf{x}) - \mathbf{y}\|_2^2 \right) + \log \exp \left(-\frac{\mu}{2} (F(\mathbf{x}) - \lambda/\mu)^2 \right) \right], \quad (7.1.13)$$

where $c \doteq \lambda/\mu$ can be viewed as a “mean” for the constraint function. As μ becomes large when enforcing the constraint via continuation¹, c becomes

¹In the same spirit of continuation in Chapter 3 where we obtained better approximations of our distribution by sending $\epsilon \rightarrow 0$, here we send $\mu \rightarrow \infty$. Larger values of μ will constrain F to take smaller and smaller values at the optimum, meaning that the optimum lies within a smaller and smaller neighborhood of the support $\mathcal{S}_{\mathbf{x}}$. Interestingly, the theory of Lagrange multipliers hints that, under certain benign conditions on F and other terms in the objective, we only need to make μ large enough in order to ensure $F(\mathbf{x}) = 0$ at the optimum, meaning that at *finite* penalty we get *perfect* approximation of the support. In general, we should have the intuition that μ plays the same role as ϵ^{-1} .

increasingly small. The above program may be interpreted in two different ways.

Firstly, one may view the first term as the conditional probability of \mathbf{y} given \mathbf{x} , and the second term as a probability density for \mathbf{x} :

$$p(\mathbf{y} | \mathbf{x}) \propto \exp\left(-\frac{1}{2}\|h(\mathbf{x}) - \mathbf{y}\|_2^2\right), \quad p(\mathbf{x}) \propto \exp\left(-\frac{\mu}{2}(F(\mathbf{x}) - c)^2\right). \quad (7.1.14)$$

Hence, solving the constrained optimization for the inverse problem is equivalent to conducting Bayes inference with the above probability densities. Hence solving the above program (7.1.13) via gradient ascent is equivalent to the above maximal likelihood estimate (7.1.7), in which the gradient takes the form:

$$\nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{x}) \quad (7.1.15)$$

$$= -(h(\mathbf{x}) - \mathbf{y})^\top \frac{\partial h}{\partial \mathbf{x}}(\mathbf{x}) - \mu(F(\mathbf{x}) - c) \frac{\partial F}{\partial \mathbf{x}}(\mathbf{x}), \quad (7.1.16)$$

where $\frac{\partial h}{\partial \mathbf{x}}(\mathbf{x})$ and $\frac{\partial F}{\partial \mathbf{x}}(\mathbf{x})$ are the Jacobian of $h(\mathbf{x})$ and $F(\mathbf{x})$, respectively.

Example 7.3 (Gradient of $F(\mathbf{x})$). For the distance function defined in Example 7.2, its gradient is given by

$$\frac{\partial F}{\partial \mathbf{x}}(\mathbf{x}) = \frac{1}{\sqrt{x_1^2 + x_2^2}} \begin{bmatrix} x_1 \\ x_2 \\ 0 \end{bmatrix}. \quad (7.1.17)$$

Notice that unlike regular smooth functions whose gradient is typically unique at a fixed point, the gradients of the distance function at a point in the support, say $(0, 0, x_3)$, can span an entire subspace complementary to the support. ■

Notice that the above (gradient)

$$-\mu(F(\mathbf{x}) - c) \frac{\partial F}{\partial \mathbf{x}}(\mathbf{x})$$

always points towards the low-dimensional support of the distribution. Hence the descent process can be viewed as a “denoising” process, studied in Chapter 3, that gradually enforces the data to be closer to the correct support. Notice that the above derivation suggests that the “step size” of the score $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ of the denoising process is

$$\mu|F(\mathbf{x}) - \lambda/\mu|$$

which for any fixed μ is nearly proportional to the distance of the point \mathbf{x} to the support.

Secondly, notice that solving the above program (7.1.13) with μ increasing is equivalent to:

$$\mathbf{x}_* = \arg \min_{\mathbf{x}} \frac{1}{2}\|h(\mathbf{x}) - \mathbf{y}\|_2^2 + \frac{\mu}{2}(F(\mathbf{x}) - \lambda/\mu)^2 \quad \text{as } \mu \uparrow \infty. \quad (7.1.18)$$

Due to the conspicuous quadratic form of the two terms in the above equation, they can also be interpreted as certain “energy” functions. Such a formulation

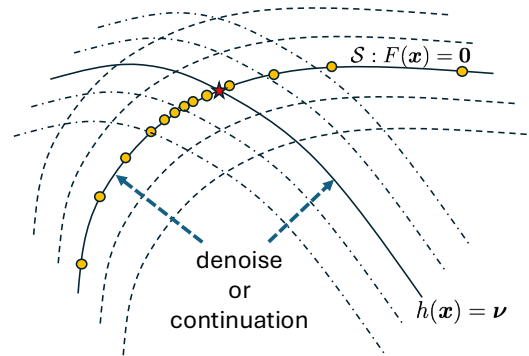


Figure 7.4: Illustration of the continuation process of enforcing the constraints $F(\mathbf{x}) = 0$ and $\mathbf{y} = h(\mathbf{x})$.

is often referred to as “Energy Minimization” in the machine learning literature, advocated by people like Yann LeCun [LCH+06]. Here the energy of a point \mathbf{x} is a quadratic function of its distance to the support of the desired distribution. Notice that in Chapter 3 and Chapter 4, we have argued that the notion of entropy is to measure the “volume” or uncertainty of the data distribution whereas here the energy depends on the “distance” of \mathbf{x} to the support of the distribution. As we minimize the above energy functions, the entropy (or uncertainty) of the feasible solutions reduces until it reaches (a feasible solution of) the optimal MAP estimate, as illustrated in Figure 7.4.

7.1.3 Representative Practical Settings for Inference

Notice that the above discussion and derivation are based on the assumption that we have perfect knowledge about the function $F(\mathbf{x})$ and the observation model $h(\mathbf{x})$. In practice, however, they may not be available at all and need to be “learned” from the data given. Hence to make the above conceptual solution truly computable, we need to deal with various situations in which information about the distributions of \mathbf{x} and the relationship between \mathbf{x} and \mathbf{y} are given or accessible in different ways and forms.

In general, they can mostly be categorized into four cases, which are, conceptually, increasingly more challenging:

- *Case 1:* Both a model for the distribution of \mathbf{x} and the observation model $\mathbf{y} = h(\mathbf{x})$ (possibly with added noise \mathbf{w}) are known, even with an analytical form. This is typically the case for many classic signal processing problems, such as signal denoising, the sparse vector recovery problem we saw in Chapter 2 and the low-rank matrix recovery problem to be introduced below.
- *Case 2:* We do not have a model for the distribution but only samples $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ of \mathbf{x} , and the observation model $\mathbf{y} = h(\mathbf{x})$ (possibly

with added noise \mathbf{w}) is known.² A model for the distribution $p(\mathbf{x})$ of \mathbf{x} needs to be learned, and subsequently the conditional distribution $p(\mathbf{x} | \mathbf{y})$. Natural image completion or natural language completion (e.g., BERT and GPT) are typical examples of this class of problems.

- *Case 3:* We only have some paired samples: $(\mathbf{X}, \mathbf{Y}) = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$ of the two variables (\mathbf{x}, \mathbf{y}) . The distributions of \mathbf{x} and \mathbf{y} and their relationship $h(\cdot)$ need to be learned from these paired sample data. For example, given many images and their captions, learning to conduct text-conditioned image generation is one such problem.
- *Case 4:* We only have the samples $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$ of the observations \mathbf{y} , and the observation model $h(\cdot)$ needs to be known, at least in some parametric family $h(\cdot, \boldsymbol{\theta})$. The distribution $p(\mathbf{x})$ and $p(\mathbf{x} | \mathbf{y})$ need to be learned from $\hat{\mathbf{x}}$, estimated from \mathbf{Y} . For example, learning to render a new view from a sequence of calibrated or uncalibrated views is one such problem.

In this chapter, we will discuss general approaches to learn the desired distributions and solve the associated conditional estimation or generation for these cases, typically with a representative practical problem. Throughout the chapter, the reader should keep Figure 7.1 in mind.

7.2 Conditional Inference with a Known Data Distribution

Notice that in the setting we have discussed in previous Chapters, the autoencoding network is trained to reconstruct a set of samples of the random vector \mathbf{x} . This would allow us to regenerate samples from the learned (low-dimensional) distribution. In practice, the low-dimensionality of the distribution, once given or learned, can be exploited for stable and robust recovery, completion, or prediction tasks. That is, under rather mild conditions, one can recover \mathbf{x} from highly compressive, partial, noisy or even corrupted measures of \mathbf{x} of the kind:

$$\mathbf{y} = h(\mathbf{x}) + \mathbf{w}, \quad (7.2.1)$$

where \mathbf{y} is typically an observation of \mathbf{x} that is of much lower dimension than \mathbf{x} and \mathbf{w} can be random noise or even sparse gross corruptions. This is a class of problems that have been extensively studied in the classical signal processing literature, for low-dimensional structures such as sparse vectors, low-rank matrices, and beyond. Interested readers may see [WM22] for a complete exposition of this topic.

Here to put the classic work in a more general modern setting, we illustrate the basic idea and facts through the arguably simplest task of data (and particularly image) completion. That is, we consider the problem of recovering a

²In the literature, this setting is sometimes referred to as the *empirical Bayesian inference*.

sample \mathbf{x} when parts of it are missing (or even corrupted). We want to recover or predict the rest of \mathbf{x} from observing only a fraction of it:

$$f : \mathcal{P}_\Omega(\mathbf{x}) \mapsto \hat{\mathbf{x}}, \quad (7.2.2)$$

where $\mathcal{P}_\Omega(\cdot)$ represents a masking operation (see Figure 7.5 for an example).

In this section and the next, we will study the completion task under two different scenarios: One is when the distribution of the data \mathbf{x} of interest is already given *a priori*, even in a certain analytical form. This is the case that prevails in classic signal processing where the structures of the signals are assumed to be known, for example, band-limited, sparse or low-rank. The other is when only raw samples of \mathbf{x} are available and we need to learn the low-dimensional distribution from the samples in order to solve the completion task well. This is the case for the tasks of natural image completion or video frame prediction. As a precursor to the rest of the chapter, we start with the simplest case of image completion: when the image to be completed can be well modeled as a low-rank matrix. We will move on to increasingly more general cases and more challenging settings later.

Low-rank matrix completion. The low-rank *matrix completion* problem is a classical problem for data completion when its distribution is low-dimensional and known. Consider a random sample of a matrix $\mathbf{X}_o = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{m \times n}$ from the space of all matrices of rank r . In general, we assume the rank of the matrix is

$$\text{rank}(\mathbf{X}_o) = r < \min\{m, n\}. \quad (7.2.3)$$

So it is clear that locally the intrinsic dimension of the space of all matrices of rank r is much lower than the ambient space mn .

Now, let Ω indicate a set of indices of observed entries of the matrix \mathbf{X}_o . Let the observed entries be:

$$\mathbf{Y} = \mathcal{P}_\Omega(\mathbf{X}_o). \quad (7.2.4)$$

The remaining entries supported on Ω^c are unobserved or missing. The problem is whether we can recover from \mathbf{Y} the missing entries of \mathbf{X} correctly and efficiently. Figure 7.5 shows one example of completing such a matrix.

Notice that the fundamental reason why such a matrix can be completed is that columns and rows of the matrix are highly correlated and they all lie on a low-dimensional subspace. For the example shown in Figure 7.5, the dimension or the rank of the matrix completed is only two. Hence the fundamental idea to recover such a matrix is to seek a matrix that has the lowest rank among all matrices that have entries agreeing with the observed ones:

$$\min_{\mathbf{X}} \text{rank}(\mathbf{X}) \quad \text{subject to} \quad \mathbf{Y} = \mathcal{P}_\Omega(\mathbf{X}). \quad (7.2.5)$$

This is known as the *low-rank matrix completion* problem. See [WM22] for a full characterization of the space of all low-rank matrices. As the rank function is discontinuous and rank minimization is in general an NP-hard problem, we would like to relax it with something easier to optimize.



Figure 7.5: Illustration of completing an image as low-rank matrix with some entries masked or corrupted. Left: the masked/corrupted image \mathbf{Y} ; middle: the mask Ω ; right: the completed image $\hat{\mathbf{X}}$.

Based on our knowledge about compression from Chapter 3, we could promote the low-rankness of the recovered matrix \mathbf{X} by enforcing the lossy coding rate (or the volume spanned by \mathbf{X}) of the data in \mathbf{X} to be small:

$$\min R_\epsilon(\mathbf{X}) = \frac{1}{2} \log \det (\mathbf{I} + \alpha \mathbf{X} \mathbf{X}^\top) \quad \text{subject to} \quad \mathbf{Y} = \mathcal{P}_\Omega(\mathbf{X}). \quad (7.2.6)$$

The problem can be viewed as a continuous relaxation of the above low-rank matrix completion problem (7.2.5) and it can be solved via gradient descent. One can show that the gradient descent operator for the logdet objective is precisely minimizing a close surrogate of the rank of the matrix $\mathbf{X} \mathbf{X}^\top$.

The rate distortion function is a nonconvex function, and its gradient descent does not always guarantee finding the globally optimal solution. Nevertheless, since the underlying structure sought for \mathbf{X} is piecewise linear, the rank function admits a rather effective convex relaxation: the nuclear norm—the sum of all singular values of the matrix \mathbf{X} . As shown in the compressive sensing literature, under fairly broad conditions,³ the matrix completion problem (7.2.5) can be effectively solved by the following convex program:

$$\min \|\mathbf{X}\|_* \quad \text{subject to} \quad \mathbf{Y} = \mathcal{P}_\Omega(\mathbf{X}), \quad (7.2.7)$$

where the nuclear norm $\|\mathbf{X}\|_*$ is the sum of singular values of \mathbf{X} . In practice, we often convert the above constrained convex optimization program to an unconstrained one:

$$\min \|\mathbf{X}\|_* + \lambda \|\mathbf{Y} - \mathcal{P}_\Omega(\mathbf{X})\|_F^2, \quad (7.2.8)$$

for some properly chosen $\lambda > 0$. Interested readers may refer to [WM22] for how to develop algorithms that can solve the above programs efficiently and effectively. Figure 7.5 shows a real example in which the matrix $\hat{\mathbf{X}}$ is actually recovered by solving the above program.

³Typically, such conditions specify the necessary and sufficient amount of entries needed for the completion to be computationally feasible. These conditions have been systematically characterized in [WM22].

Further extensions. It has been shown that images (or more accurately textures) and 3D scenes with low-rank structures can be very effectively completed via solving optimization programs of the above kind, even if there is additional corruption and distortion [LRZ+12; YZB+23; ZLG+10]:

$$\mathbf{Y} \circ \tau = \mathbf{X}_o + \mathbf{E}, \quad (7.2.9)$$

where τ is some unknown nonlinear distortion of the image and \mathbf{E} is an unknown matrix that models some (sparse) occlusion and corruption. Again, interested readers may refer to [WM22] for a more detailed account.

7.3 Conditional Inference with a Learned Data Representation

In the previous subsection, the reason we can infer \mathbf{x} from the partial observation \mathbf{y} is because (the support of) the distribution of \mathbf{X} is known or specified *a priori*, say as the set of all low-rank matrices. For many practical datasets, we do not have their distribution in an analytical form like the low-rank matrices, say the set of all natural images. Nevertheless, if we have sufficient samples of the data \mathbf{x} , we should be able to learn its low-dimensional distribution first and leverage it for future inference tasks based on an observation $\mathbf{y} = h(\mathbf{x}) + \mathbf{w}$. In this section, we assume the observation model $h(\cdot)$ is given and known. We will study the case when $h(\cdot)$ is not explicitly given in the next section.

7.3.1 Image Completion with Masked Auto-Encoding

For a general image \mathbf{X} such as the one shown on the left of Figure 7.6, we can no longer view it as a low-rank matrix. However, humans still demonstrate remarkable ability to complete a scene and recognize familiar objects despite severe occlusion. This suggests that our brain has learned the low-dimensional distribution of natural images and can use it for completion, and hence recognition. However, the distribution of all natural images is not as simple as a low-dimensional linear subspace. Hence a natural question is whether we can learn the more sophisticated distribution of natural images and use it to perform image completion?

One empirical approach to the image completion task is to find an encoding and decoding scheme by solving the following *masked autoencoding* (MAE) program that minimizes the reconstruction loss:

$$\min_{f,g} L_{\text{MAE}}(f,g) \doteq \mathbb{E}[\|(g \circ f)(\mathcal{P}_\Omega(\mathbf{X})) - \mathbf{X}\|_2^2]. \quad (7.3.1)$$

Unlike the matrix completion problem which has a simple underlying structure, we should no longer expect that the encoding and decoding mappings admit simple closed forms or the program can be solved by explicit algorithms.

For a general natural image, we can no longer assume that its columns or rows are sampled from a low-dimensional subspace or a low-rank Gaussian.

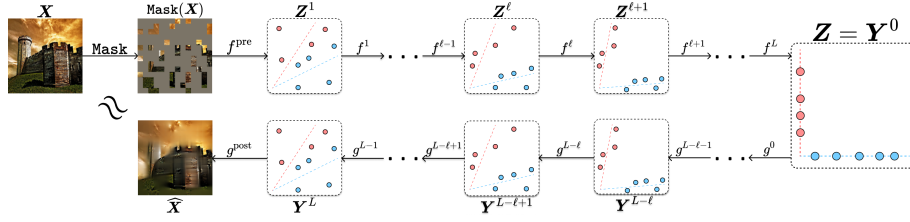


Figure 7.6: **Diagram of the overall (masked) autoencoding process.** The (image) token representations are transformed iteratively towards a parsimonious (e.g., compressed and sparse) representation by each encoder layer f^ℓ . Furthermore, such representations are transformed back to the original image by the decoder layers g^ℓ . Each encoder layer f^ℓ is meant to be (partially) inverted by a corresponding decoder layer $g^{L-\ell}$.

However, it is reasonable to assume that the image consists of multiple regions. Image patches in each region are similar and can be modeled as one (low-rank) Gaussian or subspace. Hence, to exploit the low-dimensionality of the distribution, the objective of the encoder f is to transform \mathbf{X} to a representation \mathbf{Z} :

$$f : \mathbf{X} \mapsto \mathbf{Z} \quad (7.3.2)$$

such that the distribution of \mathbf{Z} can be well modeled as a mixture of subspaces, say $\{\mathbf{U}_{[K]}\}$, such that the rate reduction is maximized while the sparsity is minimized:

$$\mathbb{E}_{\mathbf{Z}=f(\mathbf{X})}[\Delta R_\epsilon(\mathbf{Z} \mid \mathbf{U}_{[K]}) - \lambda \|\mathbf{Z}\|_0] = \mathbb{E}_{\mathbf{Z}=f(\mathbf{X})}[R_\epsilon(\mathbf{Z}) - R_\epsilon^c(\mathbf{Z} \mid \mathbf{U}_{[K]}) - \lambda \|\mathbf{Z}\|_0], \quad (7.3.3)$$

where the functions $R_\epsilon(\cdot)$ and $R_\epsilon^c(\cdot)$ are defined in (5.2.2) and (5.2.3), respectively.

As we have shown in the previous Chapter 5, the encoder f that minimizes the above objective can be constructed as a sequence of transformer-like operators. As shown in the work of [PBW+24], the decoder g can be viewed and hence constructed explicitly as the inverse process of the encoder f . Figure 7.7 illustrates the overall architectures of both the encoder and the corresponding decoder at each layer. The parameters of the encoder f and decoder g can be learned by optimizing the reconstruction loss (7.3.1) via gradient descent.

Figure 7.8 shows some representative results of the thus-designed masked auto-encoder. More implementation details and results of the masked auto-encoder for natural image completion can be found in Chapter 8 Section 8.5.

7.3.2 Conditional Sampling with Measurement Matching

The above (masked) autoencoding problem aims to generate a sample image that is consistent with certain observations or conditions. But let us examine the approach more closely: given the visual part of an image $\mathbf{X}_v = \mathcal{P}_\Omega(\mathbf{X})$, we

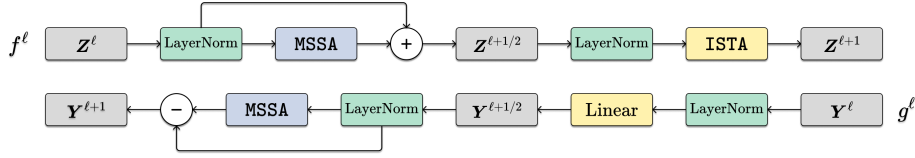


Figure 7.7: **Diagram of each encoder layer (*top*) and decoder layer (*bottom*).** Notice that the two layers are highly anti-parallel — each is constructed to do the operations of the other in reverse order. That is, in the decoder layer g^ℓ , the ISTA block of $f^{L-\ell}$ is partially inverted first using a linear layer, then the MSSA block of $f^{L-\ell}$ is reversed; this order unravels the transformation done in $f^{L-\ell}$.

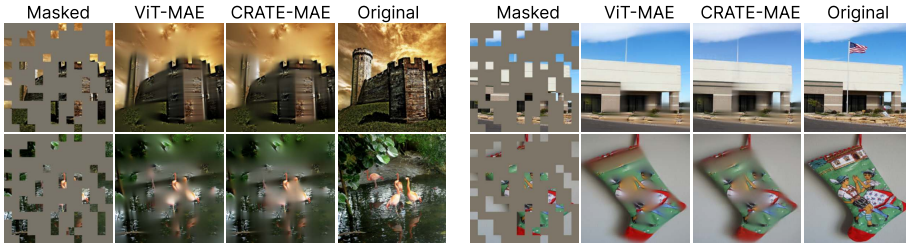


Figure 7.8: **Autoencoding visualizations of CRATE-Base and ViT-MAE-Base [HCX+22] with 75% patches masked.** We observe that the reconstructions from CRATE-Base are on par with the reconstructions from ViT-MAE-Base, despite using $< 1/3$ of the parameters.

try to estimate the masked part $\mathbf{X}_m = \mathcal{P}_{\Omega^c}(\mathbf{X})$. For realizations (Ξ_v, Ξ_m) of the random variable $\mathbf{X} = (\mathbf{X}_v, \mathbf{X}_m)$, let

$$p_{\mathbf{X}_m|\mathbf{X}_v}(\Xi_m | \Xi_v)$$

be the conditional distribution of \mathbf{X}_m given \mathbf{X}_v . It is easy to show that the optimal solution to the MAE formulation (7.3.1) is given by the conditional expectation:

$$\arg \min_{h=g \circ f} L_{\text{MAE}}(h) = \Xi_v \mapsto \Xi_v + \mathbb{E}[\mathbf{X}_m | \mathbf{X}_v = \Xi_v]. \quad (7.3.4)$$

In general, however, this expectation may not even lie on the low-dimensional distribution of natural images! This partially explains why some of the recovered patches in Figure 7.8 are a little blurry.

For many practical purposes, we would like to learn (a representation of) the conditional distribution $p_{\mathbf{X}_m|\mathbf{X}_v}$, or equivalently $p_{\mathbf{X}|\mathbf{X}_v}$, and then get a clear (most likely) sample from this distribution directly. Notice that, when the distribution of \mathbf{X} is low-dimensional, it is possible that if a sufficient part of \mathbf{X} , \mathbf{X}_v , is observed, it fully determines \mathbf{X} and hence the missing part \mathbf{X}_m . In other words, the distribution $p_{\mathbf{X}|\mathbf{X}_v}$ is a generalized function (analogous to a delta function).

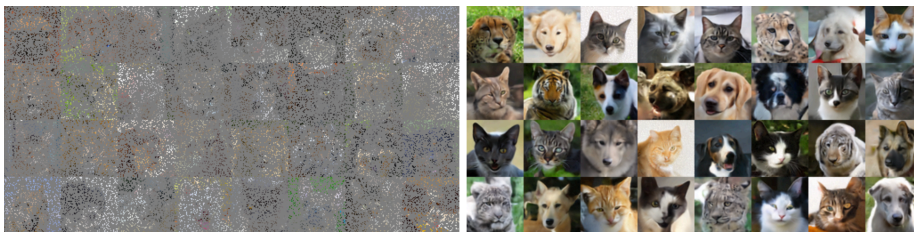


Figure 7.9: **Sampling visualizations from models trained via ambient diffusion [DSD+23b] with 80% of the pixels masked.** Using a similar ratio of masked pixels as in Figure 7.8, the ambient diffusion sampling algorithm recovers a much sharper image than the blurry image recovered by the MAE-based method. The former method samples from the distribution of natural images, while the latter approximates the conditional expectation (i.e., average) of this distribution given the observation; this averaging causes the blurriness.

Hence, instead of solving the completion task as a conditional estimation problem, we should address it as a conditional sampling problem. To that end, we should first learn the (low-dimensional) distribution of all natural images \mathbf{X} . If we have sufficient samples of natural images, we can learn the distribution via a denoising process \mathbf{X}_t described in Chapter 3. Then the problem of recovering \mathbf{X} from its partial observation $\mathbf{Y} = \mathcal{P}_\Omega(\mathbf{x}) + \mathbf{w}$ becomes a conditional generation problem—to sample the distribution conditioned on the observation.

General linear measurements. In fact, we may even consider recovering \mathbf{X} from a more general linear observation model:

$$\mathbf{Y} = \mathbf{A}\mathbf{X}_0, \quad \mathbf{X}_t = \mathbf{X}_0 + \sigma_t\mathbf{G}, \quad (7.3.5)$$

where \mathbf{A} is a linear operator on matrix space⁴ and $\mathbf{G} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The masking operator $\mathcal{P}_\Omega(\cdot)$ in the image completion task is one example of such a linear model. Then it has been shown by [DSD+23a] that

$$\hat{\mathbf{X}}_* = \arg \min_{\hat{\mathbf{X}}} \mathbb{E}[\|\mathbf{A}(\hat{\mathbf{X}}(\mathbf{A}\mathbf{X}_t, \mathbf{A}) - \mathbf{X}_0)\|^2] \quad (7.3.6)$$

satisfies the condition that:

$$\mathbf{A}\hat{\mathbf{X}}_*(\mathbf{A}(\mathbf{X}_t), \mathbf{A}) = \mathbf{A}\mathbb{E}[\mathbf{X}_0 \mid \mathbf{A}\mathbf{X}_t, \mathbf{A}]. \quad (7.3.7)$$

Notice that in the special case when \mathbf{A} is of full column rank, we have $\mathbb{E}[\mathbf{X}_0 \mid \mathbf{A}\mathbf{X}_t, \mathbf{A}] = \mathbb{E}[\mathbf{X}_0 \mid \mathbf{X}_t]$. Hence, in the more general case, it has been suggested by [DSD+23a] that one could still use the so obtained $\mathbb{E}[\mathbf{X}_0 \mid \mathbf{A}(\mathbf{X}_t), \mathbf{A}]$ to replace the $\mathbb{E}[\mathbf{X}_0 \mid \mathbf{X}_t]$ in the normal denoising process for \mathbf{X}_t :

$$\mathbf{X}_{t-s} = \gamma_t\mathbf{X}_t + (1 - \gamma_t)\mathbb{E}[\mathbf{X}_0 \mid \mathbf{A}\mathbf{X}_t, \mathbf{A}]. \quad (7.3.8)$$

⁴i.e., if we imagine unrolling \mathbf{X} into a long vector then \mathbf{A} takes the role of a matrix on \mathbf{X} -space

This usually works very well in practice, say for many image restoration tasks, as shown in [DSD+23a]. Compared to the blurry images recovered from MAE, the images recovered by the above method are much sharper as it leverages a learned distribution of natural images and samples a (sharp) image from the distribution that is consistent with the measurement, as shown in Figure 7.9 (cf. Figure 7.8).

General nonlinear measurements. To generalize the above (image) completion problems and make things more rigorous, we may consider that a random vector $\mathbf{x} \sim p$ is partially observed through a more general observation function:

$$\mathbf{y} = h(\mathbf{x}) + \mathbf{w}, \quad (7.3.9)$$

where \mathbf{w} usually stands for some random measurement noise, say of a Gaussian distribution $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. It is easy to see that, for \mathbf{x} and \mathbf{y} so related, their joint distribution $p(\mathbf{x}, \mathbf{y})$ is naturally nearly degenerate if the noise \mathbf{w} is small. To a large extent, we may view $p(\mathbf{x}, \mathbf{y})$ as a noisy version of a hypersurface defined by the function $\mathbf{y} = h(\mathbf{x})$ in the joint space (\mathbf{x}, \mathbf{y}) . Practically speaking, we will consider a setting more akin to masked autoencoding than to pure matrix completion, where we always have access to a corresponding clean sample \mathbf{x} for every observation \mathbf{y} we receive.⁵

Like image/matrix completion, we are often faced with a setting where \mathbf{y} denotes a degraded or otherwise “lossy” observation of the input \mathbf{x} . This can manifest in quite different forms. For example, in various scientific or medical imaging problems, the measured data \mathbf{y} may be a compressed and corrupted observation of the underlying data \mathbf{x} ; whereas in 3D vision tasks, \mathbf{y} may represent an image captured by a camera of a physical object with an unknown (low-dimensional) pose \mathbf{x} . Generally, by virtue of mathematical modeling (and, in some cases, co-design of the measurement system), we know h and can evaluate it on any input, and we can exploit this knowledge to help reconstruct and sample \mathbf{x} .

At a technical level, we want the learned representation of the data to facilitate us to sample the conditional distribution $p_{\mathbf{x}|\mathbf{y}}$, also known as the posterior, effectively and efficiently. More precisely, write $\boldsymbol{\nu}$ to denote a realization of the random variable \mathbf{y} . We want to generate samples $\hat{\mathbf{x}}$ such that:

$$\hat{\mathbf{x}} \sim p_{\mathbf{x}|\mathbf{y}}(\cdot \mid \mathbf{y} = \boldsymbol{\nu}). \quad (7.3.10)$$

Recall that in Section 3.2, we have developed a natural and effective way to produce *unconditional* samples of the data distribution p . The ingredients are the denoisers $\bar{\mathbf{x}}^*(t, \boldsymbol{\xi}) = \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}]$, or their learned approximations $\bar{\mathbf{x}}_\theta(t, \boldsymbol{\xi})$, for different levels of noisy observations $\mathbf{x}_t = \mathbf{x} + t\mathbf{g}$ (and $\boldsymbol{\xi}$ for their realizations) under Gaussian noise $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and $t \in [0, T]$ with a choice of times $0 = t_1 <$

⁵In some more specialized applications, in particular in scientific imaging, it is of interest to be able to learn to generate samples from the posterior $p_{\mathbf{x}|\mathbf{y}}$ without access to any clean/ground-truth samples of \mathbf{x} . We give a brief overview of methods for this setting in the end-of-chapter notes.

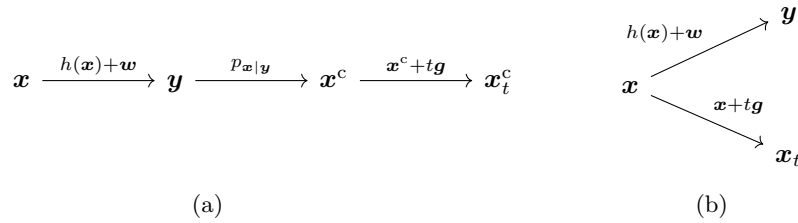


Figure 7.10: Statistical dependency diagrams for the conditional sampling process. **Left:** In a direct (conceptual) application of the diffusion-denoising scheme we have developed in Chapter 3 to conditional sampling, we use samples from the posterior $p_{\mathbf{x}|\mathbf{y}}$ to train denoisers directly on the posterior at different noise levels, then use them to generate new samples. In practice, however, we do not normally have direct samples from the posterior, but rather paired samples (\mathbf{x}, \mathbf{y}) from the joint. **Right:** It turns out that it suffices to have only noisy observations of \mathbf{x} to realize the denoisers corresponding to $p_{\mathbf{x}_i^c|\mathbf{x}^c}$: this follows from conditional independence of \mathbf{x}_t and \mathbf{y} given \mathbf{x} . It implies that $p_{\mathbf{x}_i^c|\mathbf{y}} = p_{\mathbf{x}_i|\mathbf{y}}$, which gives a score function for denoising that consists of the unconditional score function, plus a correction term that enforces measurement consistency.

$\dots < t_L = T$ at which to perform the iterative denoising, starting from $\hat{\mathbf{x}}_{t_L} \sim \mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$ (recall Equation (3.2.82)).⁶ We could directly apply this scheme to generate samples from the posterior $p_{\mathbf{x}|\mathbf{y}}$ if we had access to a dataset of samples $\mathbf{x}^c \sim p_{\mathbf{x}|\mathbf{y}}(\cdot | \boldsymbol{\nu})$ for each realization $\boldsymbol{\nu}$ of \mathbf{y} , by generating noisy observations \mathbf{x}_t^c and training denoisers to approximate $\mathbb{E}[\mathbf{x}^c | \mathbf{x}_t^c = \cdot, \mathbf{y} = \boldsymbol{\nu}]$, the mean of the posterior under the noisy observation (see Figure 7.10(a)). However, performing this resampling given only paired samples (\mathbf{x}, \mathbf{y}) from the joint distribution (say by binning the samples over values of \mathbf{y}) requires prohibitively many samples for high-dimensional data, and alternate approaches explicitly or implicitly rely on density estimation, which similarly suffers from the curse of dimensionality.

Fortunately, it turns out that this is not necessary. Consider the alternate statistical dependency diagram in Figure 7.10(b), which corresponds to the random variables in the usual denoising-diffusion process, together with the measurement \mathbf{y} . Because our assumed observation model (7.3.9) implies that

⁶Recall from our discussion in Section 3.2.2 that a few small improvements to this basic iterative denoising scheme are sufficient to bring competitive practical performance. For clarity as we develop conditional sampling, we will focus here on the simplest instantiation.

\mathbf{x}_t and \mathbf{y} are independent conditioned on \mathbf{x} , we have for any realization $\boldsymbol{\nu}$ of \mathbf{y}

$$\begin{aligned}
p_{\mathbf{x}_t|\mathbf{y}}(\cdot | \boldsymbol{\nu}) &= \int \underbrace{p_{\mathbf{x}_t|\mathbf{x}^c}(\cdot | \boldsymbol{\xi})}_{=\mathcal{N}(\boldsymbol{\xi}, t^2 \mathbf{I})} \cdot \underbrace{p_{\mathbf{x}^c|\mathbf{y}}(\boldsymbol{\xi} | \boldsymbol{\nu})}_{=p_{\mathbf{x}|\mathbf{y}}} d\boldsymbol{\xi} \\
&= \int p_{\mathbf{x}_t|\mathbf{x}, \mathbf{y}}(\cdot | \boldsymbol{\xi}, \boldsymbol{\nu}) \cdot p_{\mathbf{x}|\mathbf{y}}(\boldsymbol{\xi} | \boldsymbol{\nu}) d\boldsymbol{\xi} \quad (7.3.11) \\
&= \int p_{\mathbf{x}_t, \mathbf{x}|\mathbf{y}}(\cdot, \boldsymbol{\xi} | \boldsymbol{\nu}) d\boldsymbol{\xi} \\
&= p_{\mathbf{x}_t|\mathbf{y}}(\cdot | \boldsymbol{\nu}).
\end{aligned}$$

Above, the first line recognizes an equivalence between the distributions arising in Figure 7.10 (a,b); the second line applies this together with conditional independence of \mathbf{x}_t and \mathbf{y} given \mathbf{x} ; the third line uses the definition of conditional probability; and the final line marginalizes over \mathbf{x} . Thus, the denoisers from the conceptual posterior sampling process are equal to $\mathbb{E}[\mathbf{x} | \mathbf{x}_t = \cdot, \mathbf{y} = \boldsymbol{\nu}]$, which we can learn solely from paired samples (\mathbf{x}, \mathbf{y}) , and by Tweedie’s formula (Theorem 3.2), we can express these denoisers in terms of the score function of $p_{\mathbf{x}_t|\mathbf{y}}$, which, by Bayes’ rule, satisfies

$$p_{\mathbf{x}_t|\mathbf{y}}(\boldsymbol{\xi} | \boldsymbol{\nu}) = \frac{p_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\nu} | \boldsymbol{\xi}) p_{\mathbf{x}_t}(\boldsymbol{\xi})}{p_{\mathbf{y}}(\boldsymbol{\nu})}. \quad (7.3.12)$$

Recall that the density of \mathbf{x}_t is given by $p_t = \varphi_t * p$, where φ_t denotes the standard Gaussian density with zero mean and covariance $t^2 \mathbf{I}$ and $*$ denotes convolution. This is nothing but the *unconditional score function* obtained from the standard diffusion training that we developed in Section 3.2! The conditional score function then satisfies, for any realization $(\boldsymbol{\xi}, \boldsymbol{\nu})$ of $(\mathbf{x}_t, \mathbf{y})$,

$$\nabla_{\boldsymbol{\xi}} \log p_{\mathbf{x}_t|\mathbf{y}}(\boldsymbol{\xi} | \boldsymbol{\nu}) = \underbrace{\nabla_{\boldsymbol{\xi}} \log p_t(\boldsymbol{\xi})}_{\text{score matching}} + \underbrace{\nabla_{\boldsymbol{\xi}} \log p_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\nu} | \boldsymbol{\xi})}_{\text{measurement matching}}, \quad (7.3.13)$$

giving (by Tweedie’s formula) our proposed denoisers as

$$\begin{aligned}
\mathbb{E}[\mathbf{x} | \mathbf{x}_t = \boldsymbol{\xi}, \mathbf{y} = \boldsymbol{\nu}] &= \boldsymbol{\xi} + t^2 \nabla_{\boldsymbol{\xi}} \log p_t(\boldsymbol{\xi}) + t^2 \nabla_{\boldsymbol{\xi}} \log p_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\nu} | \boldsymbol{\xi}) \\
&= \mathbb{E}[\mathbf{x} | \mathbf{x}_t = \boldsymbol{\xi}] + t^2 \nabla_{\boldsymbol{\xi}} \log p_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\nu} | \boldsymbol{\xi}). \quad (7.3.14)
\end{aligned}$$

The resulting operators are interpretable as a *corrected* version of the unconditional denoiser for the noisy observation, where the correction term (the so-called “measurement matching” term) enforces consistency with the observations \mathbf{y} . This decomposition is the basis of *classifier guidance*, which we will develop in detail in Section 7.4.1. The reader should take care to note to which argument the gradient operators are applying in the above score functions in order to fully grasp the meaning of this operator.

The key remaining issue in making this procedure computational is to prescribe how to compute the measurement matching correction, since in general

we do not have a closed-form expression for the likelihood $p_{\mathbf{y}|\mathbf{x}_t}$ except for when $t = 0$. Before taking up this problem, we discuss an illustrative concrete example of the entire process, continuing from those we have developed in Section 3.2.

Example 7.4. Consider the case where the data distribution is Gaussian with mean $\boldsymbol{\mu} \in \mathbb{R}^D$ and covariance $\boldsymbol{\Sigma} \in \mathbb{R}^{D \times D}$, i.e., $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Assume that $\boldsymbol{\Sigma} \succeq \mathbf{0}$ is nonzero. Moreover, in the measurement model (7.3.9), suppose we obtain linear measurements of \mathbf{x} with independent Gaussian noise, where $\mathbf{A} \in \mathbb{R}^{d \times D}$ and $\mathbf{y} = \mathbf{A}\mathbf{x} + \sigma\mathbf{w}$ with $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ independent of \mathbf{x} . Below, we are going to work out the following consequences of this model for the general framework we have derived above. Specifically:

1. We will work out the precise forms of the posterior distribution $p_{\mathbf{x}|\mathbf{y}}$ that we are aiming to sample from.
2. We will work out the form of the measurement matching term from (7.3.13), which gives us the additive correction to the unconditional denoiser for $p_{\mathbf{x}|\mathbf{x}_t}$ that yields the denoiser for $p_{\mathbf{x}|\mathbf{x}_t, \mathbf{y}}$ (recall (7.3.14)).
3. We will assess the specific form of the measurement matching term we derived in the previous step, and posit a natural approximation that makes it far less sample and compute intensive to learn the measurement matching term from data. This approximation will be applied in greater generality, after the conclusion of this example.

First, we note that $\mathbf{x} \stackrel{\text{d}}{=} \boldsymbol{\Sigma}^{1/2}\mathbf{g} + \boldsymbol{\mu}$, where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is independent of \mathbf{w} and $\boldsymbol{\Sigma}^{1/2}$ is the unique positive square root of the covariance matrix $\boldsymbol{\Sigma}$, and after some algebra, we can then write

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \stackrel{\text{d}}{=} \begin{bmatrix} \boldsymbol{\Sigma}^{1/2} & \mathbf{0} \\ \mathbf{A}\boldsymbol{\Sigma}^{1/2} & \sigma\mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{g} \\ \mathbf{w} \end{bmatrix} + \begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} \end{bmatrix}.$$

By independence, we have that (\mathbf{g}, \mathbf{w}) is jointly Gaussian, which means that (\mathbf{x}, \mathbf{y}) is also jointly Gaussian, as the affine image of a jointly Gaussian vector. Its covariance matrix is given by

$$\begin{bmatrix} \boldsymbol{\Sigma}^{1/2} & \mathbf{0} \\ \mathbf{A}\boldsymbol{\Sigma}^{1/2} & \sigma\mathbf{I} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}^{1/2} & \mathbf{0} \\ \mathbf{A}\boldsymbol{\Sigma}^{1/2} & \sigma\mathbf{I} \end{bmatrix}^\top = \begin{bmatrix} \boldsymbol{\Sigma} & \boldsymbol{\Sigma}\mathbf{A}^\top \\ \mathbf{A}\boldsymbol{\Sigma} & \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top + \sigma^2\mathbf{I} \end{bmatrix}.$$

Now, we apply the fact that conditioning a random vector with joint Gaussian distribution on a subset of coordinates is again a Gaussian distribution (Exercise 3.2). By this, we obtain that

$$p_{\mathbf{x}|\mathbf{y}}(\cdot | \boldsymbol{\nu}) = \mathcal{N}\left(\underbrace{\boldsymbol{\mu} + \boldsymbol{\Sigma}\mathbf{A}^\top (\mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top + \sigma^2\mathbf{I})^{-1} (\boldsymbol{\nu} - \mathbf{A}\boldsymbol{\mu})}_{\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{\nu})}, \underbrace{\boldsymbol{\Sigma} - \boldsymbol{\Sigma}\mathbf{A}^\top (\mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top + \sigma^2\mathbf{I})^{-1} \mathbf{A}\boldsymbol{\Sigma}}_{\boldsymbol{\Sigma}_{\mathbf{x}|\mathbf{y}}}\right). \quad (7.3.15)$$

Next, since $\mathbf{x}_t = \mathbf{x} + t\mathbf{g}'$, where $\mathbf{g}' \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is independent of all other random vectors, and since $\mathbf{x} \mid \mathbf{y}$ is Gaussian, by the calculation directly above, it follows by another application of Exercise 3.2 that $\mathbf{x} \mid \mathbf{x}_t, \mathbf{y}$ is Gaussian. Its conditional expectation function is given by reading off the mean of the Gaussian distribution, which in this case is

$$\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}, \mathbf{y} = \boldsymbol{\nu}] = \boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{\nu}) + \boldsymbol{\Sigma}_{\mathbf{x}|\mathbf{y}} (\boldsymbol{\Sigma}_{\mathbf{x}|\mathbf{y}} + t^2\mathbf{I})^{-1} (\boldsymbol{\xi} - \boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}(\boldsymbol{\nu})). \quad (7.3.16)$$

The functional form of this denoiser is quite simple, but it carries an unwieldy dependence on the problem data $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}$, \mathbf{A} , and σ^2 . We can gain further insight into its behavior by comparing it with Equation (7.3.14). We have as usual

$$\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}] = \boldsymbol{\mu} + \boldsymbol{\Sigma} (\boldsymbol{\Sigma} + t^2\mathbf{I})^{-1} (\boldsymbol{\xi} - \boldsymbol{\mu}), \quad (7.3.17)$$

which is rather simple—suggesting that the measurement matching term is rather complicated. To confirm this, we can calculate the likelihood $p_{\mathbf{y}|\mathbf{x}_t}$ directly using the following expression for the joint distribution of $(\mathbf{x}, \mathbf{x}_t, \mathbf{y})$:

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{x}_t \\ \mathbf{y} \end{bmatrix} \stackrel{d}{=} \begin{bmatrix} \boldsymbol{\Sigma}^{1/2} & \mathbf{0} & \mathbf{0} \\ \boldsymbol{\Sigma}^{1/2} & t\mathbf{I} & \mathbf{0} \\ \mathbf{A}\boldsymbol{\Sigma}^{1/2} & \mathbf{0} & \sigma\mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{g} \\ \mathbf{g}' \\ \mathbf{w} \end{bmatrix} + \begin{bmatrix} \boldsymbol{\mu} \\ \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} \end{bmatrix}, \quad (7.3.18)$$

where $\mathbf{g}' \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ independent of the other Gaussians. This is again a jointly Gaussian distribution; restricting to only the final two rows, we have the covariance

$$\begin{bmatrix} \boldsymbol{\Sigma}^{1/2} & t\mathbf{I} & \mathbf{0} \\ \mathbf{A}\boldsymbol{\Sigma}^{1/2} & \mathbf{0} & \sigma\mathbf{I} \end{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}^{1/2} & t\mathbf{I} & \mathbf{0} \\ \mathbf{A}\boldsymbol{\Sigma}^{1/2} & \mathbf{0} & \sigma\mathbf{I} \end{bmatrix}^\top = \begin{bmatrix} \boldsymbol{\Sigma} + t^2\mathbf{I} & \boldsymbol{\Sigma}\mathbf{A}^\top \\ \mathbf{A}\boldsymbol{\Sigma} & \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top + \sigma^2\mathbf{I} \end{bmatrix}.$$

Another application of Exercise 3.2 then gives us

$$p_{\mathbf{y}|\mathbf{x}_t}(\cdot \mid \boldsymbol{\xi}) = \mathcal{N}\left(\underbrace{\mathbf{A}\boldsymbol{\mu} + \mathbf{A}\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + t^2\mathbf{I})^{-1} (\boldsymbol{\xi} - \boldsymbol{\mu})}_{\boldsymbol{\mu}_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\xi})}, \underbrace{\mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top + \sigma^2\mathbf{I} - \mathbf{A}\boldsymbol{\Sigma} (\boldsymbol{\Sigma} + t^2\mathbf{I})^{-1} \boldsymbol{\Sigma}\mathbf{A}^\top}_{\boldsymbol{\Sigma}_{\mathbf{y}|\mathbf{x}_t}}\right). \quad (7.3.19)$$

Now notice that $\boldsymbol{\mu}_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\xi}) = \mathbf{A}\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}]$. So, by the chain rule,

$$\begin{aligned} & t^2 \nabla_{\boldsymbol{\xi}} \log p_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\nu} \mid \boldsymbol{\xi}) \\ &= t^2 \nabla_{\boldsymbol{\xi}} \left[-\frac{1}{2} (\boldsymbol{\nu} - \mathbf{A}\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}])^\top \boldsymbol{\Sigma}_{\mathbf{y}|\mathbf{x}_t}^{-1} (\boldsymbol{\nu} - \mathbf{A}\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}]) \right] \\ &= t^2 (\boldsymbol{\Sigma} + t^2\mathbf{I})^{-1} \boldsymbol{\Sigma}\mathbf{A}^\top \boldsymbol{\Sigma}_{\mathbf{y}|\mathbf{x}_t}^{-1} (\boldsymbol{\nu} - \mathbf{A}\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}]). \end{aligned} \quad (7.3.20)$$

This gives us a more interpretable decomposition of the conditional posterior denoiser (7.3.16): following Equation (7.3.14), it is the sum of the unconditional posterior denoiser (7.3.17) and the measurement matching term (7.3.20).

We can further analyze the measurement matching term to understand cases under which it can be approximated. Notice that

$$\Sigma_{\mathbf{y}|\mathbf{x}_t} = \sigma^2 \mathbf{I} + \mathbf{A} \Sigma^{1/2} \left(\mathbf{I} - \Sigma^{1/2} (\Sigma + t^2 \mathbf{I})^{-1} \Sigma^{1/2} \right) \Sigma^{1/2} \mathbf{A}^\top. \quad (7.3.21)$$

If we let $\Sigma = \mathbf{V} \Lambda \mathbf{V}^\top$ denote an eigenvalue decomposition of Σ , where (\mathbf{v}_i) are the columns of \mathbf{V} , we can further write

$$\Sigma^{1/2} \left(\mathbf{I} - \Sigma^{1/2} (\Sigma + t^2 \mathbf{I})^{-1} \Sigma^{1/2} \right) \Sigma^{1/2} = t^2 \mathbf{V} \Lambda^{1/2} (\Lambda + t^2 \mathbf{I})^{-1} \Lambda^{1/2} \mathbf{V}^\top \quad (7.3.22)$$

$$= t^2 \sum_{i=1}^D \frac{\lambda_i}{\lambda_i + t^2} \mathbf{v}_i \mathbf{v}_i^*. \quad (7.3.23)$$

Then for any eigenvalue of Σ equal to zero, the corresponding summand is zero; and writing $\lambda_{\min}(\Sigma)$ for the smallest positive eigenvalue of Σ , we have that whenever $t \ll \sqrt{\lambda_{\min}(\Sigma)}$, it holds

$$\frac{\lambda_i t^2}{\lambda_i + t^2} \approx 0. \quad (7.3.24)$$

So, when $t \ll \sqrt{\lambda_{\min}(\Sigma)}$, we have the approximation

$$\Sigma_{\mathbf{y}|\mathbf{x}_t} \approx \sigma^2 \mathbf{I}. \quad (7.3.25)$$

The right-hand side of this approximation is equal to $\Sigma_{\mathbf{y}|\mathbf{x}}$. So we have in turn

$$\nabla_{\boldsymbol{\xi}} \log p_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\nu} | \boldsymbol{\xi}) \approx \nabla_{\boldsymbol{\xi}} \log p_{\mathbf{y}|\mathbf{x}}(\boldsymbol{\nu} | \mathbb{E}[\mathbf{x} | \mathbf{x}_t = \boldsymbol{\xi}]). \quad (7.3.26)$$

■

Let us take some time to interpret the results of Example 7.4 before we generalize it. The approximation (7.3.26) is, of course, a direct consequence of the specific modeling assumptions we have made in Example 7.4. However, notice that if we directly interpret this approximation, it is *ab initio* tractable: the likelihood $p_{\mathbf{y}|\mathbf{x}} = \mathcal{N}(\mathbf{A}\mathbf{x}, \sigma^2 \mathbf{I})$ is a simple Gaussian distribution centered at the observation, and the approximation to the measurement matching term that we arrive at can be interpreted as simply evaluating the log-likelihood at the conditional expectation $\mathbb{E}[\mathbf{x} | \mathbf{x}_t = \boldsymbol{\xi}]$, then taking gradients with respect to $\boldsymbol{\xi}$ (which involves backpropagating through the conditional expectation).

To gain insight into the effect of the convenient approximation (7.3.26), we implement and simulate a simple numerical experiment in the Gaussian setting in Figure 7.11. The sampler we implement is a direct implementation of the simple scheme (3.2.82) we have developed in Chapter 3 and recalled above, using the true conditional posterior denoiser, i.e. Equation (7.3.16) (samples are marked with blue circles), and the convenient approximation to this denoiser made through the measurement matching approximation (7.3.26) (samples are

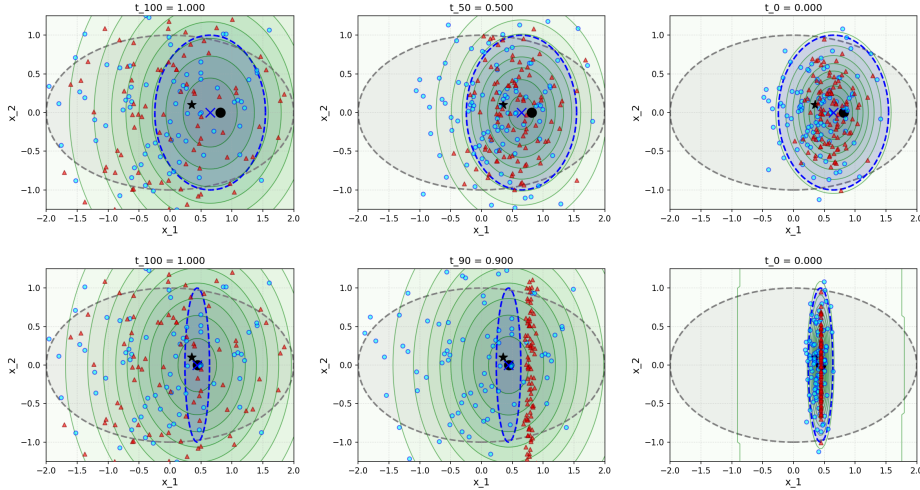


Figure 7.11: Numerical simulation of the conditional sampling setup (7.3.9), with Gaussian data, linear measurements, and Gaussian noise. We simulate $D = 2$ and $d = 1$, with $\Sigma = e_1 e_1^\top + \frac{1}{4} e_2 e_2^\top$, $\mu = \mathbf{0}$, and $\mathbf{A} = e_1^\top$. The underlying signal \mathbf{x} is marked with a black star, and the measurement \mathbf{y} is marked with a black circle. Each individual plot corresponds to a different value of sampler time t_ℓ , with different rows corresponding to different observation noise levels σ^2 . In each plot, the covariance matrix of \mathbf{x} is plotted in gray, the posterior covariance matrix and posterior mean of $p_{\mathbf{x}|\mathbf{y}}$ are plotted in blue (with the posterior mean marked by a blue “x”), and contours for $p_{\mathbf{x}_{t_\ell}|\mathbf{y}}$ are drawn in green. The sampler hyperparameters are $T = 1$, $L = 100$, and we draw 100 independent samples to initialize the samplers. Samplers are implemented with the closed-form denoisers derived in Example 7.4, with those using the approximation (7.3.26) marked with red triangles, and those using the exact conditional posterior denoiser marked with blue circles. **Top:** For large observation noise $\sigma = 0.5$, both the exact conditional posterior denoiser and the approximate one do a good job of converging to the posterior $p_{\mathbf{x}|\mathbf{y}}$. Sampling time (corresponding to time in the “forward process”, so larger times mean larger noise) decreases from left to right. The convergence dynamics for the exact and approximate measurement matching term are similar. **Bottom:** For smaller observation noise $\sigma = 0.1$, the approximate measurement matching term leads to extreme bias in the sampler (red triangles): samples rapidly converge to an affine subspace of points that are consistent, modulo some shrinkage from the posterior mean denoiser, with the measured ground truth, and later sampling iterations are unable to recover the lost posterior variance along this dimension. Note that different times t_ℓ are plotted in the bottom row, compared to the top row, to show the rapid collapse of the approximation to the posterior along the measurement dimension.

marked with red triangles). The top row of Figure 7.11 shows a setting with large measurement noise σ^2 , and the bottom with small measurement noise. The measurement matching approximation works very well in the large-noise setting, with the caveat that in the small-noise setting, it suffers from rapid collapse of the variance of the sampling distribution along directions that are parallel to the rows of the linear measurement operator \mathbf{A} , which cannot be corrected by later iterations of sampling. Our analysis in Example 7.4 precisely characterizes the level at which the noise becomes too “small” in terms of the data covariance matrix Σ .

Generalizing the approximation: Diffusion Posterior Sampling. Example 7.4 suggests a convenient approximation for the measurement matching term, i.e. (7.3.26), which can be made beyond the Gaussian setting of the example. To motivate this approximation in greater generality, notice that by conditional independence of \mathbf{y} and \mathbf{x}_t given \mathbf{x} , we can write

$$p_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\nu} | \boldsymbol{\xi}) = \int p_{\mathbf{y}|\mathbf{x}}(\boldsymbol{\nu} | \boldsymbol{\xi}') p_{\mathbf{x}|\mathbf{x}_t}(\boldsymbol{\xi}' | \boldsymbol{\xi}) d\boldsymbol{\xi}'. \quad (7.3.27)$$

Formally, when the posterior $p_{\mathbf{x}|\mathbf{x}_t}$ is a delta function centered at its mean $\mathbb{E}[\mathbf{x} | \mathbf{x}_t = \boldsymbol{\xi}]$, the approximation (7.3.26) is exact. More generally, when the posterior $p_{\mathbf{x}|\mathbf{x}_t}$ is highly concentrated around its mean, the approximation (7.3.26) is accurate. This holds, for example, for sufficiently small t , which we saw explicitly in the Gaussian setting of Example 7.4. Although the numerical simulation in Figure 7.11 suggests that this approximation is not without its caveats in certain regimes, it has proved to be a reliable baseline in practice, after being proposed by Chung et al. as “Diffusion Posterior Sampling” (DPS) [CKM+23]. In addition, there are even principled and generalizable approaches to improve it by incorporating better estimates of the posterior variance (which turn out to be exact in the Gaussian setting of Example 7.4), which we discuss further in the end-of-chapter summary.

Thus, with the DPS approximation, we arrive at the following approximation for the conditional posterior denoisers $\mathbb{E}[\mathbf{x} | \mathbf{y}, \mathbf{x}_t]$, via Equation (7.3.14):

$$\mathbb{E}[\mathbf{x} | \mathbf{x}_t = \boldsymbol{\xi}, \mathbf{y} = \boldsymbol{\nu}] \approx \mathbb{E}[\mathbf{x} | \mathbf{x}_t = \boldsymbol{\xi}] + t^2 \nabla_{\boldsymbol{\xi}} \log p_{\mathbf{y}|\mathbf{x}}(\boldsymbol{\nu} | \mathbb{E}[\mathbf{x} | \mathbf{x}_t = \boldsymbol{\xi}]). \quad (7.3.28)$$

And, for a neural network or other model $\bar{\mathbf{x}}_{\theta}(t, \boldsymbol{\xi})$ trained as in Section 3.2 to approximate the denoisers $\mathbb{E}[\mathbf{x} | \mathbf{x}_t = \boldsymbol{\xi}]$ for each $t \in [0, T]$, we arrive at the learned conditional posterior denoisers

$$\bar{\mathbf{x}}_{\theta}(t, \boldsymbol{\xi}, \boldsymbol{\nu}) = \bar{\mathbf{x}}_{\theta}(t, \boldsymbol{\xi}) + t^2 \nabla_{\boldsymbol{\xi}} \log p_{\mathbf{y}|\mathbf{x}}(\boldsymbol{\nu} | \bar{\mathbf{x}}_{\theta}(t, \boldsymbol{\xi})). \quad (7.3.29)$$

Note that the approximation (7.3.28) is valid for arbitrary forward models h in the observation model (7.3.9), including nonlinear h , and even to arbitrary noise models for which a clean expression for the likelihood $p_{\mathbf{y}|\mathbf{x}}$ is known. Indeed, in the case of Gaussian noise, we have

$$p_{\mathbf{y}|\mathbf{x}}(\boldsymbol{\nu} | \boldsymbol{\xi}) \propto \exp\left(-\frac{1}{2\sigma^2} \|h(\boldsymbol{\xi}) - \boldsymbol{\nu}\|_2^2\right). \quad (7.3.30)$$

Hence, evaluating the right-hand side of (7.3.29) requires only

1. A pretrained denoiser $\bar{\mathbf{x}}_\theta(t, \boldsymbol{\xi})$ for the data distribution p (of \mathbf{x}), learned as in Section 3.2 via Algorithm 3.2;
2. Forward and backward pass access to the forward model h for the measurements (7.3.9);
3. A forward and backward pass through $\bar{\mathbf{x}}_\theta(t, \boldsymbol{\xi})$, which can be evaluated efficiently using (say) backpropagation.

Algorithm 7.1 Conditional sampling under measurements (7.3.9), with an unconditional denoiser and DPS.

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

Input: An unconditional denoiser $\bar{\mathbf{x}}_\theta: \{t_\ell\}_{\ell=1}^L \times \mathbb{R}^D \rightarrow \mathbb{R}^D$ for $p_{\mathbf{x}}$.

Input: Measurement realization $\boldsymbol{\nu}$ of \mathbf{y} (Equation (7.3.9)) to condition on.

Input: Forward model $h: \mathbb{R}^D \rightarrow \mathbb{R}^d$ and measurement noise variance $\sigma^2 > 0$.

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 $p_{\mathbf{x}|\mathbf{y}}$.

- 1: **function** DDIMSAMPLERCONDITIONALDPS($\bar{\mathbf{x}}_\theta, \boldsymbol{\nu}, h, \sigma^2, (t_\ell)_{\ell=0}^L$)
 - 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

$$w_{t_{\ell-1}} \doteq \alpha_{t_{\ell-1}} - \frac{\sigma_{t_{\ell-1}}}{\sigma_{t_\ell}} \alpha_{t_\ell};$$

$$\hat{\mathbf{x}}_{t_{\ell-1}} \doteq \frac{\sigma_{t_{\ell-1}}}{\sigma_{t_\ell}} \hat{\mathbf{x}}_{t_\ell} + \left(\bar{\mathbf{x}}_\theta(t_\ell, \hat{\mathbf{x}}_{t_\ell}) - \frac{\sigma_{t_\ell}^2}{2\alpha_{t_\ell}\sigma^2} \nabla_{\boldsymbol{\xi}} \left[\|h(\bar{\mathbf{x}}_\theta(t_\ell, \boldsymbol{\xi})) - \boldsymbol{\nu}\|_2^2 \right] \Big|_{\boldsymbol{\xi}=\hat{\mathbf{x}}_{t_\ell}} \right)$$
 - 5: **end for**
 - 6: **return** $\hat{\mathbf{x}}_{t_0}$
 - 7: **end function**
-

Combining this scheme with the basic implementation of unconditional sampling we developed in Section 3.2, we obtain a practical algorithm for conditional sampling of the posterior $p_{\mathbf{x}|\mathbf{y}}$ given measurements following (7.3.9). Algorithm 7.1 records this scheme for the case of Gaussian observation noise with known standard deviation σ , with minor modifications to extend to a general noising process, as in Equation (3.2.85) and the surrounding discussion in Chapter 3 (our discussion above made the simplifying choices $\alpha_t = 1$, $\sigma_t = t$, and $t_\ell = T\ell/L$, as for Equation (3.2.82) in Section 3.2).

Application to medical image reconstruction. The framework that we have developed above is already powerful enough to be applied to solve various *scientific inverse problems*, where the measurements $\mathbf{y} = h(\mathbf{x}_o) + \mathbf{w}$ are the

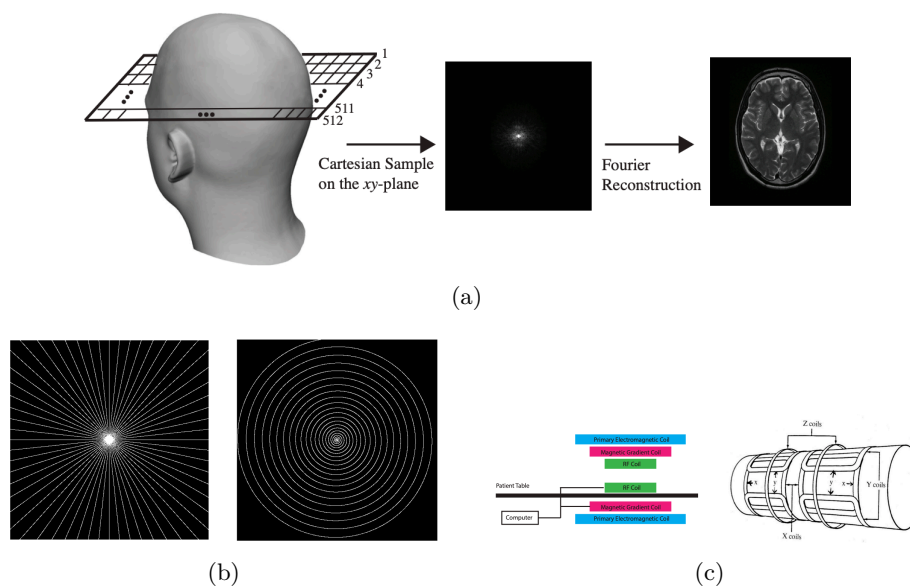


Figure 7.12: Simplified overview of magnetic resonance imaging (MRI) reconstruction. **(a)** An overview of the imaging process. A slice of the object being imaged (say, a human brain) is measured in the Fourier domain via the MRI machine. The spatial image can be reconstructed via direct inversion of the Fourier domain measurements. **(b)** In practice, the MRI machine *subsamples* the Fourier domain via a structured sampling pattern for the sake of efficiency. This creates a linear inverse problem, where priors about the structure of the underlying image must be exploited to enable unique reconstruction. **(c)** The actual measurement apparatus uses a controlled magnetic field to obtain the Fourier domain measurements. Images reproduced with permission from [WM22].

outputs of a known measurement process h on a specific signal of interest \mathbf{x}_o , and where there are infinitely many possible signals \mathbf{x} that satisfy $\mathbf{y} = h(\mathbf{x}) + \mathbf{w}$. This means that to find the *true* \mathbf{x}_o that generated \mathbf{y} , we need to exploit prior information about the structure of \mathbf{x}_o . We will consider the probabilistic setting where $\mathbf{x}_o \sim p_{\mathbf{x}}$, and where we have access to a prior on $p_{\mathbf{x}}$, say in the form of a diffusion model.

An instructive example of a scientific inverse problem is magnetic resonance imaging (MRI) (Figure 7.12), wherein the data distribution $p_{\mathbf{x}}$ is over images corresponding to an underlying object of interest (such as a 2D slice of a human subject’s brain). In MRI applications, the measurement operator h corresponds to the machine generating and modulating the magnetic field and recording the patient’s response to it (Figure 7.12(c)). It can be modeled as implementing the *Fourier transform* of the underlying spatial signal (Figure 7.12(a))—see [WM22, Chapter 10] for a detailed mathematical derivation. In addition, for efficiency’s

sake, the MRI machine typically only measures a subset of the measured frequencies, in a structured measurement pattern that can be modeled as either a radial or a spiral pattern (Figure 7.12(b)). Writing \mathcal{F} to denote the (discrete) Fourier transform operator and \mathcal{S} to denote the machine’s (known) subsampling pattern, the measured signal is then

$$\mathbf{y} = \mathcal{S} \circ \mathcal{F}(\mathbf{x}_o) + \mathbf{w}. \quad (7.3.31)$$

The measurement operator $\mathcal{S} \circ \mathcal{F}$ is *linear*, implying that it can be represented by a $m \times n$ matrix after appropriate reshaping, but it is also *compressive*, which means that $m \ll n$ and that we need to exploit prior knowledge about $p_{\mathbf{x}}$ to reconstruct \mathbf{x}_o (or more generally to sample from the posterior $p_{\mathbf{x}|\mathbf{y}}$).

The distribution $p_{\mathbf{x}}$ of MRI images is simpler than some image distributions (e.g., natural images), but it is still complex enough to benefit from using a learned model for $p_{\mathbf{x}}$ rather than an analytical one (cf. [WM22]). One direct approach is to apply the diffusion posterior sampling approximation we developed in the previous section, which led to Algorithm 7.1, with the specific MRI measurement operator $\mathbf{A} = \mathcal{S} \circ \mathcal{F}$ and a diffusion model pretrained on MRI data. This precise combination of approaches has not appeared in the literature. Instead, we highlight results from a different but related approach, which enforces measurement consistency on each diffusion step. Following the pioneering work of Song et al. [SSX+22], numerous such measurement consistency algorithms for medical image reconstruction have been developed, including the improved approaches of Song et al. [SKZ+24] and Rout et al. [RRD+23] (which utilize latent diffusion models, as we have developed for representation autoencoders in Chapter 6). Figure 7.13 shows visual comparisons of reconstructed MRI images across three different samples, comparing FISTA-TV (a classical, accelerated version of the ISTA algorithm we studied in Chapter 2), the measurement consistency diffusion modeling approach of Song et al. [SSX+22], and the ground truth. It can be seen that superior reconstruction quality is achieved by the diffusion-based method, as well as superior faithfulness to fine detail in the ground truth image. Table 7.1 shows MRI reconstruction accuracy results (measured in terms of PSNR) for the approach developed by Song et al. [SSX+22]. It compares favorably to supervised baselines (trained on many independent paired samples $(\mathbf{y}, \mathbf{x}_o)$, which is an additional requirement) and unsupervised baselines (which correspond to measurement matching approximations that came before diffusion posterior sampling).

7.4 Conditional Inference with Paired Data and Measurements

In many practical applications, we do not know either the distribution of the data \mathbf{x} of interest or the explicit relationship between the data and certain observed attributes \mathbf{y} of the data. We only have a (large) set of paired samples $(\mathbf{X}, \mathbf{Y}) = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$ from which we need to infer the data

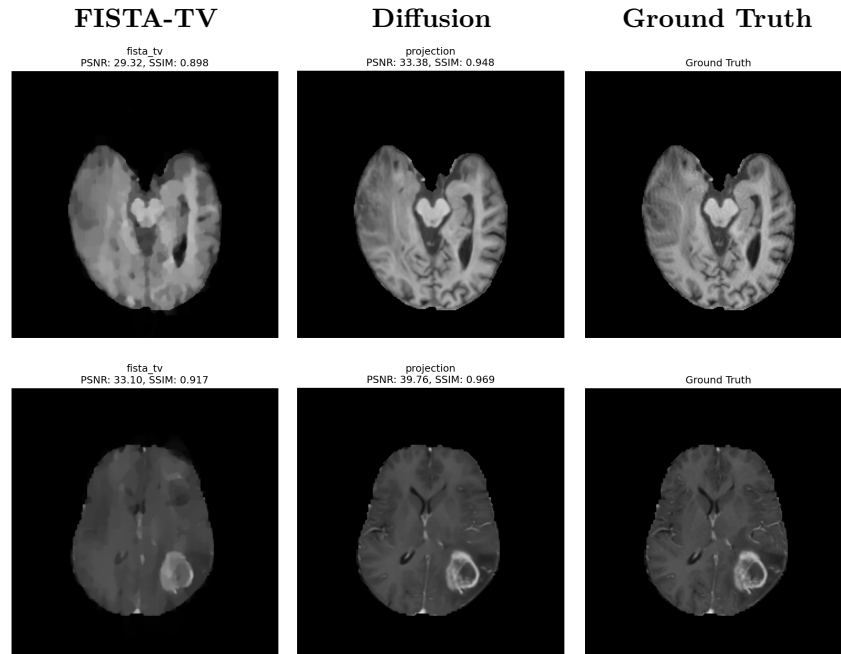


Figure 7.13: **Visual comparison of MRI reconstruction methods.** Three example brain MRI reconstructions showing FISTA-TV, the measurement consistency diffusion approach of Song et al. [SSX+22], and ground truth images on the BraTS MRI reconstruction dataset at 8x undersampling. The diffusion-based approach produces reconstructions that more closely match the ground truth structure compared to FISTA-TV, which can be seen to over-smooth the fine details present in the ground truth image. At the same time, the degree of faithfulness to the ground truth image in the diffusion-based reconstruction (arising from the proper enforcement of measurement consistency) is remarkable.

distribution and a mapping that models their relationship:

$$h : \mathbf{x} \mapsto \mathbf{y}. \quad (7.4.1)$$

The problem of image classification can be viewed as one such example. In a sense, the classification problem is to learn an (extremely lossy) compressive encoder for natural images. Say, given a random sample of an image \mathbf{x} , we would like to predict its class label \mathbf{y} that best correlates the content in \mathbf{x} . We know the distribution of natural images of objects is low-dimensional compared to the dimension of the pixel space. From the previous chapters, we have learned that given sufficient samples, in principle, we can learn a structured low-dimensional representation \mathbf{z} for \mathbf{x} through a learned compressive encoding:

$$f : \mathbf{x} \mapsto \mathbf{z}. \quad (7.4.2)$$

Method	24× Accel.	8× Accel.	4× Accel.
Cascade DenseNet	23.39 \pm 2.17	28.35 \pm 2.30	30.97 \pm 2.33
DuDoRNet	18.46 \pm 3.05	37.88\pm3.03	30.53 \pm 4.13
Song et al. [SSK+21]	27.83 \pm 2.73	35.04 \pm 2.11	37.55 \pm 2.08
Jalal et al. [JAD+21]	28.80 \pm 3.21	36.44 \pm 2.28	38.76 \pm 2.32
Song et al. [SSX+22]	29.42\pm3.03	37.63 \pm 2.70	39.91\pm2.67

Table 7.1: **MRI reconstruction results** comparing measurement matching approaches across different acceleration factors (the degree of Fourier space subsampling in \mathcal{S}) on the BraTS MRI reconstruction dataset. PSNR metrics (higher is better) with standard deviations are reported. Cascade DenseNet [ZFZ19] and DuDoRNet [ZZ20] are supervised baselines trained on paired samples at an 8x acceleration level. The last three methods are diffusion-based approaches using measurement consistency. Reproduced with permission from Song et al. [SSX+22].

The representation \mathbf{z} can also be viewed as a learned (lossy but structured) code for \mathbf{x} . It is rather reasonable to assume that if the class assignment \mathbf{y} truly depends on the low-dimensional structures of \mathbf{x} and the learned code \mathbf{z} truly reflects such structures, \mathbf{y} and \mathbf{z} can be made highly correlated and hence their joint distribution $p(\mathbf{z}, \mathbf{y})$ should be extremely low-dimensional. Therefore, we may combine the two desired codes \mathbf{y} and \mathbf{z} together and try to learn a combined encoder:

$$f : \mathbf{x} \mapsto (\mathbf{z}, \mathbf{y}) \quad (7.4.3)$$

where the joint distribution of (\mathbf{z}, \mathbf{y}) is highly low-dimensional.

From our study in previous chapters, the mapping f is usually learned as a sequence of compression or denoising operators in the same space. Hence to leverage such a family of operations, we may introduce an auxiliary vector \mathbf{w} that can be viewed as an initial random guess of the class label \mathbf{y} . In this way, we can learn a compression or denoising mapping:

$$f : (\mathbf{x}, \mathbf{w}) \mapsto (\mathbf{z}, \mathbf{y}) \quad (7.4.4)$$

within a common space. In fact, the common practice of introducing an auxiliary “class token” in the training of a transformer for classification tasks, such as in ViT, can be viewed as learning such a representation by compressing (the coding rate of) given (noisy) samples of (\mathbf{x}, \mathbf{w}) . If the distribution of the data \mathbf{x} is already a mixture of (low-dimensional) Gaussians, the work [WTL+08] has shown that classification can be done effectively by directly minimizing *the (lossy) coding length* associated with the given samples.

7.4.1 Class-Conditioned Image Generation

While a learned classifier allows us to classify a given image \mathbf{x} to its corresponding class, we often would like to generate an image of a given class, by sampling

the learned distribution of natural images. To some extent, this can be viewed as the “inverse” problem to image classification. Let $p_{\mathbf{x}}$ denote the distribution of natural images, say modeled by a diffusion-denoising process. Given a class label random variable $y \in [K]$ with realization ν , say an “Apple”, we would like to sample the conditional distribution $p_{\mathbf{x}|y}(\cdot | \nu)$ to generate an image of an apple:

$$\hat{\mathbf{x}} \sim p_{\mathbf{x}|y}(\cdot | \nu). \quad (7.4.5)$$

We call this *class-conditioned image generation*.

In Section 7.3.2, we have seen how to use the denoising-diffusion paradigm for conditional sampling from the posterior $p_{\mathbf{x}|y}$ given *model-based* measurements $\mathbf{y} = h(\mathbf{x}) + \mathbf{w}$ (Equation (7.3.9)), culminating in the DPS algorithm (Algorithm 7.1). This is a powerful framework, but it does not apply to the class (or text) conditioned image generation problem here, where an explicit generative model h for the observations/attributes y is not available due to the intractability of analytical modeling. In this section, we will present techniques for extending conditional sampling to this setting.

Thus, we now assume only that we have access to samples from the joint distribution of (\mathbf{x}, y) :

$$(\mathbf{x}, y) \sim p_{\mathbf{x},y}. \quad (7.4.6)$$

As in the previous section, we define $\mathbf{x}_t = \alpha_t \mathbf{x} + \sigma_t \mathbf{g}$ with $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ independent of (\mathbf{x}, y) , as in Equation (3.2.85) in Chapter 3, and we will repeatedly use the notation $\boldsymbol{\xi}$ to denote realizations of \mathbf{x} and \mathbf{x}_t .

This section focuses on the conceptual and algorithmic underpinnings for conditional sampling in modern generative models. For a treatment focused on implementation issues in various applications, see Chapter 8, Section 8.7 and subsequent sections.

Classifier Guidance: Consistency via a Pretrained Classifier

Recall from Section 7.3.2 that the following decomposition of the optimal conditional posterior denoiser holds, by virtue of Bayes’ rule and conditional independence of y and \mathbf{x}_t given \mathbf{x} (recall Figure 7.10):

$$\mathbb{E}[\mathbf{x} | \mathbf{x}_t = \boldsymbol{\xi}, y = \nu] = \mathbb{E}[\mathbf{x} | \mathbf{x}_t = \boldsymbol{\xi}] + \frac{\sigma_t^2}{\alpha_t} \nabla_{\boldsymbol{\xi}} \log p_{y|\mathbf{x}_t}(\nu | \boldsymbol{\xi}). \quad (7.4.7)$$

In other words, the representation (7.4.7) for the optimal conditional denoiser holds regardless of whether the explicit observation model $\mathbf{y} = h(\mathbf{x}) + \mathbf{w}$ is valid. This means that if we can learn such a denoiser (7.4.7) in the paired data setting, we can follow the denoising-diffusion methodology we have developed in Chapter 3 to perform class-conditional sampling.

A natural idea is then to directly implement the likelihood correction term in (7.4.7) using a deep network f_{θ_c} with parameters θ_c , as in Equation (7.4.4):

$$f_{\theta_c} : (t, \mathbf{x}_t) \mapsto \text{softmax}(\mathbf{W}_{\text{head}} \mathbf{z}(t, \mathbf{x}_t)). \quad (7.4.8)$$

This expression combines the final representations $\mathbf{z}(t, \mathbf{x}_t)$ (which also depend on θ_c) of the noisy inputs \mathbf{x}_t with a classification head $\mathbf{W}_{\text{head}} \in \mathbb{R}^{K \times d}$, which maps the representations to a probability distribution over the K possible classes. As is common in practice, it also takes the time t in the noising process as input. Thus, with appropriate training, it provides an approximation to the log-likelihood $\log p_{y|\mathbf{x}_t}$, and differentiating $\log f_{\theta_c}$ with respect to its input \mathbf{x}_t allows an approximation to the second term in Equation (7.4.7):

$$\bar{\mathbf{x}}_{\theta}^{\text{naive}}(t, \mathbf{x}_t, y) = \bar{\mathbf{x}}_{\theta_d}(t, \mathbf{x}_t) + \frac{\sigma_t^2}{\alpha_t} \nabla_{\mathbf{x}_t} \langle \log f_{\theta_c}(t, \mathbf{x}_t), \mathbf{e}_y \rangle \quad (7.4.9)$$

where, as usual, we approximate the first term in Equation (7.4.7) via a learned unconditional denoiser for \mathbf{x}_t with parameters θ_d , and where we write \mathbf{e}_k for $k \in [K]$ to denote the k -th canonical basis vector for \mathbb{R}^K (i.e., the vector with a one in the k -th position, and zeros elsewhere). The reader should note that the conditional denoiser $\bar{\mathbf{x}}_{\theta}$ requires two separate training runs, with separate losses: one for the classifier parameters θ_c , on a classification loss,⁷ and one for the denoiser parameters θ_d , on a denoising loss. Such an approach to conditional sampling was already recognized and exploited to perform conditional sampling in pioneering early works on diffusion models, notably those by Sohl-Dickstein et al. [SWM+15] and by Song et al. [SSK+21].

However, this straightforward methodology has two key drawbacks (which is why we label it as “naive”). The first is that, empirically, such a trained deep network classifier frequently does not provide a strong enough guidance signal (in Equation (7.4.7)) to ensure that generated samples reflect the conditioning information y . This was first emphasized by Dhariwal and Nichol [DN21b], who noted that in the setting of class-conditional ImageNet generation, the learned deep network classifier’s probability outputs for the class y being conditioned on were frequently around 0.5—large enough to be the dominant class, but not large enough to provide a strong guidance signal—and that upon inspection, generations were not consistent with the conditioning class y . Dhariwal and Nichol [DN21b] proposed to address this heuristically by incorporating an “inverse temperature” hyperparameter $\gamma > 0$ into the definition of the naive conditional denoiser (7.4.9), referring to the resulting conditional denoiser as having incorporated “**classifier guidance**” (CG):

$$\bar{\mathbf{x}}_{\theta}^{\text{CG}}(t, \mathbf{x}_t, y) = \bar{\mathbf{x}}_{\theta_d}(t, \mathbf{x}_t) + \gamma \frac{\sigma_t^2}{\alpha_t} \nabla_{\mathbf{x}_t} \langle \log f_{\theta_c}(t, \mathbf{x}_t), \mathbf{e}_y \rangle \quad (7.4.10)$$

with the case $\gamma = 1$ coinciding with (7.4.9).

Dhariwal and Nichol [DN21b] found that a setting $\gamma > 1$ performed best empirically. One possible interpretation for this is as follows: note that, in the context of the *true* likelihood term Equation (7.4.7), scaling by γ gives equivalently

$$\gamma \frac{\sigma_t^2}{\alpha_t} \nabla_{\boldsymbol{\xi}} \log p_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\nu} | \boldsymbol{\xi}) = \frac{\sigma_t^2}{\alpha_t} \nabla_{\boldsymbol{\xi}} \log (p_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\nu} | \boldsymbol{\xi})^{\gamma}), \quad (7.4.11)$$

⁷In Chapter 8, we review the process of training such a classifier in full detail.

which suggests the natural interpretation of the parameter γ performing (inverse) *temperature scaling* on the likelihood $p_{\mathbf{y}|\mathbf{x}_t}$, which is precise if we consider the renormalized distribution $p_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\nu}|\boldsymbol{\xi})^\gamma / \int p_{\mathbf{y}|\mathbf{x}_t}(\boldsymbol{\nu}'|\boldsymbol{\xi})^\gamma d\boldsymbol{\nu}'$. However, note that this *is not* a rigorous interpretation in the context of Equation (7.4.7), because the gradients are taken with respect to $\boldsymbol{\xi}$, and the normalization constant in the temperature-scaled distribution is in general a function of $\boldsymbol{\xi}$. Instead, the parameter γ should simply be understood as amplifying large values of the deep network classifier’s output probabilities $f_{\theta_c}(t, \mathbf{x}_t)$ relative to smaller ones, which effectively amplifies the guidance signal provided in cases where the deep network f assigns it the largest probability among the K classes.

Simplifications and Improvements via Classifier-Free Guidance

Nevertheless, classifier guidance does not address the second key drawback of the naive methodology: it is both cumbersome and wasteful to have to train an auxiliary classifier f_{θ_c} in addition to the unconditional denoiser $\boldsymbol{x}_{\theta_d}$, given that it is not possible to directly adapt a pretrained classifier due to the need for it to work well on noisy inputs \mathbf{x}_t and incorporate other empirically-motivated architecture modifications. In particular, Dhariwal and Nichol [DN21b] found that it was necessary to *explicitly design the architecture of the deep network implementing the classifier to match that of the denoiser*.

Moreover, from a purely practical perspective—trying to obtain the best possible performance from the resulting sampler—the best-performing configuration of classifier guidance-based sampling departs even further from the idealized and conceptually sound framework we have presented above. To obtain the best performance, Dhariwal and Nichol [DN21b] found it necessary to provide the class label y as an additional input to the denoiser $\boldsymbol{x}_{\theta_d}$. As a result, the idealized classifier-guided denoiser (7.4.10), derived by Dhariwal and Nichol [DN21b] as we have done above from the conditional posterior denoiser decomposition (7.4.7), is not exactly reflective of the best-performing denoiser in practice—such a denoiser actually combines a *conditional* denoiser for \mathbf{x}_t given y with an additional guidance signal from an auxiliary classifier!

This state of affairs, empirically motivated as it is, led Ho and Salimans [HS22a] in subsequent work to propose a more empirically pragmatic methodology, known as **classifier-free guidance (CFG)**. Instead of representing the conditional denoiser (7.4.7) as a weighted sum of an unconditional denoiser for \mathbf{x}_t with a log-likelihood correction term (with possibly modified weights, as in classifier guidance), they accept the apparent necessity of training a conditional denoiser for \mathbf{x}_t given y , as demonstrated by the experimental results of Dhariwal and Nichol [DN21b], and replace the log-likelihood gradient term with a correctly-weighted sum of this conditional denoiser with an *unconditional* denoiser for \mathbf{x} given \mathbf{x}_t .⁸

⁸That said, Ho and Salimans [HS22a] actually proposed to use a different weighting than what we present here, based on the fact that Dhariwal and Nichol [DN21b] heuristically replaced the unconditional denoiser in (7.4.7) with a conditional denoiser. In fact, the weighting we derive and present here reflects modern practice, and in particular is used in state-of-the-art

Deriving Classifier-Free Guidance. To see how this structure arises, we begin with an ‘idealized’ version of the classifier guidance denoiser $\bar{\mathbf{x}}_\theta^{\text{CG}}$ defined in (7.4.10), for which the denoiser $\bar{\mathbf{x}}_{\theta_d}$ and the classifier f_{θ_c} perfectly approximate their targets, via (7.4.7):

$$\bar{\mathbf{x}}_\theta^{\text{CG, ideal}}(t, \boldsymbol{\xi}, \nu) = \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}] + \gamma \frac{\sigma_t^2}{\alpha_t} \nabla_{\boldsymbol{\xi}} \log p_{y|\mathbf{x}_t}(\nu \mid \boldsymbol{\xi}). \quad (7.4.12)$$

We then use Bayes’ rule, in the form

$$\log p_{y|\mathbf{x}_t} = \log p_{\mathbf{x}_t|y} + \log p_y - \log p_{\mathbf{x}_t}, \quad (7.4.13)$$

together with Tweedie’s formula (Theorem 3.2, modified as in Equation (3.2.86)) to convert between score functions and denoisers, to obtain

$$\begin{aligned} \bar{\mathbf{x}}_\theta^{\text{CG, ideal}}(t, \boldsymbol{\xi}, \nu) &= \frac{1}{\alpha_t} \boldsymbol{\xi} + (1 - \gamma) \frac{\sigma_t^2}{\alpha_t} \nabla_{\boldsymbol{\xi}} \log p_{\mathbf{x}_t}(\boldsymbol{\xi}) + \gamma \frac{\sigma_t^2}{\alpha_t} \nabla_{\boldsymbol{\xi}} \log p_{\mathbf{x}_t|y}(\boldsymbol{\xi} \mid \nu) \\ &= (1 - \gamma) \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}] + \gamma \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}, y = \nu], \end{aligned} \quad (7.4.14)$$

where in the last line, we apply Equation (7.3.11). Now, Equation (7.4.14) suggests a natural approximation strategy: we combine a learned unconditional denoiser for \mathbf{x} given \mathbf{x}_t , as previously, with a learned *conditional* denoiser for \mathbf{x} given \mathbf{x}_t and y .

However, following Ho and Salimans [HS22a] and the common practice of training deep network denoisers, it is standard to use the *same* deep network to represent both the conditional and unconditional denoisers by introducing an additional label, which we will denote by \emptyset , to denote the “unconditional” case. This leads to the form of the CFG denoiser:

$$\bar{\mathbf{x}}_\theta^{\text{CFG}}(t, \mathbf{x}_t, y) = (1 - \gamma) \bar{\mathbf{x}}_\theta(t, \mathbf{x}_t, \emptyset) + \gamma \bar{\mathbf{x}}_\theta(t, \mathbf{x}_t, y). \quad (7.4.15)$$

To train a denoiser $\bar{\mathbf{x}}_\theta(t, \mathbf{x}_t, y^+)$ for use with classifier-free guidance sampling, where $y^+ \in \{1, \dots, K, \emptyset\}$, we proceed almost identically to the unconditional training procedure in Algorithm 3.2, but with two modifications:

1. When we sample from the dataset, we sample a pair (\mathbf{x}, y) rather than just a sample \mathbf{x} .
2. Every time we sample a pair from the dataset, we sample the augmented label y^+ via

$$y^+ = \begin{cases} \emptyset & \text{with probability } p_{\text{uncond}}; \\ y & \text{else.} \end{cases} \quad (7.4.16)$$

Here, $p_{\text{uncond}} \in [0, 1]$ is a new hyperparameter. This can be viewed as a form of dropout [SHK+14].

In this way, we train a conditional denoiser suitable for use in classifier-free guidance sampling. We summarize the overall sampling process for class-conditioned sampling with classifier-free guidance in Algorithm 7.2.

diffusion models such as Stable Diffusion 3.5 [EKB+24].

Algorithm 7.2 Conditional sampling with classification data, using class-conditioned denoiser.

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

Input: Class label $\nu \in \{1, \dots, K\}$ to condition on.

Input: A denoiser $\bar{\mathbf{x}}_\theta: \{t_\ell\}_{\ell=1}^L \times \mathbb{R}^D \times \{1, \dots, K, \emptyset\} \rightarrow \mathbb{R}^D$ for $p_{\mathbf{x}|y}$ and $p_{\mathbf{x}}$ (input \emptyset for $p_{\mathbf{x}}$).

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

Input: Guidance strength $\gamma \geq 0$ ($\gamma > 1$ preferred for performance).

Output: A sample $\hat{\mathbf{x}}$, approximately from $p_{\mathbf{x}|y}(\cdot | \nu)$.

1: **function** DDIMSAMPLERCONDITIONALCFG($\bar{\mathbf{x}}_\theta, \nu, \gamma, (t_\ell)_{\ell=0}^L$)

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

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

5: **end for**

6: **return** $\hat{\mathbf{x}}_{t_0}$

7: **end function**

Ho and Salimans [HS22a] reports strong empirical performance for class-conditional image generation with classifier-free guidance, and it has become a mainstay of the largest-scale practical diffusion models, such as Stable Diffusion [RBL+22] and its derivatives.

What Denoisers Does Classifier-Free Guidance Promote Learning?

As we have seen, the derivation of classifier-free guidance is rather opaque and empirically motivated, giving little insight into the mechanisms behind its strong performance. A number of theoretical works have attempted to address this [BN24b; LWQ25; WCL+24]. They provide explanations for some parts of the overall CFG methodology—itsself encompassing denoiser parameterization and training, as well as configuration of the guidance strength and performance at sampling time. Below, following our running theme in the book, we will give an interpretation in the simplifying setting of a Gaussian mixture model data distribution and denoiser. First, we consider *the effect of CFG*: what does a large weight $\gamma \gg 1$ promote?

Example 7.5. Let us recall the low-rank mixture of Gaussians data generating process we studied in Example 3.3 (and specifically, the form in Equation (3.2.69)). Given $K \in \mathbb{N}$ classes, we assume that

$$\mathbf{x} \sim \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\mathbf{0}, \mathbf{U}_k \mathbf{U}_k^\top), \quad (7.4.17)$$

where each $\mathbf{U}_k \in \mathcal{O}(D, P) \subseteq \mathbb{R}^{D \times P}$ is a matrix with orthogonal columns, and $P \ll D$. Moreover, we assume that the class label $y \in [K]$ is a deterministic function of \mathbf{x} mapping an example to its corresponding mixture component. We consider the simple noise process (3.2.1) studied in Chapter 3, so that $\mathbf{x}_t = \mathbf{x} + t\mathbf{g}$ for all $t \in [0, T]$, with $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

Applying the analysis in Example 3.3 (and the subsequent analysis of the low-rank case, culminating in Equation (3.2.80)), we obtain for the class-conditional optimal denoisers

$$\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}, y = \nu] = \frac{1}{1+t^2} \mathbf{U}_\nu \mathbf{U}_\nu^\top \boldsymbol{\xi} \quad (7.4.18)$$

for each $\nu \in [K]$, and for the optimal unconditional denoiser, we obtain

$$\mathbb{E}[\mathbf{x} \mid \mathbf{x}_t = \boldsymbol{\xi}] = \frac{1}{1+t^2} \sum_{k=1}^K \frac{\exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_k^\top \boldsymbol{\xi}\|_2^2\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_i^\top \boldsymbol{\xi}\|_2^2\right)} \mathbf{U}_k \mathbf{U}_k^\top \boldsymbol{\xi}. \quad (7.4.19)$$

As a result, we can express the CFG denoiser with guidance strength $\gamma > 1$ as

$$\begin{aligned} \bar{\mathbf{x}}^{\text{CFG, ideal}}(t, \mathbf{x}_t, y) = \frac{1}{1+t^2} \left((1-\gamma) \sum_{k=1}^K \frac{\exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_k^\top \mathbf{x}_t\|_2^2\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_i^\top \mathbf{x}_t\|_2^2\right)} \mathbf{U}_k \mathbf{U}_k^\top \right. \\ \left. + \gamma \mathbf{U}_y \mathbf{U}_y^\top \right) \mathbf{x}_t. \end{aligned} \quad (7.4.20)$$

This denoiser has a simple, interpretable form:

1. The first term, corresponding to *the unconditional denoiser*, performs denoising of the signal \mathbf{x}_t against an average of the denoisers associated with each subspace, weighted by how correlated \mathbf{x}_t is with each subspace.
2. The second term, corresponding to *the conditional denoiser*, simply performs denoising with the conditioning class's denoiser.

The CFG scheme averages these two denoisers. The effect of this averaging can be gleaned from the refactoring

$$\begin{aligned} \bar{\mathbf{x}}^{\text{CFG, ideal}}(t, \mathbf{x}_t, y) = \frac{1}{1+t^2} \left(\left[\gamma + (1-\gamma) \frac{\exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_y^\top \mathbf{x}_t\|_2^2\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_i^\top \mathbf{x}_t\|_2^2\right)} \right] \mathbf{U}_y \mathbf{U}_y^\top \right. \\ \left. + (1-\gamma) \sum_{k \neq y} \frac{\exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_k^\top \mathbf{x}_t\|_2^2\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_i^\top \mathbf{x}_t\|_2^2\right)} \mathbf{U}_k \mathbf{U}_k^\top \right) \mathbf{x}_t, \end{aligned} \quad (7.4.21)$$

where we combined the conditional denoiser with the corresponding $k = y$ term in the sum over $k \in [K]$. We have

$$\sum_{k=1}^K \frac{\exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_k^\top \mathbf{x}_t\|_2^2\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_i^\top \mathbf{x}_t\|_2^2\right)} = 1, \quad (7.4.22)$$

and each summand is nonnegative, hence also bounded above by 1. So we can conclude two regimes for the terms in Equation (7.4.21):

1. **Well-correlated regime:** In this case, suppose the noisy input \mathbf{x}_t correlates well with one of the subspaces \mathbf{U}_y . If \mathbf{x}_t correlates well with \mathbf{U}_y , then the normalized weight corresponding to the $k = y$ summand in the unconditional denoiser is near to 1. Then

$$\gamma + (1 - \gamma) \frac{\exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_y^\top \mathbf{x}_t\|_2^2\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_i^\top \mathbf{x}_t\|_2^2\right)} \approx 1, \quad (7.4.23)$$

all other weights are necessarily near to zero, and the CFG denoiser is approximately equal to the denoiser associated to the conditioning class y .

This implies: *if the CFG denoiser's input is highly correlated with one of the class subspaces, the denoiser is approximately equal to the corresponding class-conditional denoiser!* Hence, we expect the CFG denoising sampler to converge.

2. **Poorly-correlated regime:** In contrast, if \mathbf{x}_t *does not* correlate well with \mathbf{U}_y (say because t is large), then the normalized weight corresponding to the $k = y$ summand in the unconditional denoiser is near to 0. As a result,

$$\gamma + (1 - \gamma) \frac{\exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_y^\top \mathbf{x}_t\|_2^2\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_i^\top \mathbf{x}_t\|_2^2\right)} \approx \gamma, \quad (7.4.24)$$

and thus the guidance strength $\gamma \gg 1$ places a large positive weight on the denoiser associated to y .

Meanwhile, in the second term of Equation (7.4.21), any classes $k \neq y$ that are well-correlated with \mathbf{x}_t receive a large *negative* weight from the $1 - \gamma$ coefficient. This simultaneously has the effect of making the denoised signal vastly more correlated with the conditioning class y , and making it negatively correlated with the previous iterate (i.e., the iterate before denoising).

In other words, *CFG steers the iterative denoising process towards the conditioning class and away from the previous iterate*, a different dynamics from purely conditional sampling (i.e., the case $\gamma = 1$).

■

Example 7.5 shows that in the mixture of Gaussians setting, the CFG denoiser provides a strong bias towards the conditioning class, while preserving local correctness (hence convergence of sampling) when the noise level t is small. This is suggestive of why a large weight $\gamma \gg 1$ has been found essential in practice.

Next, the following example will demonstrate an insight into the *parameterization* of the denoiser in the presence of low-dimensional structure, again in the Gaussian mixture model.

Example 7.6. Consider again the mixture of Gaussians model studied in Example 7.5:

$$\mathbf{x} \sim \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\mathbf{0}, \mathbf{U}_k \mathbf{U}_k^\top), \quad (7.4.25)$$

where each $\mathbf{U}_k \in \mathcal{O}(D, P) \subseteq \mathbb{R}^{D \times P}$ is a matrix with orthogonal columns, and $P \ll D$. The class label $y \in [K]$ is a deterministic function of \mathbf{x} mapping an example to its corresponding mixture component. For this distribution, and again under the simple noise process (3.2.1) studied in Chapter 3, where $\mathbf{x}_t = \mathbf{x} + t\mathbf{g}$ for all $t \in [0, T]$, with $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, we saw in Example 7.5 that the ideal CFG denoiser takes the form

$$\begin{aligned} \bar{\mathbf{x}}^{\text{CFG, ideal}}(t, \mathbf{x}_t, y) = & \frac{1}{1+t^2} \left((1-\gamma) \sum_{k=1}^K \frac{\exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_k^\top \mathbf{x}_t\|_2^2\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)} \|\mathbf{U}_i^\top \mathbf{x}_t\|_2^2\right)} \mathbf{U}_k \mathbf{U}_k^\top \right. \\ & \left. + \gamma \mathbf{U}_y \mathbf{U}_y^\top \right) \mathbf{x}_t. \end{aligned} \quad (7.4.26)$$

Now we consider the problem of parameterizing a learnable denoiser $\mathbf{x}_\theta^{\text{CFG}}$ to represent the optimal denoiser (7.4.26). For tractability, we add an additional assumption associated to the subspaces \mathbf{U}_k being ‘distinguishable’ from one another, which is natural in practice: specifically, we assume that for any pair of indices $k, k' \in [K]$ with $k \neq k'$, we can find a set of K nonzero directions $\mathbf{v}_k \in \mathbb{R}^D$ such that

$$\mathbf{U}_k \mathbf{U}_k^\top \mathbf{v}_k = \mathbf{v}_k, \quad \mathbf{U}_{k'} \mathbf{U}_{k'}^\top \mathbf{v}_k = \mathbf{0}, \quad k' \neq k. \quad (7.4.27)$$

This is a slightly stronger assumption than simple distinguishability, but it should be noted that it is not overly restrictive: for example, it still allows the subspaces \mathbf{U}_k to have significant correlations with one another.⁹

These vectors \mathbf{v}_k can then be thought of as *embeddings* of the class label $y \in [K]$, and we can use them to define a more general operator that can represent both the unconditional and class-conditional denoisers. More precisely, consider

⁹More generally, this assumption is naturally formulated as an *incoherence condition* between the subspaces $\mathbf{U}_{[K]}$, familiar from the theory of compressive sensing [WM22].

the mapping

$$(\mathbf{x}_t, \mathbf{v}) \mapsto \sum_{k=1}^K \frac{\exp\left(\frac{1}{2t^2(1+t^2)} |\mathbf{x}_t^\top \mathbf{U}_k \mathbf{U}_k^\top \mathbf{v}|\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)} |\mathbf{x}_t^\top \mathbf{U}_i \mathbf{U}_i^\top \mathbf{v}|\right)} \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}_t. \quad (7.4.28)$$

Then if we plug in $\mathbf{v} = \mathbf{x}_t$ to the operator (7.4.28), we obtain a result proportional to the first term in (7.4.19), corresponding to the optimal unconditional denoiser for \mathbf{x}_t . On the other hand, if we substitute $\mathbf{v} = \mathbf{v}_y$ for some $y \in [K]$, we get by construction of the embedding vectors \mathbf{v}_k

$$\begin{aligned} & \sum_{k=1}^K \frac{\exp\left(\frac{1}{2t^2(1+t^2)} |\mathbf{x}_t^\top \mathbf{U}_k \mathbf{U}_k^\top \mathbf{v}_y|\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)} |\mathbf{x}_t^\top \mathbf{U}_i \mathbf{U}_i^\top \mathbf{v}_y|\right)} \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}_t \\ &= \frac{\exp\left(\frac{1}{2t^2(1+t^2)} |\mathbf{x}_t^\top \mathbf{v}_y|\right)}{\exp\left(\frac{1}{2t^2(1+t^2)} |\mathbf{x}_t^\top \mathbf{v}_y|\right) + K - 1} \mathbf{U}_y \mathbf{U}_y^\top \mathbf{x}_t \end{aligned} \quad (7.4.29)$$

$$+ \sum_{k \neq y} \frac{1}{\exp\left(\frac{1}{2t^2(1+t^2)} |\mathbf{x}_t^\top \mathbf{v}_y|\right) + K - 1} \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}_t. \quad (7.4.30)$$

We aim to simplify the preceding expression further. To this end, recall that when using this denoiser for sampling, we have $\mathbf{x}_t = \mathbf{x} + t\mathbf{g}$, with y denoting the label of \mathbf{x} . In other words, conditioned on y , we have that $\mathbf{x}_t = \mathbf{U}_y \mathbf{z} + t\mathbf{g}$, where \mathbf{z} is an independent P -dimensional Gaussian noise vector with zero mean and identity covariance. We have

$$\begin{aligned} |\mathbf{x}_t^\top \mathbf{v}_y| &= |\langle \mathbf{U}_y \mathbf{z}, \mathbf{v}_y \rangle + t\langle \mathbf{g}, \mathbf{v}_y \rangle| \\ &= |\langle \mathbf{z}, \mathbf{U}_y^\top \mathbf{v}_y \rangle + t\langle \mathbf{g}, \mathbf{v}_y \rangle|. \end{aligned}$$

Because \mathbf{z} and \mathbf{g} are independent Gaussian random variables with identity covariance matrices, it follows from *rotational invariance* of the Gaussian distribution (see [Ver18]) that

$$\begin{aligned} |\mathbf{x}_t^\top \mathbf{v}_y| &\stackrel{d}{=} | \|\mathbf{U}_y^\top \mathbf{v}_y\|_2 z + t\|\mathbf{v}_y\|_2 g | \\ &= | \|\mathbf{v}_y\|_2 z + t\|\mathbf{v}_y\|_2 g |, \end{aligned}$$

where z and g are independent *scalar* $\mathcal{N}(0, 1)$ random variables, and $\stackrel{d}{=}$ denotes equality in distribution (the random variables have the same distribution). The second line above follows from the definition of \mathbf{v}_y . Now, by independence, the random variable $z + tg \sim \mathcal{N}(0, (1+t^2))$. Then using the well-known result that $\mathbb{E}[|g|] = \sqrt{2/\pi}$ if $g \sim \mathcal{N}(0, 1)$, we conclude

$$\mathbb{E}[|\mathbf{x}_t^\top \mathbf{v}_y|] = \sqrt{(2/\pi)(1+t^2)} \|\mathbf{v}_y\|_2.$$

In practice, the quantity $|\mathbf{x}_t^\top \mathbf{v}_y|$ is close to its expectation *with overwhelming probability*, by Gaussian concentration (see, again, [Ver18]). On the other hand, we have by similar reasoning

$$\mathbb{E} \left[\|\mathbf{x}_t\|_2^2 \right] = P + t^2 D,$$

and properties of the Gaussian distribution again imply that $\|\mathbf{x}_t\|_2$ will concentrate around $\sqrt{P + t^2 D}$ with overwhelming probability. It is then sensible to renormalize the embeddings \mathbf{v}_y so that their dot-product with \mathbf{x}_t has the same scale: defining

$$\mathbf{v}_{t,y} = \frac{(P + t^2 D)}{\sqrt{(2/\pi)(1 + t^2)} \|\mathbf{v}_y\|_2} \mathbf{v}_y$$

for the normalized embeddings, our previous argument implies

$$|\mathbf{x}_t^\top \mathbf{v}_{t,y}| \approx (P + t^2 D)$$

with overwhelming probability. Returning our attention to our starting point, i.e. (7.4.28) evaluated at $(\mathbf{x}_t, \mathbf{v}_y)$ and our subsequent simplification, we have for the “softmax” argument

$$\frac{1}{2t^2(1 + t^2)} |\mathbf{x}_t^\top \mathbf{v}_y| \approx \frac{P + t^2 D}{2t^2(1 + t^2)}$$

with overwhelming probability with respect to \mathbf{x}_t . Hence if t is bounded and the subspace dimension P is sufficiently large,¹⁰ the following approximation holds with overwhelming probability:

$$\sum_{k=1}^K \frac{\exp\left(\frac{1}{2t^2(1+t^2)} \mathbf{x}_t^\top \mathbf{U}_k \mathbf{U}_k^\top \mathbf{v}_y\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)} \mathbf{x}_t^\top \mathbf{U}_i \mathbf{U}_i^\top \mathbf{v}_y\right)} \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}_t \approx \mathbf{U}_y \mathbf{U}_y^\top \mathbf{x}_t. \quad (7.4.31)$$

In words, **evaluating the operator (7.4.28) at the embedding for a class $\mathbf{v}_{t,y}$ yields the class-conditional denoiser for y !**

Thus, the family of operators (7.4.28) provides a unified way to parameterize the constituent operators in the optimal denoiser for \mathbf{x} *within a single ‘network’*. More precisely, it is enough to add the output of an instantiation of (7.4.28) with input $(\mathbf{x}_t, \mathbf{x}_t)$ to an instantiation with input $(\mathbf{x}_t, \mathbf{v}_{t,y})$. The resulting operator is a function of (t, \mathbf{x}_t, y) , and computationally, the subspaces $(\mathbf{U}_k)_{k=1}^K$ and embeddings $y \mapsto \mathbf{v}_{t,y}$ become its learnable parameters. ■

¹⁰For example, it suffices to consider a large-scale, asymptotic regime where $P, D \rightarrow \infty$ with their ratio P/D converging to a fixed constant. In such a scenario, the *aspect ratio* of the subspaces is fixed while their sizes go to infinity, which can act as a model for some discretization phenomena (e.g., a fixed scene sampled as an image with the number of pixels going to infinity).

Example 7.6 shows that in the special case of a low-rank mixture of Gaussians data distribution for \mathbf{x} with incoherent components, operators of the form

$$(\mathbf{x}_t, \mathbf{v}) \mapsto \sum_{k=1}^K \frac{\exp\left(\frac{1}{2t^2(1+t^2)} \mathbf{x}_t^\top \mathbf{U}_k \mathbf{U}_k^\top \mathbf{v}\right)}{\sum_{i=1}^K \exp\left(\frac{1}{2t^2(1+t^2)} \mathbf{x}_t^\top \mathbf{U}_i \mathbf{U}_i^\top \mathbf{v}\right)} \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}_t \quad (7.4.32)$$

provide a sufficiently rich class of operators to parameterize the ideal denoiser for noisy observations \mathbf{x}_t of \mathbf{x} when using classifier-free guidance. For such operators, the auxiliary input \mathbf{v} can be taken as either \mathbf{x}_t or a suitable embedding of the class label $y \mapsto \mathbf{v}_y$ in order to realize such a denoiser—as the example shows, the ideal denoiser is a weighted sum of such operators, with weights derived from the guidance strength $\gamma > 1$.

Based on the framework in Chapter 5, which develops deep network architectures suitable for transforming more general data distributions to structured representations using the low-rank mixture of Gaussians model as a primitive, it is natural to imagine that operators of the type (7.4.32) may be leveraged in denoisers for general data distributions \mathbf{x} with low-dimensional geometrically-structured components that are sufficiently distinguishable (say, incoherent) from one another. The next section demonstrates that this is indeed the case.

7.4.2 Caption-Conditioned Image Generation

In Section 7.4.1 we have seen, conceptually, how to set up denoisers for use with *classifier-free guidance*, the dominant practical methodology for performing conditional sampling via the diffusion-denoising paradigm. In short, such denoisers take as input the noisy data \mathbf{x}_t as well as either the conditioning signal \mathbf{y} or a “null input” \emptyset , which determines whether the denoiser should perform *unconditional denoising* or *conditional denoising* (conditioned on \mathbf{y}).

Given this design, the most important practical question is *how to realize this family of operators as a neural network?* An empirically-successful neural network layer design known as “cross attention”, which we will present in due course, provides a practical answer to this question. However, even before this, there are several more fundamental issues to sort out:

1. *Input modality mismatch:* In the most interesting practical applications, the data \mathbf{x} and the conditioning signal \mathbf{y} *do not lie in directly-comparable spaces*, although they share many underlying structural relationships. For example, we saw this in the case of a class label y in Section 7.4.1. More generally, one could consider a natural language description \mathbf{y} of an image \mathbf{x} (referred to as a “caption”), where \mathbf{y} would correspond to a sequence of characters in different alphabets.
2. *Parameterization and flexibility:* Downstream of the input modality issue, there is a fundamental question of whether the same basic architecture for the network can be used for different modalities, or whether this needs to be tuned in a very precise way per-modality.

3. *Training*: This entails the proper way to use our paired samples (\mathbf{x}, \mathbf{y}) to train and deploy the conditional denoisers. The design space is vast: one could imagine applying various modality-specific encoder-decoder architectures, as we have developed in Chapters 4 and 6, together with different denoising backbones in various combinations. The key questions here are of *expressivity*, getting the best conditional denoising performance; *efficiency*, getting the best performance per unit of training data; and *stability*, guaranteeing that the training process converges quickly without collapse.

We will cover these three issues in detail below, providing guidance on how to explore the design space and, where possible, generally-useful prescriptions.

Modality-Agnostic Input Embedding

Although we do not know the precise relationship between \mathbf{x} and \mathbf{y} in the general paired data setting, in all practical settings of interest, they have significant correlations. For example, when the conditioning signal \mathbf{y} is a “caption” describing the image \mathbf{x} in natural language, it is easy to imagine a human artist being able to produce an accurate sketch of \mathbf{x} based solely on a high-quality caption. In our probabilistic setting, one natural way to quantify these “correlations” is in terms of the mutual information, which we saw previously in Chapter 4 (Equation (4.1.12)):

$$I(\mathbf{x}; \mathbf{y}) = h(\mathbf{x}) - h(\mathbf{x} | \mathbf{y}). \quad (7.4.33)$$

The mutual information (7.4.33) is a purely information-theoretic quantity that does not depend on modality-specific aspects of \mathbf{x} and \mathbf{y} . Hence it offers a natural criterion for learning a useful representation of either input \mathbf{x}, \mathbf{y} : one simply seeks to learn, say for \mathbf{x} , a representation $\mathbf{z} = f(\mathbf{x})$ that preserves the mutual information, so that $I(\mathbf{z}; \mathbf{y}) \approx I(\mathbf{x}; \mathbf{y})$. That is:

$$\text{Learn } \mathbf{z} = f(\mathbf{x}) \text{ such that } I(\mathbf{z}; \mathbf{y}) \approx I(\mathbf{x}; \mathbf{y}). \quad (7.4.34)$$

Actually, we can be even more prescriptive here. For an encoder f which is deterministic (matching all the cases we have studied throughout Chapters 4 and 6), a fundamental result in information theory known as the *data processing inequality* tells us that processing \mathbf{x} with f to produce $\mathbf{z} = f(\mathbf{x})$ changes the mutual information in a predictable way: *it can only reduce it*.

Theorem 7.1 (Data Processing Inequality). *Consider random variables $\mathbf{y}, \mathbf{x}, \mathbf{z}$ defined such that \mathbf{z} is independent of \mathbf{y} , conditioned on \mathbf{x} .¹¹ Then the data processing inequality holds:*

$$I(\mathbf{x}; \mathbf{y}) \geq I(\mathbf{z}; \mathbf{y}). \quad (7.4.35)$$

Proof. Consider the mutual information $I(\mathbf{y}; \mathbf{x}, \mathbf{z})$ between \mathbf{y} and (\mathbf{x}, \mathbf{z}) . Because \mathbf{z} is independent of \mathbf{y} conditioned on \mathbf{x} , we have

$$h(\mathbf{y} | \mathbf{x}, \mathbf{z}) = h(\mathbf{y} | \mathbf{x}),$$

¹¹In other words, the random variables form a Markov chain $\mathbf{y} \rightarrow \mathbf{x} \rightarrow \mathbf{z}$.

so $I(\mathbf{y}; \mathbf{x}, \mathbf{z}) = I(\mathbf{y}; \mathbf{x})$. Meanwhile, we have by Bayes' rule that

$$p_{\mathbf{y}|\mathbf{x},\mathbf{z}} = p_{\mathbf{y}|\mathbf{z}} \frac{p_{\mathbf{y},\mathbf{x},\mathbf{z}}}{p_{\mathbf{y},\mathbf{z}}p_{\mathbf{x}|\mathbf{z}}},$$

which implies by the definition of conditional differential entropy that

$$I(\mathbf{y}; \mathbf{x}, \mathbf{z}) = I(\mathbf{y}; \mathbf{z}) - \mathbb{E}_{\mathbf{x},\mathbf{y},\mathbf{z}} \left[\log \left(\frac{p_{\mathbf{x}|\mathbf{z}}p_{\mathbf{y},\mathbf{z}}}{p_{\mathbf{y},\mathbf{x},\mathbf{z}}} \right) \right].$$

Applying Jensen's inequality to the expectation in the previous line then implies the claim, because log is a concave function (and the mutual information is symmetric, i.e. $I(\mathbf{x}; \mathbf{y}) = I(\mathbf{y}; \mathbf{x})$). \square

For a deterministic encoder f , $\mathbf{z} = f(\mathbf{x})$ is always independent of \mathbf{y} when conditioned on \mathbf{x} . The data processing inequality therefore suggests the following computational procedure to find the encoder f : *seek to maximize the mutual information between the representation and the conditioning signal*. Equation (7.4.34) then leads to the objective

$$\max_f I(f(\mathbf{x}); \mathbf{y}). \quad (7.4.36)$$

This is known as the *Infomax principle* [Lin88], and it dates to the earliest years of representation learning.

Computational considerations in applying the Infomax principle. There are two key conceptual issues to keep in mind when thinking about how to apply the (abstract) Infomax principle (7.4.36) for learning a representation. The first is that in general it is computationally hard to estimate the mutual information, which constitutes the objective in (7.4.36). As we discussed in the context of rate distortion and lossy coding when we introduced the mutual information in Chapter 4, the mutual information can be relaxed to remain well-defined even when \mathbf{x} and \mathbf{y} are low-dimensional and have differential entropy approaching $-\infty$ (e.g., by connection to the Gaussian rate reduction function that we have studied in Chapters 4 to 6). Moreover, empirically such relaxations have non-worst-case statistical properties in practice, as we saw in the experiments of Chapter 4; and the accuracy of these approximations can be improved by incremental processing, as we saw in Chapter 5. This makes a criterion like Equation (7.4.36) an excellent foundation for learning representations of paired data (\mathbf{x}, \mathbf{y}) .

The second issue to keep in mind is that for the objective (7.4.36) to represent a sensible criterion for learning f , it is necessary that f itself is parameterized in a specific way, and/or that some additional structural enforcement is done on the representations $f(\mathbf{x})$. Indeed, if f is completely unconstrained, there are simple trivial solutions to (7.4.36) that do not correspond to useful representations: for example, just take $f(\mathbf{x}) = \mathbf{x}$ to be the identity! We discussed principled techniques for structuring the representation space in Chapter 6 when we discussed

autoencoding and closed loop transcription; many of these techniques can be brought to bear in our present joint embedding setting, and we will detail this further later in the section, when we move to practical implementation.

Joint embedding with the Infomax principle. Now, the question of *joint* representation of (\mathbf{x}, \mathbf{y}) remains. Applying the data processing inequality once more, we obtain for another encoder g that produces a representation of \mathbf{y} that

$$I(\mathbf{x}; \mathbf{y}) \geq I(f(\mathbf{x}); \mathbf{y}) \geq I(f(\mathbf{x}); g(\mathbf{y})).$$

The joint objective

$$\max_{f, g} I(f(\mathbf{x}); g(\mathbf{y})) \quad (7.4.37)$$

is then a well-founded consequence of the Infomax principle.

To design the encoders f and g , recall the intuition that we obtained through our work in Example 7.6. There, we saw that in the case of a mixture-of-Gaussians model for \mathbf{x} , with the conditioning signal y identifying which specific mixture component generates the observation, a flexible denoiser architecture could be constructed via a suitable embedding of the class label y into the *same space* as the data \mathbf{x} . More precisely, such an embedding $\mathbf{u}_y \in \mathbb{R}^D$ allows, when correlated with the embeddings \mathbf{U}_x of $\mathbf{x} \in \mathbb{R}^D$ (generally a matrix) via $\mathbf{u}_y^\top \mathbf{U}_x$, to ‘pick out’ only the embeddings relevant to the correct mixture component of the distribution. How exactly to design these layers for use in neural networks being trained on more complex nonlinear high-dimensional data is a challenging question that we will dig into in more detail below. But a necessary condition is clear: the embeddings f and g should map their inputs to a *common embedding space*, say \mathbb{R}^d , to facilitate downstream comparison when performing conditional sampling. This leads to the objective

$$\max_{f: \mathbb{R}^{D_x} \rightarrow \mathbb{R}^d, g: \mathbb{R}^{D_y} \rightarrow \mathbb{R}^d} I(f(\mathbf{x}); g(\mathbf{y})). \quad (7.4.38)$$

The next step is to instantiate a suitable computationally-convenient relaxation of the mutual information in (7.4.38) for use with large-scale datasets.

Practical Prescription for Learning Joint Embeddings: CLIP

A concrete and highly influential realization of learning joint embeddings for paired data is the method of Contrastive Language-Image Pre-training (CLIP) [RKH+21b]. Take the simple case of image-text pairs $[(\mathbf{x}_i, \mathbf{y}_i)]_{i=1}^N$ sampled from a dataset \mathcal{D} , where \mathbf{x}_i is an image and \mathbf{y}_i is a text caption. CLIP trains two encoders—an image encoder $f_\theta(\mathbf{x})$ and a text encoder $g_\phi(\mathbf{y})$ —to map paired samples into a shared d -dimensional unit sphere (i.e. $\|f_\theta(\mathbf{x}_i)\|_2 = 1$ and $\|g_\phi(\mathbf{y}_i)\|_2 = 1$). To achieve this, we can adopt a simple symmetric cross-entropy loss that pushes the cosine similarity of matched pairs $(f_\theta(\mathbf{x}_i), g_\phi(\mathbf{y}_i))$ closer together while repelling unmatched pairs:

$$\mathcal{L}_{\text{CLIP}} = -\frac{1}{N} \mathbb{E}_{[(\mathbf{x}_i, \mathbf{y}_i)]_{i=1}^N \sim \mathcal{D}} \left[\sum_{i=1}^N \log \frac{\exp(f_\theta(\mathbf{x}_i)^\top g_\phi(\mathbf{y}_i))}{\sum_{j=1}^N \exp(f_\theta(\mathbf{x}_i)^\top g_\phi(\mathbf{y}_j))} \right]$$

$$+ \sum_{i=1}^N \log \frac{\exp(g_\phi(\mathbf{y}_i)^\top f_\theta(\mathbf{x}_i))}{\sum_{j=1}^N \exp(g_\phi(\mathbf{y}_i)^\top f_\theta(\mathbf{x}_j))} \Big], \quad (7.4.39)$$

From the perspective of mutual information maximization, CLIP is a tractable surrogate for the lower bound on the mutual information between the image and text representations:¹²

$$I(f_\theta(\mathbf{x}), g_\phi(\mathbf{y})) \geq -\mathcal{L}_{\text{CLIP}} + \log N. \quad (7.4.40)$$

Thus, CLIP provides a practical recipe for learning joint embeddings for paired data that increases the mutual information between the paired samples. Empirically, CLIP has been shown to be effective at learning joint embeddings for image-text pairs, and has been used in a variety of applications, including image captioning, image-text retrieval, image generation and beyond. We will see more specific examples in Chapter 8 (with key implementation details described in detail in Section 8.3).

Application to Caption-Conditioned Image Generation

Finally, we will describe an end-to-end system that integrates joint image-text encoders, learned with the CLIP loss as described above, with conditioning operators that are evocative of the mixture-of-Gaussians operator (7.4.32) we saw in Example 7.6. These techniques formed the basis for the original open-source Stable Diffusion implementation [RBL+22], where the embedding and subsequent conditioning is performed not on a class label, but a text prompt, which describes the desired image content (Figure 7.14). This section will give a high-level overview of StableDiffusion-style image generation. For full implementation details sufficient to re-implement such a system from scratch, we refer to Chapter 8: specifically, Section 8.11 describes the process of encoding strings of text to a sequence of integer IDs; Section 8.3 describes the training of the necessary image-text encoders; and Sections 8.6 and 8.7 describe the training of the (latent) diffusion model.

Stable Diffusion follows the conditional generation methodology we outline in Section 7.4.1, with two key modifications: (i) The conditioning signal is a tokenized text prompt $\mathbf{Y} \in \mathbb{R}^{D_{\text{text}} \times N}$, rather than a class label; (ii) Image denoising is performed in “latent” space rather than on raw pixels, using a specialized, pretrained variational autoencoder pair $f : \mathbb{R}^{D_{\text{img}}} \rightarrow \mathbb{R}^{d_{\text{img}}}$, $g : \mathbb{R}^{d_{\text{img}}} \rightarrow \mathbb{R}^{D_{\text{img}}}$ (see Sections 6.1.4 and 8.6), where f is the encoder and g is the decoder. For issue (i), in the context of the iterative conditional denoising framework we have developed in Section 7.4.1, this concerns the parameterization of the denoisers $\bar{z}_\theta(t, \mathbf{z}_t, \mathbf{Y}^+)$.¹³ Rombach et al. [RBL+22] implement text conditioning

¹²This result holds under the assumption that the pairs $(\mathbf{x}_i, \mathbf{y}_i)$ are i.i.d. within each batch of N [OLC+25, Lemma 1]. The tightness of the bound improves as N increases [OLV18]. However, for certain data distributions, it can remain loose for all N [MS19; WI20].

¹³In the setting of text conditioning, the ‘augmented’ label \mathbf{Y}^+ , which is either the encoded text prompt or \emptyset , denoting unconditional denoising, is often implemented by mapping \emptyset to the empty string “”, then encoding this text prompt with the tokenizer as usual. This gives a simple, unified way to treat conditional and unconditional denoising with text conditioning.

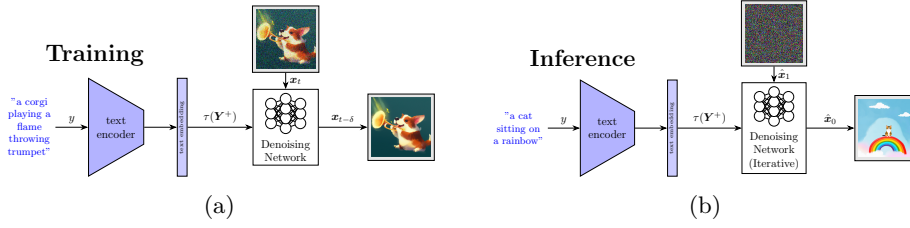


Figure 7.14: A high-level schematic of training and applying a text-to-image generative model, via conditional generation with a text prompt. **Left:** To train a text-to-image model, a large dataset of images paired with corresponding text captions is used. An encoder is used to map the captions to sequences of vectors, which are used as conditioning signals for a conditional denoiser, trained as described algorithmically in Section 7.4.1 (see Section 8.6 for implementation details). The text encoder may be pretrained and frozen, or jointly trained with the denoiser. **Right:** When applying a trained model, a desired text prompt is used as conditioning, then sampling is performed with the trained model, as in Algorithm 7.2 (*mutatis mutandis* for use with an encoded text prompt).

in the denoiser using a layer known as cross attention, inspired by the original encoder-decoder transformer architecture of Vaswani et al. [VSP+17b]. Cross attention is implemented as follows. We let $\tau : \mathbb{R}^{D_{\text{text}} \times N} \rightarrow \mathbb{R}^{d_{\text{model}} \times N_{\text{text}}}$ denote an encoding network for the text embeddings (often a causal transformer—see Section 8.11), and let $\psi : \mathbb{R}^{d_{\text{img}}} \rightarrow \mathbb{R}^{d_{\text{model}} \times N_{\text{img}}}$ denote the mapping corresponding to one of the intermediate representations in the denoiser.¹⁴ Here, N_{text} is the maximum tokenized text prompt length, and N_{img} roughly corresponds to the number of image channels (layer-dependent) in the representation, which is fixed if the input image resolution is fixed. Cross attention (with K heads, and no bias) is defined as

$$\text{MHCA}(z_t, \mathbf{Y}^+) = \mathbf{U}_{\text{out}} \begin{bmatrix} \text{SA}([\mathbf{U}_{\text{qry}}^1]^\top \psi(z_t), [\mathbf{U}_{\text{key}}^1]^\top \tau(\mathbf{Y}^+), [\mathbf{U}_{\text{val}}^1]^\top \tau(\mathbf{Y}^+)) \\ \vdots \\ \text{SA}([\mathbf{U}_{\text{qry}}^K]^\top \psi(z_t), [\mathbf{U}_{\text{key}}^K]^\top \tau(\mathbf{Y}^+), [\mathbf{U}_{\text{val}}^K]^\top \tau(\mathbf{Y}^+)) \end{bmatrix}, \quad (7.4.41)$$

where SA denotes the scaled dot-product attention operation in the transformer (which we recall in detail in Chapter 8: see Equations (8.2.17) and (8.2.18)), and $\mathbf{U}_*^k \in \mathbb{R}^{d_{\text{model}} \times d_{\text{attn}}}$ for $* \in \{\text{qry}, \text{key}, \text{val}\}$ (as well as the output projection \mathbf{U}_{out}) are the learnable parameters of the layer.

Notice that, by the definition of the self-attention operation, cross attention *outputs linear combinations of the value-projected text embeddings weighted by correlations between the image features and the text embeddings*. In the denoiser architecture used by Rombach et al. [RBL+22], self-attention residual blocks in the denoiser architecture, applied to the image representation at the current

¹⁴In practice, text-conditioned denoisers add cross attention layers at regular intervals within the forward pass of the denoiser, so ψ should be seen as layer-dependent, in contrast to τ .

layer and defined analogously to those in Equation (8.2.15) for the vision transformer, are followed by cross attention residual blocks of the form (7.4.41). Such a structure requires the text encoder τ to, in a certain sense, share some structure in its output with the image feature embedding ψ : this can be achieved via joint embedding, learned by a procedure such as CLIP as described above. Conceptually, this joint text-image embedding space and the cross attention layer itself bear a strong resemblance to the conditional mixture of Gaussians denoiser that we derived in the previous section (recall (7.4.28)), in the special case of a single token sequence. Deeper connections can be drawn in the multi-token setting following the rate reduction framework for deriving deep network architectures discussed in Chapter 5, and manifested in the derivation of the CRATE transformer-like architecture.

This same basic design has been further scaled to even larger model and dataset sizes, in particular in modern instantiations of Stable Diffusion [EKB+24], as well as in competing models such as FLUX.1 [LBB+25], Imagen [SCS+22], and DALL-E [RDN+22]. The conditioning mechanism of cross attention has also become ubiquitous in other applications, as in EgoAllo (Section 8.10) for conditioned pose generation and in Michelangelo (Section 8.8) for conditional 3D shape generation based on images or texts.

7.5 Conditional Inference with Measurement Self-Consistency

In this last section, we consider the more extreme, but actually ubiquitous, case for distribution learning in which we only have a set of observed samples $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$ of the data \mathbf{x} , but no samples of \mathbf{x} directly! In general, the observation $\mathbf{y} \in \mathbb{R}^d$ is of lower dimension than $\mathbf{x} \in \mathbb{R}^D$. To make the problem well-defined, we do assume that the observation model between \mathbf{y} and \mathbf{x} is known to belong to a certain family of analytical models, denoted as $\mathbf{y} = h(\mathbf{x}, \theta) + \mathbf{w}$, with θ either known or not known.

Let us first try to understand the problem conceptually with the simple case when the measurement function h is known and the observed $\mathbf{y} = h(\mathbf{x}) + \mathbf{w}$ is informative about \mathbf{x} . That is, we assume that h is surjective from the space of \mathbf{x} to that of \mathbf{y} and the support of the distribution $\mathbf{y}_0 = h(\mathbf{x}_0)$ is low-dimensional. This typically requires that the *extrinsic* dimension d of \mathbf{y} is higher than the *intrinsic* dimension of the support of the distribution of \mathbf{x} . Without loss of generality, we may assume that there exist functions:

$$F(\mathbf{x}) = \mathbf{0}, \quad G(\mathbf{y}) = \mathbf{0}. \quad (7.5.1)$$

Notice that here we may assume that we know $G(\mathbf{y})$ but not $F(\mathbf{x})$. Let $\mathcal{S}_{\mathbf{y}} \doteq \{\mathbf{y} \mid G(\mathbf{y}) = \mathbf{0}\}$ be the support of $p(\mathbf{y})$. In general, $h^{-1}(\mathcal{S}_{\mathbf{y}}) = \{\mathbf{x} \mid G(h(\mathbf{x})) = \mathbf{0}\}$ is a superset of $\mathcal{S}_{\mathbf{x}} \doteq \{\mathbf{x} \mid F(\mathbf{x}) = \mathbf{0}\}$. That is, we have $h(\mathcal{S}_{\mathbf{x}}) \subseteq \mathcal{S}_{\mathbf{y}}$.

7.5.1 Linear Measurement Models

First, for simplicity, let us consider that the measurement is a linear function of the data \mathbf{x} of interest:

$$\mathbf{y} = \mathbf{A}\mathbf{x}. \quad (7.5.2)$$

Here the matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is of full row rank and m is typically smaller than n . We assume \mathbf{A} is known for now. We are interested in how to learn the distribution of \mathbf{x} from such measurements. Since we no longer have direct samples of \mathbf{x} , we wonder whether we can still develop a denoiser for \mathbf{x} with observations \mathbf{y} . Let us consider the following diffusion process:

$$\mathbf{y}_t = \mathbf{y}_0 + t\mathbf{g}, \quad \mathbf{y}_0 = \mathbf{A}(\mathbf{x}_0), \quad (7.5.3)$$

where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

Without loss of generality, we assume \mathbf{A} is of full row rank, i.e., under-determined. Let us define the corresponding process \mathbf{x}_t as one that satisfies:

$$\mathbf{y}_t = \mathbf{A}\mathbf{x}_t. \quad (7.5.4)$$

From the denoising process of \mathbf{y}_t , we have

$$\mathbf{y}_{t-s} \approx \mathbf{y}_t + st\nabla \log p_t(\mathbf{y}_t). \quad (7.5.5)$$

Then we have:

$$\mathbf{A}\mathbf{x}_{t-s} \approx \mathbf{A}\mathbf{x}_t + st\nabla \log p_t(\mathbf{A}\mathbf{x}_t), \quad (7.5.6)$$

for a small $s > 0$. So \mathbf{x}_{t-s} and \mathbf{x}_t need to satisfy:

$$\mathbf{A}(\mathbf{x}_{t-s} - \mathbf{x}_t) \approx st\nabla \log p_t(\mathbf{A}\mathbf{x}_t). \quad (7.5.7)$$

Among all \mathbf{x}_{t-s} that satisfy the above constraint, we arbitrarily choose the one that minimizes the distance $\|\mathbf{x}_{t-s} - \mathbf{x}_t\|_2^2$. Therefore, we obtain a “denoising” process for \mathbf{x}_t :

$$\mathbf{x}_{t-s} \approx \mathbf{x}_t + st\mathbf{A}^\dagger \nabla \log p_t(\mathbf{A}\mathbf{x}_t). \quad (7.5.8)$$

Notice that this process does not sample from the distribution of \mathbf{x}_t . In particular, there are components of \mathbf{x} in the null space/kernel of \mathbf{A} that can never be recovered from observations. Thus more information is needed to recover the full distribution of \mathbf{x} , strictly speaking. But this recovers the component of \mathbf{x} that is orthogonal to the null space of \mathbf{A} .

7.5.2 Learning 3D World Model from 2D Images

In practice, the measurement model is often nonlinear or only partially known. A typical problem of this kind is actually behind how we can learn a working model of the external world from the images perceived, say through our eyes, telescopes or microscopes. In particular, humans and animals are able to build a model of the 3D world (or 4D for a dynamical world) through a sequence

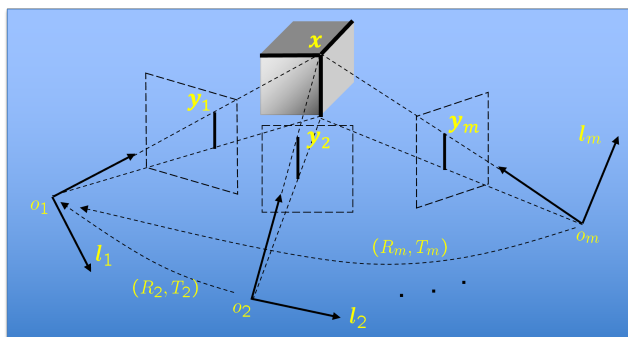


Figure 7.15: Relationship between a 3D object/scene and its 2D projections. Here we illustrate the projection of a point \mathbf{x} and a line intersecting the point.

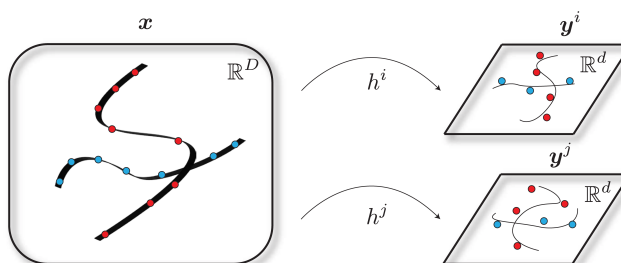


Figure 7.16: **Inference with distributed measurements.** We have a low-dimensional distribution \mathbf{x} (here, similarly to Figure 7.1, depicted as a union of two 2-dimensional manifolds in \mathbb{R}^3) and a measurement model $\mathbf{y}^i = h^i(\mathbf{x}) + \mathbf{w}^i$. As before, we want to infer various properties of the conditional distribution of \mathbf{x} given \mathbf{y} , where \mathbf{y} is the collection of all the measurements \mathbf{y}^i .

of its 2D projections—a sequence of 2D images (or stereo image pairs). The mathematical or geometric model of the projection is generally known:

$$\mathbf{y}^i = h(\mathbf{x}, \theta^i) + \mathbf{w}^i, \quad (7.5.9)$$

where $h(\cdot)$ represents a (perspective) projection of the 3D (or 4D) scene from a certain camera view at time t_i to a 2D image (or a stereo pair) and \mathbf{w} is some possibly additive small measurement noise. Figure 7.15 illustrates this relationship concretely, while Figure 7.16 illustrates the model problem in the abstract. A full exposition of geometry related to multiple 2D views of a 3D scene is beyond the scope of this book. Interested readers may refer to the book [MKS+04]. For now, all we need to proceed is that such projections are well understood and multiple images of a scene contain sufficient information about the scene.

In general, we would like to learn the distribution $p(\mathbf{x})$ of the 3D (or 4D)

world scene \mathbf{x} ¹⁵ from the perceived 2D images of the world so far. The primary function of such a (visual) world model is to allow us to recognize places where we had been before or predict what the current scene would look like in a future time at a new viewpoint.

Let us first examine the special but important case of stereo vision. In this case, we have two calibrated views of the 3D scene \mathbf{x} :

$$\mathbf{y}^0 = h(\mathbf{x}, \theta^0) + \mathbf{w}^0, \quad \mathbf{y}^1 = h(\mathbf{x}, \theta^1) + \mathbf{w}^1, \quad (7.5.10)$$

where parameters θ^0 and θ^1 for the view poses can be assumed to be known. \mathbf{y}^0 and \mathbf{y}^1 are two 2D-projections of the 3D scene \mathbf{x} . We may also assume that they have the same marginal distribution $p(\mathbf{y})$ and we have learned a diffusion and denoising model for it. That is, we know the denoiser:

$$\mathbb{E}[\mathbf{y} \mid \mathbf{y}_t = \boldsymbol{\nu}] = \boldsymbol{\nu} + t^2 \nabla_{\boldsymbol{\nu}} \log p_t(\boldsymbol{\nu}). \quad (7.5.11)$$

Or, furthermore, we may assume that we have a sufficient number of samples of stereo pairs $(\mathbf{y}^0, \mathbf{y}^1)$ and have also learned the joint distribution of the pairs. By a little abuse of notation, we also use $\mathbf{y} = h(\mathbf{x})$ to indicate the pair $\mathbf{y} = (\mathbf{y}^0, \mathbf{y}^1)$ and $p(\mathbf{y})$ as the learned probability distribution of the pair (say via a denoiser as above).

The main question now is: How to learn (a representation for) the distribution of the 3D scene \mathbf{x} from its two projections with known relationships? People might question the rationale for doing this: why is this necessary if the function $h(\cdot)$ is largely invertible? That is, the observation \mathbf{y} can largely determine the unknown \mathbf{x} , which is kind of the case for stereo—in general, two (calibrated) images contain sufficient information about the scene depth, from the given vantage point. However, 2D images are far from the most compact representation of the 3D scene as the same scene can produce infinitely many (highly correlated) 2D images or image pairs. In fact, a good representation of a 3D scene should be invariant to the viewpoint. Hence, a correct representation of the distribution of 3D scenes should be much more compact and structured than the distribution of 2D images, stereo pairs, or image-depth pairs.

Consider the (inverse) denoising process for the diffusion: $\mathbf{y}_t = \mathbf{y} + t\mathbf{g}$ in (7.5.11), where \mathbf{g} is standard Gaussian. From the denoising process of (7.5.11), we have

$$\mathbf{y}_{t-s} = \mathbf{y}_t + st \nabla_{\mathbf{y}} \log p_t(\mathbf{y}_t). \quad (7.5.12)$$

We try to find a corresponding “denoising” process of \mathbf{x}_t such that \mathbf{x} is related to \mathbf{y} as:

$$\mathbf{y} = h(\mathbf{x}). \quad (7.5.13)$$

Then we have:

$$h(\mathbf{x}_{t-s}) \approx h(\mathbf{x}_t) + st \nabla_{\mathbf{y}} \log p_t(h(\mathbf{x}_t)), \quad (7.5.14)$$

¹⁵Here by abuse of notation, we use \mathbf{x} to represent either a point in 3D or a sample of an entire 3D object or a scene that consists of many points.

for a small $s > 0$. Suppose $\mathbf{x}_{t-s} = \mathbf{x}_t + s\mathbf{v}$ for some vector \mathbf{v} and small increment s . We have

$$h(\mathbf{x}_{t-s}) \approx h(\mathbf{x}_t) + \frac{\partial h}{\partial \mathbf{x}}(\mathbf{x}_t) \cdot \mathbf{v}s \doteq h(\mathbf{x}_t) + \mathbf{A}(\mathbf{x}_t)\mathbf{v}s. \quad (7.5.15)$$

Hence, we have

$$\mathbf{A}(\mathbf{x}_t)\mathbf{v} = t\nabla_{\mathbf{y}} \log p_t(h(\mathbf{x}_t)). \quad (7.5.16)$$

Geometrically the vector \mathbf{v} in the domain of \mathbf{x} can be viewed as the pullback of the vector field $t\nabla \log p_t(\mathbf{y})$ under the map $\mathbf{y} = h(\mathbf{x})$. In general, as before, we may (arbitrarily) choose \mathbf{v} to be the minimum 2-norm vector that satisfies the pullback relationship. Hence, we can express $\hat{\mathbf{x}}_{t-s}$ approximately as:

$$\hat{\mathbf{x}}_{t-s} \approx \mathbf{x}_t + st\mathbf{A}(\mathbf{x}_t)^\dagger \nabla_{\mathbf{y}} \log p_t(h(\mathbf{x}_t)). \quad (7.5.17)$$

Remark 7.2 (Parallel Sensing and Distributed Denoising.). There is something very interesting about the above equation (7.5.17). It seems to suggest we could try to learn the distribution of \mathbf{x} through a process that is coupled with (many of) its (partial) observations:

$$\mathbf{y}^i = h^i(\mathbf{x}) + \mathbf{w}^i, i = 1, \dots, K. \quad (7.5.18)$$

In this case, we obtain a set of equations that the vector field \mathbf{v} in the domain of \mathbf{x} should satisfy:

$$\mathbf{A}^i(\mathbf{x}_t)\mathbf{v} = t\nabla_{\mathbf{y}^i} \log p_t(h^i(\mathbf{x}_t)), \quad (7.5.19)$$

where $\mathbf{A}^i(\mathbf{x}_t) = \frac{\partial h^i}{\partial \mathbf{x}}(\mathbf{x}_t)$. The final \mathbf{v} can be chosen as a “centralized” solution that satisfies all the above equations, or it could be chosen as a certain (stochastically) “aggregated” version of all \mathbf{v}^i :

$$\mathbf{v}^i = t\mathbf{A}^i(\mathbf{x}_t)^\dagger [\nabla_{\mathbf{y}^i} \log p_t(h^i(\mathbf{x}_t))], \quad i = 1, \dots, K, \quad (7.5.20)$$

that are computed in a parallel and distributed fashion? An open question here is exactly what the so-defined “denoising” process for \mathbf{x}_t converges to, even in the linear measurement model case. When would it converge to a distribution that has the same low-dimensional support as the original \mathbf{x}_0 , as \mathbf{y}_t converges to $\mathbf{y} = h(\mathbf{x}_0)$?

Visual World Model from Uncalibrated Image Sequences. In the above derivation, we have assumed that the measurement model $h(\cdot)$ is fully known. In the case of stereo vision, this is rather reasonable as the relative pose (and calibration) of the two camera views (or two eyes¹⁶) is usually known in advance. Hence, through the stereo image pairs, in principle we should be able to learn the distribution of 3D scenes, at least the ego-centric distribution of 3D scenes. However, the low-dimensional structures of the so-called learned distribution contain variation caused by changing the viewpoints. That is, the appearance

¹⁶The relative pose of our two eyes is well known to our brain.

of the stereo images varies when we change our viewpoints with respect to the same 3D scene. For many practical vision tasks (such as localization and navigation), it is important that we can decouple this variation of viewpoints from an invariant representation of (the distribution of) 3D scenes.

Remark 7.3. Note that the above goal aligns well with Klein’s Erlangen Program for modern geometry, which is to study invariants of a manifold under a group of transformations. Here, we may view the manifold of interest as the distribution of ego-centric representations of 3D scenes. We have learned that it admits a group of three-dimensional rigid-body motion acting on it. It is remarkable that our brain has learned to effectively decouple such transformations from the observed 3D world.

Notice that we have studied learning representations that are invariant to translation and rotation in a limited setting in Chapter 5. We know that the associated compression operators take the necessary form of (multi-channel) convolutions, hence leading to the (deep) convolutional neural networks. Nevertheless, operators that are associated with compression or denoising that are invariant to more general transformation groups remain elusive to characterize [CW16b]. For the 3D Vision problem in its most general setting, we know the change of our viewpoints can be well modeled as a rigid-body motion. However, the exact relative motion of our eyes between different viewpoints is usually not known. More generally, there could also be objects (e.g., cars, humans, hands) moving in the scene and we normally do not know their motion either. How can we generalize the problem of learning the distribution of 3D scenes with calibrated stereo pairs to such more general settings? More precisely, we want to learn a compact representation \mathbf{x} of the 3D scenes that is invariant to the camera/eye motions. Once such a representation is learned, we could sample and generate a 3D scene and render images or stereo pairs from arbitrary poses.

To this end, note that we can model a sequence of stereo pairs as:

$$\mathbf{y}^k = h(\mathbf{x}^k, \theta^k), \quad k = 1, \dots, K, \quad (7.5.21)$$

where $h(\cdot)$ represents the projection map from 3D to 2D. θ^k denotes the rigid-body motion parameters of the k th view, with respect to some canonical frame in the world. \mathbf{x}^k represents the 3D scene at time k . If the scene is static, \mathbf{x}^k should all be the same $\mathbf{x}^k = \mathbf{x}$. To simplify the notation, we may denote the set of k equations as one:

$$\mathbf{Y} = H(\mathbf{x}, \Theta). \quad (7.5.22)$$

We may assume that we are given many samples of such stereo image sequences $\{\mathbf{Y}_i\}$. The problem is how to recover the associated motion sequence $\{\Theta_i\}$ and learn the distribution of the scene \mathbf{x} (that is invariant to the motion). To the best of our knowledge, this remains an open challenging problem, probably as the final frontier for the 3D Vision problem.

7.6 Summary and Notes

Measurement matching without clean samples. In our development of conditional sampling, we considered measurement matching under an observation model (7.3.9), where we assume that we have paired data (\mathbf{x}, \mathbf{y}) —i.e., ground truth for each observation \mathbf{y} . In many practically relevant inverse problems, this is not the case: one of the most fundamental examples is in the context of compressed sensing, which we recalled in Chapter 2, where we need to reconstruct \mathbf{x} from \mathbf{y} using prior knowledge about \mathbf{x} (i.e., sparsity). In the setting of denoising-diffusion, we have access to an implicit prior for \mathbf{x} via the learned denoisers $\bar{\mathbf{x}}_\theta(t, \boldsymbol{\xi})$. Can we still perform conditional sampling without access to ground truth samples \mathbf{x} ?

For intuition as to why this might be possible, we recall a classical example from statistics known as Stein’s unbiased risk estimator (SURE). Under an observation model $\mathbf{x}_t = \mathbf{x} + t\mathbf{g}$ with $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $t > 0$, it turns out that for any weakly differentiable $f : \mathbb{R}^D \rightarrow \mathbb{R}^D$,

$$\mathbb{E}_{\mathbf{g}} \left[\|\mathbf{x} - f(\mathbf{x} + t\mathbf{g})\|_2^2 \right] = \mathbb{E}_{\mathbf{g}} \left[\|\mathbf{x} + t\mathbf{g} - f(\mathbf{x} + t\mathbf{g})\|_2^2 + 2t^2 \nabla \cdot f(\mathbf{x} + t\mathbf{g}) \right] - t^2 D, \quad (7.6.1)$$

where $\nabla \cdot$ denotes the divergence operator:

$$\nabla \cdot f = \sum_{i=1}^D \partial_i f_i.$$

The \mathbf{x} -dependent part of the RHS of Equation (7.6.1) is called Stein’s unbiased risk estimator (SURE). If we take expectations over \mathbf{x} in Equation (7.6.1), note that the RHS can be written as an expectation with respect to \mathbf{x}_t —in particular, the mean-squared error of *any denoiser* f can be estimated *solely from noisy samples!* This remarkable fact, in refined forms, constitutes the basis for many practical techniques for performing image restoration, denoising-diffusion, etc. using only noisy data: notable examples include the “noise2noise” paradigm [LMH+18] and Ambient Diffusion [DSD+23a].

As a fun aside, we point out that Equation (7.6.1) leads to an alternate proof of Tweedie’s formula (Theorem 3.2). At a high level, one takes expectations over \mathbf{x} and expresses the main part of the RHS of Equation (7.6.1) equivalently, via integration by parts, as

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}_t} \left[\|\mathbf{x}_t - f(\mathbf{x}_t)\|_2^2 + 2t^2 \nabla \cdot f(\mathbf{x}_t) \right] \\ = & \mathbb{E}_{\mathbf{x}_t} \left[\|\mathbf{x}_t - f(\mathbf{x}_t)\|_2^2 \right] - 2t^2 \int \langle \nabla p_{\mathbf{x}_t}(\boldsymbol{\xi}), f(\boldsymbol{\xi}) \rangle d\boldsymbol{\xi}. \end{aligned} \quad (7.6.2)$$

This is a quadratic function of f , and formally taking derivatives gives that the optimal f satisfies Tweedie’s formula (Theorem 3.2). This argument can be made rigorous using basic ideas from the calculus of variations.

Corrections to the Diffusion Posterior Sampling (DPS) approximation. In Example 7.4 and in particular in Figure 7.11, we pointed out a limitation of the DPS approximation Equation (7.3.26) at small levels of measurement noise. This limitation is well-understood, and a principled approach to ameliorating it has been proposed by Rozet et al. [RAL+24]. The approach involves incorporating an additional estimate for the variance of the noisy posterior $p_{\mathbf{x}|\mathbf{x}_t}$ to Equation (7.3.26)—we refer to the paper for details. Natural estimates for the posterior variance are slightly less scalable than DPS itself due to the need to invert an affine transformation of the Jacobian of the posterior denoiser $\mathbb{E}[\mathbf{x} | \mathbf{x}_t = \boldsymbol{\xi}]$ (a large matrix). This is done relatively efficiently by Rozet et al. [RAL+24] using automatic differentiation and an approximation for the inverse based on conjugate gradients. It seems that it should be possible to improve further over this approach (say, using classical ideas from second-order optimization).

More about measurement matching and diffusion models for inverse problems. Diffusion models have become an extremely popular tool for solving inverse problems arising in scientific applications. Many more methods beyond the simple DPS algorithm we have presented in Algorithm 7.1 have been developed and continue to be developed, as the area is evolving rapidly. Popular and performant classes of approaches beyond DPS, which we have presented due to its generality, include variable splitting approaches like DAPS [ZCB+24], which allow for specific measurement constraints to be enforced much more strongly than in DPS, and exact approaches that can avoid the use of approximations as in DPS, such as TDS [WTN+23]. For more on this area, we recommend [ZCZ+25], which functions simultaneously as a survey and a benchmark of several popular methods on specific scientific inverse problem datasets.

History of the Infomax principle. The Infomax principle [Lin88] that we introduced in our discussion of joint embedding learning and CLIP in Section 7.4.2 is a well-known and important foundation for much of modern research into unsupervised representation learning. For further reading, we refer to [OLV18], which our conceptual discussion in Section 7.4.2 expands upon, and [OLC+25], which details further important theoretical issues associated with optimizing Infomax-type objectives from finite samples.

7.7 Exercises and Extensions

- Exercise 7.1* (Posterior Variance Correction to DPS). 1. Using the code provided in the book GitHub for implementing Figure 7.11, implement the posterior variance correction proposed by Rozet et al. [RAL+24].
2. Verify that it ameliorates the posterior collapse at low noise variance issue observed in Figure 7.11.

3. Discuss any issues of sampling correctness that are retained or introduced by the corrected method, as well as its efficiency, relative to diffusion posterior sampling (DPS).

Exercise 7.2 (Conditional Sampling on MNIST). 1. Train a simple classifier for the MNIST dataset, using an architecture of your choice. Additionally train a denoiser suitable for use in conditional sampling (Algorithm 7.2, since this denoiser can be used for unconditional denoising as well).

2. Integrate the classifier into a conditional sampler based on classifier guidance, as described in the first part of Section 7.4.1. Evaluate the resulting samples in terms of faithfulness to the conditioning class (visually; in terms of nearest neighbor; in terms of the output of the classifier).
3. Integrate the classifier into a conditional sampler based on classifier-free guidance, as described in Section 7.4.1 and Algorithm 7.2. Perform the same evaluation as in the previous step, and compare the results.
4. Repeat the experiment on the CIFAR-10 dataset.