

Chapter 2

Learning Linear and Independent Structures

“The art of doing mathematics consists in finding that special case which contains all the germs of generality.”

– David Hilbert

Real data has low-dimensional structure. To see why this is true, let us consider the unassuming case of static on a TV when the satellite isn’t working. At each frame (approximately every $\frac{1}{30}$ seconds), the RGB static on a screen of size $H \times W$ is, roughly, sampled independently from a uniform distribution on $[0, 1]^{3 \times H \times W}$. In theory, the static *could* resolve to a natural image on any given frame, but even if you spend a thousand years looking at the TV screen, it will not. This discrepancy is explained by the fact that the set of $H \times W$ natural images takes up a vanishingly small fraction of the hypercube $[0, 1]^{3 \times H \times W}$. In particular, it is extremely low-dimensional, compared to the ambient space dimension. Similar phenomena occur for all other types of natural data, such as text, audio, and video. Thus, when we design systems and methods to process natural data and learn its structure or distribution, this is a central property of natural data which we need to take into account.

Therefore, our central task is to learn a distribution that has low intrinsic dimension in a high-dimensional space. In the remainder of this chapter, we discuss several *classical* methods to perform this task for several somewhat *idealistic* models for the distribution, namely models that are geometrically linear or statistically independent. While these models and methods are important and useful in their own right, we discuss them here as they motivate, inspire, and serve as a predecessor or analogue to more modern methods for more general distributions that involve deep (representation) learning.

Our main approach (and general problem formulation) can be summarized as:

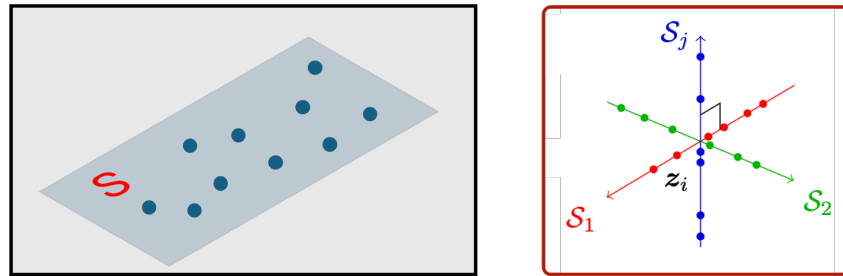


Figure 2.1: Data on a single low-dimensional subspace (left), say $\mathcal{S} = \text{col}(\mathbf{U})$, or a mixture of low-dimensional subspaces (right), say $\mathcal{S}_j = \text{col}(\mathbf{U}_j)$.

Problem: *Given one or several (noisy or incomplete) observations of a ground truth sample from the data distribution, obtain an estimate of this sample.*

This approach underpins several classical methods for data processing, which we discuss in this chapter.

- Section 2.1 — Principal Components Analysis (PCA): Given noisy samples from a distribution supported on *one low-dimensional subspace*, obtain an estimate of the true sample that lies on this subspace.
- Section 2.2 — Complete Dictionary Learning and Independent Components Analysis (ICA): Given noisy samples from a distribution supported on *a union (not the span) of a few low-dimensional subspaces*, obtain an estimate of the true samples.
- Section 2.3 — Sparse Coding and Dictionary Learning: Given noisy samples from *a distribution supported on combinations of a few incoherent vectors*, such as the coordinate axes, obtain an estimate of the true sample, which also has this property.

In this chapter, as described above, we make simplifying modeling assumptions that essentially assume (the support of) the data distribution has geometrically (nearly, piece-wise) linear structures and statistically independent components, as illustrated in Figure 2.1. In Chapter 1, we have referred to such models of the data as “analytical models”. These modeling assumptions allow us to derive efficient algorithms with provable efficiency guarantees, both in terms of data and computational complexity, for processing data at scale. However, they are imperfect models for often-complex real-world data distributions, and so their underlying assumptions only approximately hold. This means that the guarantees afforded by detailed analysis of these algorithms also only approximately hold in the case of real data.

Nonetheless, the techniques discussed in this chapter are useful in their own right, and beyond that, serve as the “special case with the germ of generality”,

so to speak, in that they present a guiding motivation and intuition for the more general paradigms within (deep) learning of more general distributions that we later address. As we will soon reveal in later chapters, in the deep learning era, modern “data-driven” or “non-parametric” approaches to learn (low-dimensional) data distributions essentially adopt the same concepts and ideas to learn.

A Note on Notation. We encourage the reader to consult the [Notation](#) chapter, particularly the “Modeling and Optimization Notation” section, as they begin to go through this chapter. We use the same notational conventions for modeling data and optimizing the parameters of these models throughout the book, and the simple, explicit models we study in this chapter help make the ‘semantic’ content of these conventions clear.

2.1 A Low-Dimensional Subspace

2.1.1 Principal Components Analysis (PCA)

Perhaps the simplest modeling assumption possible for low-dimensional structure is the so-called *low-rank* assumption. Letting D be the dimension of our data space, we assume that our data belong to a low-dimensional subspace of dimension $d \ll D$, possibly plus some small disturbances.¹ This ends up being a nearly valid assumption for some surprisingly complex data, such as images of handwritten digits and face data [RD03] as shown in Figure 2.2, yet as we will see, it will lend itself extremely well to comprehensive analysis.

Problem Formulation. To write this in mathematical notation, we represent a subspace $\mathcal{S} \subseteq \mathbb{R}^D$ of dimension d by an orthonormal matrix $\mathbf{U} \in \mathcal{O}(D, d) \subseteq \mathbb{R}^{D \times d}$ such that the columns of \mathbf{U} span \mathcal{S} . Then, we say that our data $\{\mathbf{x}_i\}_{i=1}^N \subseteq \mathbb{R}^D$ have (approximate) low-rank structure if there exists an orthonormal matrix $\mathbf{U} \in \mathcal{O}(D, d)$, vectors $\{\mathbf{z}_i\}_{i=1}^N \subseteq \mathbb{R}^d$, and *small* vectors $\{\boldsymbol{\varepsilon}_i\}_{i=1}^N \subseteq \mathbb{R}^D$ such that

$$\mathbf{x}_i = \mathbf{U}\mathbf{z}_i + \boldsymbol{\varepsilon}_i, \quad \forall i \in [N]. \quad (2.1.1)$$

Here $\boldsymbol{\varepsilon}_i$ are disturbances that prevent the data from being perfectly low rank; their presence in our model allows us to quantify the degree to which our analysis remains relevant in the presence of deviations from our model. The true supporting subspace is $\mathcal{S} \doteq \text{col}(\mathbf{U})$, the span of the columns of \mathbf{U} . To process all we can from these data, we need to recover \mathcal{S} ; to do this it is sufficient to recover \mathbf{U} , also called the *principal components*. Fortunately, this is a computationally tractable task named **Principal Component Analysis**, and we discuss now how to solve it.

¹Plus, potentially, a non-trivial offset from $\mathbf{0}$; in practice we account for this offset by subtracting the empirical mean over data samples from each sample, and in the following Chapter we pay no more heed to it.



Figure 2.2: Images of aligned human faces and handwritten digits. Despite the seemingly large variety in their appearances, each set of these data spans (approximately) a very low-dimensional (nearly) linear subspace. This can be made technically precise in the notation of (2.1.1): for a relatively small subspace dimension d (say, ≈ 20 , compared to the dimension $D = 784$ of MNIST digits), the error vectors ε_i are small for every $i \in [N]$ (where $N = 6 \cdot 10^4$ for MNIST is the dataset size). We can *learn* the optimal such subspace from data.

Given data $\{\mathbf{x}_i\}_{i=1}^N$ distributed as in (2.1.1), we aim to recover the model \mathbf{U} . A natural approach is to find the subspace \mathbf{U}^* and latent vectors $\{\mathbf{z}_i^*\}_{i=1}^N$ which yield the best approximation $\mathbf{x}_i \approx \mathbf{U}^* \mathbf{z}_i^*$. Namely, we aim to solve the problem

$$\min_{\tilde{\mathbf{U}}, \{\tilde{\mathbf{z}}_i\}_{i=1}^N} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{U}} \tilde{\mathbf{z}}_i\|_2^2, \quad (2.1.2)$$

where $\tilde{\mathbf{U}}$ is constrained to be an orthonormal matrix, as above. We will omit this constraint in similar statements below for the sake of concision.

Subspace Encoding-Decoding via Denoising. To simplify this problem, for a fixed $\tilde{\mathbf{U}}$, we have (proof as exercise):

$$\min_{\{\tilde{\mathbf{z}}_i\}_{i=1}^N} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{U}} \tilde{\mathbf{z}}_i\|_2^2 = \frac{1}{N} \sum_{i=1}^N \min_{\tilde{\mathbf{z}}_i} \|\mathbf{x}_i - \tilde{\mathbf{U}} \tilde{\mathbf{z}}_i\|_2^2 \quad (2.1.3)$$

$$= \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top \mathbf{x}_i\|_2^2. \quad (2.1.4)$$

That is, the optimal solution $(\mathbf{U}^*, \{\mathbf{z}_i^*\}_{i=1}^N)$ to the above optimization problem has $\mathbf{z}_i^* = (\mathbf{U}^*)^\top \mathbf{x}_i$.

Now, we can write the original optimization problem in $\tilde{\mathbf{U}}$ and $\{\tilde{\mathbf{z}}_i\}_{i=1}^N$ as an optimization problem just over $\tilde{\mathbf{U}}$, i.e., to obtain the basis \mathbf{U}^* and compact

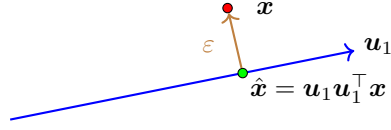


Figure 2.3: **Geometry of PCA.** A data point \mathbf{x} (red) is projected onto the one-dimensional learned subspace spanned by the unit basis vector \mathbf{u}_1 (blue arrow). The projection $\mathbf{U}\mathbf{U}^\top \mathbf{x} = \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}$ (green) is the denoised version of \mathbf{x} using the low-dimensional structure given by \mathbf{u}_1 , and ε (brown arrow) represents the projection residual or noise.

codes $\{\mathbf{z}_i^*\}_{i=1}^N$ it suffices to solve either of the two following equivalent problems:

$$\min_{\tilde{\mathbf{U}}, \{\tilde{\mathbf{z}}_i\}_{i=1}^N} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{U}} \tilde{\mathbf{z}}_i\|_2^2 = \min_{\tilde{\mathbf{U}}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top \mathbf{x}_i\|_2^2. \quad (2.1.5)$$

Note that the problem on the right-hand side of (2.1.5) is a *denoising* problem: given noisy observations \mathbf{x}_i of low-rank data, we aim to find the *noise-free* copy of \mathbf{x}_i , which under the model (2.1.1) is $\mathbf{U}\mathbf{z}_i$. That is, the denoised input $\hat{\mathbf{x}}_i = \mathbf{U}\mathbf{U}^\top \mathbf{x}_i$. Notice that this is the point on the subspace that is closest to \mathbf{x}_i , as visualized in Figure 2.3. Here by solving the equivalent problem of finding the best subspace, parameterized by the learned basis \mathbf{U}^* , we learn an approximation to the *denoiser*, i.e., the projection matrix $\mathbf{U}^*(\mathbf{U}^*)^\top \approx \mathbf{U}\mathbf{U}^\top$ that projects the noisy data point \mathbf{x}_i onto the subspace \mathcal{S} .

Putting the above process together, we essentially obtain a simple encoding-decoding scheme that maps a data point \mathbf{x} in \mathbb{R}^D to a lower-dimensional (latent) space \mathbb{R}^d and then back to \mathbb{R}^D :

$$\mathbf{x} \xrightarrow{\mathcal{E}=(\mathbf{U}^*)^\top} \mathbf{z} \xrightarrow{\mathcal{D}=\mathbf{U}^*} \hat{\mathbf{x}}. \quad (2.1.6)$$

Here, $\mathbf{z} \in \mathbb{R}^d$ can be viewed as the low-dimensional compact code (or a latent representation) of a data point $\mathbf{x} \in \mathbb{R}^D$ and the learned subspace basis \mathbf{U}^* as the associated codebook whose columns are the (learned) optimal code words. The process achieves the function of denoising \mathbf{x} by projecting it onto the subspace spanned by \mathbf{U}^* .

Computing the Subspace Basis. For now, we continue our calculation. Let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ be the matrix whose columns are the observations \mathbf{x}_i . We have (proof as exercise)

$$\arg \min_{\tilde{\mathbf{U}}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top \mathbf{x}_i\|_2^2 = \arg \max_{\tilde{\mathbf{U}}} \frac{1}{N} \sum_{i=1}^N \|\tilde{\mathbf{U}}^\top \mathbf{x}_i\|_2^2 \quad (2.1.7)$$

$$= \arg \max_{\tilde{\mathbf{U}}} \text{tr} \left\{ \tilde{\mathbf{U}}^\top \left(\frac{\mathbf{X}\mathbf{X}^\top}{N} \right) \tilde{\mathbf{U}} \right\}. \quad (2.1.8)$$

Thus, to compute the principal components, we find the orthogonal matrix $\tilde{\mathbf{U}}$ which maximizes the term $\text{tr}(\tilde{\mathbf{U}}^\top(\mathbf{X}\mathbf{X}^\top/N)\tilde{\mathbf{U}})$. We can prove via induction that this matrix \mathbf{U}^* has columns which are the *top d unit eigenvectors* of $\mathbf{X}\mathbf{X}^\top/N$. We leave the whole proof to the reader in Exercise 2.1, but we handle the base case of the induction here. Suppose that $d = 1$. Then we only have a single unit vector \mathbf{u} to recover, so the above problem reduces to

$$\max_{\tilde{\mathbf{u}}: \|\tilde{\mathbf{u}}\|_2=1} \tilde{\mathbf{u}}^\top(\mathbf{X}\mathbf{X}^\top/N)\tilde{\mathbf{u}}. \quad (2.1.9)$$

This is the so-called *Rayleigh quotient* of $\mathbf{X}\mathbf{X}^\top/N$. By invoking the spectral theorem we diagonalize $\mathbf{X}\mathbf{X}^\top/N = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$, where \mathbf{Q} is orthogonal and $\mathbf{\Lambda}$ is diagonal with non-negative entries. Hence

$$\tilde{\mathbf{u}}^\top(\mathbf{X}\mathbf{X}^\top/N)\tilde{\mathbf{u}} = \tilde{\mathbf{u}}^\top\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top\tilde{\mathbf{u}} = (\mathbf{Q}^\top\tilde{\mathbf{u}})^\top\mathbf{\Lambda}(\mathbf{Q}^\top\tilde{\mathbf{u}}). \quad (2.1.10)$$

Since \mathbf{Q} is an invertible orthogonal transformation, $\mathbf{Q}^\top\tilde{\mathbf{u}}$ is a unit vector, and optimizing over $\tilde{\mathbf{u}}$ is equivalent to optimizing over $\tilde{\mathbf{w}} \doteq \mathbf{Q}^\top\tilde{\mathbf{u}}$. Hence, we can write

$$\tilde{\mathbf{u}}^\top(\mathbf{X}\mathbf{X}^\top/N)\tilde{\mathbf{u}} = \tilde{\mathbf{w}}^\top\mathbf{\Lambda}\tilde{\mathbf{w}}, \quad (2.1.11)$$

whose optimal solutions \mathbf{w}^* among unit vectors are one-hot vectors whose only nonzero (hence unit) entry is in one of the indices corresponding to the largest eigenvalue of $\mathbf{X}\mathbf{X}^\top/N$. This means that $\mathbf{u}^* = \mathbf{Q}\mathbf{w}^*$, the optimal solution to the original problem, corresponds to a unit eigenvector of $\mathbf{X}\mathbf{X}^\top/N$ (i.e., column of \mathbf{Q}) which corresponds to the largest eigenvalue. Suitably generalizing this to the case $d > 1$, and summarizing all the previous discussion, we have the following informal Theorem.

Theorem 2.1. *Suppose that our dataset $\{\mathbf{x}_i\}_{i=1}^N \subseteq \mathbb{R}^D$ can be written as*

$$\mathbf{x}_i = \mathbf{U}\mathbf{z}_i + \boldsymbol{\varepsilon}_i, \quad \forall i \in [N], \quad (2.1.12)$$

where $\mathbf{U} \in \mathcal{O}(D, d)$ captures the low-rank structure, $\{\mathbf{z}_i\}_{i=1}^N \subseteq \mathbb{R}^d$ are the compact codes of the data, and $\{\boldsymbol{\varepsilon}_i\}_{i=1}^N \subseteq \mathbb{R}^D$ are small vectors indicating the deviation of our data from the low-rank model. Then the principal components $\mathbf{U}^* \in \mathcal{O}(D, d)$ of our dataset are given by the top d eigenvectors of $\mathbf{X}\mathbf{X}^\top/N$, where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$, and approximately correspond to the optimal linear denoiser: $\mathbf{U}^*(\mathbf{U}^*)^\top \approx \mathbf{U}\mathbf{U}^\top$.

We do not give explicit rates of approximation here as they can become rather technical. In the special case that $\boldsymbol{\varepsilon}_i = \mathbf{0}$ for all i , the learned \mathbf{U}^* spans the support of the samples $\{\mathbf{x}_i\}_{i=1}^N$. If in addition the \mathbf{z}_i are sufficiently diverse (say, spanning all of \mathbb{R}^d) then we would have perfect recovery: $\mathbf{U}^*(\mathbf{U}^*)^\top = \mathbf{U}\mathbf{U}^\top$.

Remark 2.1. In some data analysis tasks, the data matrix \mathbf{X} is formatted such that each data point is a *row* rather than a *column* as is presented here. In this case the principal components are the top d eigenvectors of $\mathbf{X}^\top\mathbf{X}/N$.

Remark 2.2 (Basis Selection via Denoising Eigenvalues). In many cases, either our data will not truly be distributed according to a subspace-plus-noise model, or we will not know the true underlying dimension d of the subspace. In this case, we have to choose d ; this problem is called *model selection*. In the restricted case of PCA, one way to perform model selection is to compute $\mathbf{X}\mathbf{X}^\top/N$ and look for instances where adjacent eigenvalues sharply decrease; this is one indicator that the index of the larger eigenvalue is the “true dimension d ”, and the rest of the eigenvalues of $\mathbf{X}\mathbf{X}^\top/N$ are contributed by the noise or disturbances ε_i . Model selection is a difficult problem and, nowadays in the era of deep learning where it is called “hyperparameter optimization”, is usually done via brute force or Bayesian optimization, often using the loss on a holdout (“validation”) dataset as an objective.

Remark 2.3 (Denoising Samples). The expression on the right-hand side of (2.1.5), that is,

$$\min_{\tilde{\mathbf{U}}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{U}}\tilde{\mathbf{U}}^\top \mathbf{x}_i\|_2^2, \quad (2.1.13)$$

is what’s known as a *denoising problem*, thus named because it is an optimization problem whose solution *removes the noise from the samples so it fits on the subspace*. Denoising—learning a map which removes noise from noisy samples so that it fits on the data structure (such as in (2.1.1), but maybe more complicated)—is a common method for learning distributions that will be discussed in the sequel and throughout the manuscript. Note that we have already discussed this notion, but it bears repeating due to its central importance in later chapters.

Remark 2.4 (Neural Network Interpretation). If we do a PCA, we approximately recover the support of the distribution encoded by the parameter \mathbf{U}^* . The learned denoising map then takes the form $\mathbf{U}^*(\mathbf{U}^*)^\top$. On top of being a denoiser, this can be viewed as a *simple two-layer weight-tied linear neural network*: the first layer multiplies by $(\mathbf{U}^*)^\top$, and the second layer multiplies by \mathbf{U}^* , namely

$$\text{denoise}(\mathbf{x}) = \mathbf{U}^* \circ \underbrace{\text{id} \circ \underbrace{(\mathbf{U}^*)^\top \mathbf{x}}_{\text{first “layer”}}}_{\text{post-activation of first “layer”}}_{\text{output of “NN”}} \quad (2.1.14)$$

Contrasting this to a standard two-layer neural network, we see a structural similarity:

$$\text{NN}(\mathbf{x}) = \mathbf{W}^* \circ \underbrace{\text{ReLU} \circ \underbrace{(\mathbf{U}^*)^\top \mathbf{x}}_{\text{first layer}}}_{\text{post-activation of first layer}}_{\text{output of NN}} \quad (2.1.15)$$

In particular, PCA can be interpreted as *learning a simple two-layer denoising*

autoencoder,² one of the simplest examples of a non-trivial neural network. In this framework, the *learned representations* would just be $(\mathbf{U}^*)^\top \mathbf{x} \approx \mathbf{z}$. In this way, PCA serves as a model problem for (deep) representation learning, which we will draw upon further in the monograph. Notice that in this analogy, the representations reflect, or are projections of, the input data towards a learned low-dimensional structure. This property will be particularly relevant in the future.

2.1.2 Pursuing Low-Rank Structure via Power Iteration

There is a computationally efficient way to estimate the top eigenvectors of $\mathbf{X}\mathbf{X}^\top/N$ or any symmetric positive semidefinite matrix \mathbf{M} , called *power iteration*. This method is the building block of several algorithmic approaches to high-dimensional data analysis that we discuss later in the chapter, so we discuss it here.

Let \mathbf{M} be a symmetric positive semidefinite matrix. There exists an orthonormal basis for \mathbb{R}^D consisting of eigenvectors $(\mathbf{w}_i)_{i=1}^D$ of \mathbf{M} , with corresponding eigenvalues $\lambda_1 \geq \dots \geq \lambda_D \geq 0$. By definition, any eigenvector \mathbf{w}_i satisfies $\lambda_i \mathbf{w}_i = \mathbf{M}\mathbf{w}_i$. Therefore, for any $\lambda_i > 0$, \mathbf{w}_i is a “fixed point” to the following equation:

$$\mathbf{w} = \frac{\mathbf{M}\mathbf{w}}{\|\mathbf{M}\mathbf{w}\|_2}. \quad (2.1.16)$$

Theorem 2.2 (Power Iteration). *Assume that $\lambda_1 > \lambda_i$ for all $i > 1$. If we compute the fixed point of (2.1.16) using the following iteration:*

$$\mathbf{v}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{1}), \quad \mathbf{v}_{t+1} \leftarrow \frac{\mathbf{M}\mathbf{v}_t}{\|\mathbf{M}\mathbf{v}_t\|_2}, \quad (2.1.17)$$

then, in the limit, \mathbf{v}_t will converge to a top unit-norm eigenvector of \mathbf{M} .

Proof. First, note that for all t , we have

$$\mathbf{v}_t = \frac{\mathbf{M}\mathbf{v}_{t-1}}{\|\mathbf{M}\mathbf{v}_{t-1}\|_2} = \frac{\mathbf{M}^2\mathbf{v}_{t-2}}{\|\mathbf{M}\mathbf{v}_{t-1}\|_2\|\mathbf{M}\mathbf{v}_{t-2}\|_2} = \dots = \frac{\mathbf{M}^t\mathbf{v}_0}{\prod_{s=1}^t \|\mathbf{M}\mathbf{v}_s\|_2}. \quad (2.1.18)$$

Thus, \mathbf{v}_t has the same direction as $\mathbf{M}^t\mathbf{v}_0$ and is unit norm, so we can write

$$\mathbf{v}_t = \frac{\mathbf{M}^t\mathbf{v}_0}{\|\mathbf{M}^t\mathbf{v}_0\|_2}. \quad (2.1.19)$$

Because all the eigenvectors \mathbf{w}_i of \mathbf{M} form an orthonormal basis for \mathbb{R}^D , we can write

$$\mathbf{v}_0 = \sum_{i=1}^D \alpha_i \mathbf{w}_i, \quad (2.1.20)$$

²In fact, as we have mentioned in the previous chapter, PCA was one of the first problems that neural networks were used to solve [BH89; Oja82].

where because \mathbf{v}_0 is Gaussian the α_i are all nonzero with probability 1. Thus, we can use our earlier expression for \mathbf{v}_t to write

$$\mathbf{v}_t = \frac{\mathbf{M}^t \mathbf{v}_0}{\|\mathbf{M}^t \mathbf{v}_0\|_2} = \frac{\sum_{i=1}^D \lambda_i^t \alpha_i \mathbf{w}_i}{\left\| \sum_{i=1}^D \lambda_i^t \alpha_i \mathbf{w}_i \right\|_2} = \frac{\sum_{i=1}^D \lambda_i^t \alpha_i \mathbf{w}_i}{\sqrt{\sum_{i=1}^D \lambda_i^{2t} \alpha_i^2}}. \quad (2.1.21)$$

Now, let us consider the case where $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_D \geq 0$. (The case with repeated top eigenvalues goes similarly.) Then we can write

$$\mathbf{v}_t = \frac{\alpha_1 \mathbf{w}_1 + \sum_{i=2}^D (\lambda_i / \lambda_1)^t \alpha_i \mathbf{w}_i}{\sqrt{\alpha_1^2 + \sum_{i=2}^D (\lambda_i / \lambda_1)^{2t} \alpha_i^2}}. \quad (2.1.22)$$

Because $\lambda_1 > \lambda_i$ for all $i > 1$, the terms inside the summation go to 0 exponentially fast, and the remainder is the limit

$$\lim_{t \rightarrow \infty} \mathbf{v}_t = \frac{\alpha_1}{|\alpha_1|} \mathbf{w}_1 = \text{sign}(\alpha_1) \mathbf{w}_1, \quad (2.1.23)$$

which is a top unit eigenvector of \mathbf{M} . The top eigenvalue λ_1 of \mathbf{M} can be estimated by $\mathbf{v}_t^\top \mathbf{M} \mathbf{v}_t$, which converges similarly rapidly to λ_1 . \square

Remark 2.5 (Power Iteration as a Denoising Process). Notice that computing the principal component via power iteration can be interpreted as a “denoising” process: starting from random noise, it converges iteratively onto a one-dimensional subspace. Indeed, as we will show later in Chapter 3, if we assume *the support of the data distribution is a single subspace*, the associated denoising operator (or the score function) is precisely the same as the above power iteration. In Chapter 3, we will show how to generalize such a denoising process to distributions of more general low-dimensional structures. This is essentially a universal and most efficient approach to learn any low-dimensional distributions.

To find the second top eigenvector, we apply the power iteration algorithm to $\mathbf{M} - \lambda_1 \mathbf{v}_1 \mathbf{v}_1^\top$, which has eigenvectors $(\mathbf{w}_i)_{i=2}^D$ and corresponding eigenvalues $(\lambda_i)_{i=2}^D$. By repeating this procedure d times in sequence, we can very efficiently estimate the top d eigenvectors of \mathbf{M} for any symmetric positive semidefinite matrix \mathbf{M} . Thus we can also apply it to $\mathbf{X} \mathbf{X}^\top / N$ to recover the top d principal components, which is what we were after in the first place. Notice that this approach recovers one principal component at a time; we will contrast this to other algorithmic approaches, such as gradient descent on all eigenvectors, in future sections.

2.1.3 Local and Incremental Learning via Online PCA

Notice that in the above construction, the linear transform \mathbf{U} (i.e., the matrix of principal components) used for encoding and decoding is computed “offline” from all the input data beforehand. One question is whether this transform can be learned “online” as the input data comes in order. This question was answered by the work of Oja in 1982 [Oja82].

Example 2.1 (Normalized Hebbian learning scheme for PCA). Consider a sequence of i.i.d. random vectors $\mathbf{x}_1, \dots, \mathbf{x}_i, \dots \in \mathbb{R}^n$ with covariance $\mathbf{\Sigma} \in \mathbb{R}^{n \times n}$. Let $\mathbf{u}_0 \in \mathbb{R}^n$ and define the response of an input vector \mathbf{x}_i against a weight vector \mathbf{u}_i to be their inner product:

$$\eta_i = \mathbf{u}_i^T \mathbf{x}_i \quad (2.1.24)$$

and we update the weight vector according to the following scheme:

$$\mathbf{u}_{i+1} = \frac{\mathbf{u}_i + \gamma \eta_i \mathbf{x}_i}{\|\mathbf{u}_i + \gamma \eta_i \mathbf{x}_i\|_2} \quad (2.1.25)$$

for some small gain $\gamma > 0$. This update scheme can be viewed as a normalized Hebbian scheme, in which the weights of connections between neurons become stronger if (products of) both the input \mathbf{x} and output η are strong. One may view the vector of weights \mathbf{u} as “learned” based on a form of feedback from the output η . Then, under reasonable assumptions, Oja [Oja82] has shown that \mathbf{u}_i converges to the eigenvector associated with the largest eigenvalue of $\mathbf{\Sigma}$. ■

To interpret the normalized Hebbian scheme (2.1.25), let us think of the data sample PCA objective (2.1.13) as an approximation to the *population* version of the objective, defined via an expectation over the distribution of samples \mathbf{x} :

$$\min_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \mathbb{E}_{\mathbf{x}} \left[\|\mathbf{x} - \mathbf{U} \mathbf{U}^T \mathbf{x}\|_2^2 \right] \quad (2.1.26)$$

We will dig into this objective further in the following section. Then the normalized Hebbian scheme can be interpreted as a first-order approximation to a *stochastic* projected gradient descent scheme on Equation (2.1.26) (with batch size 1, and with the number of columns of \mathbf{U} equal to 1) as long as $\|\mathbf{u}\|_2 = 1$, which is maintained by the projection operation in (2.1.25). Such an *online algorithm* is a useful primitive for more complex algorithms that process sequentially-structured data with causal structure, as we will study more in Section 5.3.3 and Section 8.11. Interestingly, and more generally, the simple online PCA algorithm is suggestive of the fact that *simpler algorithms than (end-to-end) back propagation can succeed in learning encoder-decoder pairs*, which we will briefly revisit in Chapter 9 as an important goal for the future study of intelligence.

2.1.4 The Statistical View: Probabilistic PCA

Notice that the above formulation makes no statistical assumptions on the data-generating process. However, it is common to include statistical elements within a given data model, as it may add further enlightening interpretations about the result of the analysis. As such, we ask the natural question: *what is the statistical analogue to low-dimensional structure?* Our answer is that a low-dimensional *distribution* is one whose support is concentrated around a low-dimensional geometric structure.

To illustrate this point, we discuss *probabilistic principal component analysis* (PPCA) [Row97; TB99]. This formulation can be viewed as a statistical variant of regular PCA. Mathematically, we now consider our data as samples from a random variable \mathbf{x} taking values in \mathbb{R}^D (also sometimes called a *random vector*). We say that \mathbf{x} has (approximate) low-rank statistical structure if and only if there exists an orthonormal matrix $\mathbf{U} \in \mathcal{O}(D, d)$, a random variable \mathbf{z} taking values in \mathbb{R}^d , and a *small* random variable $\boldsymbol{\varepsilon}$ taking values in \mathbb{R}^D such that \mathbf{z} and $\boldsymbol{\varepsilon}$ are independent, and

$$\mathbf{x} = \mathbf{U}\mathbf{z} + \boldsymbol{\varepsilon}. \quad (2.1.27)$$

Our goal is again to recover \mathbf{U} . Towards this end, we set up the analogous problem as in Section 2.1.1, i.e., optimizing over subspace supports $\tilde{\mathbf{U}}$ and random variables \mathbf{z} to solve the following problem:

$$\min_{\tilde{\mathbf{U}}, \tilde{\mathbf{z}}} \mathbb{E} \|\mathbf{x} - \tilde{\mathbf{U}}\tilde{\mathbf{z}}\|_2^2. \quad (2.1.28)$$

Since we are finding the best such random variable \mathbf{z} we can find its realization $\mathbf{z}(\mathbf{x})$ separately for each value of \mathbf{x} . Performing the same calculations as in Section 2.1.1, we obtain

$$\min_{\tilde{\mathbf{U}}, \tilde{\mathbf{z}}} \mathbb{E} \|\mathbf{x} - \tilde{\mathbf{U}}\tilde{\mathbf{z}}\|_2^2 = \min_{\tilde{\mathbf{U}}} \mathbb{E} \min_{\tilde{\mathbf{z}}(\mathbf{x})} \|\mathbf{x} - \tilde{\mathbf{U}}\tilde{\mathbf{z}}(\mathbf{x})\|_2^2 = \min_{\tilde{\mathbf{U}}} \mathbb{E} \|\mathbf{x} - \tilde{\mathbf{U}}\tilde{\mathbf{U}}^\top \mathbf{x}\|_2^2, \quad (2.1.29)$$

again re-emphasizing the fact that the estimated subspace with principal components \mathbf{U}^* corresponds to a denoiser $\mathbf{U}^*(\mathbf{U}^*)^\top$ which projects onto that subspace. As before, we obtain

$$\arg \min_{\tilde{\mathbf{U}}} \mathbb{E} \|\mathbf{x} - \tilde{\mathbf{U}}\tilde{\mathbf{U}}^\top \mathbf{x}\|_2^2 = \arg \max_{\tilde{\mathbf{U}}} \mathbb{E} \|\tilde{\mathbf{U}}^\top \mathbf{x}\|_2^2 \quad (2.1.30)$$

$$= \arg \max_{\tilde{\mathbf{U}}} \text{tr}(\tilde{\mathbf{U}}^\top \mathbb{E}[\mathbf{x}\mathbf{x}^\top] \tilde{\mathbf{U}}), \quad (2.1.31)$$

and the solution to the latter problem is just the top d principal components of the second moment matrix $\mathbb{E}[\mathbf{x}\mathbf{x}^\top]$. Actually, the above problems are visually very similar to the equations for computing the principal components in the previous subsection, except with $\mathbb{E}[\mathbf{x}\mathbf{x}^\top]$ replacing $\mathbf{X}\mathbf{X}^\top/N$. In fact, the latter quantity is an estimate for the former. Both formulations effectively do the same thing, and have the same practical solution—compute the left singular vectors of the data matrix \mathbf{X} , or equivalently the top eigenvectors of the estimated covariance matrix $\mathbf{X}\mathbf{X}^\top/N$. The statistical formulation, however, has an additional interpretation. Suppose that $\mathbb{E}[\mathbf{z}] = \mathbf{0}$ and $\mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}$. We have

$$\mathbb{E}[\mathbf{x}] = \mathbf{U} \mathbb{E}[\mathbf{z}] + \mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}, \quad (2.1.32)$$

so that $\text{Cov}[\mathbf{x}] = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$. Now working out $\text{Cov}[\mathbf{x}]$ we have

$$\text{Cov}[\mathbf{x}] = \mathbf{U} \text{Cov}[\mathbf{z}]\mathbf{U}^\top + \text{Cov}[\boldsymbol{\varepsilon}] = \mathbf{U} \mathbb{E}[\mathbf{z}\mathbf{z}^\top]\mathbf{U}^\top + \text{Cov}[\boldsymbol{\varepsilon}]. \quad (2.1.33)$$

In particular, if $\text{Cov}[\boldsymbol{\varepsilon}]$ is small, it holds that $\text{Cov}[\mathbf{x}] = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$ is approximately a low-rank matrix, in particular rank d . Thus the top d eigenvectors of $\mathbb{E}[\mathbf{x}\mathbf{x}^\top]$

essentially summarize the whole covariance matrix. But they are also the principal components, so we can interpret principal component analysis as performing a low-rank decomposition of $\text{Cov}[\mathbf{x}]$.

Remark 2.6. By using the probabilistic viewpoint of PCA, we achieve a clearer and more quantitative understanding of how it relates to denoising. First, consider the denoising problem in (2.1.29), namely

$$\min_{\tilde{\mathbf{U}}} \mathbb{E} \|\mathbf{x} - \tilde{\mathbf{U}}\tilde{\mathbf{U}}^\top \mathbf{x}\|_2^2. \quad (2.1.34)$$

It is not too hard to prove that if $\tilde{\mathbf{U}}$ has d columns and if $\boldsymbol{\varepsilon}$ is an isotropic Gaussian random variable, i.e., with distribution $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$,³ then for *any* optimal solution \mathbf{U}^* to this problem, we have

$$\mathbf{U}^*(\mathbf{U}^*)^\top = \mathbf{U}\mathbf{U}^\top \quad (2.1.35)$$

and so the true supporting subspace, say $\mathcal{S} \doteq \text{col}(\mathbf{U})$, is recovered as the span of columns of \mathbf{U}^* , since

$$\mathcal{S} = \text{col}(\mathbf{U}) = \text{col}(\mathbf{U}\mathbf{U}^\top) = \text{col}(\mathbf{U}^*(\mathbf{U}^*)^\top) = \text{col}(\mathbf{U}^*). \quad (2.1.36)$$

In particular, the learned *denoising map* $\mathbf{U}^*(\mathbf{U}^*)^\top$ is an orthogonal projection onto \mathcal{S} , pushing noisy points onto the ground truth supporting subspace. We can establish a similar technical result in the case where we only have finite samples, as in Theorem 2.1, but this takes more effort and technicality. Summarizing this discussion, we have the following informal Theorem.

Theorem 2.3. *Suppose that the random variable \mathbf{x} can be written as*

$$\mathbf{x} = \mathbf{U}\mathbf{z} + \boldsymbol{\varepsilon} \quad (2.1.37)$$

where $\mathbf{U} \in \mathbf{O}(D, d)$ captures the low-rank structure, \mathbf{z} is a random vector taking values in \mathbb{R}^d , and $\boldsymbol{\varepsilon}$ is a random vector taking values in \mathbb{R}^D such that \mathbf{z} and $\boldsymbol{\varepsilon}$ are independent, and $\boldsymbol{\varepsilon}$ is small. Then the principal components $\mathbf{U}^* \in \mathbf{O}(D, d)$ of our dataset are given by the top d eigenvectors of $\mathbb{E}[\mathbf{x}\mathbf{x}^\top]$, and approximately correspond to the optimal linear denoiser: $\mathbf{U}^*(\mathbf{U}^*)^\top \approx \mathbf{U}\mathbf{U}^\top$.

2.1.5 Masked Low-Rank Matrix Completion

In the previous subsections, we discussed the problem of *learning a low-rank geometric or statistical distribution*, where the data were sampled from a subspace with additive noise. But this is not the only type of disturbance from a low-dimensional distribution that is worthwhile to study. In this subsection, we introduce one more class of non-additive errors which become increasingly important in deep learning. Let us consider the case where we have some data $\{\mathbf{x}_i\}_{i=1}^N$ generated according to (2.1.1). Now we arrange them into a matrix

³Other distributions work so long as they support all of \mathbb{R}^D , but the Gaussian is the easiest to work with here.

$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$. Unlike before, we do not observe \mathbf{X} directly; we instead imagine that our observation was corrupted en route and we obtained

$$\mathbf{Y} = \mathbf{M} \odot \mathbf{X}, \quad (2.1.38)$$

where $\mathbf{M} \in \{0, 1\}^{D \times N}$ is a *mask* that is known by us, and \odot is element-wise multiplication. In this case, our goal is to recover \mathbf{X} (from which point we can use PCA to recover \mathbf{U} , etc), given only the corrupted observation \mathbf{Y} , the mask \mathbf{M} , and the knowledge that \mathbf{X} is (approximately) low-rank. This is called *low-rank matrix completion*.

This and similar generalizations of this low-rank structure recovery problem are referred to as “Robust PCA” problems. There have been many excellent resources discussing algorithms and approaches to solve these problems [WM22]. Hence we will not go into the solution method here. Instead, we will discuss under what conditions this problem is *plausible* to solve. On one hand, in the most absurd case, suppose that each entry of the matrix \mathbf{X} were chosen independently from all the others. Then there would be no hope of recovering \mathbf{X} exactly even if only one entry were missing and we had $DN - 1$ entries. On the other hand, suppose that we knew that \mathbf{X} has rank 1 exactly, which is an extremely strong condition on the low-dimensional structure of the data, and we were dealing with the mask

$$\mathbf{M} = \begin{bmatrix} \mathbf{1}_{(D-1) \times 1} & \mathbf{0}_{(D-1) \times (N-1)} \\ 1 & \mathbf{1}_{1 \times (N-1)} \end{bmatrix}. \quad (2.1.39)$$

Then we know that the data are distributed on a line, and we know a vector on that line—it is just the first column of the matrix $\mathbf{Y} = \mathbf{M} \odot \mathbf{X}$. From the last coordinate of each column, also revealed to us by the mask, we can solve for each column since for each final coordinate there is only one vector on the line with this coordinate. Thus we can reconstruct \mathbf{X} with perfect accuracy, and we only need a linear number of observations $D + N - 1$.

In the real world, actual problems are somewhere in between the two limiting cases discussed above. Yet the differences between these two extremes, as well as the earlier discussion of PCA, reveal a general kernel of truth:

The lower-dimensional and more structured the data distribution is, the easier it is to process, and the fewer observations are needed—provided that the algorithm effectively utilizes this low-dimensional structure.

As is perhaps predictable, we will encounter this motif repeatedly in the remainder of the manuscript, starting in the very next section.⁴

⁴In later chapters, we will see that one can complete real-world data such as images with only a fraction of unmasked observations, e.g., the masked autoencoding of images in Section 8.5. This is also because the distribution of natural images has rather low-dimensional (albeit nonlinear) structures in the space of all images.

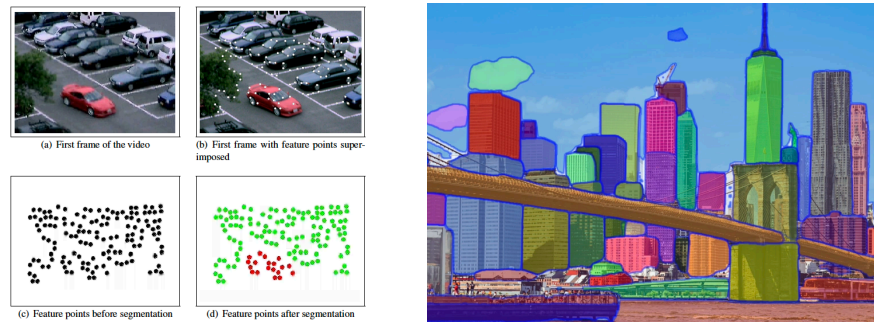


Figure 2.4: **Left:** features tracked on independently moving objects in a scene. **Right:** image patches associated with different regions of an image.

2.2 A Mixture of Complete Low-Dimensional Subspaces

As we have seen, low-rank signal models are rich enough to provide a full picture of the interplay between low-dimensionality in data and efficient and scalable computational algorithms for representation and recovery under errors. These models imply a *linear* and symmetric representation learning pipeline (2.1.6):

$$\mathbf{z} = \mathcal{E}(\mathbf{x}) = \mathbf{U}^\top \mathbf{x}, \quad \hat{\mathbf{x}} = \mathcal{D}(\mathbf{z}) = \mathbf{U}\mathbf{z},$$

which can be provably learned from finite samples of \mathbf{x} with principal component analysis (solved efficiently, say, with the power method) whenever the distribution of \mathbf{x} truly is linear. This is a restrictive assumption—for as Harold Hotelling, the distinguished 20th century statistician,⁵ objected following George Dantzig’s presentation of his theory of linear programming for the first time [Dan02],

“...we all know the world is nonlinear.”

Even accounting for its elegance and simplicity, the low-rank assumption is too restrictive to be broadly applicable to modeling real-world data. A key limitation is the assumption of a *single* linear subspace that is responsible for generating the structured observations. In many practical applications, structure generated by a *mixture* of distinct low-dimensional subspaces provides a more realistic model. For example, consider a video sequence that captures the motion of several distinct objects, each subject to its own independent displacement (Figure 2.4 left). Under suitable assumptions on the individual motions, each object becomes responsible for an independent low-dimensional subspace in the concatenated sequence of video frames [VM04]. As another example, consider modeling natural images via learning a model for the distribution of

⁵By coincidence, also famous for his contributions to the development and naming of Principal Component Analysis [Hot33].

patches, spatially-contiguous collections of pixels, within an image (Figure 2.4 right). Unlike in the simplistic setting of Figure 2.2, where images of faces with matched poses can be well-approximated by a single low-dimensional subspace, the patch at a specific location in a natural image can correspond to objects with very different properties—for example, distinct color or shape due to occlusion boundaries. Therefore, modeling the distribution of patches with a single subspace is futile, but a *mixture* of subspaces, one for each region, performs surprisingly well in practice, say for segmenting or compressing purposes [MRY+11].⁶ We will see a concrete example in a later chapter (Example 4.10).

In this section, we will begin by discussing the conceptual and algorithmic foundations for structured representation learning when the data distribution can be modeled by a *mixture of low-dimensional subspaces*, as illustrated in Figure 2.1. In this setting, the decoder mapping will be almost as simple as the case of a single subspace: we simply represent via

$$\hat{\mathbf{x}} = \mathcal{D}(\mathbf{z}) = \left(\sum_{k=1}^K \pi_k(\mathbf{z}) \mathbf{U}_k \right) \mathbf{z}, \quad (2.2.1)$$

where $\pi_k : \mathbb{R}^d \rightarrow \{0, 1\}$ are a set of *sparse* weighting random variables, such that only a single subspace $\mathcal{S}_k = \text{col}(\mathbf{U}_k)$ is selected in the sum. However, the task of encoding such data \mathbf{x} to suitable representations \mathbf{z} , and learning such an encoder-decoder pair from data, will prove to be more involved.

We will see how ideas from the rich literature on *sparse representation* and *independent component analysis* lead to a natural reformulation of the above decoder architecture through the lens of sparsity, the corresponding encoder architecture (obtained through a power-method-like algorithm analogous to that of principal component analysis), and strong guarantees of correctness and efficiency for learning such encoder-decoder pairs from data. In this sense, the case of mixed linear low-dimensional structure already leads to many of the key conceptual components of structured representation learning that we will develop in far greater generality in this book.

2.2.1 Mixtures of Subspaces and Sparsely-Used Dictionaries

Let $\mathbf{U}_1, \dots, \mathbf{U}_K$, each of size $D \times d$, denote a collection of orthonormal bases for K subspaces of dimension d in \mathbb{R}^D . To say that \mathbf{x} follows a mixture-of-subspaces distribution parameterized by $\mathbf{U}_1, \dots, \mathbf{U}_K$ means, geometrically speaking, that

$$\mathbf{x} = \mathbf{U}_k \mathbf{z} \quad \text{for some } k \in [K], \quad \mathbf{z} \in \mathbb{R}^d. \quad (2.2.2)$$

The statistical analogue of this geometric model, as we saw for the case of PCA and linear structure, is that \mathbf{x} follows a *mixture of Gaussians* distribution: that

⁶We will return to this observation in Chapter 5, where we will show it can be significantly generalized to learn representations for large-scale modern datasets.

is,

$$\mathbf{x} \sim \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{0}, \mathbf{U}_k \mathbf{U}_k^\top), \quad \text{for some } \pi_k \geq 0, \quad \sum_{k=1}^K \pi_k = 1. \quad (2.2.3)$$

In other words, for each $k \in [K]$, \mathbf{x} is Gaussian on the low-dimensional subspace $\text{col}(\mathbf{U}_k)$ with probability π_k .

Remark 2.7 (A Mixture of Gaussians v.s. A Superposition of Gaussians). One should be aware that the above model (2.2.3) is a *mixture* of Gaussian distributions, not to be confused with a mixture of Gaussian variables by superposition, say

$$\mathbf{x} = \sum_{i=1}^n w_i \mathbf{x}_i, \quad \mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{U}_i \mathbf{U}_i^\top), \quad (2.2.4)$$

where \mathbf{x}_i are independent random Gaussian vectors and w_i are a set of fixed weights. As we know from the properties of Gaussian vectors, such a superposition \mathbf{x} will remain a Gaussian distribution.

For now, we focus on the geometric perspective offered by (2.2.2). There is an algebraically convenient alternative to this conditional representation. Consider a *lifted* representation vector $\mathbf{z} = [\mathbf{z}_1^\top, \dots, \mathbf{z}_K^\top]^\top \in \mathbb{R}^{dK}$, such that \mathbf{z} is *d-sparse* with support on one of the K consecutive non-overlapping blocks of d coordinates out of dK . Then (2.2.2) can be written equivalently as

$$\mathbf{x} = \underbrace{\begin{bmatrix} | & \dots & | \\ \mathbf{U}_1 & \dots & \mathbf{U}_K \\ | & \dots & | \end{bmatrix}}_{\mathbf{U}} \underbrace{\begin{bmatrix} \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_K \end{bmatrix}}_{\mathbf{z}}, \quad \left\| \begin{bmatrix} \|\mathbf{z}_1\|_2 \\ \vdots \\ \|\mathbf{z}_K\|_2 \end{bmatrix} \right\|_0 = 1. \quad (2.2.5)$$

Here, the ℓ^0 “norm” $\|\cdot\|_0$ measures sparsity by counting the number of nonzero entries:

$$\|\mathbf{z}\|_0 = |\{i \mid z_i \neq 0\}|, \quad (2.2.6)$$

and the matrix $\mathbf{U} \in \mathbb{R}^{D \times Kd}$ is called a *dictionary* with all the $\{\mathbf{U}_i\}_{i=1}^K$ as code words. In general, if the number of subspaces in the mixture K is large enough, there is no bound on the number of columns contained in the dictionary \mathbf{U} . In the case where $Kd < D$, \mathbf{U} is called *undercomplete*; when $Kd = D$, it is called *complete*; and when $Kd > D$, it is called *overcomplete*.

Now, (2.2.5) suggests a convenient relaxation for tractability of analysis: rather than modeling \mathbf{x} as coming from a mixture of K *specific* subspaces $\mathbf{U}_1, \dots, \mathbf{U}_K$, we may instead start with a dictionary $\mathbf{U} \in \mathbb{R}^{D \times m}$, where m may be smaller or larger than D , and simply seek to represent $\mathbf{x} = \mathbf{U}\mathbf{z}$ with the sparsity $\|\mathbf{z}\|_0$ sufficiently small. This leads to the *sparse dictionary model* for \mathbf{x} :

$$\mathbf{x} = \mathbf{U}\mathbf{z} + \boldsymbol{\varepsilon}, \quad \|\mathbf{z}\|_0 \ll d, \quad (2.2.7)$$

where $\boldsymbol{\varepsilon}$ represents an unknown noise vector. Geometrically, this implies that \mathbf{x} lies close to the span of a subset of $\|\mathbf{z}\|_0$ columns of \mathbf{U} , making this an instantiation of the mixture-of-subspaces model (2.2.2) with a very large value of K , and specific correlations between the subspaces \mathbf{U}_k .

Orthogonal dictionary for sparse coding. Now we can formulate the structured representation learning problem for mixtures of low-dimensional subspaces that we will study in this section. We assume that we have samples $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ from an unknown sparse dictionary model (2.2.7), possibly with added noises ε_i . Let us begin from the assumption that the dictionary \mathbf{U} in the sparse dictionary model (2.2.7) is complete and orthogonal,⁷ and that each coefficient vector \mathbf{z} is d -sparse, with $d \ll D$. In this setting, representation learning amounts to correctly learning the orthogonal dictionary \mathbf{U} via optimization: we can then take $\mathcal{E}(\mathbf{x}) = \mathbf{U}^\top \mathbf{x}$ as the encoder and $\mathcal{D}(\mathbf{z}) = \mathbf{U}\mathbf{z}$ as the decoder, and $\mathcal{D} = \mathcal{E}^{-1}$. In diagram form:

$$\mathbf{x} \xrightarrow{\mathcal{E}=\mathbf{U}^\top} \mathbf{z} \xrightarrow{\mathcal{D}=\mathbf{U}} \hat{\mathbf{x}}. \quad (2.2.8)$$

We see that the autoencoding pair $(\mathcal{E}, \mathcal{D})$ for complete dictionary learning is symmetric, as in the case of a single linear subspace, making the computational task of encoding and decoding no more difficult than in the linear case. On the other hand, the task of learning the dictionary \mathbf{U} is strictly more difficult than learning a single linear subspace by PCA. To see why we cannot simply use PCA to learn the orthogonal dictionary \mathbf{U} correctly, note that the loss function that gave rise to PCA, namely (2.1.5), is completely invariant to rotations of the rows of the matrix \mathbf{U} : that is, if \mathbf{Q} is any $d \times d$ orthogonal matrix, then \mathbf{U} and $\mathbf{U}\mathbf{Q}$ are both feasible and have an identical loss for (2.1.5). The sparse dictionary model is decidedly not invariant to such transformations: if we replaced \mathbf{U} by $\mathbf{U}\mathbf{Q}$ and made a corresponding rotation $\mathbf{Q}^\top \mathbf{z}$ of the representation coefficients \mathbf{z} , we would in general destroy the sparsity structure of \mathbf{z} , violating the modeling assumption. Thus, we need new algorithms for learning orthogonal dictionaries.

2.2.2 Complete Dictionary Learning

In this section, we will derive algorithms for solving the orthogonal dictionary learning problem. To be more precise, we assume that the observed vector $\mathbf{x} \in \mathbb{R}^D$ follows a statistical model

$$\mathbf{x} = \mathbf{U}\mathbf{z} + \varepsilon, \quad (2.2.9)$$

where $\mathbf{U} \in \mathbb{R}^{D \times D}$ is an unknown orthogonal dictionary, \mathbf{z} is a random vector with statistically independent components z_i , each with zero mean, and $\varepsilon \in \mathbb{R}^D$ is an independent random vector of small (Gaussian) noises. The goal is to recover \mathbf{U} (and hence \mathbf{z}) from samples of \mathbf{x} .

Here we assume that each independent component z_i is distributed as

$$z_i \sim \text{Ber}(\theta) \cdot \mathcal{N}(0, 1/\theta).$$

That is, it is the product of a Bernoulli random variable with probability θ of being 1 and $1 - \theta$ of being 0, and an independent Gaussian random variable with variance $1/\theta$. This distribution is formally known as the *Bernoulli-Gaussian* distribution. The normalization is chosen so that $\text{Var}(z_i) = 1$ and

⁷It can be shown that for the complete case, we do not lose any generality by making the orthogonal assumption (Exercise 2.4).

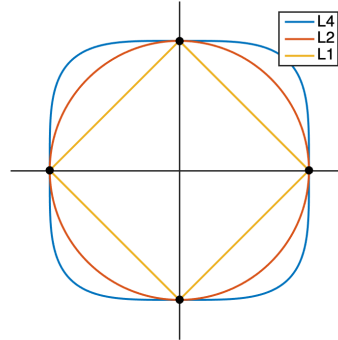


Figure 2.5: Maximizing ℓ^4 norm or minimizing ℓ^1 norm promotes sparsity (for vectors on the sphere).

hence $\mathbb{E}[\|\mathbf{z}\|_2^2] = d$. This modeling assumption implies that the vector of independent components \mathbf{z} is typically very sparse: we calculate $\mathbb{E}[\|\mathbf{z}\|_0] = d\theta$, which is small when θ is inversely proportional to d .

Remark 2.8 (The Orthogonal Assumption). At first sight, the assumption that the dictionary \mathbf{U} is orthogonal might seem to be somewhat restrictive. But there is actually no loss of generality. One may consider a complete dictionary to be any square invertible matrix \mathbf{U} . With samples generated from this dictionary: $\mathbf{X} = \mathbf{U}\mathbf{Z} \in \mathbb{R}^{D \times N}$, it is easy to show that with some preconditioning of the data matrix \mathbf{X} :

$$\bar{\mathbf{X}} = \left(\frac{1}{N\theta} \mathbf{X}\mathbf{X}^\top \right)^{-\frac{1}{2}} \mathbf{X}, \quad (2.2.10)$$

then there exists an orthogonal matrix $\mathbf{U}_o \in \mathcal{O}(D)$ such that

$$\bar{\mathbf{X}} = \mathbf{U}_o \mathbf{Z}. \quad (2.2.11)$$

See Exercise 2.4 or [SQW17b] for more details.

Dictionary learning via the MSP algorithm. Now suppose that we are given a set of observations:

$$\mathbf{x}_i = \mathbf{U}\mathbf{z}_i + \boldsymbol{\varepsilon}_i, \quad \forall i \in [N]. \quad (2.2.12)$$

Let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ and $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N]$. The goal is to recover \mathbf{U} from the data \mathbf{X} . Therefore, given any orthogonal matrix $\mathbf{A} \in \mathcal{O}(D)$,

$$\mathbf{A}\mathbf{x}_i = \mathbf{A}\mathbf{U}\mathbf{z}_i + \mathbf{A}\boldsymbol{\varepsilon}_i \quad (2.2.13)$$

would be nearly sparse if $\mathbf{A} = \mathbf{U}^\top$ (as by assumption the noise $\boldsymbol{\varepsilon}_i$ is of small magnitude).

Also, given \mathbf{U} is orthogonal and the fact ε is small, the vector \mathbf{x} has a predictable expected norm, i.e., $\mathbb{E}[\|\mathbf{x}\|_2^2] \approx \mathbb{E}[\|\mathbf{z}\|_2^2] = d$. It is a known fact that for vectors on a sphere, maximizing the ℓ^4 norm is equivalent to minimizing the ℓ^0 norm (for promoting sparsity),

$$\arg \max_{\mathbf{z} \in \mathbb{S}^n} \|\mathbf{z}\|_4 \Leftrightarrow \arg \min_{\mathbf{z} \in \mathbb{S}^n} \|\mathbf{z}\|_0. \quad (2.2.14)$$

This is illustrated in Figure 2.5.

An orthogonal matrix \mathbf{A} preserves the Euclidean (ℓ^2) norm: $\|\mathbf{A}\mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2$. Therefore, to find the correct orthogonal dictionary \mathbf{U} from \mathbf{X} , we may try to solve the following optimization program:

$$\max_{\tilde{\mathbf{A}} \in \mathcal{O}(D)} \frac{1}{4} \left\| \tilde{\mathbf{A}}\mathbf{X} \right\|_4^4 = \frac{1}{4} \sum_{i=1}^N \left\| \tilde{\mathbf{A}}\mathbf{x}_i \right\|_4^4 \quad (2.2.15)$$

This is known as the ℓ^4 maximization problem [ZMZ+20]. After we find the solution \mathbf{A}^* , we can take the transpose $\mathbf{U}^* = (\mathbf{A}^*)^\top$.

Remark 2.9. It is also known that for vectors on a sphere, minimizing the ℓ^1 norm is equivalent to minimizing the ℓ^0 norm (for promoting sparsity),

$$\arg \min_{\mathbf{z} \in \mathbb{S}^n} \|\mathbf{z}\|_1 \Leftrightarrow \arg \min_{\mathbf{z} \in \mathbb{S}^n} \|\mathbf{z}\|_0,$$

which is also illustrated in Figure 2.5. This fact can also be exploited to learn the dictionary \mathbf{A} effectively and efficiently. This was actually explored earlier than the ℓ^4 norm used here. Interested readers may refer to the work of [QZL+20b].

Note that the above problem is equivalent to the following constrained optimization problem:

$$\min -\frac{1}{4} \left\| \tilde{\mathbf{A}}\mathbf{X} \right\|_4^4 \quad \text{subject to} \quad \tilde{\mathbf{A}}^\top \tilde{\mathbf{A}} = \mathbf{I}. \quad (2.2.16)$$

As shown in [WM22], using the Lagrange multiplier method, one can derive that the optimal solution to the problem should satisfy the following “fixed point” condition:

$$\mathbf{A}^* = \mathcal{P}_{\mathcal{O}(D)}[(\mathbf{A}^*\mathbf{X})^{\odot 3}\mathbf{X}^\top], \quad (2.2.17)$$

where $\mathcal{P}_{\mathcal{O}(D)}[\cdot]$ is a projection onto the space of orthogonal matrices $\mathcal{O}(D)$.⁸

To compute the fixed point for the above equation, similar to how we computed eigenvectors for PCA (2.1.16), we may take the following power iteration:

$$\tilde{\mathbf{A}}_{t+1} = \mathcal{P}_{\mathcal{O}(D)}[(\tilde{\mathbf{A}}_t\mathbf{X})^{\odot 3}\mathbf{X}^\top]. \quad (2.2.18)$$

This is known as the *matching, stretching, and projection* (MSP) algorithm proposed by [ZMZ+20]. It was shown that under broad conditions such a greedy algorithm indeed converges to the correct solution at a superlinear rate.

⁸For any matrix \mathbf{M} with SVD $\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, $\mathcal{P}_{\mathcal{O}(D)}[\mathbf{M}] = \mathbf{U}\mathbf{V}^\top$. We leave this as an exercise for the reader.

Remark 2.10 (Global Optimality of ℓ^4 Maximization). The constrained ℓ^4 maximization problem is a nonconvex program. In general one should *not* expect that any greedy (say gradient-descent type) algorithms would converge to the globally optimal solution. Surprisingly, one can show that, unlike general non-convex programs, the landscape of ℓ^4 maximization over a sphere

$$\min -\frac{1}{4} \|\tilde{\mathbf{q}}^\top \mathbf{X}\|_4^4 \quad \text{subject to} \quad \tilde{\mathbf{q}}^\top \tilde{\mathbf{q}} = 1. \quad (2.2.19)$$

is very benign: All local minima are close to the global optima and all critical points are saddle points with a direction of negative curvature. Hence, any descent method with the ability of escaping strict saddle points provably finds global optimal solutions. For more precise statements, interested readers may refer to [QZL+20a].

Remark 2.11 (Stable Deep Linear Network). The above iterative process of computing the dictionary has a natural incremental “deep learning” interpretation. Let us define $\delta \mathbf{A}_{t+1} = \tilde{\mathbf{A}}_{t+1} \tilde{\mathbf{A}}_t^\top$ and $\mathbf{Z}_t = \tilde{\mathbf{A}}_t \mathbf{X}$, then it is easy to show that

$$\delta \mathbf{A}_{t+1} = \mathcal{P}_{\mathcal{O}(D)}[(\mathbf{Z}_t)^{\odot 3} \mathbf{Z}_t^\top].$$

If $\tilde{\mathbf{A}}_t$ converges to the correct dictionary \mathbf{U} , then the above iterative encoding process is essentially equivalent to a “deep linear network”:

$$\mathbf{Z} \longleftarrow \mathbf{Z}_{t+1} = \underbrace{\delta \mathbf{A}_{t+1} \delta \mathbf{A}_t \dots \delta \mathbf{A}_1}_{\text{forward constructed layers}} \mathbf{X}.$$

Note that the computation of the increment transforms $\delta \mathbf{A}_{t+1}$ at each layer depends only on the feature output from the previous layer \mathbf{Z}_t . The network is naturally stable as each layer is a norm-preserving orthogonal transform. Despite its resemblance to a linear deep network, backpropagation is unnecessary to learn each layer. All layers are learned in one forward pass!

2.2.3 Connection to ICA and Kurtosis

With the Bernoulli-Gaussian model, the variables z_i are independent and non-Gaussian. Then, there is a clear correspondence between the dictionary learning and the classic independent component analysis (ICA), to the extent that algorithms to solve one problem can be used to solve the other.⁹

Towards deriving an algorithm based on ICA, we focus on an objective function known as *kurtosis*, which is used in ICA as a direct consequence of the non-Gaussianity of the components. The *kurtosis*, or fourth-order cumulant, of a zero-mean random variable X is defined as

$$\text{kurt}(X) = \mathbb{E} X^4 - 3(\mathbb{E} X^2)^2. \quad (2.2.20)$$

⁹We explore this issue in more depth in Exercise 2.3, where a connection between non-Gaussianity of the independent components and the purely geometric notion of symmetry is made. This issue is related to our observation above that PCA does not work for recovering sparsely-used orthogonal dictionaries: in the statistical setting, it can be related to rotational invariance of the Gaussian distribution (Exercise 2.2).

If we have only finite samples from the random variable X arranged into a vector $\mathbf{x} = [x_1, \dots, x_N]$, we define kurtosis through their empirical average, which yields

$$\text{kurt}(\mathbf{x}) = \frac{1}{N} \|\mathbf{x}\|_4^4 - \frac{3}{N^2} \|\mathbf{x}\|_2^4. \quad (2.2.21)$$

Finally, for random vectors, we define their kurtosis as the sum of each component's scalar kurtosis. Kurtosis is a natural loss function for ICA because for Gaussian X , kurtosis is zero; the reader can verify further that the Bernoulli-Gaussian distribution has positive kurtosis. Thus a natural procedure for seeking non-Gaussian independent components is to search for a set of mutually-orthogonal directions $\tilde{\mathbf{U}} \in \mathbb{R}^{d \times k}$, such that $\tilde{\mathbf{U}}^\top \mathbf{X}$ has maximal kurtosis, where $\mathbf{X} = \mathbf{U}\mathbf{Z} \in \mathbb{R}^{D \times N}$ is the Bernoulli-Gaussian ICA data matrix (such that $\tilde{\mathbf{U}}$ approximates \mathbf{U} up to permutations of the columns). Formally, we seek to solve the problem

$$\max_{\tilde{\mathbf{U}}^\top \tilde{\mathbf{U}} = \mathbf{I}} \text{kurt}(\tilde{\mathbf{U}}^\top \mathbf{X}). \quad (2.2.22)$$

At one extreme, we may set $k = D$ and seek to recover the entire dictionary \mathbf{U} in a single shot. It can be shown that this problem can be solved with the MSP algorithm we have seen previously. At the other extreme, we may set $k = 1$ and seek to recover a single direction (column of \mathbf{U}) at a time, performing *deflation*, i.e., replacing the data matrix \mathbf{X} by $(\mathbf{I} - \mathbf{u}\mathbf{u}^\top)\mathbf{X}$, after each step before finding another direction. There is a natural tradeoff between the scalability of the $k = 1$ incremental approach and the efficiency and robustness of the $k = D$ approach.

Incremental ICA: correctness and FastICA algorithm. The FastICA algorithm, advanced by Hyvärinen and Oja [HO97], is a fast fixed-point algorithm for solving the $k = 1$ kurtosis maximization scheme for ICA. The problem at hand is

$$\max_{\|\tilde{\mathbf{u}}\|_2=1} \text{kurt}(\mathbf{X}^\top \tilde{\mathbf{u}}). \quad (2.2.23)$$

First, we will perform some very basic analysis of this objective to verify its correctness. Notice by the change of variables $\tilde{\mathbf{w}} = \mathbf{U}^\top \tilde{\mathbf{u}}$ that this problem is equivalent to

$$\max_{\|\tilde{\mathbf{w}}\|_2=1} \text{kurt}(\mathbf{Z}^\top \tilde{\mathbf{w}}).$$

This objective is simple enough that we can make strong statements about its correctness as a formulation for recovering the dictionary \mathbf{U} . For example, in the population setting where $N \rightarrow \infty$, we may use additivity properties of the kurtosis (Exercise 2.5) and our assumed normalization on the independent components to write the previous problem equivalently as

$$\max_{\|\tilde{\mathbf{w}}\|_2=1} \sum_{i=1}^d \text{kurt}(z_i) \tilde{w}_i^4. \quad (2.2.24)$$

It can be shown that under the Bernoulli-Gaussian assumption, the optimization landscape of this problem is “benign” (Exercise 2.7)—meaning that all local maxima of the objective function correspond to the recovery of one of the independent components.

In practice, we cannot directly compute the change of variables sending $\tilde{\mathbf{u}}$ to $\tilde{\mathbf{w}}$ because it requires the ground-truth dictionary \mathbf{U} , so we develop algorithms directly for the $\tilde{\mathbf{u}}$ parameterization and the data matrix \mathbf{X} . Then one efficient and scalable way to compute one of these maxima is via first-order optimization algorithms, which iteratively follow the gradient of the objective function and project onto the constraint set $\{\tilde{\mathbf{u}} \mid \|\tilde{\mathbf{u}}\|_2^2 = 1\}$. Since we have assumed that each z_i satisfies $\text{Var}(z_i) = 1$, we have for any $\tilde{\mathbf{u}}$, for large enough N

$$\text{kurt}(\mathbf{X}^\top \tilde{\mathbf{u}}) \approx \frac{1}{N} \|\mathbf{X}^\top \tilde{\mathbf{u}}\|_4^4 - 3\|\tilde{\mathbf{u}}\|_2^4. \quad (2.2.25)$$

We can then derive a corresponding approximation to the gradient:

$$\nabla_{\tilde{\mathbf{u}}} \text{kurt}(\mathbf{X}^\top \tilde{\mathbf{u}}) \approx \frac{4}{N} \mathbf{X}(\mathbf{X}^\top \tilde{\mathbf{u}})^{\odot 3} - 12\|\tilde{\mathbf{u}}\|_2^2 \tilde{\mathbf{u}}.$$

The FastICA algorithm uses a fixed-point method to compute directions of maximum kurtosis. It starts from the first-order optimality conditions for the kurtosis maximization problem, given the preceding gradient approximation and the constraint set, which read

$$\mathbf{X}(\mathbf{X}^\top \tilde{\mathbf{u}})^{\odot 3} = \underbrace{\langle \tilde{\mathbf{u}}, \mathbf{X}(\mathbf{X}^\top \tilde{\mathbf{u}})^{\odot 3} \rangle}_{\lambda} \tilde{\mathbf{u}}, \quad (2.2.26)$$

where the specific value of λ is determined using the unit norm constraint on $\tilde{\mathbf{u}}$. Exercise 2.6 describes the mathematical background necessary to derive these optimality conditions from first principles. Equation (2.2.26) is satisfied by *any* critical point of the kurtosis maximization problem; we want to derive an equation satisfied by only the maximizers. After noticing that $\lambda = \|\mathbf{X}^\top \tilde{\mathbf{u}}\|_4^4$, we equivalently re-express (2.2.26) as the modified equation

$$\frac{1}{N} \mathbf{X}(\mathbf{X}^\top \tilde{\mathbf{u}})^{\odot 3} - 3\tilde{\mathbf{u}} = \left(\frac{\lambda}{N} - 3 \right) \tilde{\mathbf{u}}, \quad (2.2.27)$$

and realize that any maximizer of (2.2.23) must satisfy $\lambda/N - 3 > 0$, assuming that N is sufficiently large. Hence, we may *normalize* both sides of (2.2.27), giving the following fixed-point equation satisfied by every maximizer of (2.2.23):

$$\frac{\frac{1}{N} \mathbf{X}(\mathbf{X}^\top \tilde{\mathbf{u}})^{\odot 3} - 3\tilde{\mathbf{u}}}{\left\| \frac{1}{N} \mathbf{X}(\mathbf{X}^\top \tilde{\mathbf{u}})^{\odot 3} - 3\tilde{\mathbf{u}} \right\|_2} = \tilde{\mathbf{u}}. \quad (2.2.28)$$

Iterating the mapping defined by the lefthand side of this fixed point expression then yields the FastICA algorithm of Hyvärinen and Oja [HO97]:

$$\tilde{\mathbf{v}}^+ = \frac{1}{N} \mathbf{X}(\mathbf{X}^\top \tilde{\mathbf{u}})^{\odot 3} - 3\tilde{\mathbf{u}}, \quad (2.2.29)$$

$$\tilde{\mathbf{u}}^+ = \tilde{\mathbf{v}}^+ / \|\tilde{\mathbf{v}}^+\|_2. \quad (2.2.30)$$

It turns out that the FastICA algorithm converges extremely rapidly (actually at a *cubic* rate) to columns of the dictionary \mathbf{U} ; interested readers may consult [HO97] for details. Comparing the FastICA algorithm (2.2.29) to the power method studied in Section 2.1.1 for the PCA problem and the MSP algorithm (2.2.18), we notice a striking similarity. Indeed, FastICA is essentially a modified power method, involving the gradient of the empirical kurtosis rather than the simpler linear gradient of the PCA objective.

2.3 A Mixture of Overcomplete Low-Dimensional Subspaces

As we have seen, complete dictionary learning enjoys an elegant computational theory in which we maintain a symmetric autoencoding structure $\mathcal{E}(\mathbf{x}) = \mathbf{U}^\top \mathbf{x}$, $\mathcal{D}(\mathbf{z}) = \mathbf{U}\mathbf{z}$, with a scalable power-method-like algorithm (the MSP algorithm) for learning an orthogonal dictionary/codebook \mathbf{U} from data. Nevertheless, for learning representations of general high-dimensional data distributions, one must expand the size of the codebook beyond the orthogonality requirement—meaning that we must have $\mathbf{A} \in \mathbb{R}^{D \times m}$, with $m \gg D$, corresponding to the case of an *overcomplete* dictionary/codebook,¹⁰ and the signal model

$$\mathbf{x} = \mathbf{A}\mathbf{z} + \boldsymbol{\varepsilon}, \quad \|\mathbf{z}\|_0 = d \ll m. \quad (2.3.1)$$

There are both geometric and physical/modeling motivations for passing to the overcomplete case. Geometrically, recall that in our original reduction from the mixture of subspace data model to the sparse dictionary model, a mixture of K subspaces in \mathbb{R}^D , each of dimension d , led to a dictionary of shape $\mathbf{A} \in \mathbb{R}^{D \times Kd}$. In other words, overcomplete dictionaries correspond to *richer* mixtures of subspaces, with more distinct modes of variability for modeling the high-dimensional data distribution. On the modeling side, we may run a computational experiment on real-world data that reveals the additional modeling power conferred by an overcomplete representation.

Example 2.2. Given sampled images of hand-written digits, Figure 2.6(a) shows the result of fitting an orthogonal dictionary to the dataset. In contrast, Figure 2.6(b) shows the result of running an optimization algorithm for learning overcomplete dictionaries (which we will present in detail later in the Chapter) on these samples. Notice that the representations become far sparser and the codebooks far more interpretable—they consist of fundamental primitives for the strokes composing the digits, i.e. oriented edges. ■

In fact, overcomplete dictionary learning was originally proposed as a biologically plausible algorithm for image representation based on empirical evidence of how early stages of the visual cortex represent visual stimuli [OF97; OF96].

In the remainder of this section, we will overview the conceptual and computational foundations of overcomplete dictionary learning. Supposing that the

¹⁰We change the notation here from \mathbf{U} to \mathbf{A} in order to emphasize the non-orthogonality and non-square-shape of the overcomplete dictionary \mathbf{A} .

model (2.3.1) is satisfied with sparse codes \mathbf{z} , overcomplete dictionary \mathbf{A} , and sparsity level d , and given samples $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ of \mathbf{x} , we want to learn an encoder $\mathcal{E} : \mathbb{R}^D \rightarrow \mathbb{R}^m$ mapping each \mathbf{x} to its sparse code \mathbf{z} , and a decoder $\mathcal{D}(\mathbf{z}) = \mathbf{A}\mathbf{z}$ reconstructing each \mathbf{x} from its sparse code. In diagram form:

$$\mathbf{x} \xrightarrow{\mathcal{E}} \mathbf{z} \xrightarrow{\mathcal{D}=\mathbf{A}} \hat{\mathbf{x}}. \quad (2.3.2)$$

We will start from the construction of the encoder \mathcal{E} . We will work incrementally: first, *given the true dictionary \mathbf{A}* , we will show how the problem of *sparse coding* gives an elegant, scalable, and provably-correct algorithm for recovering the sparse code \mathbf{z} of \mathbf{x} . Although this problem is NP-hard in the worst case, it can be solved efficiently and scalably for dictionaries \mathbf{A} which are *incoherent*, i.e. having columns that are not too correlated. The encoder architecture encompassed by this solution will no longer be symmetric: we will see it has the form of a primitive deep network, which depends on the dictionary \mathbf{A} .

Then we will proceed to the task of learning the decoder \mathcal{D} , or equivalently the overcomplete dictionary \mathbf{A} . We will derive an algorithm for overcomplete dictionary learning that allows us to simultaneously learn the codebook \mathbf{A} and the sparse codes \mathbf{z} , using ideas from sparse coding. Finally, we will discuss a more modern perspective on learnable sparse coding that leads us to a fully asymmetric encoder-decoder structure, as an alternative to (2.3.2). Here, the decoder will correspond to an incremental solution to the sparse dictionary learning problem, and yield a pair of deep network-like encoder decoders for sparse dictionary learning. This structure will foreshadow many developments to come in the remainder of the monograph, as we progress from analytical models to modern neural networks.

2.3.1 Sparse Coding with an Overcomplete Dictionary

In this section, we will consider the data model (2.3.1), which accommodates sparse linear combinations of many motifs, or *atoms*. Given data $\{\mathbf{x}_i\}_{i=1}^N \subseteq \mathbb{R}^D$ satisfying this model, i.e. expressible as

$$\mathbf{x}_i = \mathbf{A}\mathbf{z}_i + \boldsymbol{\varepsilon}_i, \quad \forall i \in [N] \quad (2.3.3)$$

for some dictionary $\mathbf{A} \in \mathbb{R}^{D \times m}$ with m atoms, sparse codes \mathbf{z}_i such that $\|\mathbf{z}_i\|_0 \leq d$, and small errors $\boldsymbol{\varepsilon}_i$, the sparse coding problem is to recover the codes \mathbf{z}_i as accurately as possible from the data \mathbf{x}_i , given the dictionary \mathbf{A} . Efficient algorithms to solve this problem succeed when the dictionary \mathbf{A} is *incoherent* in the sense that the inner products $\mathbf{a}_i^\top \mathbf{a}_j$ are uniformly small, hence the atoms are nearly orthogonal.¹¹

Note that we can collect the \mathbf{x}_i into $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$, collect the \mathbf{z}_i into $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N] \in \mathbb{R}^{d \times N}$, and collect the $\boldsymbol{\varepsilon}_i$ into $\mathbf{E} = [\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_N] \in \mathbb{R}^{D \times N}$, to rewrite (2.3.3) as

$$\mathbf{X} = \mathbf{AZ} + \mathbf{E}. \quad (2.3.4)$$

¹¹As it turns out, in a high-dimensional space, it is rather easy to pack a number of nearly orthogonal vectors that is much larger than the ambient dimension [WM22].

A natural approach to solving the sparse coding problem is to seek the sparsest signals that are consistent with our observations, and this naturally leads to the following optimization problem:

$$\min_{\tilde{\mathbf{Z}} \in \mathbb{R}^{d \times N}} \left\{ \|\mathbf{X} - \mathbf{A}\tilde{\mathbf{Z}}\|_F^2 + \lambda \|\tilde{\mathbf{Z}}\|_1 \right\}, \quad (2.3.5)$$

where the ℓ^1 norm $\|\tilde{\mathbf{Z}}\|_1$, taken elementwise, is known to promote sparsity of the solution [WM22]. This optimization problem is known as the LASSO problem.

However, unlike PCA or the complete dictionary learning problem, there is no clear power iteration-type algorithm to compute the optimal \mathbf{Z}^* . A natural alternative is to consider solving the above optimization problem with gradient descent type algorithms. Let $\mathcal{L}(\tilde{\mathbf{Z}}) = \|\mathbf{X} - \mathbf{A}\tilde{\mathbf{Z}}\|_2^2 + \lambda \|\tilde{\mathbf{Z}}\|_1$. Conceptually, we could try to compute \mathbf{Z}^* with the following iterations:

$$\tilde{\mathbf{Z}}_{t+1} \leftarrow \tilde{\mathbf{Z}}_t + \eta \nabla \mathcal{L}(\tilde{\mathbf{Z}}_t). \quad (2.3.6)$$

However, because the term associated with the ℓ^1 norm $\|\tilde{\mathbf{Z}}\|_1$ is non-smooth, we cannot just run gradient descent. For this type of function, we need to replace the gradient $\nabla \mathcal{L}(\tilde{\mathbf{Z}})$ with something that generalizes the notion of gradient, known as the subgradient $\partial \mathcal{L}(\tilde{\mathbf{Z}})$:

$$\tilde{\mathbf{Z}}_{t+1} \leftarrow \tilde{\mathbf{Z}}_t + \eta \partial \mathcal{L}(\tilde{\mathbf{Z}}_t). \quad (2.3.7)$$

However, it is known that the convergence of subgradient descent is usually very slow. Hence, for this type of optimization problems, it is conventional to adopt a so-called *proximal gradient descent*-type algorithm. We give a technical overview of this method in Section A.1.3.

Applying proximal gradient to the LASSO objective function (2.3.5) leads to the classic *iterative shrinkage-thresholding algorithm* (ISTA), which implements the iteration

$$\tilde{\mathbf{Z}}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (2.3.8)$$

$$\tilde{\mathbf{Z}}_{t+1} = S_{\eta\lambda} \left(\tilde{\mathbf{Z}}_t - 2\eta \mathbf{A}^\top (\mathbf{A}\tilde{\mathbf{Z}}_t - \mathbf{X}) \right), \quad \forall t \geq 1, \quad (2.3.9)$$

with step size $\eta \geq 0$, and the soft-thresholding operator S_α defined on scalars by

$$S_\alpha(x) \doteq \begin{cases} x - \alpha, & x \geq \alpha, \\ 0, & -\alpha < x < \alpha, \\ x + \alpha, & x \leq -\alpha \end{cases} \quad (2.3.10)$$

$$= \text{sign}(x) \max\{|x| - \alpha, 0\}, \quad (2.3.11)$$

and applied element-wise to the input matrix. As a proximal gradient descent algorithm applied to a convex problem, convergence to a global optimum is assured, and a precise convergence rate can be derived straightforwardly [WM22].

We now take a moment to remark on the functional form of the update operator in (2.3.9). It takes the form

$$\tilde{\mathbf{Z}}_{t+1} = \text{nonlinearity}(\tilde{\mathbf{Z}}_t + \text{linear}^\top(\text{linear}(\tilde{\mathbf{Z}}_t) + \text{bias})). \quad (2.3.12)$$

This functional form is very similar to that of a residual network layer, namely,

$$\tilde{\mathbf{Z}}_{t+1} = \tilde{\mathbf{Z}}_t + \text{linear}_1^\top(\text{nonlinearity}(\text{linear}_2(\tilde{\mathbf{Z}}_t) + \text{bias}_1) + \text{bias}_2), \quad (2.3.13)$$

only decoupling the linear maps (i.e., matrix multiplications), adding a bias, and moving the nonlinearity. The nonlinearity in (2.3.9) is notably similar to the commonly-used ReLU activation function $x \mapsto \max\{x, 0\}$ in deep learning—as visualized in Figure 2.7, the soft-thresholding operator is like a “signed” ReLU activation, which is also shifted by a bias. The ISTA, then, can be viewed as a forward pass of a primitive (recurrent one-layer) neural network. We argue in Chapter 5 that such operations are essential to deep representation learning.

2.3.2 Overcomplete Dictionary Learning

Recall that we have the data model

$$\mathbf{X} = \mathbf{AZ} + \mathbf{E}, \quad (2.3.14)$$

where \mathbf{Z} is sparse, and our goal previously was to estimate \mathbf{Z} given knowledge of the data \mathbf{X} and the dictionary atoms \mathbf{A} . Now we turn to the more practical and more difficult case where we do not know either \mathbf{A} or \mathbf{Z} and seek to learn them from a large dataset.

A direct generalization of Equation (2.3.5) suggests solving the problem

$$\min_{\tilde{\mathbf{A}}, \tilde{\mathbf{Z}}} \left\{ \|\mathbf{X} - \tilde{\mathbf{A}}\tilde{\mathbf{Z}}\|_F^2 + \lambda \|\tilde{\mathbf{Z}}\|_1 \right\}. \quad (2.3.15)$$

However, the bilinear term in Equation (2.3.15) introduces a scale ambiguity that hinders convergence: given any point $(\tilde{\mathbf{A}}, \tilde{\mathbf{Z}})$ and any constant $c > 0$, the substitution $(c^{-1}\tilde{\mathbf{A}}, c\tilde{\mathbf{Z}})$ gives loss value

$$\|\mathbf{X} - \tilde{\mathbf{A}}\tilde{\mathbf{Z}}\|_F^2 + c\lambda \|\tilde{\mathbf{Z}}\|_1. \quad (2.3.16)$$

This loss is evidently minimized over c by taking $c \rightarrow 0$, which corresponds to making the rescaled dictionary $c^{-1}\tilde{\mathbf{A}}$ go ‘to infinity’. In particular, the iterates of any optimization algorithm solving (2.3.15) will not converge.

This issue with (2.3.15) is dealt with by adding a constraint on the scale of the columns of the dictionary $\tilde{\mathbf{A}}$. For example, it is common to assume that each column $\tilde{\mathbf{A}}_j$ of the dictionary $\tilde{\mathbf{A}}$ in (2.3.14) has bounded ℓ^2 norm—say, without loss of generality, by 1. We then enforce this as a constraint, giving the objective

$$\min_{\tilde{\mathbf{Z}}, \tilde{\mathbf{A}}: \|\tilde{\mathbf{A}}_j\|_2 \leq 1} \left\{ \|\mathbf{X} - \tilde{\mathbf{A}}\tilde{\mathbf{Z}}\|_F^2 + \lambda \|\tilde{\mathbf{Z}}\|_1 \right\}. \quad (2.3.17)$$

Constrained optimization problems such as (2.3.17) can be solved by a host of sophisticated algorithms [NW06]. However, a simple and scalable one is actually furnished by the same proximal gradient descent algorithm that we used to solve the sparse coding problem in the previous section. We can encode each constraint as additional regularization term, via the characteristic function for the constraint set—details are given in Example A.1. Applying proximal gradient descent to the resulting regularized problem is equivalent to *projected gradient descent*, in which, at each iteration, the iterates after taking a gradient descent step are projected onto the constraint set.

Remark 2.12 (ℓ^4 maximization versus ℓ^1 minimization). Note that the above problem formulation follows naturally from the LASSO formulation (2.3.5) for sparse coding. We promote the sparsity of the solution via the ℓ^1 norm. Nevertheless, if we are only interested in recovering the over-complete dictionary \mathbf{A} , the ℓ^4 maximization scheme introduced in Section 2.2.2 also generalizes to the over-complete case, without any significant modification. Interested readers may refer to the work of [QZL+20a].

The above problem (2.3.17), which we call *overcomplete dictionary learning*, is nonconvex as here both \mathbf{A} and \mathbf{Z} are unknown. It cannot be solved easily by the standard convex optimization toolkit. Nevertheless, because it is interesting, simple to state, and practically important, there have been many important works dedicated to different algorithms and theoretical analysis for this problem. Here, for the interest of this manuscript, we present an idiomatic approach to solve this problem which is closer to the spirit of deep learning.

From our experience with the LASSO problem above, it is easy to see that, for the two unknowns \mathbf{A} and \mathbf{Z} , if we fix one and optimize the other, each subproblem is in fact convex and easy to solve. This naturally suggests that we could attempt to solve the above program (2.3.17) by minimizing against \mathbf{Z} or \mathbf{A} alternatively, say using gradient descent. Coupled with a natural choice of initialization, this leads to the following iterative scheme:

$$\tilde{\mathbf{Z}}^{\ell+1} = S_{\eta\lambda} \left(\tilde{\mathbf{Z}}^\ell - 2\eta \mathbf{A}_+^\top (\mathbf{A}_+ \tilde{\mathbf{Z}}^\ell - \mathbf{X}) \right), \quad \tilde{\mathbf{Z}}^1 = \mathbf{0}, \quad \forall \ell \in [L] \quad (2.3.18)$$

$$\mathbf{Z}^+ = \tilde{\mathbf{Z}}^L, \quad (2.3.19)$$

$$\tilde{\mathbf{A}}_{t+1} = \text{proj}_{\|(\cdot)_j\|_2 \leq 1, \forall j} \left(\tilde{\mathbf{A}}_t - 2\nu (\tilde{\mathbf{A}}_t \mathbf{Z}^+ - \mathbf{X})(\mathbf{Z}^+)^\top \right),$$

where $(\tilde{\mathbf{A}}_1)_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{1}{D} \mathbf{I}), \quad \forall j \in [m], \quad \forall t \in [T], \quad (2.3.20)$

$$\mathbf{A}_+ = \tilde{\mathbf{A}}_T, \quad (2.3.21)$$

where the projection operation in the update for \mathbf{A} ensures each column has at most unit ℓ^2 norm, via $\mathbf{A}_j \mapsto \mathbf{A}_j / \max\{\|\mathbf{A}_j\|_2, 1\}$, and where \mathbf{A}_+ is initialized with each column i.i.d. $\mathcal{N}(\mathbf{0}, \frac{1}{D} \mathbf{I})$. The above consists of one ‘block’ of alternating minimization, and we repeatedly perform such blocks, each with independent initializations, until convergence. Above, we have used two separate indices $\{t\}$ and $\{\ell\}$ to indicate the iterations. As we will see later, this allows us to interpret the two updates separately in the context of deep learning.

Despite the dictionary learning problem being a nonconvex problem, it has been shown that alternating minimization type algorithms indeed converge to the correct solution, at least locally. See, for example, [AAJ+16]. As a practical demonstration, the above algorithm (with $L = T = 1$) was used to generate the results for overcomplete dictionary learning in Figure 2.6.

2.3.3 Learned Deep Sparse Coding

The main insight from the alternating minimization algorithm for overcomplete dictionary learning in the previous section (Equations (2.3.18) and (2.3.20)) is to notice that *when we fix \mathbf{A} , the ISTA update for \mathbf{Z}^ℓ (2.3.18) looks like the forward pass of a deep neural network with weights given by \mathbf{A} (and \mathbf{A}^\top)*. But in general, we do not know the true \mathbf{A} , and the current estimate \mathbf{A}_+ could be erroneous. Hence it needs to be further updated using (2.3.20) based on the residual of using the current estimate of the sparse codes \mathbf{Z}^+ to reconstruct \mathbf{X} . The alternating minimization algorithm iterates these two procedures until convergence. But we can instead extrapolate, and design other learning procedures by combining these insights with techniques from deep learning. This leads to more interpretable network architectures, which will be a recurring theme throughout this manuscript.

Learned ISTA. The above deep-network interpretation of the alternating minimization is more conceptual than practical, as the process could be rather inefficient and take many layers or iterations to converge. But this is mainly because we try to infer both \mathbf{Z} and \mathbf{A} from \mathbf{X} . The problem can be significantly simplified and the above iterations can be made much more efficient in the *supervised* setting, where we have a dataset of input and output pairs (\mathbf{X}, \mathbf{Z}) distributed according to (2.3.14) and we only seek to learn \mathbf{A}^ℓ for the layerwise learnable sparse coding iterations (2.3.29):¹²

$$\mathbf{Z}^{\ell+1} = S_{\eta\lambda}(\mathbf{Z}^\ell - 2\eta(\mathbf{A}^\ell)^\top(\mathbf{A}^\ell \mathbf{Z}^\ell - \mathbf{X})), \quad \forall \ell \in [L]. \quad (2.3.22)$$

If we denote the operator for each iteration as $\mathbf{Z}^{\ell+1} = f(\mathbf{A}^\ell, \mathbf{Z}^\ell)$, the above iteration can be illustrated in terms of a diagram:

$$\mathbf{X}, \mathbf{Z}^1 \xrightarrow{f(\mathbf{A}^1, \cdot)} \mathbf{Z}^2 \xrightarrow{f(\mathbf{A}^2, \cdot)} \mathbf{Z}^3 \xrightarrow{f(\mathbf{A}^3, \cdot)} \dots \mathbf{Z}^L \xrightarrow{f(\mathbf{A}^L, \cdot)} \mathbf{Z}^{L+1} \approx \mathbf{Z}.$$

Thus, given the sequential architecture, to learn the operator \mathbf{A}^ℓ at each layer, it is completely natural to learn it, say via back propagation (BP),¹³ by minimizing the error between the final code \mathbf{Z}^L and the ground truth \mathbf{Z} :

$$\min_{\{\mathbf{A}^\ell\}} \|\mathbf{Z}^L(\mathbf{A}^1, \dots, \mathbf{A}^L) - \mathbf{Z}\|_2^2. \quad (2.3.23)$$

¹²We drop the “tilde” accents here for concision, as the indexing makes the optimization-dependence clear.

¹³See Section A.2 for a brief description of BP.

This is the basis of the Learned ISTA (LISTA) algorithm [GL10], which can be viewed as the learning algorithm for a deep neural network, which tries to emulate the sparse encoding process from \mathbf{X} to \mathbf{Z} . In particular, it can be viewed as a *simple representation learning algorithm*. In fact, this same methodology can be used as a basis to understand the representations computed in more powerful network architectures, such as transformers. We develop these implications in detail in Chapter 5.

Sparse Autoencoders. The original motivation for overcomplete dictionary learning was to provide a simple generative model for high-dimensional data. We have seen with LISTA that, in addition, iterative algorithms for learning sparsely-used overcomplete dictionaries provide an interpretation for ReLU-like deep networks, which we will generalize in later chapters to more complex data distributions than (2.3.14). But it is also worth noting that even in the modern era of large models, the data generating model (2.3.14) provides a useful practical basis for *interpreting features in pretrained large-scale deep networks*, such as transformers, following the hypothesis that the (non-interpretable, *a priori*) features in these networks consist of sparse “superpositions” of underlying features, which are themselves interpretable [EHO+22b]. These *unsupervised* learning paradigms are generally more data friendly than LISTA, as well, which requires large amounts of labeled (\mathbf{X}, \mathbf{Z}) pairs for supervised training.

We can use our development of the LISTA algorithm above to understand common practices in this field of research. In the most straightforward instantiation (see [GTT+25; HCS+24]), a large number of features from a pretrained deep network h are collected from different inputs \mathbf{x}_i , which themselves are chosen based on a desired interpretation task.¹⁴ For simplicity, we will use h to denote the pre-selected feature map in question, with D -dimensional features; given N sample inputs, let $\mathbf{H} \in \mathbb{R}^{D \times N}$ denote the full matrix of features of h . Then a so-called sparse autoencoder $f : \mathbb{R}^D \rightarrow \mathbb{R}^d$, with decoder $g : \mathbb{R}^d \rightarrow \mathbb{R}^D$, is trained via the LASSO loss (2.3.5):

$$\min_{f,g} \|\mathbf{H} - g(f(\mathbf{H}))\|_F^2 + \lambda \|\mathbf{H}\|_1, \quad (2.3.24)$$

where the sparse autoencoder f takes the form of a one-layer neural network, i.e. $f(\mathbf{h}_i) = \sigma(\mathbf{W}_{\text{enc}}(\mathbf{h}_i - \mathbf{b}) + \mathbf{b}_{\text{enc}})$, where $\sigma(x) = \max\{x, 0\}$ is the ReLU activation function, and the decoder g is linear, so that $g(\mathbf{z}_i) = \mathbf{W}_{\text{dec}}\mathbf{z}_i + \mathbf{b}$.¹⁵

¹⁴For example, the inputs \mathbf{x}_i could correspond to texts containing samples of computer code in different programming languages, with our task being to try to identify interpretable features in a transformer feature map h corresponding to different salient aspects of the input, such as the specific programming language (distinct across input “classes”) or the need to insert a matching parenthesis at the current position (common across input “classes”). We discuss the use of deep networks, and in particular transformers, for text representation learning in greater detail in Chapter 8.

¹⁵Since f is a sparsifying operator, we denote it analogously to the previously used f , say in LISTA where f denotes an optimization algorithm to sparsely encode the input. Throughout the book, f will be used to denote such sparse and otherwise structured transformations, including those of *deep network encoders* (Chapters 4 and 5).

The parameterization and training procedure (2.3.24) may initially seem to be an arbitrary application of deep learning to the sparse coding problem, but it is actually highly aligned with the algorithms we have studied above for layerwise sparse coding with a learned dictionary. In particular, recall the LISTA architecture $\mathbf{Z}^L = f(\mathbf{A}^L, f(\mathbf{A}^{L-1}, \dots, f(\mathbf{A}^1, \mathbf{X}) \dots))$. In the special case $L = 2$, we have

$$\mathbf{Z}^2 = f(\mathbf{A}^1, \mathbf{X}) = S_{\eta\lambda}(\mathbf{Z}^1 - 2\eta(\mathbf{A}^1)^\top(\mathbf{A}^1\mathbf{Z}^1 - \mathbf{X})). \quad (2.3.25)$$

Let us assume that the sparse codes \mathbf{Z} in question are nonnegative, i.e., that $\mathbf{Z} \geq \mathbf{0}$.¹⁶ Then (see Example A.3), we can consider an equivalent LISTA architecture obtained from the sparse coding objective with an additional non-negativity constraint on \mathbf{Z} as

$$\mathbf{Z}^2 = f(\mathbf{A}^1, \mathbf{X}) = \max\{\mathbf{Z}^1 - 2\eta(\mathbf{A}^1)^\top(\mathbf{A}^1\mathbf{Z}^1 - \mathbf{X}) - \lambda\eta\mathbf{1}, 0\}, \quad (2.3.26)$$

and after some algebra, express this as

$$\mathbf{Z}^2 = f(\mathbf{A}^1, \mathbf{X}) = \max\{2\eta(\mathbf{A}^1)^\top + (\mathbf{Z}^1 - 2\eta(\mathbf{A}^1)^\top\mathbf{A}^1\mathbf{Z}^1 - \lambda\eta\mathbf{1}), 0\}. \quad (2.3.27)$$

Given the ability to change the sparse code initialization \mathbf{Z}^1 as a learnable parameter (which, in the current framework, must have all columns equal to the same learnable vector), this has the form of a ReLU neural network with learnable bias—identical to the sparse autoencoder f ! Moreover, to *decode* the learned sparse codes \mathbf{Z}^2 , it is natural to apply the learned dictionary $\mathbf{Z}^2 \mapsto \mathbf{A}^1\mathbf{Z}^2$. Then the only difference between this and the SAE decoder g is the additional bias \mathbf{b} , which can technically be absorbed into \mathbf{H} and f in the training objective (2.3.24).

Thus, the SAE parameterization and training procedure coincides with LISTA training with $L = 1$, and a modified training objective—using the LASSO objective (2.3.5), which remains *unsupervised*, instead of the supervised reconstruction loss (2.3.23) used in vanilla LISTA. In particular, we can understand the SAE architecture in terms of our interpretation of the LISTA architecture in terms of layerwise sparse coding in (2.3.29). This connection is suggestive of a host of new design strategies for improving practical interpretability methodology, many of which remain tantalizingly unexplored. We begin to lay out some connections to broader autoencoding methodology in Chapter 6.

Layerwise learned sparse coding? In the supervised setting, LISTA provides a deep neural network analogue of the sparse coding iteration, with layerwise learned dictionaries, inspired by alternating minimization; even in the unsupervised setting, the same methodology can be applied to learning, as with sparse autoencoders. But the connection between low-dimensional-structure-seeking optimization algorithms and deep network architectures goes much deeper than

¹⁶In the data generating model (2.3.14), an arbitrary dictionary-and-sparse-code pair (\mathbf{A}, \mathbf{Z}) can be replaced by one in which $\mathbf{Z} \geq \mathbf{0}$ simply by doubling the number of columns in \mathbf{A} , so from a modeling perspective, this is a very mild assumption.

this, and suggests an array of scalable and natural neural learning architectures which may even be usable without backpropagation.

As a simple illustration, we return the alternating minimization iterations (2.3.18) and (2.3.20). This scheme randomly re-initializes the dictionary \mathbf{A}_1 on every such update. An improvement uses instead *warm starting*, where the residual is generated using the previous estimate \mathbf{A}_+ for the dictionary. If we then view each ISTA update (2.3.18) as a layer and allow the associated dictionary, now coupled with the sparse code updates as \mathbf{A}_ℓ , to update in time, this leads to a “layerwise-learnable” sparse coding scheme:

$$\mathbf{Z}^1 = \mathbf{0}, \quad (\mathbf{A}_1)_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{1}{D}\mathbf{I}), \quad \forall j \in [m], \quad (2.3.28)$$

$$\mathbf{Z}^{\ell+1} = S_{\eta\lambda}(\mathbf{Z}^\ell - 2\eta(\mathbf{A}_\ell)^\top(\mathbf{A}_\ell\mathbf{Z}^\ell - \mathbf{X})), \quad (2.3.29)$$

$$\mathbf{A}_{\ell+1} = \mathbf{A}_\ell - 2\nu(\mathbf{A}_\ell\mathbf{Z}^{\ell+1} - \mathbf{X})(\mathbf{Z}^{\ell+1})^\top. \quad (2.3.30)$$

Note that this iteration corresponds to a relabeling of (2.3.18) and (2.3.20) for $T = L = 1$, over infinitely many blocks. Each of the ‘inner’ steps updating \mathbf{Z} can be considered as a one-layer forward pass, while each of the ‘outer’ steps updating \mathbf{A} can be considered as a one-layer backward pass, of a primitive deep neural network. In particular, this algorithm is the simplest case in which a clear divide between forward optimization and backward learning manifests. This distinction between forward updates to the ephemeral “features” \mathbf{Z} of data \mathbf{X} versus backward updates to the persistent parameters \mathbf{A} is still observed in current neural networks and autoencoders—we will have much more to say about it in Chapter 5 and in Chapter 6.

Notice that the above layer-wise scheme also suggests a plausible alternative to the current end-to-end optimization strategy that primarily relies on back propagation (BP) detailed in Section A.2.3. Freeing training large networks from BP would be one of the biggest challenges and opportunities in the future, as we will discuss more at the end of the book in Chapter 9.

2.4 Summary and Notes

Key Messages. The idealistic models we have presented in this chapter—PCA, ICA, and dictionary learning—were developed over the course of the twentieth century. Many books have been written solely about each method, so we only attempted here to give a broad overview of the key concepts and results. As we will see in coming chapters, although these are arguably rather idealistic simple models for low-dimensional data distributions, the basic ideas and even resulting algorithms to identify such distributions can be naturally generalized or extended to arbitrary low-dimensional distributions – the thesis of this book. In addition, as we will see in coming chapters, these basic (linear) models, to a large extent, can also serve as local prototypes to effectively approximately model general distributions with nonlinear low-dimensional structures.

More History and References. Jolliffe [Jol86] attributes principal component analysis to Pearson [Pea01], and independently Hotelling [Hot33]. In mathematics, the main result on the related problem of low-rank approximation in unitarily invariant norms is attributed to Eckart and Young [EY36], and to Mirsky for full generality [Mir60]. PCA continues to play an important role in research as perhaps the simplest model problem for unsupervised representation learning: as early as the 1980s, works such as Oja [Oja82] and Baldi and Hornik [BH89] used the problem to understand learning in primitive neural networks, and more recently, it has served as a tool for understanding more complex representation learning frameworks, such as diffusion models [WZZ+24].

Independent component analysis was proposed by B. Ans et al. [BJC85] and pioneered by Aapo Hyvärinen in the 1990s and early 2000s in a series of influential works: see Hyvärinen and Oja [HO00b] for a summary. As a simple model for structure that arises in practical data, it initially saw significant use in applications such as blind source separation, where each independent component z_i represents an independent source (such as sound associated to a distinct instrument in a musical recording) that is superimposed to produce the observation $\mathbf{x} = \mathbf{U}\mathbf{z}$.

The problem of dictionary learning can, in the complete or orthogonal case, be seen as one of the foundational problems of twentieth-century signal processing, particularly in linear systems theory, where the Fourier basis plays the key role; from the 1980s onward, the field of computational harmonic analysis crystallized around the study of alternate such dictionaries for classes of signals in which optimal approximation could only be realized in a basis other than Fourier (e.g., wavelets) [DVD+98]. However, the importance of the case of redundant bases, or overcomplete dictionaries, was only highlighted following the pioneering work of Olshausen and Field [OF97; OF96]. Early subsequent work established the conceptual and algorithmic foundations for learning sparsely-used overcomplete dictionaries, often aimed at representing natural images [AEB06; DM03; Don01; EA06; GJB15; MBP14; MK07]. Later, a significant amount of theoretical interest in the problem, as an important and nontrivial model problem for unsupervised representation learning, led to its study by the signal processing, theoretical machine learning, and theoretical computer science communities, in particular focused on conditions under which the problem could be provably and efficiently solved. A non-exhaustive list of notable works in this line include those of Spielman et al. [SWW12] and Sun et al. [SQW17a] on the complete case; Arora et al. [AGM+15]; Barak et al. [BKS15]; and Qu et al. [QZL+20a]. Many deep theoretical questions about this simple-to-state problem remain open, perhaps in part due to a tension with the problem’s worst-case NP-hardness (e.g., see Tillmann [Til15]).

One point that we wish to highlight from the study of these classical analytical models for low-dimensional structure is the common role played by various *generalized power methods*—algorithms that very rapidly converge, at least locally, to various types of low-dimensional structures. The terminology for this class of algorithms follows the work of M. Journée et al. [MYP+10]. At a high level, modeled on the classical power iteration for computation of the top

Table 2.1: Summary of (generalized) power methods presented in this chapter.

Problem	Algorithm	Iteration	Type of Structure Enforced
PCA	Power Method	$\mathbf{u}_t = \frac{\mathbf{X}\mathbf{X}^\top \mathbf{u}_t}{\ \mathbf{X}\mathbf{X}^\top \mathbf{u}_t\ _2}$	1-dim. subspace (unit vector)
ICA	FastICA	$\mathbf{u}_{t+1} = \frac{\frac{1}{N}\mathbf{X}(\mathbf{X}^\top \mathbf{u}_t)^{\odot 3} - 3\mathbf{u}_t}{\ \frac{1}{N}\mathbf{X}(\mathbf{X}^\top \mathbf{u}_t)^{\odot 3} - 3\mathbf{u}_t\ _2}$	1-dim. subspace (unit vector)
Complete DL	MSP Algorithm	$\mathbf{U}_{t+1} = \mathcal{P}_{\mathcal{O}(D)}[(\mathbf{U}_t \mathbf{X})^{\odot 3} \mathbf{X}^\top]$	D -dim. subspace (ortho. matrix)

eigenvector of a semidefinite matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, that is

$$\mathbf{u}_{t+1} = \frac{\mathbf{A}\mathbf{u}_t}{\|\mathbf{A}\mathbf{u}_t\|_2}, \quad (2.4.1)$$

this class of algorithms consists of a “powering” operation involving a matrix \mathbf{A} associated to the data, along with a “projection” operation that enforces a desired type of structure. Table 2.1 presents a summary of the algorithms we have studied in this chapter that follow this structure. The reader may appreciate the applicability of this methodology to different types of low-dimensional structure, and different losses (i.e., both the quadratic loss from PCA, and the kurtosis-type losses from ICA), as well as the lack of such an algorithm for overcomplete dictionary learning, despite the breadth of the literature on these algorithms. We see the development of power methods for further families of low-dimensional structures, particularly those relevant to applications where deep learning is prevalent, as one of the more important (and open) research questions suggested by this chapter.

The connection we make in Section 2.2.1 between the geometric mixture-of-subspaces distributional assumption and the more analytically-convenient sparse dictionary assumption has been mentioned in prior work, especially by those focused on generalized principal component analysis and applications such as subspace clustering, e.g., work of Vidal et al. [VMS16]. The mixture of subspaces assumption will continue to play a significant role throughout this manuscript, both as an analytical test case for different algorithmic paradigms, and as a foundation for deriving different deep network architectures, as with LISTA in Section 2.3.3, but which can scale to more complex data distributions.

2.5 Exercises and Extensions

Exercise 2.1. Prove that, for any symmetric matrix \mathbf{A} , the solution to the problem $\max_{\mathbf{U} \in \mathcal{O}(D,d)} \text{tr}(\mathbf{U}^\top \mathbf{A} \mathbf{U})$ is the matrix \mathbf{U}^* whose columns are the top d unit eigenvectors of \mathbf{A} .

Exercise 2.2. Let $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ be a Gaussian random variable with independent components, each with variance σ^2 . Prove that for any orthogonal matrix \mathbf{Q} (i.e., $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$), the random variable $\mathbf{Q}\mathbf{z}$ is distributed identically to \mathbf{z} . (Hint: recall the formula for the Gaussian probability density function, and the formula for the density of a linear function of a random variable.)

Exercise 2.3. The notion of statistical identifiability discussed above can be related to *symmetries* of the model class, allowing estimation to be understood in a purely deterministic fashion without any statistical assumptions.

Consider the model $\mathbf{X} = \mathbf{U}\mathbf{Z}$ for matrices $\mathbf{X}, \mathbf{U}, \mathbf{Z}$ of compatible sizes.

1. Show that if \mathbf{A} is any square invertible matrix of compatible size, then the pair $(\mathbf{U}\mathbf{A}, \mathbf{A}^{-1}\mathbf{Z})$ also equals \mathbf{X} under the model. We call this a *GL(d) symmetry*.
2. Suppose each column of \mathbf{Z} is an independent and identically distributed observation from a common statistical model \mathbf{z} , which moreover has zero mean and independent components z_i with positive variance. Show that for any square invertible matrix \mathbf{A} , if $\mathbf{A}\mathbf{z}$ has uncorrelated components, then \mathbf{A} can be written as $\mathbf{D}_1\mathbf{Q}\mathbf{D}_2$, where \mathbf{Q} is an orthogonal matrix and $\mathbf{D}_1, \mathbf{D}_2$ are diagonal matrices. *This links the “independence” assumption in ICA to a “symmetry breaking” effect, which only allows scale and rotational symmetries.*

Exercise 2.4. Consider the model $\mathbf{x} = \mathbf{U}\mathbf{z}$, where $\mathbf{U} \in \mathbb{R}^{D \times d}$ with $D \geq d$ is fixed and has rank d , and \mathbf{z} is a zero-mean random variable. Let $\mathbf{x}_1, \dots, \mathbf{x}_N$ denote i.i.d. observations from this model.

1. Show that the matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ has rank no larger than d , and therefore there is an orthonormal matrix $\mathbf{V} \in \mathbb{R}^{D \times d}$ so that $\mathbf{X} = \mathbf{V}\mathbf{Y}$, where $\mathbf{Y} \in \mathbb{R}^{d \times N}$. (*Hint: use PCA.*)
2. Show that the *whitened matrix* $(\mathbf{Y}\mathbf{Y}^\top)^{-1/2}\mathbf{Y}$ exists in expectation whenever $\text{Cov}(\mathbf{z})$ is nonsingular, and that it has identity empirical covariance.¹⁷
3. Show, by using the singular value decomposition of \mathbf{U} , that the matrix \mathbf{V} can be chosen so that the whitened matrix satisfies $(\mathbf{Y}\mathbf{Y}^\top)^{-1/2}\mathbf{Y} = \mathbf{W}[\mathbf{z}_1, \dots, \mathbf{z}_N]$, where \mathbf{W} is an orthonormal matrix.

Exercise 2.5. Let X and Y be zero-mean independent random variables.

1. Show that $\text{kurt}(X + Y) = \text{kurt}(X) + \text{kurt}(Y)$.
2. For any $\alpha \in \mathbb{R}$, show that $\text{kurt}(\alpha X) = \alpha^4 \text{kurt}(X)$.

Exercise 2.6. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a given twice-continuously-differentiable objective function. Consider the spherically-constrained optimization problem

$$\max_{\|\mathbf{u}\|_2=1} f(\mathbf{u}). \quad (2.5.1)$$

In this exercise, we will derive the expressions we gave in the FastICA derivation for maximizing kurtosis over the sphere via a gradient ascent algorithm. These

¹⁷In particular, it can be proved mathematically that this is enough to guarantee that the whitened matrix exists with high probability whenever \mathbf{z} satisfies a suitable concentration inequality and N is sufficiently large.

expressions are special cases of a rich theory of calculus and optimization on manifolds, of which the sphere is a particular example. A deep technical study of this field is out-of-scope for our purposes, so we only mention two key references for the interested reader: the pioneering textbook by Absil, Mahony, and Sepulchre [AMS09], and a more recent introductory treatise by Boumal [Bou23].

1. For any constraint set \mathcal{M} that is a differentiable submanifold of \mathbb{R}^d , the *tangent space* at a point $\mathbf{u} \in \mathcal{M}$ is, informally, the best local linear approximation to the manifold \mathcal{M} at the point \mathbf{u} . In the important special case where \mathcal{M} is defined locally at \mathbf{u} as a level set of a function $F : \mathbb{R}^d \rightarrow \mathbb{R}$, that is

$$U \cap \mathcal{M} = F^{-1}(\{0\})$$

for some open set $U \subset \mathcal{M}$ with $\mathbf{u} \in U$, the tangent space to \mathcal{M} at \mathbf{u} can be calculated via differentiation:

$$T_{\mathbf{u}}\mathcal{M} = \text{Ker}(DF_{\mathbf{u}}).$$

It is easily seen that the sphere has the defining equation $F(\mathbf{u}) = \|\mathbf{u}\|_2^2 - 1$. Show, using these facts, that the tangent space to the sphere at \mathbf{u} is given by

$$T_{\mathbf{u}}\mathbb{S}^{d-1} = \{\mathbf{v} \in \mathbb{R}^d \mid \langle \mathbf{v}, \mathbf{u} \rangle = 0\},$$

and that the orthogonal projection onto this subspace is $\mathbf{P}_{\mathbf{u}}^{\perp} = \mathbf{I} - \mathbf{u}\mathbf{u}^{\top}$.

2. The vector field

$$\text{grad } f(\mathbf{u}) = \mathbf{P}_{\mathbf{u}}^{\perp} \nabla f \tag{2.5.2}$$

is known as the *Riemannian gradient* of the function f restricted to the sphere. The *first order optimality conditions* for the optimization problem (2.5.1) can be expressed in terms of the Riemann gradient:

$$\text{grad } f(\mathbf{u}) = \mathbf{0}.$$

Geometrically, this says that the Euclidean gradient of f at \mathbf{u} must be orthogonal to the tangent space to the sphere at \mathbf{u} . Now suppose $\mathbf{v} \in \mathbb{R}^d$ is nonzero. Show that

$$\text{proj}_{\mathbb{S}^{d-1}}(\mathbf{v}) \doteq \min_{\|\mathbf{u}\|_2=1} \|\mathbf{u} - \mathbf{v}\|_2 = \frac{\mathbf{v}}{\|\mathbf{v}\|_2},$$

using the first-order optimality conditions.

3. In optimization over \mathbb{R}^d , one checks the second-order optimality conditions (to determine whether a critical point is a maximizer, a minimizer, or a saddle point) using the *Hessian matrix* $\nabla^2 f(\mathbf{u})$. Show, by differentiating the Riemann gradient $\text{grad } f(\mathbf{u})$ for the sphere with respect to \mathbf{u} as in the first part of this exercise, that the corresponding object for determining second-order optimality conditions for sphere-constrained optimization is the *Riemannian Hessian*, defined as

$$\text{Hess } f(\mathbf{u}) = \mathbf{P}_{\mathbf{u}}^{\perp} (\nabla^2 f(\mathbf{u}) - \langle \nabla f(\mathbf{u}), \mathbf{u} \rangle \mathbf{I}) \mathbf{P}_{\mathbf{u}}^{\perp}. \tag{2.5.3}$$

Exercise 2.7. In this exercise, we sketch an argument referred to in the literature as a *landscape analysis* for the spherically-constrained population kurtosis maximization problem (2.2.24). We will show that when there is at least one independent component with positive kurtosis, its global maximizers indeed lead to the recovery of one column of the dictionary \mathbf{U} . For simplicity, we will assume that $\text{kurt}(z_i) \neq 0$ for each $i = 1, \dots, d$.

- Using the results of Part 1 of Exercise 2.6, show that the first-order optimality condition for (2.2.24) is

$$\left(\sum_{i=1}^d \text{kurt}(z_i) w_i^4 \right) \mathbf{w} = \text{kurt}(\mathbf{z}) \odot \mathbf{w}^{\odot 3}, \quad (2.5.4)$$

where the kurtosis is calculated elementwise, \odot denotes elementwise multiplication of vectors and $\mathbf{w}^{\odot 3}$ denotes the elementwise cube of its argument.

- Show that the vectors \mathbf{w} with unit norm that also satisfy (2.5.4) all take the following form. Let $S^+ = \{i \in [d] \mid \text{kurt}(z_i) > 0\}$, and $S^- = \{i \in [d] \mid \text{kurt}(z_i) < 0\}$. Let S be a subset either of S^+ or S^- . Then

$$\mathbf{w}_S = \sum_{i \in S} \pm \sqrt{\frac{1}{\text{kurt}(z_i) \sum_{j \in S} \frac{1}{\text{kurt}(z_j)}}} \mathbf{e}_i \quad (2.5.5)$$

satisfies (2.5.4), where \mathbf{e}_i is the vector with a 1 in the i -th position and 0s elsewhere, and the \pm sign denotes the choice of either a positive or negative sign.

- Assume that there is at least one i such that $\text{kurt}(z_i) > 0$. Using the results of Part 2 of Exercise 2.6, show that the only local maxima of the objective of (2.2.24) are the signed one-sparse vectors $\pm \mathbf{e}_i$ with $i \in S^+$. Conclude that the global maximizers of (2.2.24) are the signed one-sparse vectors corresponding to components with maximum kurtosis. (*Hint: count the number of positive and negative eigenvalues of the Riemannian Hessian (2.5.3) at each critical point.*)
- Now assume that $\text{kurt}(z_j) < 0$ for every $j = 1, \dots, d$. This corresponds to an “over-deflated” instantiation of the kurtosis maximization problem. Using again the results of Part 2 of Exercise 2.6, show that the only local maxima of the objective of (2.2.24) are the signed dense vectors $\sum_{i=1}^d \pm \mathbf{e}_i$. This shows that the optimization formulation (2.2.24) cannot be applied naively.

Exercise 2.8. This exercise follows the structure and formalism introduced in Exercise 2.6, but applies it instead to the orthogonal group $\mathbf{O}(d) = \{\mathbf{U} \in \mathbb{R}^{d \times d} \mid \mathbf{U}^\top \mathbf{U} = \mathbf{I}\}$. Consult the description of Exercise 2.6 for the necessary conceptual background; the formalism applies identically to the case where the ambient space is the set of $d \times d$ matrices as long as one recalls that the relevant inner product on matrices is $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{tr}(\mathbf{X}^\top \mathbf{Y})$. An excellent general reference for

facts about optimization on the orthogonal group is Edelman, Arias, and Smith [EAS98].

Let $f : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$ be a given twice-continuously-differentiable objective function. Consider the orthogonally-constrained optimization problem

$$\max_{\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}} f(\mathbf{Q}). \quad (2.5.6)$$

1. It is easily seen that the orthogonal group has the defining equation $F(\mathbf{Q}) = \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$. Show, using this fact, that the tangent space to the orthogonal group at \mathbf{Q} is given by

$$T_{\mathbf{Q}}\mathbf{O}(d) = \{\mathbf{Q}\boldsymbol{\Omega} \in \mathbb{R}^{d \times d} \mid \boldsymbol{\Omega}^\top = -\boldsymbol{\Omega}\},$$

and that the orthogonal projection onto this subspace is

$$\mathcal{P}_{T_{\mathbf{Q}}\mathbf{O}(d)}(\boldsymbol{\Delta}) = \mathbf{Q} \operatorname{skew}(\mathbf{Q}^\top \boldsymbol{\Delta}),$$

where $\operatorname{skew}(\boldsymbol{\Delta}) = \frac{1}{2}(\boldsymbol{\Delta} - \boldsymbol{\Delta}^\top)$ is the orthogonal projection onto the set of skew-symmetric matrices. The vector field

$$\operatorname{grad} f(\mathbf{Q}) = \mathcal{P}_{T_{\mathbf{Q}}\mathbf{O}(d)}(\nabla f(\mathbf{Q})) \quad (2.5.7)$$

is known as the *Riemannian gradient* of the function f restricted to the orthogonal group. The *first order optimality conditions* for the optimization problem (2.5.6) can be expressed in terms of the Riemann gradient:

$$\operatorname{grad} f(\mathbf{Q}) = \mathbf{0}.$$

2. Show, by differentiating the Riemann gradient $\operatorname{grad} f(\mathbf{Q})$ for the orthogonal group with respect to \mathbf{Q} as in the first part of this exercise, that the *Riemannian Hessian* is given by

$$\operatorname{Hess} f(\mathbf{Q}) = \mathcal{P}_{T_{\mathbf{Q}}\mathbf{O}(d)}(\nabla^2 f(\mathbf{Q}) - \operatorname{sym}(\mathbf{Q}^\top \nabla f(\mathbf{Q})) \otimes \mathbf{I}) \mathcal{P}_{T_{\mathbf{Q}}\mathbf{O}(d)}, \quad (2.5.8)$$

where $\operatorname{sym}(\boldsymbol{\Delta}) = \frac{1}{2}(\boldsymbol{\Delta} + \boldsymbol{\Delta}^\top)$ denotes the orthogonal projection onto the set of symmetric matrices, and \otimes denotes the Kronecker product of matrices. Take care to interpret the operators appearing in the previous expression as *linear transformations on $d \times d$ matrices*, **not** as $d \times d$ matrices themselves. The *second-order optimality conditions* for the optimization problem (2.5.6) can be expressed in terms of the Riemann Hessian:

$$\operatorname{Hess} f(\mathbf{Q}) \preceq \mathbf{0}.$$

For a minimization problem, the sign is reversed.

(Hint: The key is to manipulate one's calculations to obtain the form (2.5.8), which is as compact as possible. To this end, make use of the following

isomorphism of the Kronecker product: if \mathbf{A} , \mathbf{X} , and \mathbf{B} are matrices of compatible sizes, then one has

$$(\mathbf{B}^\top \otimes \mathbf{A}) \operatorname{vec}(\mathbf{X}) = \operatorname{vec}(\mathbf{A}\mathbf{X}\mathbf{B}),$$

where $\operatorname{vec}(\mathbf{X})$ denotes the “left-to-right” stacking of the columns of the matrix argument into a vector. We use this isomorphism in (2.5.8) in order to define the Kronecker product of two matrices as an operator on matrices in a canonical way.)

3. Now suppose $\mathbf{X} \in \mathbb{R}^{d \times d}$ is full-rank. In this and the next part of the exercise, we consider the projection onto the orthogonal group of \mathbf{X} :

$$\operatorname{proj}_{\mathbf{O}(d)}(\mathbf{X}) \doteq \min_{\mathbf{Q} \in \mathbf{O}(d)} \|\mathbf{Q} - \mathbf{X}\|_F^2. \quad (2.5.9)$$

We will prove that the solution to this problem is given by

$$\operatorname{proj}_{\mathbf{O}(d)}(\mathbf{X}) = \mathbf{U}\mathbf{V}^\top,$$

where $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ is a singular value decomposition of \mathbf{X} .

1. Using the first and second-order optimality conditions, show that every local minimizer \mathbf{Q} of (2.5.9) satisfies

$$\begin{aligned} (\mathbf{Q}^\top \mathbf{X})^\top &= \mathbf{Q}^\top \mathbf{X}, \\ \mathbf{Q}^\top \mathbf{X} &\succeq \mathbf{0}. \end{aligned}$$

(Hint: use linearity of the Kronecker product in either of its two arguments when the other is fixed.)

2. Using these conditions, argue that at every local minimizer \mathbf{Q} of (2.5.9), one has $\mathbf{Q}^\top \mathbf{X} = (\mathbf{X}^\top \mathbf{X})^{1/2}$. (Hint: Use the fact from linear algebra that if $\mathbf{S} \succeq \mathbf{0}$ is a symmetric positive semidefinite matrix, then $(\mathbf{S}^\top \mathbf{S})^{1/2} = \mathbf{S}$.)
3. Using the singular value decomposition $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$, conclude that

$$\mathbf{U}\mathbf{V}^\top = \operatorname{proj}_{\mathbf{O}(d)}(\mathbf{X}).$$

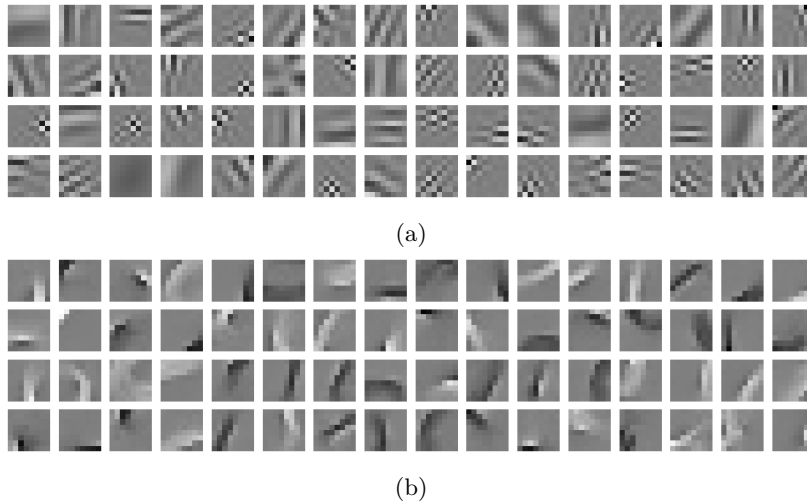


Figure 2.6: Comparison of learned dictionary atoms for complete (orthogonal) and overcomplete dictionaries, trained to reconstruct 8 by 8 patches taken from MNIST digits. Both dictionaries are trained for 6000 epochs on 10^4 random patches with nontrivial content, and sparse codes are computed with the LASSO objective and $\lambda = 0.1$ (see (2.3.5)). Colormaps have black for negative values, and white for positive values. **Top:** An orthogonal dictionary learned with the MSP algorithm (2.2.18) is constrained to have no more than 64 atoms; the learned atoms roughly correspond to a “spike and slab” dictionary, and achieve relatively poor reconstruction sparsity levels on held-out test data (codes are approximately 17-sparse on average, with respect to a threshold of 10^{-1}). **Bottom:** In contrast, an overcomplete dictionary (here, with 8^3 atoms; we visualize a random subset of 64) learns semantically-meaningful dictionary atoms corresponding to signed oriented edges, which can be pieced together to create digit patches and achieve superior reconstruction and sparsity levels. Codes are approximately 20-sparse on average, while being 8 times larger than those of the orthogonal dictionary. To compute the dictionary, we use an optimizer based on proximal alternating linearized minimization on a suitably-regularized version of (2.3.17).

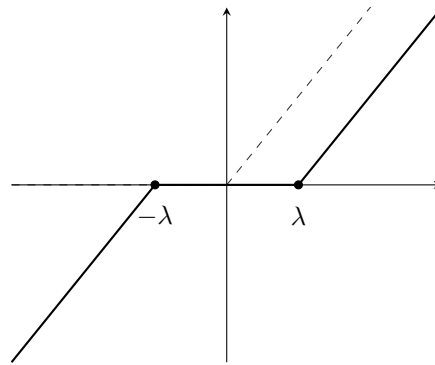


Figure 2.7: The graph of the soft thresholding function S_λ (solid) compared to the ReLU activation $x \mapsto \max\{x, 0\}$ (dashed). Between $-\lambda$ and $+\lambda$, the value of S_λ is clipped to 0, while large inputs' (absolute) value is simply shrunk by λ . Unlike ReLU, S_λ is symmetric and acts on both positive and negative inputs.