

Chapter 2

Sparse Signal Models

This book is about modeling and exploiting *simple structure* in signals, images, and data. In this chapter, we take our first steps in this direction. We study a class of models known as *sparse models*, in which the signal of interest is a superposition of a few basic signals (called “atoms”) selected from a large “dictionary.” This basic model arises in a surprisingly large number of applications. It also illustrates fundamental tradeoffs in modeling and computation that will recur throughout the book.

2.1 Applications of Sparse Signal Modeling

Why do we need signal models at all? We give a pragmatic answer. Many problems arising in modern signal processing and data analysis are intrinsically *ill-posed*. Often, the number of unknowns vastly exceeds the number of observations. In this situation, prior knowledge is absolutely essential to solving the problem correctly.

To describe this phenomenon mathematically, consider the simple equation

$$\underset{\text{Observation}}{\mathbf{y}} = \mathbf{A} \underset{\text{Unknown}}{\mathbf{x}}. \tag{2.1.1}$$

Here, $\mathbf{y} \in \mathbb{R}^m$ is our observation, while $\mathbf{x} \in \mathbb{R}^n$ is unknown. The matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ represents the data generation process: the observed data \mathbf{y} is a linear function of the unknown (or hidden) signal \mathbf{x} . This is a simple

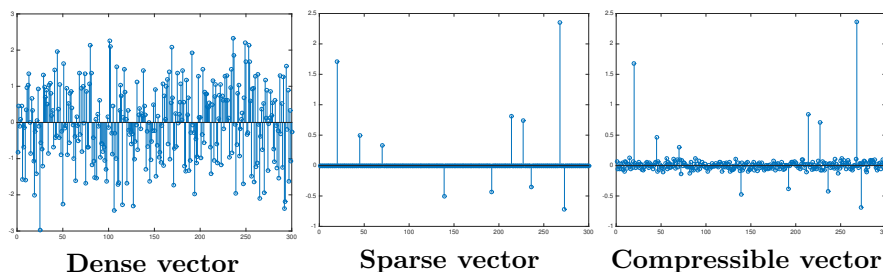


Figure 2.1. **Dense vs. sparse vectors.** Left: a *generic* vector $\mathbf{x} \in \mathbb{R}^n$, with entries independent standard normal random variables. Center: a *sparse* vector, with only a few nonzero entries. Right: a *compressible* vector, with only a few significant entries.

models; however, we will see that it is rich enough to bear on a vast array of practical applications.

Recovering the unknown \mathbf{x} from observation \mathbf{y} may appear trivial: we simply have to solve a linear system of equations! However, many practical applications raise a substantial challenge: the number of observations, m , can be significantly smaller than the number of elements n in the signal to be recovered. From linear algebra,¹ we know that when $m < n$, the system of equations $\mathbf{y} = \mathbf{A}\mathbf{x}$ does not necessarily have any solution, but if it has any solution at all, then the solution space has dimension $n - m$. Hence, either there is no solution, or there are infinitely many solutions. Only one of them is the one we wish to recover! To make progress, we need to leverage some additional properties of the target solution.

Sparsity is one such property, which has strong implications on our ability to solve underdetermined systems. A vector $\mathbf{x} \in \mathbb{R}^n$ is considered *sparse* if only a few of its elements are nonzero. Figure 2.1 (center) shows an example of such a vector. Some form of sparsity arises naturally in almost every type of high-dimensional signal or data that we encounter in practical applications. Below, we illustrate with a few representative examples.

2.1.1 An Example from Medical Imaging

Figure 2.2 shows a *magnetic resonance* (MR) image of the brain. This is a digital image $I \in \mathbb{R}^{N \times N}$. Each entry $I(\mathbf{v})$ (here, $\mathbf{v} \in \mathbb{R}^2$) corresponds to the density of protons at a given spatial location inside the brain. This essentially indicates where water is in the brain, and can reveal many biological structures that are important for disease monitoring and diagnosis.

¹Appendix A provides a detailed review of linear algebra and matrix analysis. In particular, Appendix A.6 reviews the existence and uniqueness of solutions to linear systems, which we use here to motivate our study of sparse approximation.

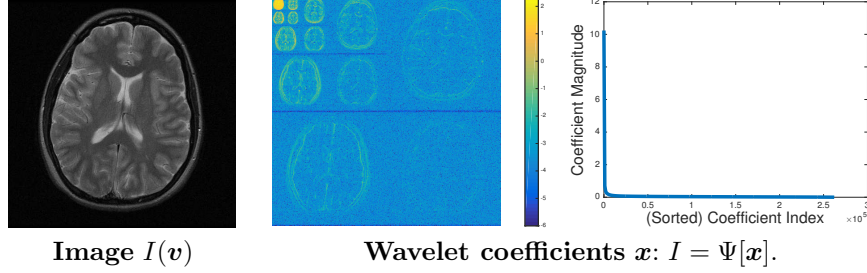


Figure 2.2. **A Magnetic Resonance image.** Left: target image of a human brain. Right: coefficients in the wavelet decomposition $I = \sum_i \psi_i x_i$, and their magnitudes, sorted in descending order. The large wavelet coefficients concentrate around sharp edges in the image; wavelet coefficients corresponding to smooth regions are much smaller. The wavelet coefficients are highly compressible: their magnitude decays rapidly. Data courtesy M. Lustig [Lustig, 2013].

To caricature the MRI problem a bit, our goal is to estimate I , without opening up the brain! This is possible, if we subject the patient to a large, spatially and temporally varying magnetic field. The magnetic field causes the protons to oscillate at a frequency that depends on position their locations and energy states. Each proton essentially acts as its own radio transmitter, and in aggregate they create a signal we can measure.

As we will see from a more detailed derivation of the physical model for MRI in Chapter 11, it turns out that the signal we observe is simply a sample of the two-dimensional Fourier transform of I :

$$y = \int_{\mathbf{v}} I(\mathbf{v}) \exp(-i 2\pi \mathbf{u}^* \mathbf{v}) d\mathbf{v}. \quad (2.1.2)$$

Here, i is the imaginary unit, and \cdot^* denotes the (complex conjugate) transpose of a vector. The two-dimensional frequency vector $\mathbf{u} = (u_1, u_2) \in \mathbb{R}^2$ depends how the magnetic field we applied varies over space. Here, letting \mathcal{F} denote the 2D Fourier transform, the above expression is

$$y = \mathcal{F}[I](\mathbf{u}). \quad (2.1.3)$$

By changing the applied magnetic field, we can vary \mathbf{u} , and collect m samples of the Fourier transform, corresponding to different applied magnetic fields, parameterized by $U = \{\mathbf{u}_1, \dots, \mathbf{u}_m\}$. We can concatenate all of our observations into a vector $\mathbf{y} \in \mathbb{C}^m$, given by

$$\mathbf{y} = \begin{bmatrix} \mathcal{F}[I](\mathbf{u}_1) \\ \vdots \\ \mathcal{F}[I](\mathbf{u}_m) \end{bmatrix} \doteq \mathcal{F}_U[I]. \quad (2.1.4)$$

Here, \mathcal{F}_U is simply the operator that obtains from I the Fourier samples indexed by U . If you imagine the Fourier transform as acting by matrix

multiplication, \mathcal{F}_U is simply the matrix we get if we discard all the rows of \mathcal{F} that are not indexed by U .

One very basic property of the integral (2.1.2), and hence of the operator \mathcal{F}_U , is that it is *linear* in its input I . This means that for any pair of inputs I and J and complex scalars α, β ,

$$\mathcal{F}_U[\alpha I + \beta J] = \alpha \mathcal{F}_U[I] + \beta \mathcal{F}_U[J]. \quad (2.1.5)$$

Because \mathcal{F}_U is a linear operator, the problem of finding I from \mathbf{y} using the observation equation (2.1.4) “just” consists of solving a large linear system of equations.

There is a substantial catch, though. In this system of equations, there are typically far more unknowns (here $n = N^2$) than observations m . This is necessary: it is generally too time and energy intensive to simply measure all N^2 Fourier coefficients. This is even more pressing of a concern in *dynamic MRI*, where the object being imaged is changing over time, and so acquisition needs to be time-efficient. So, in general, we need m to be as small as is just necessary to guarantee accurate reconstruction – and certainly significantly smaller than n .

This leaves us with a seemingly impossible situation: we have n unknowns and $m \ll n$ equations. Unless we can make some additional assumptions on the structure of I , the problem is ill-posed. Fortunately, real signals are not completely unstructured.² Figure 2.2 (right) shows a *wavelet transform* of I . The wavelet transform expresses I as a superposition of a collection of basis functions $\Psi = \{\psi_1, \dots, \psi_{N^2}\}$:

$$\underset{\text{image}}{I} = \sum_{i=1}^{N^2} \underset{\text{i-th basis signal}}{\psi_i} \times \underset{\text{i-th coefficient}}{x_i}. \quad (2.1.6)$$

Here, $x_1, \dots, x_{N^2} \in \mathbb{R}$ are coefficients of the image I with respect to the basis Ψ . The entries in Figure 2.2 (right) are the magnitudes $|x_i|$ for the N^2 coefficients x_i . The important point is that many of these coefficients are extremely small. If let $J = \{i_1, \dots, i_k\}$ denote the k largest coefficients, we can approximate I as

$$\underset{\text{Target image}}{I} \approx \underset{\text{Superposition of } k \text{ basis functions}}{\tilde{I}_k} = \sum_{i \in J} \psi_i x_i. \quad (2.1.7)$$

Figure 2.3 visualizes the reconstruction and reconstruction error $I - \tilde{I}_k$. It seems that even if we retain only a relatively small fraction of the coefficients, we still obtain an accurate approximation, and most of what remains

²Indeed, we can construct a “generic” element I_{generic} of $\mathbb{R}^{N \times N}$, by choosing its entries at random – say from a standard Gaussian distribution $\mathcal{N}(0, 1)$. With very high probability, I_{generic} will simply look like noise. The target magnetic resonance image in Figure 2.2 certainly does not look like noise!

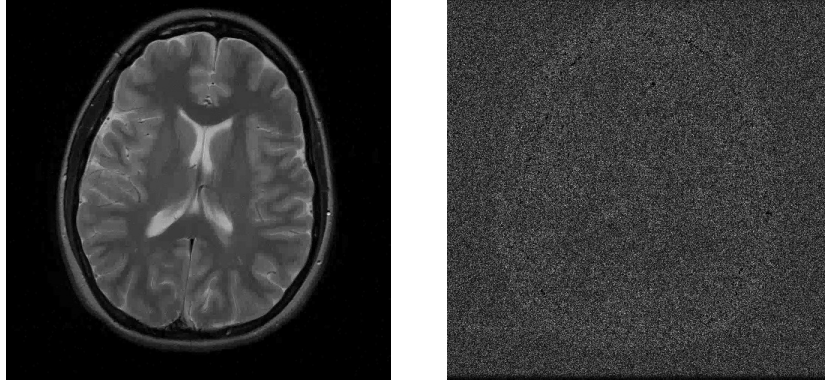


Figure 2.3. **Wavelet approximation \hat{I} to I , and approximation error.** Left: approximation to the image in Figure 2.2 using the most significant 7% of the wavelet coefficients. Right: approximation error $|I - \hat{I}|$. The error contains mostly noise, suggesting that most of the important structure of the image is captured in the wavelet approximation \hat{I} .

is noise. This suggests that the sequence \mathbf{x} is *compressible* – it is very close to a sparse vector.

In order to recover I , we can first try to reconstruct the sparse vector \mathbf{x} , using the observation equation

$$\begin{aligned}
 \mathbf{y} &= \mathcal{F}_U[I], \\
 \text{Observed Fourier coefficients} & \\
 &= \mathcal{F}_U[\psi_1 x_1 + \cdots + \psi_{N^2} x_{N^2}], \\
 &= \mathcal{F}_U[\psi_1] x_1 + \cdots + \mathcal{F}_U[\psi_{N^2}] x_{N^2}, \\
 &= [\mathcal{F}_U[\psi_1] \mid \cdots \mid \mathcal{F}_U[\psi_{N^2}]] \mathbf{x}, \\
 &\quad \text{Matrix } \mathbf{A} \in \mathbb{R}^{m \times N^2}, m \ll N^2. \\
 &= \mathbf{A} \mathbf{x}.
 \end{aligned} \tag{2.1.8}$$

After these manipulations, we end up with a system of equations $\mathbf{y} = \mathbf{A} \mathbf{x}$. The vector \mathbf{x} contains the coefficients of the target image I in the wavelet basis. The i -th column of the matrix \mathbf{A} contains a subset U of the Fourier coefficients of the i -th basis signal ψ_i . To reconstruct I , we can look for a solution $\hat{\mathbf{x}}$ to this system, and then set

$$\hat{I} = \sum_{i=1}^{N^2} \psi_i \hat{x}_i. \tag{2.1.9}$$

Because \mathbf{x} has N^2 entries, but we only have $m \ll N^2$ observations, the system $\mathbf{y} = \mathbf{A} \mathbf{x}$ is underdetermined. Nevertheless, because the wavelet coefficients of I are (nearly) sparse – say only its k largest coefficients are

significant and others are negligible, the desired solution \mathbf{x} to this system is sparse. To reconstruct I we need to find a sparse solution to an under-determined system! In Chapter 11, we will illustrate how to actually apply such a “compressive sampling” scheme to real MRI images under more realistic conditions.

2.1.2 An Example from Image Processing

In the previous example, we used the fact that the image I had a good sparse approximation in terms of a “dictionary” of basic elements $\psi_1, \dots, \psi_{N^2}$:

$$I \approx \sum_{i \in J} \psi_i x_i = \underbrace{\Psi}_{N^2 \times N^2 \text{ matrix}} \underbrace{\mathbf{x}}_{\text{sparse vector}}, \quad (2.1.10)$$

where $x_i = 0$ for $i \notin J$, and $k = |J| \ll N^2$. Expressions of this form play a central role in lossy data compression. Image compression standards such as JPEG [Wallace, 1991] and JPEG 2000 [Taubman and Marcellin, 2001] leverage sparse approximations (in the discrete cosine transform (DCT) [Ahmed et al., 1974] and wavelet bases, respectively). Generally speaking, the sparser the representation is, the more input image can be compressed. However, sparse representations of images are not *just* useful for compression: they can be used for solving inverse problems, in which we try to reconstruct I from noisy, corrupted or incomplete observations. We already saw an example of this in the previous section, in which we used sparsity in the wavelet domain to reconstruct MR images. To facilitate all of these tasks, we can seek representations of I that are as sparse as possible, by replacing Ψ with more general dictionaries \mathbf{A} . For example, we might consider *overcomplete dictionaries* $\mathbf{A} \in \mathbb{R}^{m \times n}$, $n > m$, which consist of several orthobases (e.g., DCT and wavelets together). The idea is that each individual representation may capture a particular type of signal well – say, DCT for smooth variations and wavelets for signals with sharp edges. Together, they can represent a broader class of signals.

An even more aggressive idea is to simply *learn* \mathbf{A} from data, rather than designing it by hand. This approach tends to produce better sparsity-accuracy tradeoffs for representing images I , and is also useful for a wealth of other problems, including denoising, inpainting, and superresolution that involve reconstructing I from incomplete or corrupted observations. Each of these problems leads to an underdetermined linear system of equations; the goal is to use the prior knowledge that the target signal I has a compact representation in some dictionary \mathbf{A} to make the problem well-posed. Figure 2.4 shows an example of this for the problem of color image denoising, from [Mairal et al., 2008a]. We observe a noisy image

$$I_{\text{noisy}} = \underbrace{I_{\text{clean}}}_{\text{target image}} + \underbrace{\mathbf{z}}_{\text{noise}}. \quad (2.1.11)$$

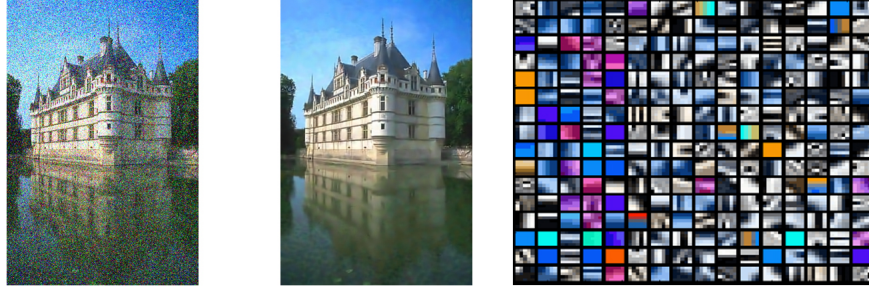


Figure 2.4. **Image denoising by sparse approximation.** Left: A noisy input image. The image is broken into patches $\mathbf{y}_1, \dots, \mathbf{y}_p$. A dictionary $\mathbf{A} = [\mathbf{a}_1 \mid \dots \mid \mathbf{a}_n]$ is learned such that each input patch can be approximated as $\mathbf{y}_i \approx \mathbf{A}\mathbf{x}_i$, with \mathbf{x}_i sparse. Right: dictionary patches $\mathbf{a}_1, \dots, \mathbf{a}_n$. Center: denoised image, reconstructed from the approximations $\hat{\mathbf{y}}_i = \mathbf{A}\mathbf{x}_i$. Results of Mairal, Saprio and Elad [Mairal et al., 2008a]. Figure from [Wright et al., 2010].

We assume³ that patches of the clean image have an accurate sparse approximation in some dictionary \mathbf{A} : if we break I_{clean} into patches $\mathbf{y}_{1\text{clean}}, \dots, \mathbf{y}_{p\text{clean}}$

$$\mathbf{y}_{i\text{clean}} \approx \underset{\text{patch dictionary}}{\mathbf{A}} \times \underset{\text{sparse coefficient vector}}{\mathbf{x}_i}. \quad (2.1.12)$$

In denoising, we do not actually observe $\mathbf{y}_{i\text{clean}}$. Rather, we observe *noisy* patches $\mathbf{y}_i = \mathbf{y}_{i\text{clean}} + \mathbf{z}_i$. Based on these patches $\mathbf{y}_1, \dots, \mathbf{y}_p$, we learn⁴ a dictionary $\hat{\mathbf{A}}$ such that

$$\underset{i\text{-th image patch}}{\mathbf{y}_i} \approx \underset{\text{learned dictionary}}{\hat{\mathbf{A}}} \times \underset{\text{sparse coefficient vector}}{\hat{\mathbf{x}}_i} = \underset{\text{denoised patch}}{\hat{\mathbf{y}}_i}.$$

We take $\hat{\mathbf{y}}_i = \hat{\mathbf{A}}\hat{\mathbf{x}}_i$ as an estimate of $\mathbf{y}_{i\text{clean}}$.

Figure 2.4 (left) shows the noisy input image; Figure 2.4 (center) shows a denoised image constructed from $\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_p$. Figure 2.4 (right) shows the dictionary $\hat{\mathbf{A}}$ learned from the noisy patches. Although the sparse dictionary prior is relatively simple, and does not capture all of the global geometric structure of the image, it leads to surprisingly good performance on many low-level image processing tasks. We discuss modeling and computational aspects of this area in detail in Chapter ???. For now, the key point is that the problem of reconstructing the clean image from noisy

³Of course, this assumption needs to be justified! See Exercise 2.14 and the notes and references to this chapter.

⁴We learn the dictionary $\hat{\mathbf{A}}$ by solving an optimization problem, which attempts to strike an optimal balance between sparsity of the coefficients $\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_p$ and the accuracy of the approximation $\mathbf{y}_i \approx \hat{\mathbf{A}}\hat{\mathbf{x}}_i$. We will explain this in more detail in Chapter ??.

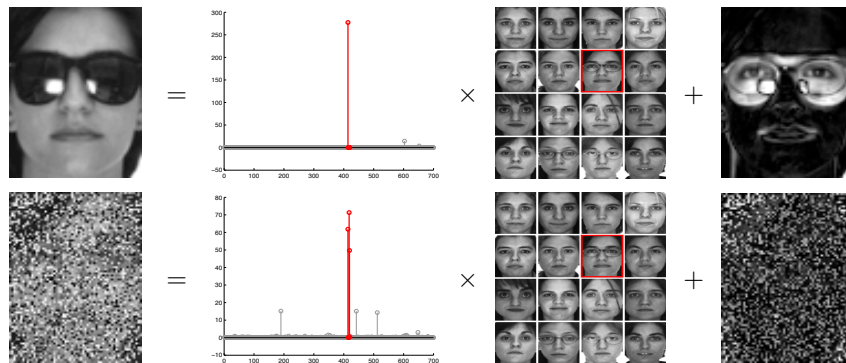


Figure 2.5. **Face recognition via sparse representation.** Top: input face image \mathbf{y} is wearing a sunglasses; Bottom: input face image \mathbf{y} is with 50% pixels arbitrarily corrupted. Each test image \mathbf{y} is approximated as a sparse combination $\mathbf{B}\mathbf{x}$ of the training images, plus a sparse error \mathbf{e} due to occlusion. In this example, red coefficients correspond to images of the correct subject.

patches again leads us to an underdetermined linear systems of equations, $\mathbf{y}_i \approx \mathbf{A}\mathbf{x}_i$.

2.1.3 An Example from Face Recognition

Sparsity also arises naturally in problems in which we wish to perform reliable inference from unreliable measurements. For example, due to sensor errors or malicious tampering, a vector-valued observation $\mathbf{y} \in \mathbb{R}^m$ might be grossly corrupted in a few of its entries:

$$\underset{\text{observation}}{\mathbf{y}} = \underset{\text{clean data}}{\mathbf{y}_o} + \underset{\text{sparse error}}{\mathbf{e}}. \quad (2.1.13)$$

We illustrate this more concretely using an example from automatic face recognition. Imagine that we have a database consisting of a number of subjects. For each subject i , we collect grayscale training images $I_{i,1}, \dots, I_{i,n_i} \in \mathbb{R}^{W \times H}$, and vectorize them to form a base matrix $\mathbf{B}_i \in \mathbb{R}^{m \times n_i}$, with $m = W \times H$. We can further concatenate these matrices to form a large training “dictionary”

$$\mathbf{B} = [\underset{\text{all training images}}{\mathbf{B}_1} \mid \mathbf{B}_2 \mid \dots \mid \mathbf{B}_n] \in \mathbb{R}^{m \times n}, \quad n = \sum_i n_i. \quad (2.1.14)$$

Suppose our system is confronted with a new image $\mathbf{y} \in \mathbb{R}^m$, taken under some new lighting condition, and possibly occluded – see Figure 2.5. For now, we can assume that the input \mathbf{y} is well-aligned to the training images (i.e., the faces occur at the same position in the training and test

images).⁵ There is a beautiful physical argument [Basri and Jacobs, 2003a] that shows that in an average case sense, images of “nice” objects taken under varying lighting conditions lie very close to low-dimensional linear subspaces of the high-dimensional image space \mathbb{R}^m . This suggests that if we have seen enough training examples, we can approximate the input sample \mathbf{y} as a linear combination of the training samples from the same class:

$$\begin{array}{ccc} \mathbf{y} & \approx & \mathbf{B}_{i_*} \mathbf{x}_{i_*} \\ \text{observed image} & & \text{linear combination of training images from } i_*\text{-th class} \end{array} \quad (2.1.15)$$

Unfortunately, in practice, this equation is violated in at least two ways: first, we don’t know the true identity i_* ahead of time. Second, nuisance factors such as occlusion cause the equation to be badly violated on a portion of the image pixels (those that are occluded). For the first problem, we note that we can *still* write down an expression for \mathbf{y} as a linear combination of elements of the database \mathbf{B} as a whole: $\mathbf{y} \approx \mathbf{B}\mathbf{x}$. To deal with occlusion, we need to introduce an additional term \mathbf{e} , giving

$$\mathbf{y} \approx \mathbf{B}\mathbf{x} + \mathbf{e}. \quad (2.1.16)$$

Because the errors caused by occlusion are large in magnitude, this error \mathbf{e} cannot simply be ignored, or treated with techniques designed for small noise. Unfortunately, this means that the system is underdetermined: we have m equations, but $m + n$ unknowns $\bar{\mathbf{x}} = (\mathbf{x}, \mathbf{e})$. Writing $\mathbf{A} = [\mathbf{B} \mid \mathbf{I}]$, we again have a very large underdetermined system

$$\mathbf{y} = \mathbf{A}\bar{\mathbf{x}}. \quad (2.1.17)$$

If we did not have prior information about $\bar{\mathbf{x}}$, there would be no hope of recovering it from this observation. Fortunately, both \mathbf{x} and \mathbf{e} are very structured. The nonzero values of \mathbf{x} should be concentrated only on those images of the true subject, i_* , and so it should be a *sparse vector*. The nonzero values of the error \mathbf{e} should be concentrated only on those pixels that are occluded or corrupted, and so it should also be sparse.⁶

Figure 2.5 shows two examples of a sparse solution to this system of equations for a given input image \mathbf{y} . Notice that the coefficients in the estimated $\hat{\mathbf{x}}$ are concentrated on images of the correct subject (red) and that the error indeed corresponds to the physical occlusion. The setting

⁵Relaxing this assumption is essential to building systems that work with unconstrained input images. We will talk about how to relax this assumption in Chapter 10.

⁶Of course, the goal is to correct as many errors as possible. One of the surprises of high dimensions is it is indeed possible to correct large fractions of errors using simple, efficient algorithms. Understanding precisely how many errors we can correct (and how dense the vector $\bar{\mathbf{x}}$ can be before our methods break down) will be a major theoretical thrust of this book. In Chapter 10, we will give a more precise characterization about how large a fraction of errors can be corrected for a system of linear equations similar to those arise in the robust face recognition setting.

we have described so far is somewhat idealized – we will discuss both the modeling and system building aspects of this problem in more detail later in the book. For our purposes here, it is enough to note that *if* we can somehow obtain a sparse (\mathbf{x}, \mathbf{e}) , it should suffice to identify the subject, despite nuisances such as illumination, occlusion, and corruption.

2.2 Recovering a Sparse Solution

Suppose, as in the above examples, that we know the ground truth signal \mathbf{x}_o is sparse. How powerful is this knowledge? Can it render ill-posed problems such as MR image acquisition or occluded face recognition well-posed? To answer these questions, we need a formal notion of sparsity. In the next two subsections, we begin by introducing the concept of a norm of a vector, which generalizes the concept of *length*. We then introduce an “ ℓ^0 norm”, which counts the number of nonzero entries in a vector, a basic measure of how dense (or sparse) that vector is.

2.2.1 Norms on Vector Spaces

A *vector space* \mathbb{V} consists of a collection of elements (vectors), field such as the real numbers \mathbb{R} or complex numbers \mathbb{C} (scalars) and operations (adding vectors and multiplying vectors with scalars) that work in ways that conform to our intuitions from \mathbb{R}^3 . Appendix A reviews the formal definition of a vector space, and gives examples. In the above application examples, our signals of interest consisted of collections of real or complex numbers – e.g., in MR imaging, the target image I was an element of $\mathbb{R}^{N \times N}$. We can view $\mathbb{R}^{N \times N}$ as a vector space, with scalar field \mathbb{R} (write $\mathbb{V} = (\mathbb{R}^{N \times N}, \mathbb{R})$). In the other examples as well, the signals of interest reside in vector spaces.

A *norm* on a vector space \mathbb{V} gives a way of measuring lengths of vectors, that conforms in important ways to our intuition from lengths in \mathbb{R}^3 . Formally:

Definition 2.2.1. A **norm** on a vector space \mathbb{V} over \mathbb{R} is a function $\|\cdot\| : \mathbb{V} \rightarrow \mathbb{R}$ that is

1. *Nonnegatively homogeneous:* $\|\alpha \mathbf{x}\| = |\alpha| \|\mathbf{x}\|$ for all vectors $\mathbf{x} \in \mathbb{V}$, scalars $\alpha \in \mathbb{R}$,
2. *Positive definite:* $\|\mathbf{x}\| \geq 0$, and $\|\mathbf{x}\| = 0$ if and only if $\mathbf{x} = \mathbf{0}$,
3. *Subadditive:* $\|\cdot\|$ satisfies the triangle inequality $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{V}$.

For our purposes, the most important family of norms are the ℓ^p norms. We will use norms from this family to derive practical algorithms for find-

ing sparse solutions to linear systems of equations, and for studying their properties. If we take $\mathbb{V} = (\mathbb{R}^n, \mathbb{R})$, and $p \in (0, \infty)$, we can write

$$\|\mathbf{x}\|_p = \left(\sum_i |x_i|^p \right)^{1/p}. \quad (2.2.1)$$

The function $\|\mathbf{x}\|_p$ is a norm for any $p \geq 1$. The most familiar example is the ℓ^2 norm or “Euclidean norm”

$$\|\mathbf{x}\|_2 = \sqrt{\sum_i |x_i|^2} = \sqrt{\mathbf{x}^* \mathbf{x}}$$

which coincides with our usual way of measuring lengths. Two other cases are of almost equal importance: $p = 1$, and $p \rightarrow \infty$. Setting $p = 1$ in (A.9.1), we obtain

$$\|\mathbf{x}\|_1 = \sum_i |x_i|, \quad (2.2.2)$$

which will play a very large role in this book.⁷ Finally, as p becomes larger, the expression in (A.9.1) accentuates large $|x_i|$. As $p \rightarrow \infty$, $\|\mathbf{x}\|_p \rightarrow \max_i |x_i|$. We extend the definition of the ℓ^p norm to $p = \infty$ by defining

$$\|\mathbf{x}\|_\infty = \max_i |x_i|. \quad (2.2.3)$$

To appreciate the distinction between the various ℓ^p norms, we can visualize their unit balls $\mathbf{B}_{\|\cdot\|_p}$, which consist of all vectors \mathbf{x} whose norm is at most one:

$$\mathbf{B}_{\|\cdot\|_p} = \left\{ \mathbf{x} \mid \|\mathbf{x}\|_p \leq 1 \right\}. \quad (2.2.4)$$

The ℓ^2 ball is a (solid) sphere, the ℓ^∞ ball is a cube, and the ℓ^1 ball a kind of diamond shape – see Figure 2.6.⁸

Notice that for $p \leq p'$, $\mathbf{B}_{\|\cdot\|_p} \subseteq \mathbf{B}_{\|\cdot\|_{p'}}$. This is because when $p \leq p'$, $\|\mathbf{x}\|_p \geq \|\mathbf{x}\|_{p'}$ for all \mathbf{x} .

Remark 2.2.2. *This containment becomes even more striking in higher dimensions: in \mathbb{R}^n , $\text{vol}(\mathbf{B}_\infty) = 2^n$, while $\text{vol}(\mathbf{B}_1) = 2^n/n!$ (see e.g., [Maltousek, 2002a]). So, in $n = 2$ dimensions $\text{vol}(\mathbf{B}_1) = (1/2) \times \text{vol}(\mathbf{B}_\infty)$, while in $n = 1,000$ dimensions $\text{vol}(\mathbf{B}_1) \approx 10^{-2,568} \times \text{vol}(\mathbf{B}_\infty)$ – a truly negligible fraction!*

⁷Anyone who has traveled in Manhattan should have good appreciation for the distinction between ℓ^1 and ℓ^2 – in fact, the ℓ^1 norm is sometimes called the Manhattan norm! This example illustrates a simple, but important point – the proper choice of norm depends quite a bit on the properties of the problem and design goals. Unless you can leap tall buildings in a single bound, measuring distance using the ℓ^2 norm would underestimate how much travel you need to reach your destination.

⁸To see this in action, you can run `Chapter_2_Illustrate_Lp_Balls.m`.

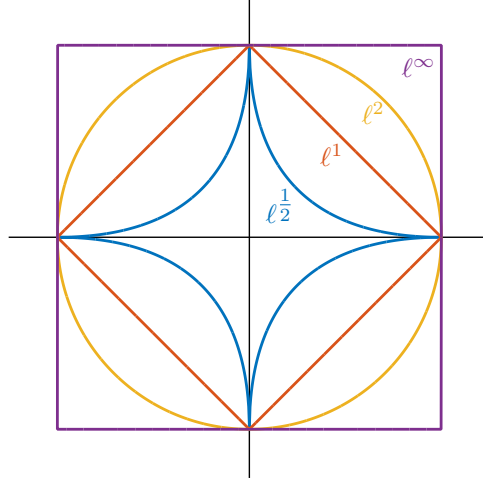


Figure 2.6. The ℓ^p balls $\mathbf{B}_p = \{\mathbf{x} \mid \|\mathbf{x}\|_p \leq 1\}$ for $0 < p \leq \infty$. For $p \geq 1$, \mathbf{B}_p is a convex set, and $\|\cdot\|_p$ is a norm. For $p < 1$, $\|\cdot\|_p$ is not a norm, in the formal sense.

Remark 2.2.3. This may seem to be in contrast to the mathematical fact that “in finite dimensions, all norms are equivalent” (see, e.g., Appendix A). Formally, this statement means that in a finite dimensional vector space \mathbb{V} , such as \mathbb{R}^n , for any pair of norms $\|\cdot\|_\diamond$ and $\|\cdot\|_\square$ there exist numbers $0 < \alpha, \beta < \infty$ such that for every $\mathbf{x} \in \mathbb{V}$,

$$\alpha \|\mathbf{x}\|_\square \leq \|\mathbf{x}\|_\diamond \leq \beta \|\mathbf{x}\|_\square. \quad (2.2.5)$$

So, the norms $\|\cdot\|_\square$ and $\|\cdot\|_\diamond$ can be compared in size. However, as the example in Remark 2.2.2 shows, in high dimensions, the unit balls of the various ℓ^p norms can be very different. In applications involving high-dimensional signals, the different choices of norm can lead to radically different solutions.

2.2.2 The ℓ^0 Norm

With the notion of a norm in hand, we are prepared to define a formal notion of sparsity. For this, we introduce a function, called the “ ℓ^0 norm” (read “ell zero norm”), which is simply the number of nonzero entries in a vector \mathbf{x} :

$$\|\mathbf{x}\|_0 = \#\{i \mid \mathbf{x}(i) \neq 0\}. \quad (2.2.6)$$

Loosely speaking, \mathbf{x} is sparse whenever $\|\mathbf{x}\|_0$ is small.

The ℓ^0 norm $\|\cdot\|_0$ is *not* a norm, in the formal sense of Definition A.9.1: since for $\alpha \neq 0$, $\|\alpha\mathbf{x}\|_0 = \|\mathbf{x}\|_0$, it does not have the property of nonnegative

homogeneity. It *does* have the other two properties, however. In particular, $\|\cdot\|_0$ is subadditive:

$$\forall \mathbf{x}, \mathbf{x}', \quad \|\mathbf{x} + \mathbf{x}'\|_0 \leq \|\mathbf{x}\|_0 + \|\mathbf{x}'\|_0. \quad (2.2.7)$$

This is easily checked by noting that the set of nonzero entries for $\mathbf{x} + \mathbf{x}'$ is contained union of the set of nonzero entries of \mathbf{x} and the set of nonzero entries of \mathbf{x}' .

Although the ℓ^0 norm is not a norm in the mathematical sense, it is closely related to the ℓ^p norms. To understand this relationship, note that for every $\mathbf{x} \in \mathbb{R}^n$,

$$\begin{aligned} \lim_{p \searrow 0} \|\mathbf{x}\|_p^p &= \sum_{i=1}^n \lim_{p \searrow 0} |\mathbf{x}(i)|^p \\ &= \sum_{i=1}^n \mathbb{1}_{\mathbf{x}(i) \neq 0} \\ &= \|\mathbf{x}\|_0. \end{aligned} \quad (2.2.8)$$

In this sense, the ℓ^0 norm can be considered to be generated from the ℓ^p norms, by taking p (infinitesimally) small. In the context of Figure 2.6, this can be understood as follows: in \mathbb{R}^2 , the sparse vectors correspond to the coordinate axes. As p drops towards zero, the unit ball of the ℓ^p norm becomes more concentrated around the coordinate axes, i.e., around the sparse vectors.

The geometric relationship between the ℓ^0 and ℓ^p norms is useful for deriving algorithms, and for understanding why small p tends to favor sparse solutions. With this said, the formal notation $\|\mathbf{x}\|_0$ has a very simple meaning: *it counts the number of nonzero entries in \mathbf{x}* . In all of the applications discussed above, our goal is to recover a vector \mathbf{x}_{true} with $\|\mathbf{x}_{\text{true}}\|_0$ small. In this book, we often use \mathbf{x}_o as a shorthand for $\mathbf{x}_{\text{true}} = \mathbf{x}_o$.

2.2.3 The Sparsest Solution: Minimizing the ℓ^0 Norm

Suppose we observe $\mathbf{y} \in \mathbb{R}^m$, with $\mathbf{y} = \mathbf{A}\mathbf{x}_o$, and that our goal is to recover \mathbf{x}_o . If we know that \mathbf{x}_o is sparse, it seems reasonable to form an estimate $\hat{\mathbf{x}}$ by choosing the *sparsest* vector \mathbf{x} that satisfies the equation $\mathbf{y} = \mathbf{A}\mathbf{x}$. That is, we choose the sparsest \mathbf{x} that could have generated our observation. We can write this as an optimization problem

$$\begin{aligned} &\text{minimize} && \|\mathbf{x}\|_0 \\ &\text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{y}. \end{aligned} \quad (2.2.9)$$

How might we solve this problem numerically? Call

$$\text{supp}(\mathbf{x}) = \{i \mid \mathbf{x}_i \neq 0\} \subset \{1 \dots n\} \quad (2.2.10)$$

the *support* of the vector \mathbf{x} – this set contains the indices of the nonzero entries. The ℓ^0 minimization problem (2.2.9) asks us to find a vector \mathbf{x}

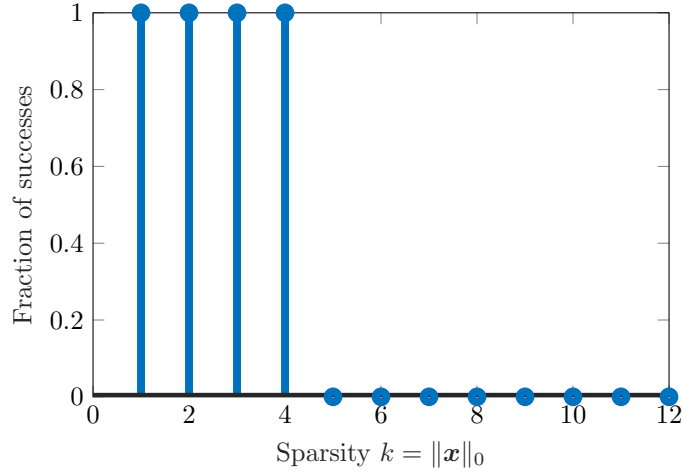


Figure 2.7. **Transitions in ℓ^0 recovery.** Fraction of correct recoveries across 100 trials, as a function of sparsity of the target solution \mathbf{x}_o . The system is of size 5×12 . In this experiment, ℓ^0 minimization successfully recovers all \mathbf{x}_o with $k \leq 4$ nonzeros.

of smallest support that agrees with the observation \mathbf{y} . One approach to finding such an \mathbf{x} is to simply try every possible support set $I \subseteq \{1, \dots, n\}$. For each such set I , we can form a system of equations

$$\mathbf{A}_I \mathbf{x}_I = \mathbf{y}, \quad (2.2.11)$$

where $\mathbf{A}_I \in \mathbb{R}^{m \times |I|}$ is the column submatrix of \mathbf{A} formed by keeping only those columns indexed by I , and $\mathbf{x}_I \in \mathbb{R}^{|I|}$. We can attempt to solve (2.2.11) for \mathbf{x}_I . If such an \mathbf{x}_I exists, we can obtain a solution \mathbf{x} to $\mathbf{A}\mathbf{x} = \mathbf{y}$ by filling in the remaining entries of \mathbf{x} with zeros. This exhaustive search procedure is spelled out formally as Algorithm 2.1.

Algorithm 2.1: ℓ^0 -minimization by exhaustive search

- 1: **Input:** a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and a vector $\mathbf{y} \in \mathbb{R}^m$.
 - 2: **for** $k = 0, 1, 2, \dots, n$,
 - 3: **for** each $I \subseteq \{1, \dots, n\}$ of size k ,
 - 4: **if** the system of equations $\mathbf{A}_I \mathbf{z} = \mathbf{y}$ has a solution \mathbf{z} ,
 - 5: set $\mathbf{x}(I) = \mathbf{z}$, $\mathbf{x}(I^c) = \mathbf{0}$.
 - 6: **return** \mathbf{x} .
 - 7: **end if**
 - 8: **end for**
 - 9: **end for**
-

Example 2.2.4. *Let us examine how the algorithm behaves numerically, using `Chapter_2_L0_recovery.m` and `Chapter_2_L0_transition.m`. These examples generate random underdetermined linear systems $\mathbf{y} = \mathbf{A}\mathbf{x}$, with $\mathbf{y} = \mathbf{A}\mathbf{x}_o$, and \mathbf{x}_o sparse. They then apply Algorithm 2.1 (`minimize_L0.m`) to recover a vector $\hat{\mathbf{x}}$, and ask whether $\hat{\mathbf{x}}$ is equal to \mathbf{x}_o up to machine precision. Fixing the system parameters (m, n) , varying the sparsity $k = 0, 1, \dots$, and performing many random trials, we produce Figure 2.7, which shows that as long as the k is not too large, the algorithm almost always succeeds.*

Is there any mathematical explanation for this phenomenon? To understand why ℓ^0 minimization succeeds, it is worth first thinking about when it would fail. Suppose that there is a k -sparse vector $\mathbf{x}_o \in \text{null}(\mathbf{A})$. Then

$$\mathbf{A}\mathbf{x}_o = \mathbf{0} = \mathbf{A}\mathbf{0}. \quad (2.2.12)$$

Hence, for this \mathbf{x}_o , when $\mathbf{y} = \mathbf{A}\mathbf{x}_o$, the ℓ^0 minimizer is simply $\hat{\mathbf{x}} = \mathbf{0}$, and \mathbf{x}_o is not recovered. Put simply: if the null space of \mathbf{A} contains sparse vectors (aside from $\mathbf{0}$), ℓ^0 minimization may fail to recover the desired sparse vector \mathbf{x}_o .

In fact, the converse statement is also true: when the null space of \mathbf{A} *does not* contain sparse vectors (aside from $\mathbf{0}$), ℓ^0 minimization *does* recover any sufficiently sparse vector \mathbf{x}_o . To state the argument simply, let us suppose that $\|\mathbf{x}_o\|_0 \leq k$, and

$$(\star) \text{ The only } \boldsymbol{\delta} \in \text{null}(\mathbf{A}) \text{ with } \|\boldsymbol{\delta}\|_0 \leq 2k \text{ is } \boldsymbol{\delta} = \mathbf{0}.$$

Let $\hat{\mathbf{x}}$ denote the solution to the ℓ^0 minimization problem, so $\|\hat{\mathbf{x}}\|_0 \leq \|\mathbf{x}_o\|_0 \leq k$. If we define the *estimation error*

$$\boldsymbol{\delta} = \hat{\mathbf{x}} - \mathbf{x}_o, \quad (2.2.13)$$

then

$$\begin{aligned} \|\boldsymbol{\delta}\|_0 &= \|\hat{\mathbf{x}} - \mathbf{x}_o\|_0 \\ &\leq \|\hat{\mathbf{x}}\|_0 + \|\mathbf{x}_o\|_0 \\ &\leq 2k. \end{aligned} \quad (2.2.14)$$

So, $\boldsymbol{\delta}$ is a sparse vector. Moreover,

$$\begin{aligned} \mathbf{A}\boldsymbol{\delta} &= \mathbf{A}(\hat{\mathbf{x}} - \mathbf{x}_o) \\ &= \mathbf{A}\hat{\mathbf{x}} - \mathbf{A}\mathbf{x}_o \\ &= \mathbf{y} - \mathbf{y} = \mathbf{0}. \end{aligned} \quad (2.2.15)$$

So, $\boldsymbol{\delta}$ is a sparse vector *in the null space of \mathbf{A}* . Property (\star) states that the only sparse vector in $\text{null}(\mathbf{A})$ is $\mathbf{0}$. So, if (\star) holds, $\boldsymbol{\delta} = \mathbf{0}$, and so $\hat{\mathbf{x}} = \mathbf{x}_o$: ℓ^0 minimization indeed recovers \mathbf{x}_o .

Property (\star) is a property of the matrix \mathbf{A} . The above reasoning suggests a slogan: *the “good” \mathbf{A} for recovering sparse vectors \mathbf{x}_o are those \mathbf{A} that have no sparse vectors in their null space.* We can restate property (\star) more

conveniently in terms of the columns of \mathbf{A} : property (\star) holds if and only if every set of $2k$ columns of \mathbf{A} is linearly independent.

Definition 2.2.5 (Kruskal rank [Kruskal, 1977]). *The Kruskal rank of a matrix \mathbf{A} , written $\text{krank}(\mathbf{A})$ is the largest number r such that every subset of r columns of \mathbf{A} is linearly independent.*

From the above reasoning, if $\|\mathbf{x}_o\|_0$ is at most half of $\text{krank}(\mathbf{A})$, ℓ^0 minimization will recover \mathbf{x}_o :

Theorem 2.2.6 (ℓ^0 recovery). *Suppose that $\mathbf{y} = \mathbf{A}\mathbf{x}_o$, with*

$$\|\mathbf{x}_o\|_0 \leq \frac{1}{2} \text{krank}(\mathbf{A}). \quad (2.2.16)$$

Then \mathbf{x}_o is the unique optimal solution to the ℓ^0 minimization problem

$$\begin{aligned} &\text{minimize} && \|\mathbf{x}\|_0 \\ &\text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{y}. \end{aligned} \quad (2.2.17)$$

Notice that Theorem 2.2.6 agrees with the behavior in Figure 2.7.⁹ Theorem 2.2.6 predicts that as long as \mathbf{x}_o is *sufficiently sparse*, it will be recovered by ℓ^0 minimization. The level of allowable sparsity depends on the Kruskal rank of the matrix \mathbf{A} . It is not hard to see that in general,

$$0 \leq \text{krank}(\mathbf{A}) \leq \text{rank}(\mathbf{A}). \quad (2.2.18)$$

For “generic” \mathbf{A} , the Kruskal rank is quite large:

Proposition 2.2.7. *Let $\mathbf{A} \in \mathbb{R}^{m \times n}$, $n \geq m$, with \mathbf{A}_{ij} independent identically distributed $\mathcal{N}(0, 1)$ random variables. Then, with probability one, $\text{krank}(\mathbf{A}) = m$.*

Proof. Exercise 2.6 guides the interested reader through the proof. \square

The intuition is that to have $\text{krank}(\mathbf{A}) < m$, there must be some subset of m columns of \mathbf{A} which are linearly dependent, i.e., there is some subset $\mathbf{a}_{i_1}, \mathbf{a}_{i_2}, \dots, \mathbf{a}_{i_m}$ which lie on a linear subspace of dimension $m - 1$. For a Gaussian random matrix \mathbf{A} , the probability that this happens is zero. This is true of many other random matrices.¹⁰ We can interpret this as saying that under generic circumstances, knowing that the target \mathbf{x}_o is sparse turns an ill-posed problem into a well-posed one. The ℓ^0 minimization problem recovers vectors \mathbf{x}_o whose number of nonzeros is as large as $\frac{m}{2}$. This level of sparsity is well beyond what is needed for most applications.

⁹Actually, the behavior in Figure 2.7 is slightly better than what Theorem 2.2.6 predicts – with probability one the Kruskal rank of \mathbf{A} is m , and so the theorem shows that ℓ^0 minimization succeeds when $k \leq \frac{m}{2} = 2$. However, in the experiment, success always occurs when $k \leq 4$. Exercise 2.7 asks you to explain this discrepancy, by proving a modified version of Theorem 2.2.6.

¹⁰For example, $\text{krank}(\mathbf{A}) = m$ with probability one whenever \mathbf{A} is distributed according to any absolutely continuous measure, i.e., there is a probability density function.

2.2.4 Computational Complexity of ℓ^0 Minimization

The theoretical results in the previous section show the power of sparsity: knowing that the target solution \mathbf{x}_o is even moderately sparse can render the problem of recovering \mathbf{x}_o well-posed. Unfortunately, Algorithm 2.1 is not very useful in practice. Its worst-case running time is on the order of n^k , where $k = \|\mathbf{x}_o\|_0$ is the number of nonzero entries we wish to recover. For example, at the time of writing this book, to solve a problem with $m = 50$, $n = 200$, and $k = 10$, on a standard laptop, Algorithm 2.1 would require ≈ 140 centuries. This is still a very small problem by the standard of most modern-day applications!

Exhaustively searching all possible supports I may not seem like a particularly intelligent strategy for solving (2.2.9). However, no significantly better algorithm is currently known. There is some reason to believe that none will ever be found:

Theorem 2.2.8 (Hardness of ℓ^0 minimization). *The ℓ^0 -minimization problem (2.2.9) is NP-hard.*

NP-hard problems are considered highly unlikely to be efficiently solvable (i.e., solvable on standard (model) computers in time polynomial in the size of the problem). This class of problems includes notoriously difficult examples, such as the Traveling Salesman Problem. Fully appreciating the mathematical content of Theorem 2.2.8 requires formal modeling of computation (Turing machines, complexity theory for problem classes P/NP, etc.) that is beyond our scope. For readers unfamiliar with these details, Appendix 2.A gives a brief, high-level description of the issues that arise. For practical purposes, we should interpret the result as indicating that we cannot expect to give an efficient algorithm that guarantees to solve all instances of (2.2.9).

Hardness results such as Theorem 2.2.8 are typically proved by reduction: we show that if we can solve the problem of interest efficiently, this would allow us to also efficiently solve some other problem, which is already known to be hard. For the ℓ^0 minimization problem, we do this by showing that ℓ^0 minimization can be used to solve certain (hard) set covering problems:

Proof of Theorem 2.2.8. Consider the following problem:

Exact 3-Set Cover (E3C): Given a set $S = \{1, 2, \dots, m\}$ and a collection $\mathcal{C} = \{U_1, \dots, U_n\}$ of subsets $U_j \subseteq S$ each of which has size $|U_j| = 3$, does there exist a subcollection $\mathcal{C}' \subseteq \mathcal{C}$ that exactly covers S , i.e., $\forall i \in S$ there is exactly one $U \in \mathcal{C}'$ with $i \in U$.

This problem is known to be NP-hard [Garey and Johnson, 1979, Karp, 1972]. To reduce it to ℓ^0 minimization, suppose that we are given an instance of E3C: Form a $m \times n$ matrix $\mathbf{A} \in \{0, 1\}^{m \times n}$ by letting $\mathbf{A}_{ij} = 1$ if

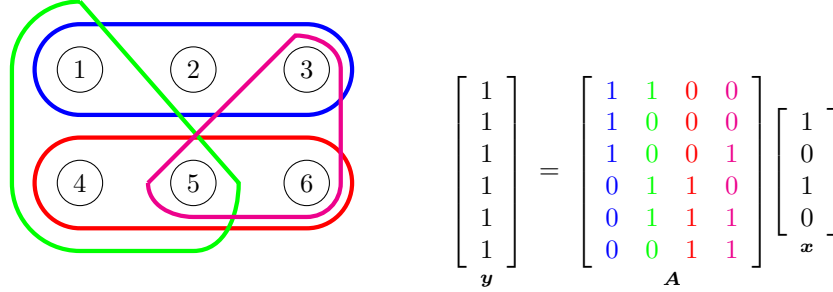


Figure 2.8. **Exact 3-Set Cover as a Sparse Representation problem.** Left: a universe $S = \{1, \dots, 6\}$ and four subsets $U_1, \dots, U_4 \subseteq S$. $\{U_1, U_3\}$ is an exact 3-set cover. Right: the same problem as a linear system of equations. The columns of \mathbf{A} are the incidence vectors for the sets U_1, U_2, U_3, U_4 . The Exact 3-Cover $\{U_1, U_3\}$ corresponds to a solution \mathbf{x} to the system $\mathbf{Ax} = \mathbf{y}$ with only $m/3 = 2$ nonzero entries.

$i \in U_j$, and $\mathbf{A}_{ij} = 0$ otherwise. Set $\mathbf{y} = \mathbf{1} \in \mathbb{R}^m$ (i.e., an m -dimensional vector of ones). Figure 2.8 illustrates this construction. We show:

Claim: The system $\mathbf{Ax} = \mathbf{y}$ has a solution \mathbf{x}_o with $\|\mathbf{x}_o\|_0 \leq m/3$ if and only if there exists an exact 3-set cover.

(\Leftarrow) Suppose there exists an exact 3-set cover \mathcal{C}' . Clearly, $|\mathcal{C}'| = m/3$. Set

$$x_j = \begin{cases} 1 & U_j \in \mathcal{C}' \\ 0 & \text{else} \end{cases}.$$

Then $\|\mathbf{x}\|_0 = m/3$, and $\mathbf{y} = \mathbf{Ax}$.

(\Rightarrow) Let \mathbf{x}_o be a solution to $\mathbf{y} = \mathbf{Ax}$ with at most $m/3$ nonzero entries. Set $\mathcal{C}' = \{U_j \mid x_o(j) \neq 0\}$. We claim \mathcal{C}' is the desired cover. Let $I = \text{supp}(\mathbf{x}_o)$. Since each column of \mathbf{A} has exactly 3 nonzero entries, and \mathbf{A}_I has at most $m/3$ columns, the matrix \mathbf{A}_I has at most m nonzero entries. Since $\mathbf{A}_I \mathbf{x}_o(I) = \mathbf{y}$, each row of \mathbf{A}_I has at least one nonzero entry. Hence, each row of \mathbf{A}_I has *exactly* one nonzero entry, and the set \mathcal{C}' gives an exact cover. \square

In fact, the truth is even worse than Theorem 2.2.8 suggests: The ℓ^0 minimization problem remains NP-hard even if we only demand that $\mathbf{Ax} \approx \mathbf{y}$, in an appropriate sense. It is also NP-hard to find an \mathbf{x} whose number of nonzero entries is within a constant factor of the smallest possible! Based on our current understanding of complexity theory, it is extraordinarily unlikely that anyone will ever discover an efficient algorithm that solves any interesting variant of the ℓ^0 minimization problem for all possible inputs (\mathbf{A}, \mathbf{y}) .

2.3 Relaxing the Sparse Recovery Problem

The rather bleak worst-case picture for ℓ^0 -minimization has not stopped engineers from searching for efficient heuristics for finding sparse solutions to linear systems. There is always some possibility for optimism:

“Although the *worst* sparse recovery problem may be impossible to solve efficiently, perhaps my *particular* instance of interest is not so hard.”

This optimism is occasionally rewarded in a rather striking fashion. In the next few chapters, we will see that many sparse recovery problems that matter for engineering practice *are* solvable efficiently. Our first step is to find a proper replacement for the ℓ^0 norm which still encourages sparsity, but can be optimized efficiently.

2.3.1 Convex Functions

If our goal is efficient optimization, perhaps the most natural class of objective functions to consider is the *convex* functions. Smooth convex functions often appear “bowl shaped” – as in Figure 2.9 (left). Indeed, a necessary and sufficient condition for a smooth function $f : \mathbb{R} \rightarrow \mathbb{R}$ to be convex is that it exhibit nonnegative curvature – its second derivative $\frac{d^2 f}{dx^2}(x) \geq 0$ at every point x .

Iterative methods for optimization seek a minimizer of an objective function $f(\mathbf{x})$, by starting from some initial point \mathbf{x}_0 ,¹¹ and then generating a new point \mathbf{x}_1 based on the local shape of the objective function in the vicinity of \mathbf{x}_0 . For a smooth function $f(\mathbf{x})$, the negative gradient $-\nabla f(\mathbf{x})$ defines the direction in which the objective function decreases most rapidly. A natural strategy for choosing \mathbf{x}_1 is to set

$$\mathbf{x}_1 = \mathbf{x}_0 - t \nabla f(\mathbf{x}_0). \quad (2.3.1)$$

where t is a step size. Continuing in this manner to produce points $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots$, we obtain the *gradient descent* method, a natural and intuitive algorithm for minimizing a smooth function $f(\mathbf{x})$. For the function f in Figure 2.9 (left), assuming we choose the step size t appropriately, the iterates $\mathbf{x}_0, \mathbf{x}_1, \dots$ will converge to the global minimizer \mathbf{x}_* . For the nonconvex function to the right, this strategy only guarantees a local minimizer.

Convex functions such as Figure 2.9 (left) have the property that every local minimizer is a global minimizer.¹² Moreover, many convex functions

¹¹In this book, we will use \mathbf{x}_0 to indicate the initial point of an iterative algorithm, which is not to be confused with \mathbf{x}_o , the desired ground truth.

¹²It is worth noting that for many of the problems we will later discuss (e.g., MRI, spectrum sensing, face recognition), global optimality is very important – there is a *true*

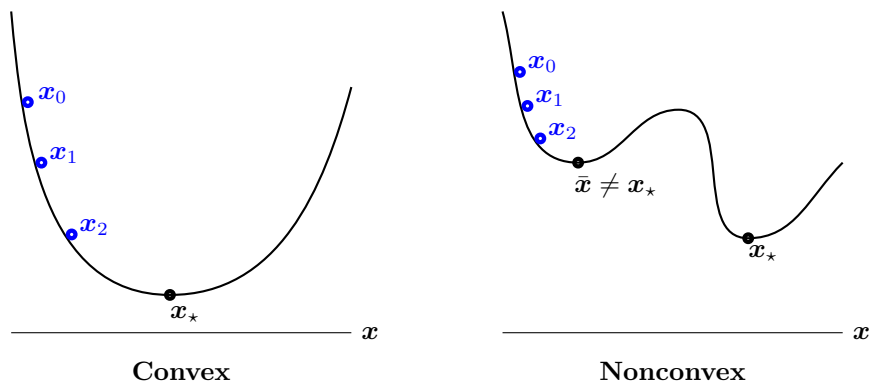


Figure 2.9. **Convex and nonconvex functions.** Left: a convex function. Local descent methods such as gradient descent produce a sequence of points $\mathbf{x}_0, \mathbf{x}_1, \dots$ which approach the global minimizer \mathbf{x}_* . Right: A nonconvex function. For this particular function, depending on the initial point \mathbf{x}_0 , local descent methods may produce the suboptimal local minimum $\bar{\mathbf{x}}$. Motivated by their good properties for optimization, in the first part of this book, we will seek convex formulations for recovering sparse (and otherwise structured) signals.

arising in practice can be optimized efficiently using variants of gradient descent. Indeed, in Chapter 8, we will see that the particular convex functions that we encounter in computing with sparse signals (and their generalizations) *can* be efficiently optimized, even on a large scale and in high dimensions.

We review the properties of convex functions more formally in Appendix C. Here, we briefly remind the reader of the general definition of convex function:¹³

Definition 2.3.1 (Convex function on \mathbb{R}^n). *A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if for every pair of points $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^n$ and $\alpha \in [0, 1]$,*

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{x}') \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{x}'). \quad (2.3.2)$$

signal that we are trying to recover, and it is important to build algorithms that can do this reliably. In our simulated example of ℓ^0 minimization, we declared the solution $\hat{\mathbf{x}}$ correct, because it coincided with the true \mathbf{x}_o that generated the observation \mathbf{y} . This is in contrast to some applications of optimization (e.g., in finance) where the objective function measures the goodness of the solution (say the expected rate of return on an investment), and locally improving the solution is meaningful, or even desirable, if the objective corresponds to dollars earned/lost!

¹³On the surface, this definition appears much more complicated than simply asking the second derivative to be positive. The reason for this complication is that we will need to work with convex functions that are not smooth; the general definition in Definition 2.3.1 handles this situation as well.

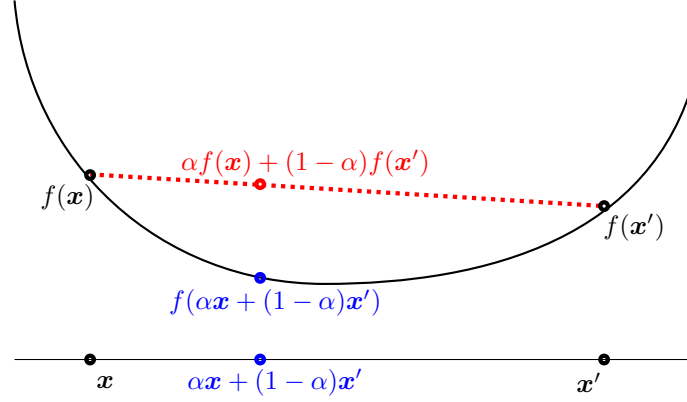


Figure 2.10. **Definition of convexity.** A convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is one which satisfies the inequality $f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{x}') \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{x}')$ for all $\alpha \in [0, 1]$ and $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^n$. Geometrically, this means that if we take the points $(\mathbf{x}, f(\mathbf{x}))$ and $(\mathbf{x}', f(\mathbf{x}'))$ on the graph of f , and then draw a **line** joining them, the graph of the function falls below this line segment.

This inequality can be visualized as follows. Consider two points $(\mathbf{x}, f(\mathbf{x}))$ and $(\mathbf{x}', f(\mathbf{x}'))$ on the graph of f . If we form the line segment joining these two points, this line segment lies above the graph of f . Figure 2.10 visualizes this inequality for several different convex functions.

A *convex combination* of a collection of points $\mathbf{x}_1, \dots, \mathbf{x}_k$ is an expression of the form $\sum_{i=1}^k \lambda_i \mathbf{x}_i$, where the weights λ_i are nonnegative and $\sum_{i=1}^k \lambda_i = 1$. For example, for $\alpha \in [0, 1]$, the expression $\mathbf{z} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{x}'$ is a convex combination of the points \mathbf{x} and \mathbf{x}' . The definition (2.3.2) states that at the point \mathbf{z} , the function f is no larger than the corresponding combination $\alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{x}')$ of the function values at the points \mathbf{x} and \mathbf{x}' .

This property of convex functions generalizes, to give Jensen's inequality, which states that the value of a convex function f at a convex combination of points is no greater than the corresponding convex combination of the function values:

Proposition 2.3.2 (Jensen's inequality). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function. Then for any k , any collection of points $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^n$ and any nonnegative scalars $\lambda_1, \dots, \lambda_k$ satisfying $\sum_{i=1}^k \lambda_i = 1$,*

$$f\left(\sum_{i=1}^k \lambda_i \mathbf{x}_i\right) \leq \sum_{i=1}^k \lambda_i f(\mathbf{x}_i). \quad (2.3.3)$$

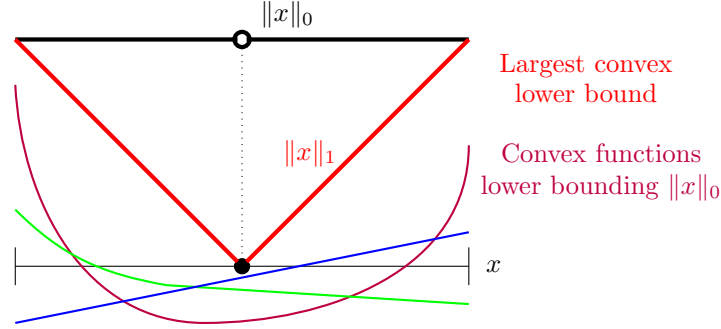


Figure 2.11. **A convex surrogate for the ℓ^0 norm.** In black, we plot the graph of the ℓ^0 norm of a scalar x , over the interval $x \in [-1, 1]$. This function takes on value 0 at $x = 0$, and +1 everywhere else. In purple, green and blue, we plot various convex function examples $f(x)$ which underestimate $\|x\|_0$ on $[-1, 1]$, in the sense that $f(x) \leq \|x\|_0$ for all $x \in [-1, 1]$. In red, we plot the function $f(x) = |x|$. This is the largest convex function which underestimates $\|x\|_0$ on $[-1, 1]$. We call $|x|$ the *convex envelope* of $\|x\|_0$ on $[-1, 1]$.

2.3.2 A Convex Surrogate for the ℓ^0 Norm: the ℓ^1 Norm

With the good properties of convex functions in mind, let us try to find a convex “surrogate” for the ℓ^0 norm. In one dimension, x is a scalar, and $\|x\|_0 = \mathbb{1}_{x \neq 0}$ is simply the indicator function for nonzero x . From Figure 2.11, it is clear that if we restrict our attention to the interval $x \in [-1, 1]$, the largest convex function which does not exceed $\|\cdot\|_0$ on this interval is simply the absolute value $|x|$. In the language of convex analysis, $|x|$ is the *convex envelope* of the function $\|x\|_0$ over the set $[-1, 1]$. This means that $|x|$ is the largest convex function f which satisfies $f(x) \leq \|x\|_0$ for every $x \in [-1, 1]$, i.e., it is the largest convex underestimator of $\|x\|_0$ over this set. Thus, in one dimension, we might consider the absolute value of x as a plausible replacement for $\|x\|_0$.

For higher-dimensional \mathbf{x} (i.e., $\mathbf{x} \in \mathbb{R}^n$), the ℓ^0 norm is

$$\|\mathbf{x}\|_0 = \sum_{i=1}^n \mathbb{1}_{\mathbf{x}(i) \neq 0}. \quad (2.3.4)$$

Applying our above reasoning to each of the coordinates $\mathbf{x}(i)$, we obtain the ℓ^1 norm

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |\mathbf{x}(i)|. \quad (2.3.5)$$

As in the scalar case, this function is the tightest convex underestimator of $\|\cdot\|_0$, over an appropriate set of vectors \mathbf{x} :

Theorem 2.3.3. *The function $\|\cdot\|_1$ is the convex envelope of $\|\cdot\|_0$, over the set $\mathfrak{B}_\infty = \{\mathbf{x} \mid \|\mathbf{x}\|_\infty \leq 1\}$ of vectors whose elements all have magnitude at most one.*

Proof. Let f be a convex function satisfying $f \leq \|\cdot\|_0$ on \mathfrak{B}_∞ . We prove that $f \leq \|\cdot\|_1$ on \mathfrak{B}_∞ as well. Consider the cube $C = [0, 1]^n$. Its vertices are the vectors $\boldsymbol{\sigma} \in \{0, 1\}^n$. Any $\mathbf{x} \in C$ can be written as a convex combination of these vertices:

$$\mathbf{x} = \sum_i \lambda_i \boldsymbol{\sigma}_i \quad (2.3.6)$$

Because $f \leq \|\cdot\|_0$, $f(\boldsymbol{\sigma}_i) \leq \|\boldsymbol{\sigma}_i\|_0 = \|\boldsymbol{\sigma}_i\|_1$. Because f is convex,

$$\begin{aligned} f(\mathbf{x}) &= f\left(\sum_i \lambda_i \boldsymbol{\sigma}_i\right) \leq \sum_i \lambda_i f(\boldsymbol{\sigma}_i) && \text{[Jensen's inequality]} \\ &\leq \sum_i \lambda_i \|\boldsymbol{\sigma}_i\|_0 = \sum_i \lambda_i \|\boldsymbol{\sigma}_i\|_1 && [\boldsymbol{\sigma}_i \text{ are binary}] \\ &= \|\mathbf{x}\|_1. \end{aligned} \quad (2.3.7)$$

Hence, $f \leq \|\cdot\|_1$ on the intersection of \mathfrak{B}_∞ with the nonnegative orthant. Repeating the argument for each of the orthants, we obtain that $f \leq \|\cdot\|_1$ on \mathfrak{B}_∞ , and hence $\|\cdot\|_1$ is the convex envelope of $\|\cdot\|_0$ over \mathfrak{B}_∞ . \square

So, at least in the sense of convex envelopes, the ℓ^1 norm provides a good replacement for the ℓ^0 norm. Replacing the ℓ^0 norm in (2.2.9) with the ℓ^1 norm, we obtain a convex optimization problem,

$$\begin{aligned} &\text{minimize} && \|\mathbf{x}\|_1 \\ &\text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{y}. \end{aligned} \quad (2.3.8)$$

In contrast to the ℓ^0 problem, this problem *can* be solved efficiently.

2.3.3 A Simple Test of ℓ^1 Minimization

Theorem 2.3.3 is a strong initial motivation for considering (2.3.8) – it says that in a certain sense, the ℓ^1 norm is the canonical convex surrogate for the ℓ^0 norm. Some care is in order, though. Theorem 2.3.3 does not say anything at all about the *correctness* of (2.3.8) – whether the solution to (2.3.8) is actually the desired sparse vector \mathbf{x}_o .

The easiest way to get some insight into this question is to do an experiment! For this, we will need to solve the problem (2.3.8) computationally and see how well it works. How do we solve the optimization problem (2.3.8)? Appendix C gives a quick introduction to some general optimization techniques that may help us solve problems of this kind. More specifically, since the objective function is convex, the geometry of a convex function in Figure 2.12 (left) suggests that we should do quite well just using local information about the slope of the objective function. Indeed, if

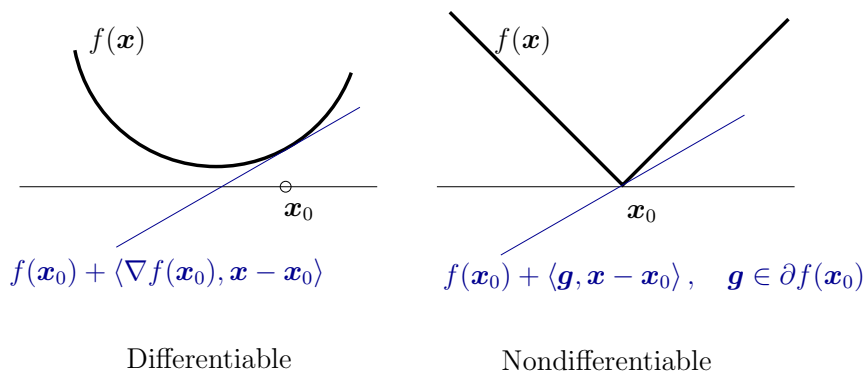


Figure 2.12. **Subgradients of convex functions.** Left: for a differentiable convex function, the best linear approximation at any point \mathbf{x}_0 is a *global* lower bound on the function. Right: for a nondifferentiable function, we say that \mathbf{g} is a subgradient of f at \mathbf{x}_0 (and write $\mathbf{g} \in \partial f(\mathbf{x}_0)$ if \mathbf{g} is the slope of a linear function that takes on value $f(\mathbf{x}_0)$ at \mathbf{x}_0 , and globally lower bounds f).

our objective function were differentiable, this would very naturally suggest the classical *gradient descent* method for solving problems of the form

$$\text{minimize } f(\mathbf{x}). \quad (2.3.9)$$

This algorithm starts at some initial point \mathbf{x}_0 , and then generates a sequence of points $(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k, \dots)$ by iteratively moving in the direction of greatest decrease of $f(\cdot)$:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \nabla f(\mathbf{x}_k). \quad (2.3.10)$$

Here, $t_k \geq 0$ is a properly chosen step size.

There are two main difficulties that prevent us from applying the gradient descent iteration (2.3.10) to the ℓ^1 minimization problem (2.3.8):

- **Nontrivial constraints:** Unlike the general unconstrained problem (2.3.9), in the problem (2.3.8) we are only interested in \mathbf{x} that satisfy $\mathbf{A}\mathbf{x} = \mathbf{y}$.
- **Nondifferentiable objective:** The objective function in (2.3.8) is not differentiable, and so at certain points the gradient does not exist. Figure 2.12 (right) shows this: the function is pointed at zero! Since zero is sparse, this is precisely one of the points we are most interested in.

Constraints.

One approach to handling the first problem is to replace the gradient descent iteration with *projected gradient descent*. This algorithm aims at

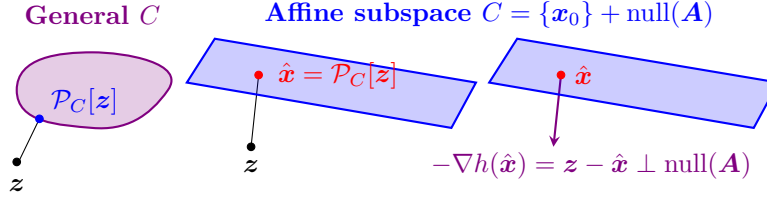


Figure 2.13. **Projection onto convex sets.** Left: projection onto a general convex set. Middle: projection onto an affine subspace. Right: the projection onto the affine subspace can be characterized as the point $\hat{\mathbf{x}}$ at which the gradient $\nabla h(\hat{\mathbf{x}})$ is orthogonal to $\text{null}(\mathbf{A})$.

general problems of the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in C, \end{aligned} \quad (2.3.11)$$

where C is some constraint set. This algorithm is exactly the same as gradient descent, except that at each iteration it *projects* the result $\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)$ onto the set C . The projection of a point \mathbf{z} onto the set C is simply the nearest point to \mathbf{z} in C :

$$\mathcal{P}_C[\mathbf{z}] = \arg \min_{\mathbf{x} \in C} \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|_2^2 \equiv h(\mathbf{x}). \quad (2.3.12)$$

For general C , the projection may not exist, or may not be unique (think about how this could happen). However, for closed, convex sets, the projection is well-defined, and satisfies a wealth of useful properties. If \mathbf{A} has full row rank, the projection onto the convex set $C = \{\mathbf{x} \mid \mathbf{A}\mathbf{x} = \mathbf{y}\}$ has an especially simple form:

$$\mathcal{P}_{\{\mathbf{x} \mid \mathbf{A}\mathbf{x} = \mathbf{y}\}}[\mathbf{z}] = \mathbf{z} - \mathbf{A}^* (\mathbf{A}\mathbf{A}^*)^{-1} [\mathbf{A}\mathbf{z} - \mathbf{y}]. \quad (2.3.13)$$

Figure 2.13 visualizes the projection onto this particular C . This formula can be derived by noting two properties of the projection $\hat{\mathbf{x}} = \mathcal{P}_C[\mathbf{z}]$:

- (i) **Feasibility:** $\hat{\mathbf{x}} \in C$, i.e., $\mathbf{A}\hat{\mathbf{x}} = \mathbf{y}$.
- (ii) **Error is orthogonal:** $\mathbf{z} - \hat{\mathbf{x}} \perp \text{null}(\mathbf{A})$. Since $\mathbf{z} - \hat{\mathbf{x}} = -\nabla h(\hat{\mathbf{x}})$, this condition can be stated as

$$-\nabla h(\hat{\mathbf{x}}) \text{ is normal to } C \text{ at } \hat{\mathbf{x}}.$$

Exercise 2.9 guides the interested reader through the derivation of this expression. For the general problem (2.3.11), with differentiable objective f , the projected gradient algorithm simply repeats the iteration

$$\mathbf{x}_{k+1} = \mathcal{P}_C[\mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)]. \quad (2.3.14)$$

Nondifferentiability.

The problem of nondifferentiability is slightly trickier. To handle it properly, we need to generalize the notion of derivative to include functions that are not differentiable. For this, we draw inspiration from geometry. Consider Figure 2.12 (left). It displays a convex, differentiable function $f(\mathbf{x})$, as well as a linear approximation $\hat{f}(\mathbf{x})$, taken at a point \mathbf{x}_0 :

$$\hat{f}(\mathbf{x}) = f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle. \quad (2.3.15)$$

The salient point here is that the graph of f lies entirely above the graph of the approximation \hat{f} :

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle \quad \forall \mathbf{x} \in \mathbb{R}^n. \quad (2.3.16)$$

It is not too difficult to prove that this property holds for *every* convex differentiable function and every point \mathbf{x}_0 , simply by using calculus and the definition of convexity.

This geometry opens the door for generalizing the notion of the gradient to nonsmooth functions. For nonsmooth functions such as $f(\mathbf{x}) = \|\mathbf{x}\|_1$, at a point of nonsmoothness \mathbf{x}_0 , the gradient does not exist, but we can still make a linear under-estimator

$$\hat{f}(\mathbf{x}) = f(\mathbf{x}_0) + \langle \mathbf{u}, \mathbf{x} - \mathbf{x}_0 \rangle, \quad (2.3.17)$$

as in Figure 2.12 (right). Here, \mathbf{u} replaces ∇f in the previous expression, and plays the role of the “slope” of the approximation. We say that \mathbf{u} is a *subgradient* of f at \mathbf{x}_0 if the linear approximation defined by \mathbf{u} is indeed an under-estimator of f (i.e., it lower bounds $f(\mathbf{x})$ at all points \mathbf{x}):

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \mathbf{u}, \mathbf{x} - \mathbf{x}_0 \rangle, \quad \forall \mathbf{x}. \quad (2.3.18)$$

Let us consider our function of interest – the ℓ^1 norm. For $\mathbf{x} \in \mathbb{R}$ (one dimension), $\|\mathbf{x}\|_1 = |x|$ is simply the absolute value. For $x < 0$, the slope of the graph of $|x|$ is -1 , while for $|x| > 0$, it is $+1$. Convince yourself that if we take $\mathbf{x}_0 \neq 0$, then the only u satisfying the above definition is $u = \text{sign}(x)$.

However, at 0 the function $|x|$ is “pointy,” namely, nondifferentiable, and something different happens: at $x_0 = 0$, every $u \in [-1, 1]$ defines a linear approximation that underestimates f . So, in fact, every $u \in [-1, 1]$ is a subgradient. Thus, at points of nondifferentiability there may exist multiple subgradients. We call the collection of all subgradients of f at a point \mathbf{x}_0 the *subdifferential* of f at \mathbf{x}_0 , and denote it by $\partial f(\mathbf{x}_0)$. Formally:

Definition 2.3.4 (Subgradient and subdifferential). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function. A subgradient of f at \mathbf{x}_0 is any \mathbf{u} satisfying*

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \mathbf{u}, \mathbf{x} - \mathbf{x}_0 \rangle, \quad \forall \mathbf{x}. \quad (2.3.19)$$

The subdifferential of f at \mathbf{x}_0 is the set of all subgradients of f at \mathbf{x}_0 :

$$\partial f(\mathbf{x}_0) = \{\mathbf{u} \mid \forall \mathbf{x} \in \mathbb{R}^n, f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \mathbf{u}, \mathbf{x} - \mathbf{x}_0 \rangle\}. \quad (2.3.20)$$

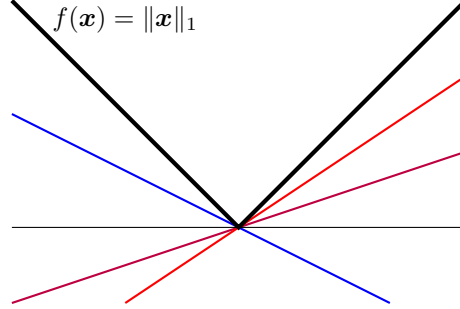


Figure 2.14. **Subdifferential of the ℓ^1 norm.** In **black**, $f(\mathbf{x}) = \|\mathbf{x}\|_1$. In **blue**, **purple**, and **red**, three linear lower bounds of the form $g(\mathbf{x}) = f(\mathbf{x}_0) + \langle \mathbf{u}, \mathbf{x} - \mathbf{x}_0 \rangle$, taken at $\mathbf{x}_0 = \mathbf{0}$, with slope $\mathbf{u} = -\frac{1}{2}$, $\frac{1}{3}$, and $\frac{2}{3}$, respectively. It should be clear that any slope $\mathbf{u} \in [-1, 1]$ defines a linear lower bound on $f(\mathbf{x})$ around $\mathbf{x}_0 = \mathbf{0}$. So, $\partial|\cdot|(0) = [-1, 1]$. For $\mathbf{x}_0 > 0$, the only linear lower bound has slope $\mathbf{u} = 1$; for $\mathbf{x}_0 < 0$, the only linear lower bound has slope $\mathbf{u} = -1$. So, $\partial|\cdot|(x) = \{-1\}$ for $x < 0$ and $\partial|\cdot|(x) = \{1\}$ for $x > 0$. Lemma 2.3.5 proves this formally, and extends to higher-dimensional $\mathbf{x} \in \mathbb{R}^n$.

With these definitions in mind, we might imagine that in the nonsmooth case, a suitable replacement for the gradient algorithm might be the *subgradient method*, which chooses (somehow) $\mathbf{g}_k \in \partial f(\mathbf{x}_k)$, and then proceeds in the direction of $-\mathbf{g}_k$: $\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \mathbf{g}_k$. Incorporating projection onto the feasible set, we arrive at the following *projected subgradient algorithm*:

$$\mathbf{x}_{k+1} = \mathcal{P}_C[\mathbf{x}_k - t_k \mathbf{g}_k], \quad \mathbf{g}_k \in \partial f(\mathbf{x}_k). \quad (2.3.21)$$

To apply the projected subgradient method, we need an expression for the subdifferential of the ℓ^1 norm. Figure 2.14 visualizes this. In one dimension, $\|\mathbf{x}\|_1 = |x|$; this function is differentiable away from $x = 0$. For $x > 0$, $\partial|\cdot|(x) = \{1\}$, while for $x < 0$, $\partial|\cdot|(x) = \{-1\}$. At $x = 0$, $|x|$ is not differentiable, and there are multiple possible linear lower bounds. Figure 2.14 visualizes three of these lower bounds. It is not difficult to see that lower bounds at $x = 0$ can have any slope from -1 to 1 ; hence, $\partial|\cdot|(x) = [-1, 1]$. The following lemma extends this observation to higher-dimensional $\mathbf{x} \in \mathbb{R}^n$:

Lemma 2.3.5 (Subdifferential of $\|\cdot\|_1$). *Let $\mathbf{x} \in \mathbb{R}^n$, with $I = \text{supp}(\mathbf{x})$,*

$$\partial \|\cdot\|_1(\mathbf{x}) = \{\mathbf{v} \in \mathbb{R}^n \mid \mathbf{P}_I \mathbf{v} = \text{sign}(\mathbf{x}), \|\mathbf{v}\|_\infty \leq 1\}. \quad (2.3.22)$$

Here, $\mathbf{P}_I \in \mathbb{R}^{n \times n}$ is the orthoprojector onto coordinates I :

$$[\mathbf{P}_I \mathbf{v}](j) = \begin{cases} v_j & j \in I \\ 0 & j \notin I \end{cases}. \quad (2.3.23)$$

Proof. The subdifferential $\partial \|\cdot\|_1(\mathbf{x})$ consists of all vectors \mathbf{v} that satisfy

$$\sum_{i=1}^n |\mathbf{x}'(i)| \geq \sum_{i=1}^n |\mathbf{x}(i)| + \mathbf{v}(i) (\mathbf{x}'(i) - \mathbf{x}(i)) \quad (2.3.24)$$

for every \mathbf{x} and \mathbf{x}' . A sufficient condition is that for every index i and every scalar z ,

$$|z| \geq |\mathbf{x}(i)| + \mathbf{v}(i)(z - \mathbf{x}(i)). \quad (2.3.25)$$

Taking $\mathbf{x}' = \mathbf{x} + (z - \mathbf{x}(i))\mathbf{e}_i$ in (2.3.24) shows that (2.3.25) is also necessary. If $\mathbf{x}(i) = 0$, (2.3.25) becomes $|z| \geq \mathbf{v}(i)z$, which holds for all z if and only if $|\mathbf{v}(i)| \leq 1$. If $\mathbf{x}(i) \neq 0$, the inequality is satisfied if and only if $\mathbf{v}(i) = \text{sign}(\mathbf{x}(i))$. Hence, $\mathbf{v} \in \partial \|\cdot\|_1$ if and only if for all $i \in I$, $\mathbf{v}(i) = \text{sign}(\mathbf{x}(i))$, and for all i , $|\mathbf{v}(i)| \leq 1$. This conclusion is summarized as (2.3.22). \square

The projected subgradient method alternates between subgradient steps, which move in the direction of $-\text{sign}(\mathbf{x})$, and orthogonal projections onto the feasible set $\{\mathbf{x} \mid \mathbf{A}\mathbf{x} = \mathbf{y}\}$ according to equation (2.3.13). We obtain a very simple algorithm that solves (2.3.8), which we spell out in detail as Algorithm 2.2.

Algorithm 2.2: ℓ^1 -minimization by projected subgradient

- 1: **Input:** a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and a vector $\mathbf{y} \in \mathbb{R}^m$.
 - 2: Compute $\mathbf{\Gamma} \leftarrow \mathbf{I} - \mathbf{A}^*(\mathbf{A}\mathbf{A}^*)^{-1}\mathbf{A}$, and $\tilde{\mathbf{x}} \leftarrow \mathbf{A}^\dagger \mathbf{y} = \mathbf{A}^*(\mathbf{A}\mathbf{A}^*)^{-1}\mathbf{y}$.
 - 3: $\mathbf{x}_0 \leftarrow \mathbf{0}$.
 - 4: $t \leftarrow 0$.
 - 5: **repeat many times**
 - 6: $t \leftarrow t + 1$.
 - 7: $\mathbf{x}_t \leftarrow \tilde{\mathbf{x}} + \mathbf{\Gamma} \left(\mathbf{x}_{t-1} - \frac{1}{t} \text{sign}(\mathbf{x}_{t-1}) \right)$.
 - 8: **end while**
-

Remark 2.3.6 (Projected Subgradient and Better Alternatives). *In many respects, this is a bad method for solving the ℓ^1 problem. It is correct, but it converges very slowly compared to methods that exploit a certain piece of problem-specific structure, which we will describe in later chapters. The main virtue of Algorithm 2.2 is that it is simple and intuitive, and also serves our exposition by introducing or reminding us of subgradients and projection operators.¹⁴ The projected subgradient method for ℓ^1 minimization can be implemented in just a few lines of Matlab code. In the later*

¹⁴Also, we would like you to have a feel for at least one *very* simple way for implementing ℓ^1 minimization in code and play with it. Our experience is that this helps to think more concretely about the optimization problem and its applications, rather than leaving it as a mathematical abstraction.

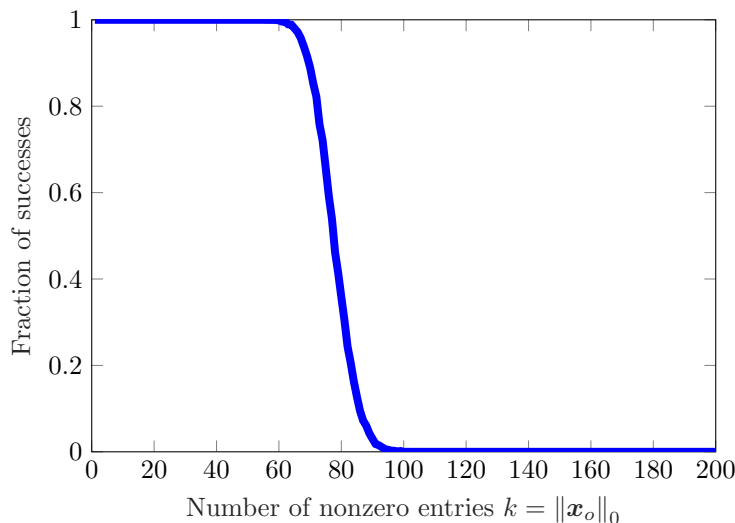


Figure 2.15. **Phase transition in ℓ^1 minimization.** We consider the problem of recovering a sparse vector \mathbf{x}_o from measurements $\mathbf{y} = \mathbf{A}\mathbf{x}_o$, where $\mathbf{A} \in \mathbb{R}^{100 \times 200}$ is a Gaussian matrix. We vary the number of nonzero entries $k = \|\mathbf{x}_o\|_0$ across $k = 0, 1, \dots, 200$, and plot the fraction of instances where ℓ^1 minimization successfully recovers \mathbf{x}_o , over 50 independent experiments for each value of k . Notice that this probability of success exhibits a (rather sharp) transition from 1 (guaranteed success) to 0 (guaranteed failure) as k increases. Notice moreover, that for sufficiently well-structured problems (k small), ℓ^1 minimization always succeeds.

chapters, we will see a number of more sophisticated and better-performing approaches that better utilize the structure in this problem.

To see how well ℓ^1 minimization (as implemented through the projected subgradient method) performs, run `Chapter_2_L1_recovery.m`. You may see an interesting phenomenon! Although the method does not *always* succeed, it *does* succeed whenever the target solution \mathbf{x}_o is *sufficiently sparse*! Figure 2.15 illustrates this in a more systematic way. In the figure, we generate random matrices \mathbf{A} of size 200×400 and random vectors \mathbf{x}_o with k nonzero entries. We vary k from 1 to 200. For each k , we run 50 experiments and plot the fraction of trials in which that ℓ^1 minimization correctly recovers \mathbf{x}_o , up to numerical error. Notice that indeed, ℓ^1 minimization succeeds whenever \mathbf{x}_o is sufficiently sparse.

Given the examples we have seen thus far of how sparsity arises in application problems, this phenomenon is certainly intriguing. In the coming chapters, we will study it first from a mathematical perspective, to understand *why* it occurs and what its limitations are; we will then investigate its implications for practical applications.

2.4 Summary

Let us briefly recap what we have learned in this chapter. In many modern data analysis and signal processing applications, we need to solve very large, underdetermined systems of linear equations:

$$\mathbf{y} = \mathbf{A}\mathbf{x}, \quad \mathbf{A} \in \mathbb{R}^{m \times n}, \quad m < n.$$

Such problems are inherently ill-posed: they admit infinitely many solutions.

To make such problems well-posed, or to make the solution unique, we need to leverage additional properties of the solution that we wish to recover. One important property, which obtains in many practical applications, is sparsity (or compressibility). This is a powerful piece of information: although the signals themselves reside in a very high-dimensional space, they have only a few intrinsic degrees of freedom – they can be represented as a linear superposition of just a few atoms from a properly chosen dictionary. As Theorem 2.2.6 shows, under fairly general conditions, imposing sparsity on \mathbf{x} can indeed make the problem of solving

$$\min \|\mathbf{x}\|_0 \quad \text{subject to} \quad \mathbf{y} = \mathbf{A}\mathbf{x}$$

well conditioned: As long as the target solution \mathbf{x}_o is sufficiently sparse w.r.t. the Kruskal rank of \mathbf{A} , the sparsest solution to $\mathbf{y} = \mathbf{A}\mathbf{x}$ is unique and is the correct solution.

Computationally, however, finding the sparsest solution to a linear system is in general intractable (Theorem 2.2.8). To alleviate the computational difficulty, we relax the ℓ^0 minimization problem and replace the ℓ^0 norm of \mathbf{x} with its convex envelope, the ℓ^1 norm:

$$\min \|\mathbf{x}\|_1 \quad \text{subject to} \quad \mathbf{y} = \mathbf{A}\mathbf{x}.$$

We have introduced a very basic subgradient descent algorithm (Algorithm 2.2) that solves the convex ℓ^1 minimization problem. From the results of the algorithm, we observe a striking phenomenon that ℓ^1 minimization can effectively recover the sparse solution under fairly broad conditions. We will explain why this is the case in the next chapter after we carefully characterize exact conditions under which ℓ^1 minimization gives the correct sparse solution.

2.5 Notes and References

Sparsity.

Historically, the notion of sparsity has arisen in many different guises, across different fields and applications.

Statistics – in statistics, sparsity arises in the context of *variable selection*. For example, in medical research, we may wish to select a few genes (predictors) whose expression levels predict the presence or absence of a certain type of disease. It is often important that we select a small (even minimal) set of predictors, so that the model we fit can be interpreted and further studied by humans. Classical variable selection methods, such as *stepwise regression* [Efroymson, 1960] are closely related to greedy methods for ℓ^0 minimization – see Chapter 8 for more on these methods. In recent years, ℓ^1 minimization (LASSO) has been widely used for variable selection, thanks to the seminal work of Tibshirani [Tibshirani, 1996]. We will examine the properties of this method in depth, beginning in the next chapter. Sparsity is also related to the principle of minimum description length [Rissanen, 1978], which argues that in choosing between various models, we should prefer models which can be encoded efficiently – see, e.g., [Hansen and Yu, 2001].

Error correction and robustness – the robust face recognition problem discussed in Section 2.1.3 illustrates the tight connection between sparsity and error correction. See Exercise 2.10 and the work of Candes and Tao [Candes and Tao, 2005a] for a more general view on this connection. Robustness to erroneous or outlying observations is an important issue for statistical inference [Huber, 1981, Hampel et al., 1986], and for communications, where we wish to ensure reliable communications over unreliable channels. The first uses of ℓ^1 minimization that we are aware of arise in the context of statistical estimation with errors, in the work of Boscovitch [Boscovitch, 1750] and later Laplace [Laplace, 1774].

Neuroscience – the human brain is believed to encode sensory phenomena in a sparse fashion, in which only a small fraction of the available neurons are active at any given time (see, e.g., Field [Field, 1987]). Seminal work of Olshausen and Field [Olshausen and Field, 1997] observes that when applied to large families of natural images, optimization methods that pursue a dictionary which represents the images as *sparsely* as possible naturally produce dictionaries that reproduce qualitative properties of the receptive fields in the human visual cortex. The problem of producing a good dictionary for sparsely encoding a data set is known as *dictionary learning*. We discuss algorithmic aspects of this problem in Chapters 9 and 15. For a contemporary discussion of the neuroscientific implications of sparsity, we refer the interested reader to Ganguli and Sompolinsky [Ganguli and Sompolinsky, 2012] and the references therein.

Data compression – there is a natural connection between sparsity and data compression.

Geophysics – amongst the earliest applications of sparse recovery are problems in geophysical exploration.

Application vignettes.

The three applications described in this chapter illustrate various aspects of sparse modeling and sparse recovery. The medical imaging application is described in the work of Lustig et. al. [Lustig et al., 2007b, Lustig et al., 2008a]. The denoising results shown in Section 2.1.2 are due to Mairal et. al. [Mairal et al., 2008a]. The face recognition formulation in Section 2.1.3 is described in [Wright et al., 2009]. The discussion in this chapter only touches the surface of these problems; we will revisit medical imaging in Chapter 11, dictionaries and image processing in Chapter 15, and face recognition in Chapter 10. Please see these chapters and their references for broader context and related work on each of these problems. These are just a few of the vast array of applications of sparse methods; a few of these are highlighted in Chapters 12 - 16.

NP hardness of ℓ^0 minimization and related problems.

The hardness result for ℓ^0 minimization, Theorem 2.2.8, is due to Natarajan [Natarajan, 1995]; see also Davis, Mallat and Avellaneda [Davis et al., 1997]. Results of Amaldi and Kann [Amaldi and Kann, 1995, Amaldi and Kann, 1998] and Arora, Babai, Stern and Sweedyk [Arora et al., 1993] show that that ℓ^0 minimization problems are also NP-hard to approximate. Delineating the boundaries between tractable and intractable instances of sparse approximation remains an active topic of research: see, e.g., Zhang, Wainwright and Jordan [Zhang et al., 2014] and Foster, Karloff and Thaler [Foster et al., 2015] for more recent results. There are hardness results for a number of problems that relate closely to sparse approximation. The above results also have implications for sparse error correction. There are also hardness results around the problem of *matrix sparsification* in numerical analysis, which seeks to replace a given matrix \mathbf{A} with a sparse matrix $\hat{\mathbf{A}}$ such that $\text{range}(\mathbf{A}) \approx \text{range}(\hat{\mathbf{A}})$: see McCormick [McCormick, 1983], Coleman and Pothén [Coleman and Pothén, 1986], and Gottlieb and Neylon [Gottlieb and Neylon,] for discussions of the hardness of this and related problems.

Exercises

2.1 (Convexity of ℓ^p norms). *Show that*

$$\|\mathbf{x}\|_p = \left(\sum_i |x_i|^p \right)^{1/p} \quad (2.5.1)$$

is convex for $p \geq 1$, and nonconvex for $0 < p < 1$.

2.2 (Relationship between ℓ^p norms). *Show that for $p < p'$,*

$$\|\mathbf{x}\|_p \geq \|\mathbf{x}\|_{p'} \quad (2.5.2)$$

for every \mathbf{x} . For what \mathbf{x} is equality obtained (i.e., $\|\mathbf{x}\|_p = \|\mathbf{x}\|_{p'}$)?

2.3 (Computing the Kruskal rank). *Write a Matlab function that takes as an input a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, and outputs the Kruskal rank $\text{krank}(\mathbf{A})$. There is no known way to efficiently compute the Kruskal rank. It is fine if your code takes time exponential in n . Corroborate the conclusion of Theorem 2.2.6, by generating a 4×8 Gaussian matrix \mathbf{A} , via $\mathbf{A} = \text{randn}(4,8)$, and computing its Kruskal rank.*

2.4 (A structured matrix with small Kruskal rank). *Consider a 4×8 dimensional complex matrix generated as*

$$\mathbf{A} = [\mathbf{I} \mid \mathbf{F}], \quad (2.5.3)$$

where \mathbf{I} is the 4×4 identity matrix, and \mathbf{F} is a 4×4 Discrete Fourier Transform (DFT) matrix: in Matlab, $\mathbf{A} = [\text{eye}(4), \text{dftmtx}(4)]$. Either using your code from Exercise 2.3, or hand calculations, determine the Kruskal rank of \mathbf{A} . You should find that it is smaller than 4! A general version of this phenomenon can be observed with the Dirac comb, which is sparse in both time and frequency.

2.5 (The spark). *Results on ℓ^0 uniqueness are sometimes described in terms of the spark of a matrix, which is the number of nonzero entries in the sparsest nonzero element of the null space of \mathbf{A} :*

$$\text{spark}(\mathbf{A}) = \min_{\mathbf{d} \neq \mathbf{0}, \mathbf{A}\mathbf{d}=\mathbf{0}} \|\mathbf{d}\|_0.$$

What is the relationship between $\text{spark}(\mathbf{A})$ and $\text{krank}(\mathbf{A})$?

2.6 (Kruskal rank of random matrices). *In this exercise we prove that for a generic $m \times n$ matrix $\mathbf{A} \sim_{\text{iid}} \mathcal{N}(0,1)$, $\text{krank}(\mathbf{A}) = 1$ with probability one.*

(i) Argue that for any $m \times n$ matrix \mathbf{A} , $\text{krank}(\mathbf{A}) \leq m$.

(ii) Let $\mathbf{A} = [\mathbf{a}_1 \mid \dots \mid \mathbf{a}_n]$. Let span denote the linear span of a collection of vectors. Argue that

$$\mathbb{P}[\mathbf{a}_m \in \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_{m-1})] = 0. \quad (2.5.4)$$

(iii) Argue that $\text{krank}(\mathbf{A}) < m$ if and only if there exist some indices i_1, \dots, i_m such that

$$\mathbf{a}_{i_m} \in \text{span}(\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_{m-1}}) \quad (2.5.5)$$

(iv) Conclude that $\text{krank}(\mathbf{A}) = m$ with probability one, by noting that

$$\begin{aligned}
 & \mathbb{P} [\exists i_1, \dots, i_m : \mathbf{a}_{i_m} \in \text{span}(\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_{m-1}})] \\
 & \leq \sum_{i_1, \dots, i_m} \mathbb{P} [\mathbf{a}_{i_m} \in \text{span}(\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_{m-1}})] \\
 & \leq m^n \times \underbrace{\mathbb{P} [\mathbf{a}_m \in \text{span}(\mathbf{a}_1, \dots, \mathbf{a}_{m-1})]}_{= 0} \\
 & = 0.
 \end{aligned}$$

2.7 (ℓ^0 minimization and typical examples). We showed that there is a worst case phase transition in ℓ^0 minimization at $\frac{\text{krank}(\mathbf{A})}{2}$. This means that ℓ^0 minimization recovers every \mathbf{x}_o satisfying $\|\mathbf{x}_o\|_0 < \frac{\text{krank}(\mathbf{A})}{2}$. We also know that for a Gaussian matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\text{krank}(\mathbf{A}) = m$, with probability one.

Using code for ℓ^0 minimization provided via CourseWorks (or write your own!), please do the following: generate a 5×12 Gaussian matrix $\mathbf{A} = \text{randn}(5, 12)$. What is $\text{rank}(\mathbf{A})$? Generate a sparse vector \mathbf{x}_o , with four nonzero entries, via $\mathbf{x}_o = \text{zeros}(12, 1); \mathbf{x}_o(1:4) = \text{randn}(4, 1)$. Now, set $\mathbf{y} = \mathbf{A} \mathbf{x}_o$. Solve the ℓ^0 minimization problem, to find the sparsest vector \mathbf{x} satisfying $\mathbf{A} \mathbf{x} = \mathbf{y}$. Is it the same as \mathbf{x}_o ? Check that $\text{norm}(\mathbf{x} - \mathbf{x}_o)$ is small, where \mathbf{x} is the solution produced by your code.

Notice that the worst case theory for ℓ^0 predicts that we can only recover vectors with at most 2 nonzero entries. But we have observed ℓ^0 succeeding with 4 nonzero entries! This is an example of a *typical case* performance which is better than the worst case.

Please explain this! Argue that if \mathbf{x}_o is a fixed vector supported on some set I of size $< m$, then the probability that there exists a subset $I' \neq I$ of size $< m$ satisfying $\mathbf{A} \mathbf{x}_o \in \text{range}(\mathbf{A}_{I'})$ is zero.

Does your argument imply that the worst case theory based on rank can be improved? Why or why not?

2.8 (Subdifferentials). Compute the subdifferentials for the following functions:

1. The subdifferential for $f(\mathbf{x}) = \|\mathbf{x}\|_\infty$ with $\mathbf{x} \in \mathbb{R}^n$.
2. The subdifferential for $f(\mathbf{X}) = \sum_{j=1}^n \|\mathbf{X} \mathbf{e}_j\|_2$ with \mathbf{X} a matrix in $\mathbb{R}^{n \times n}$.
3. The subdifferential for $f(\mathbf{x}) = \|\mathbf{X}\|_*$ with \mathbf{X} a matrix in $\mathbb{R}^{n \times n}$.

2.9 (Projection onto an affine subspace). In deriving the projected subgradient method for ℓ^1 minimization, we used the fact that for an affine subspace

$$C = \{\mathbf{x} \mid \mathbf{A} \mathbf{x} = \mathbf{y}\}, \quad (2.5.6)$$

where \mathbf{A} is a matrix with full row rank, and $\mathbf{y} \in \text{range}(\mathbf{A})$, the Euclidean projection on C is given by

$$\mathcal{P}_C[\mathbf{z}] = \arg \min_{\mathbf{Ax}=\mathbf{y}} \|\mathbf{x} - \mathbf{z}\|_2^2 \quad (2.5.7)$$

$$= \mathbf{z} - \mathbf{A}^* (\mathbf{AA}^*)^{-1} [\mathbf{Az} - \mathbf{y}]. \quad (2.5.8)$$

Prove that this formula is correct. You may use the following geometric characterization of $\mathcal{P}_C[\mathbf{z}]$: $\mathbf{x} = \mathcal{P}_C[\mathbf{z}]$ if and only if (i) $\mathbf{Ax} = \mathbf{y}$ and (ii) for any $\tilde{\mathbf{x}}$ satisfying $\mathbf{A}\tilde{\mathbf{x}} = \mathbf{y}$, we have

$$\langle \mathbf{z} - \mathbf{x}, \tilde{\mathbf{x}} - \mathbf{x} \rangle \leq 0. \quad (2.5.9)$$

2.10 (Sparse error correction). In coding theory and statistics, we often encounter the following situation: we have an observation \mathbf{z} , which should be expressible as \mathbf{Bx} , except that some of the entries are corrupted. We can express our corrupted observation as

$$\underset{\text{observation}}{\mathbf{z}} = \underset{\text{encoded message}}{\mathbf{Bx}} + \underset{\text{sparse corruption}}{\mathbf{e}} \quad (2.5.10)$$

Here $\mathbf{z} \in \mathbb{R}^n$ is the observation $\mathbf{x} \in \mathbb{R}^r$ is a message of interest; $\mathbf{B} \in \mathbb{R}^{n \times r}$ ($n > r$) is a tall matrix with full column rank r , and $\mathbf{e} \in \mathbb{R}^n$ represents any corruption of the message. In many applications, the observation may be subject to corruption which is large in magnitude, but affects only a few of the observations, i.e., \mathbf{e} is sparse vector. Let $\mathbf{A} \in \mathbb{R}^{(n-r) \times n}$ be a matrix whose rows span the left null space of \mathbf{B} , i.e., $\text{rank}(\mathbf{A}) = n - r$, and $\mathbf{AB} = \mathbf{0}$. Prove that for any k , (2.5.10) has a solution (\mathbf{x}, \mathbf{e}) with $\|\mathbf{e}\|_0 = k$ if and only if the underdetermined system

$$\mathbf{Ae} = \mathbf{Az} \quad (2.5.11)$$

has a solution \mathbf{e} with $\|\mathbf{e}\|_0 = k$. Argue that the optimization problems

$$\min_{\mathbf{x}} \|\mathbf{Bx} - \mathbf{z}\|_1 \quad (2.5.12)$$

and

$$\min_{\mathbf{e}} \|\mathbf{e}\|_1 \quad \text{subject to} \quad \mathbf{Ae} = \mathbf{Az} \quad (2.5.13)$$

are equivalent, in the sense that for every solution $\hat{\mathbf{x}}$ of (2.5.12), $\hat{\mathbf{e}} = \mathbf{B}\hat{\mathbf{x}} - \mathbf{z}$ is a solution to (2.5.13); and for every solution $\hat{\mathbf{e}}$ of (2.5.13), there is a solution $\hat{\mathbf{x}}$ of (2.5.12) such that $\hat{\mathbf{e}} = \mathbf{B}\hat{\mathbf{x}} - \mathbf{z}$.

It is sometimes observed that “sparse representation and sparse error correction are equivalent”. In what sense is this true?

2.11 (ℓ^1 vs. ℓ^p minimization). We have studied the ℓ^1 minimization problem

$$\min \|\mathbf{x}\|_1 \quad \text{subject to} \quad \mathbf{Ax} = \mathbf{y} \quad (2.5.14)$$

for recovering sparse \mathbf{x}_o . We can obtain other convex optimization problems by replace $\|\cdot\|_1$ with $\|\cdot\|_p$ for $p \in (1, \infty]$. For what kind of \mathbf{x}_o would you

expect ℓ^∞ minimization to outperform ℓ^1 minimization (in the sense of recovering \mathbf{x}_o more accurately)?

2.12 (Faces and linear subspaces). Download `face_intro_demo.zip` from the CourseWorks. Run `load_eyb_recognition` to load a collection of images under varying illumination into memory. The training images (under different lighting) will be stored in `A_train`, the identities of the subjects in `label_train`. Form a matrix `B` by selecting those columns of `A_train` that correspond to Subject 1. We will use the singular value decomposition to investigate how well-approximated the columns of `B` are by a linear subspace.

Compute the singular values of `B` using `sigma = svd(B)`. How many singular values r are needed to capture 95% of the energy of `B`? I.e., to ensure that

$$\sum_{i=1}^r \sigma_i^2 > .95 \times \sum_{i=1}^n \sigma_i^2 \quad ? \quad (2.5.15)$$

What about 99% of the energy? Repeat this calculation for several subjects.

2.13 (Sparsity of MR images).

2.14 (Sparsity of natural image patches).

2.A Complexity Classes and NP-Hardness

If you don't have any background in complexity theory, you can loosely think of the situation as follows. The problem class `P` consists of the problems that we can solve in time polynomial in the size of the problem. The problem class `NP` consists of those problems for which, if we are given a “certificate” describing the optimal solution, we can check that it is correct in polynomial time. That is, `P` contains problems for which *finding* the right answer is “easy”, while `NP` contains problems for which *checking* the right answer is easy. Anyone who has ever struggled with a problem for days, only to have a colleague or teacher easily demonstrate an obviously correct solution can appreciate the difference between finding the right answer and checking the right answer!

It turns out that amongst the `NP` problems, there are certain “`NP` complete” problems to which *every* problem in `NP` can be reduced, in polynomial time. So, solving one of these problems efficiently would enable you to solve every problem in `NP` efficiently! It is remarkable that this class of problems exists, and that it is quite large. It includes famous examples such as the Traveling Salesman Problem or the Multiway Cut Problem.

To understand the phrase “`NP` hard”, we have to appreciate one technicality regarding the above definitions of `P` and `NP`: they pertain only to *decision* problems, in which the goal is to produce a YES/NO answer. For

example, the decision version of the Traveling Salesman Problem asks: “Is it possible to visit all of the nodes of a given graph (cities) while traveling a distance at most d_* ?” The decision version of the ℓ^0 problem asks: “Does the system $\mathbf{y} = \mathbf{A}\mathbf{x}$ have a solution with at most k nonzero entries?”

Often in practice we care much more about *optimization problems* than *decision problems* – we do not just want to know whether a solution exists ... we want to know the way to find it! Optimization problems can’t be “NP-complete” – in the formal definition of NP, we only include decision problems. We call an optimization problem NP-*hard* if an efficient solution to that problem can be used to efficiently solve NP-complete problems. For example, the optimization version of the Traveling Salesman Problem asks: “Find the shortest path that visits all of the nodes in a given graph.” If I can solve this problem efficiently, I can clearly also solve the decision version efficiently, just by checking whether the optimal path that I found has length at most d_* .