

# Chapter 4

## Convex Methods for Low-Rank Matrix Recovery

In this chapter, we will branch out from sparse signals to a broader class of models: the low-rank matrices. Similar to the problem of recovering sparse signals, we consider how to recover a matrix  $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$  from linear measurements  $\mathbf{y} = \mathcal{A}[\mathbf{X}] \in \mathbb{R}^m$ . This problem can be phrased as searching for a solution  $\mathbf{X}$  to a linear system of equations

$$\mathcal{A} \left[ \begin{array}{c} \mathbf{X} \\ \text{unknown} \end{array} \right] = \begin{array}{c} \mathbf{y} \\ \text{observation} \end{array}. \quad (4.0.1)$$

Here,  $\mathcal{A} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^m$  is a linear map.

We will see that much of the mathematical structure in the sparse vector recovery problem carries over in a very natural way to this more general setting. In particular, in many interesting instances, we need to recover  $\mathbf{X}$  from far fewer measurements than the number of entries in the matrix, i.e.,  $m \ll n_1 \times n_2$ . Unless we can leverage some additional prior information about  $\mathbf{X}$ , the problem recovering  $\mathbf{X}$  from the linear measurements  $\mathbf{y}$  is ill-posed.

We will consider applications in which we can leverage the following powerful piece of structural information: the target matrix  $\mathbf{X}$  is *low-rank*. Recall that the rank of a matrix  $\mathbf{X}$  is the dimension of the linear subspace  $\text{col}(\mathbf{X})$  spanned by the columns of  $\mathbf{X}$ . If  $\mathbf{X} = [\mathbf{x}_1 \mid \cdots \mid \mathbf{x}_{n_2}] \in \mathbb{R}^{n_1 \times n_2}$  is a data matrix whose columns are  $n_1$ -dimensional vectors, then  $\text{rank}(\mathbf{X}) = r \ll n_1$  if and only if the columns of  $\mathbf{X}$  lie on an  $r$ -dimensional linear subspace of the data space  $\mathbb{R}^{n_1}$  – see Figure 4.1. Low-rank matrix recovery

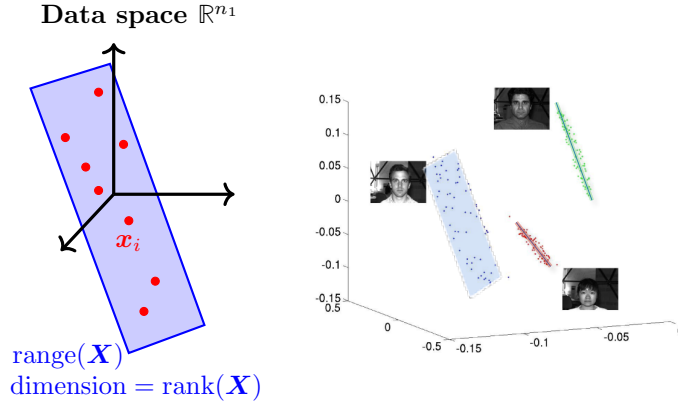


Figure 4.1. **Low-rank data matrices.** If matrix  $\mathbf{X}$  with columns  $\mathbf{x}_1, \dots, \mathbf{x}_{n_2}$  has rank  $r$ , its columns lie on an  $r$ -dimensional subspace  $\text{range}(\mathbf{X})$ . Many naturally occurring data matrices approximately satisfy this property. Right: low-dimensional approximations to images of faces under different lighting conditions.

problems arise in a broad range of application areas. We sketch a few of these below.

## 4.1 Motivating Examples of Low-Rank Modeling

### 4.1.1 Latent Semantic Analysis

Low-dimensional models are very popular in document analysis. Consider an idealized problem in search or document retrieval. The system has access to  $n_2$  documents (say, news articles), each of which is viewed as a collection of words in a dictionary of size  $n_1$ . For the  $j$ -th document, we compute a histogram of word occurrences, giving an  $n_1$ -dimensional vector  $\mathbf{y}_j$  whose  $i$ -th entry is the fraction of occurrences of word  $i$  in document  $j$ . Set

$$\underset{\text{Word occurrences}}{\mathbf{Y}} = \underset{\text{Documents}}{\text{Words}} \left[ \mathbf{y}_1 \mid \dots \mid \mathbf{y}_{n_2} \right].$$

We model these observations as follows. We imagine that there exist a set of “topics”  $\mathbf{t}_1 \dots \mathbf{t}_r$ . Each topic is a probability distribution on  $\{1, 2, \dots, n_1\}$ . We may imagine that the  $\mathbf{t}_i$  correspond loosely with our informal notion of what a topic is – say, architecture or New York city. An article on architecture in New York would involve multiple topics. We

$$\begin{array}{c} \text{Users} \end{array} \begin{bmatrix} 5 & 3 & \dots & ? \\ ? & 2 & \dots & 4 \\ \vdots & \vdots & \ddots & \vdots \\ 5 & ? & \dots & ? \end{bmatrix} = \mathcal{P}_{\Omega} \left( \begin{bmatrix} 5 & 3 & \dots & 5 \\ 4 & 2 & \dots & 4 \\ \vdots & \vdots & \ddots & \vdots \\ 5 & 5 & \dots & 3 \end{bmatrix} \right)$$

Items  
Observed (Incomplete) Ratings  $\mathbf{Y}$ 
Complete Ratings  $\mathbf{X}$

Figure 4.2. **Collaborative filtering as a matrix completion problem.** Consider a universe of  $n_1$  users and  $n_2$  items. Users experience items, and then rate their experience. Our observation  $\mathbf{Y}$  consists of those ratings that user have provided:  $Y_{ij}$  is user  $i$ 's rating of item  $j$ . We wish to predict users ratings of items that they have not yet rated. This can be viewed as attempting to recover a large matrix  $\mathbf{X}$  from a subset  $\mathbf{Y} = \mathcal{P}_{\Omega}(\mathbf{X})$  of its entries.

model this as a mixture distribution, writing

$$\begin{array}{c} \mathbf{p}_j \\ \text{Word distribution for document } j \end{array} = \sum_{l=1}^r \begin{array}{c} t_l \\ \text{topic abundance} \end{array} \alpha_{l,j},$$

where  $\alpha_{1,j} + \alpha_{2,j} + \dots + \alpha_{r,j} = 1$ . We imagine that  $\mathbf{y}_j$  is generated by sampling words independently at random from the mixture distribution  $\mathbf{p}_j$  and computing a histogram.<sup>1</sup> If the number of words sampled is large, we can imagine  $\mathbf{y}_j \approx \mathbf{p}_j$ . So, if we write  $\mathbf{T} = [\mathbf{t}_1 \dots \mathbf{t}_r]$  and  $\mathbf{A} = [\alpha_1 \dots \alpha_n]$ , then we have

$$\begin{array}{c} \mathbf{Y} \\ \text{Word occurrences} \end{array} \approx \begin{array}{c} \mathbf{T} \\ \text{Topics} \end{array} \begin{array}{c} \mathbf{A} \\ \text{Abundances} \end{array} \quad (4.1.1)$$

Notice that  $\text{rank}(\mathbf{T}\mathbf{A}) \leq r$ : *the rank is bounded by the number of topics*. *Latent semantic analysis* computes a best low-rank approximation to  $\mathbf{Y}$  and then uses it for search and indexing [Deerwester et al., 1990].

#### 4.1.2 Recommendation Systems

Here, we imagine that we have  $n_2$  products of interest, and  $n_1$  users. Users consume products and rate them based on the quality of their experience. Our goal is to use the information of all the users' ratings to predict which

<sup>1</sup>In practice, researchers have observed that more complicated methods of constructing  $\mathbf{Y}$  (say, using the *TF-IDF* weighting) improves performance compared to just using the histogram.

products will appeal to a given user. Formally, our object of interest is a large, unknown matrix

$$\mathbf{X} \in \mathbb{R}^{n_1 \times n_2},$$

whose  $(i, j)$  entry contains user  $i$ 's degree of preference for item  $j$ . If we let

$$\Omega = \{(i, j) \mid \text{user } i \text{ has rated product } j\},$$

then we observe

$$\begin{array}{c} \mathbf{Y} \\ \text{Observed ratings} \end{array} = \mathcal{P}_\Omega \begin{array}{c} \mathbf{X} \\ \text{Complete ratings} \end{array}.$$

Here,  $\mathcal{P}_\Omega$  is the projection operator onto the subset  $\Omega$ :

$$\mathcal{P}_\Omega[\mathbf{X}](i, j) = \begin{cases} X_{ij} & (i, j) \in \Omega, \\ 0 & \text{else.} \end{cases}$$

See Figure 4.2 for a schematic representation of this scenario. Our goal is to fill in the missing entries of  $\mathbf{X}$ . This problem is encountered in on-line recommendation systems – the most famous recent instance being the “Netflix Prize” competition conducted between 2006 and 2009. Obviously, with no additional assumptions, the problem of filling in the missing entries of  $\mathbf{X}$  is ill-posed. One popular assumption is that the ratings of distinct users (or distinct products) are correlated, and hence the target matrix  $\mathbf{X}$  is low-rank, or approximately so. The relevant mathematical problem then becomes filling in the missing entries of a low-rank matrix, or, somewhat equivalently, looking for the matrix  $\mathbf{X}$  of minimum rank that is consistent with our given observations:

$$\begin{aligned} & \text{minimize} && \text{rank}(\mathbf{X}), \\ & \text{subject to} && \mathcal{P}_\Omega[\mathbf{X}] = \mathbf{Y}. \end{aligned} \tag{4.1.2}$$

This problem is often referred to as *matrix completion* [Candes and Recht, 2009].

#### 4.1.3 3D Shape from Photometric Measurements.

There are also situations in which low-rank models arise due to the physical processes that generate the data. For example, in computer vision, rank constraints arise in a number of problems in reconstructing the three-dimensional shape of a scene from two-dimensional measurements (images). In *photometric stereo* [Woodham, 1980], we obtain images  $\mathbf{y}_1, \dots, \mathbf{y}_{n_2} \in \mathbb{R}^{n_1}$  of an object illuminated by different distant point light sources. Write  $\mathbf{Y} = [\mathbf{y}_1 \mid \dots \mid \mathbf{y}_{n_2}] \in \mathbb{R}^{n_1 \times n_2}$ . Let  $\mathbf{l}_1, \dots, \mathbf{l}_{n_2} \in \mathbb{S}^2$  denote the directions of these light sources. The *Lambertian model* for reflectance models the reflected light intensity as

$$Y_{ij} = \alpha_i [\langle \boldsymbol{\nu}_i, \mathbf{l}_j \rangle]_+,$$

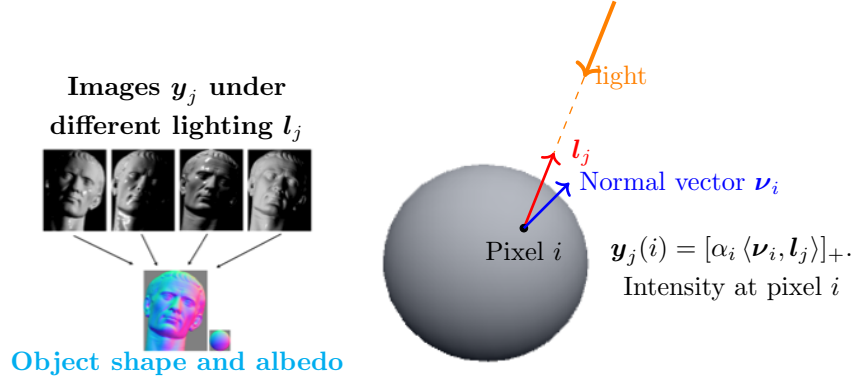


Figure 4.3. **Photometric stereo as low-rank recovery.** Photometric stereo (left) seeks to recover object shape from images taken under different illuminations. Under a diffuse reflective (Lambertian) model (right), this leads directly to a low-rank recovery problem.

where  $\boldsymbol{\nu}_i \in \mathbb{S}^2$  is the surface normal at the  $i$ -th pixel,  $\alpha_i$  is a nonnegative scalar known as the *albedo*, and  $[\cdot]_+$  takes the positive part of its argument. This model is appropriate for matte objects. See Figure 4.3 for a visualization of this model.

Under this model, if we let

$$\mathbf{N} = \begin{bmatrix} \alpha_1 \boldsymbol{\nu}_1^* \\ \vdots \\ \alpha_m \boldsymbol{\nu}_m^* \end{bmatrix} \in \mathbb{R}^{n_1 \times 3}, \quad \text{and} \quad \mathbf{L} = [\mathbf{l}_1 \mid \cdots \mid \mathbf{l}_n] \in \mathbb{R}^{3 \times n_2},$$

then we have

$$\mathbf{Y} = \mathcal{P}_\Omega[\mathbf{NL}],$$

where  $\Omega = \{(i, j) \mid \langle \boldsymbol{\nu}_i, \mathbf{l}_j \rangle \geq 0\}$ . If we can recover the low-rank matrix  $\mathbf{X} = \mathbf{NL}$ , we can then recover information about the shape and reflectance of the object. Again, a useful heuristic is to look for a solution of minimum rank consistent with the observations [Wu et al., 2010]:

$$\begin{aligned} & \text{minimize} && \text{rank}(\mathbf{X}), \\ & \text{subject to} && \mathcal{P}_\Omega[\mathbf{X}] = \mathbf{Y}. \end{aligned} \tag{4.1.3}$$

The reader can obtain an open-source implementation of this example from: <https://github.com/yasumat/RobustPhotometricStereo>. More detailed discussion will be covered in Part III, Chapter 13.

#### 4.1.4 Euclidean Distance Matrix Embedding

The abstract problem can be stated as follows: first imagine that we have  $n$  points  $\mathbf{X} = [\mathbf{x}_1 \mid \cdots \mid \mathbf{x}_n]$  living in  $\mathbb{R}^d$ . We can define a matrix  $\mathbf{D}$  via

$$D_{ij} = d^2(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2.$$

Such a  $\mathbf{D}$  is known as a *Euclidean distance matrix*. Now imagine the following scenario: rather than observing the  $\mathbf{x}_i$  themselves, we instead see their pairwise distances  $d(\mathbf{x}_i, \mathbf{x}_j)$ . How can we tell if these distances were generated by some configuration of points living in  $\mathbb{R}^d$ ? A necessary and sufficient condition is given by the following classical result:

**Theorem 4.1.1** (Schoenberg).  *$\mathbf{D} \in \mathbb{R}^{n \times n}$  is a Euclidean distance matrix for some set of  $n$  points in  $\mathbb{R}^d$  if and only if the following conditions hold:*

- $\mathbf{D}$  is symmetric.
- $D_{ii} = 0$  for all  $i$ .
- $\Phi \mathbf{D} \Phi^* \preceq \mathbf{0}$ , where  $\Phi = \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^*$  is the centering matrix (here  $\mathbf{1}$  is the vector whose entries are all one).
- $\text{rank}(\Phi \mathbf{D} \Phi^*) \leq d$ .

Now imagine we only know  $D_{ij}$  for some subset  $\Omega \subset \{1 \dots n\} \times \{1 \dots n\}$ , i.e., we observe  $\mathbf{Y} = \mathcal{P}_\Omega[\mathbf{D}]$ . We can cast the problem of looking for a Euclidean distance matrix that agrees with our observations as a rank minimization problem:

$$\begin{aligned} & \text{minimize} && \text{rank}(\Phi \mathbf{D} \Phi^*), \\ & \text{subject to} && \Phi \mathbf{D} \Phi^* \preceq \mathbf{0}, \mathbf{D} = \mathbf{D}^*, \mathcal{P}_\Omega[\mathbf{D}] = \mathbf{Y}, \forall i D_{ii} = 0. \end{aligned} \tag{4.1.4}$$

Many additional examples arise, for example in solving positioning problems, problems in system identification, quantum state tomography, image and video alignment, ect. We will survey more of these in the coming application chapters.

## 4.2 Representing Low-rank Matrix via SVD

In all of the applications described above, our goal is to recover an unknown  $\mathbf{X}$  whose columns live on an  $r$ -dimensional linear subspace of the data space  $\mathbb{R}^{n_1}$ . This subspace can be characterized via the *singular value decomposition* (SVD) of  $\mathbf{X}$  (see Appendix A.8 for a more detailed review):

**Theorem 4.2.1.** *Let  $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$  be a matrix, and  $r = \text{rank}(\mathbf{X})$ . Then there exist numbers  $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > 0$  and matrices  $\mathbf{U} \in \mathbb{R}^{n_1 \times r}$ ,  $\mathbf{V} \in \mathbb{R}^{n_2 \times r}$ , such that  $\mathbf{U}^* \mathbf{U} = \mathbf{I}$ ,  $\mathbf{V}^* \mathbf{V} = \mathbf{I}$  and if we set  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r)$ ,*

and

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^* = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^*. \quad (4.2.1)$$

Exercise 4.2 gives a guided proof of this result. This construction turns out to be a very versatile tool both for theory and for numerical computation. The *full* singular value decomposition extends the matrices  $\mathbf{U}$  and  $\mathbf{V}$  to complete orthonormal bases for  $\mathbb{R}^{n_1}$  and  $\mathbb{R}^{n_2}$ , respectively, by adding bases for the left and right null spaces of  $\mathbf{X}$ :

**Theorem 4.2.2.** *Let  $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$  be a matrix. Then there exist orthogonal matrices  $\mathbf{U} \in O(n_1)$  and  $\mathbf{V} \in O(n_2)$ , and numbers*

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{\min\{n_1, n_2\}}$$

*such that if we let  $\mathbf{\Sigma} \in \mathbb{R}^{n_1 \times n_2}$  with  $\Sigma_{ii} = \sigma_i$  and  $\Sigma_{ij} = 0$  for  $i \neq j$ ,*

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*. \quad (4.2.2)$$

**Fact 4.2.3** (Properties of the SVD). *We note the following properties of the construction in Theorem A.8.1:*

- *The right singular vectors  $\mathbf{v}_i$  are the eigenvectors of  $\mathbf{X}^* \mathbf{X}$ .*
- *The nonzero singular values  $\sigma_i$  are the positive square roots of the positive eigenvalues  $\lambda_i$  of  $\mathbf{X}^* \mathbf{X}$ .*
- *The left singular vectors  $\mathbf{u}_i$  are the eigenvectors of  $\mathbf{X} \mathbf{X}^*$  (check this!).*
- *The nonzero singular values  $\sigma_i$  are also the positive square roots of the positive eigenvalues  $\lambda_i$  of  $\mathbf{X} \mathbf{X}^*$ .*

Notice that since  $\mathbf{U}$  and  $\mathbf{V}$  are nonsingular, the  $\text{rank}(\mathbf{X}) = \text{rank}(\mathbf{\Sigma})$ . Since  $\mathbf{\Sigma}$  is diagonal, this quantity is especially simple – it is simply the number of nonzero entries  $\sigma_i$ ! Here, and below, we will let  $\boldsymbol{\sigma}(\mathbf{X}) = (\sigma_1, \dots, \sigma_{\min\{n_1, n_2\}}) \in \mathbb{R}^{\min\{n_1, n_2\}}$  denote the vector of singular values of  $\mathbf{X}$ . Then, in the language that we’ve been developing thus far,

$$\text{rank}(\mathbf{X}) = \|\boldsymbol{\sigma}(\mathbf{X})\|_0. \quad (4.2.3)$$

Hence any problem that minimizes the rank of an unknown matrix  $\mathbf{X}$  is essentially minimizing the number of nonzero singular values of  $\mathbf{X}$  – the “sparsity” of singular values, subject to data constraints.

#### 4.2.1 Singular Vectors via Nonconvex Optimization

The SVD can be computed in time  $O(\max\{n_1, n_2\} \min\{n_1, n_2\}^2)$ . The first  $r$  singular value/vector triples can be computed in time  $O(n_1 n_2 r)$ . Hence, the problem of finding a linear subspace that best fits a given set of data can be solved in polynomial time. On the surface this is quite remarkable – the

problem of computing singular vectors is nonconvex. We briefly describe why this nonconvex problem can be solved globally in an efficient manner.

We give a brief indication of *why* it is possible to efficiently compute singular vectors of a matrix  $\mathbf{Y}$ . Consider the matrix  $\mathbf{\Gamma} \doteq \mathbf{Y}\mathbf{Y}^*$ . Let  $\mathbf{\Gamma} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$  be the eigenvalue decomposition of  $\mathbf{\Gamma}$  and  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_{n_1})$  be the eigenvalues. It is obvious that the left singular vectors  $\mathbf{u}_i$  of  $\mathbf{Y}$  are the eigenvectors of  $\mathbf{\Gamma}$ . Because our goal in this paragraph is merely to convey intuition, we make the simplifying assumption that  $\mathbf{\Gamma}$  has no repeated eigenvalues and  $\lambda_1$  is the largest. We show how to use nonconvex optimization to compute the leading eigenvector  $\mathbf{u}_1$  – see Exercise 4.5 for extensions.

Consider the optimization problem

$$\begin{aligned} & \text{minimize} && \varphi(\mathbf{q}) \equiv -\frac{1}{2}\mathbf{q}^*\mathbf{\Gamma}\mathbf{q}, \\ & \text{subject to} && \|\mathbf{q}\|_2^2 = 1. \end{aligned} \quad (4.2.4)$$

The gradient and Hessian of the function  $\varphi(\mathbf{q})$  are

$$\nabla\varphi(\mathbf{q}) = -\mathbf{\Gamma}\mathbf{q} \quad \text{and} \quad \nabla^2\varphi(\mathbf{q}) = -\mathbf{\Gamma}, \quad (4.2.5)$$

respectively. A point  $\mathbf{q}$  is a *critical point* of the function  $\varphi$  over  $\mathbb{S}^{n-1} = \{\mathbf{q} \mid \|\mathbf{q}\|_2^2 = 1\}$  if there is no direction  $\mathbf{v} \perp \mathbf{q}$  (i.e., no direction that is tangent to the sphere at  $\mathbf{q}$ ) along which the function decreases. Equivalently,  $\mathbf{q}$  is a critical point of  $\varphi$  over the sphere if and only if the gradient is proportional to  $\mathbf{q}$ :

$$\nabla\varphi(\mathbf{q}) \propto \mathbf{q}. \quad (4.2.6)$$

Using our expression for  $\nabla\varphi$ , this is true if and only if  $\mathbf{\Gamma}\mathbf{q} = \lambda\mathbf{q}$  for some  $\lambda$ : *The critical points of  $\varphi$  over  $\mathbb{S}^{n-1}$  are precisely the eigenvectors  $\pm\mathbf{u}_i$  of  $\mathbf{\Gamma}$ .*

Which critical points  $\pm\mathbf{u}_i$  are actual local or global minimizers (instead of saddle points)? To answer this question, we need to study the curvature of the function  $\varphi(\mathbf{q})$  around a critical point  $\bar{\mathbf{q}}$ . In Euclidean space, the correct tool for studying curvature is the Hessian, as is justified by the second order Taylor expansion of the function along the curve  $\mathbf{x}(t) = \mathbf{x} + t\mathbf{v}$ :

$$f(\mathbf{x} + t\mathbf{v}) = f(\mathbf{x}) + t\langle \nabla f(\mathbf{x}), \mathbf{v} \rangle + \frac{1}{2}t^2\mathbf{v}^*\nabla^2 f(\mathbf{x})\mathbf{v} + o(t^2).$$

$= 0 \text{ at any critical point}$

In Euclidean space, a critical point  $\bar{\mathbf{x}}$  is a local minimizer if  $\nabla^2 f(\bar{\mathbf{x}}) \succ \mathbf{0}$ . Conversely, if  $\nabla^2 f(\bar{\mathbf{x}})$  has a negative eigenvalue, the point is not a local minimizer.



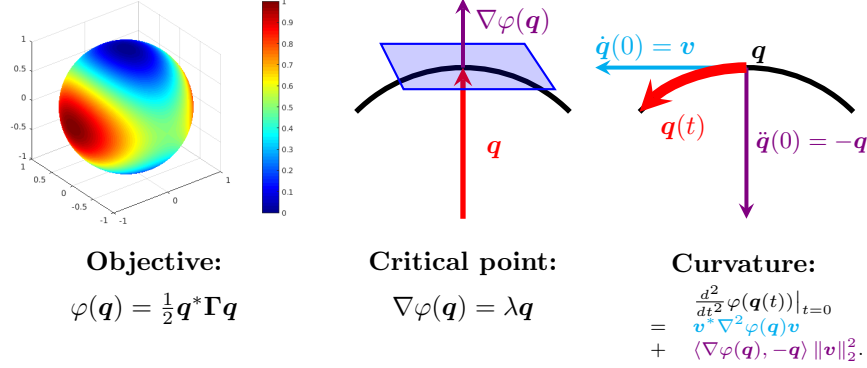


Figure 4.4. **Eigenvector Computation as Nonconvex Optimization over the Sphere.** We plot  $\varphi(\mathbf{q}) = -\frac{1}{2} \mathbf{q}^* \mathbf{\Gamma} \mathbf{q}$  over the sphere, for one particular  $\mathbf{\Gamma}$ . Critical points (middle) are points  $\mathbf{q}$  for which  $\nabla \varphi(\mathbf{q})$  is proportional to  $\mathbf{q}$ . Every critical point is an eigenvector of  $\mathbf{\Gamma}$ ; the only local minimizers are eigenvectors that correspond to the largest eigenvalue  $\lambda_1(\mathbf{\Gamma})$ . Right: curvature of the  $\varphi$  over the sphere comes from both curvature  $\nabla^2 \varphi$  of  $\varphi$  and the curvature of the sphere.

Over the sphere, we can perform a similar Taylor expansion, but we need to replace the straight line  $\mathbf{x}(t) = \mathbf{x} + t\mathbf{v}$  with a great circle<sup>2</sup>

$$\mathbf{q}(t) = \mathbf{q} \cos(t) + \mathbf{v} \sin(t), \quad (4.2.7)$$

where  $\mathbf{v} \perp \mathbf{q}$  and  $\|\mathbf{v}\|_2 = 1$ . Calculus shows that the second derivative of  $\varphi(\mathbf{q}(t))$  is given by

$$\frac{d^2}{dt^2} \varphi(\mathbf{q}(t)) \Big|_{t=0} = \underbrace{\mathbf{v}^* \nabla^2 \varphi(\mathbf{q}) \mathbf{v}}_{\text{Curvature of } \varphi} - \underbrace{\langle \nabla \varphi(\mathbf{q}), \mathbf{q} \rangle}_{\text{Curvature of the sphere}} \mathbf{v}^* \mathbf{v}. \quad (4.2.8)$$

This formula contains two terms, which combine the usual hessian of  $\varphi$  (accounting for the curvature of  $\varphi$ ) and a second correction term involving  $-\langle \nabla \varphi(\mathbf{q}), \mathbf{q} \rangle$  which accounts for the fact that the curve  $\mathbf{q}(t)$  curves in the  $-\mathbf{q}$  direction in order to stay on the sphere.

Noting that  $\nabla^2 \varphi(\mathbf{q}) = -\mathbf{\Gamma}$ . So we have  $\langle \nabla \varphi(\mathbf{u}_i), \mathbf{u}_i \rangle = -\mathbf{u}_i^* \mathbf{\Gamma} \mathbf{u}_i = -\lambda_i$ , we observe that at a critical point  $\bar{\mathbf{q}} = \pm \mathbf{u}_i$ , the second derivative in the  $\mathbf{v}$  direction is

$$\frac{d^2}{dt^2} \varphi(\mathbf{q}(t)) \Big|_{t=0} = \mathbf{v}^* (-\mathbf{\Gamma} + \lambda_i \mathbf{I}) \mathbf{v}. \quad (4.2.9)$$

The eigenvalues of the operator  $-\mathbf{\Gamma} + \lambda_i$  take the form  $-\lambda_j + \lambda_i$ ; there is a strictly negative eigenvalue if  $\mathbf{u}_i$  is an eigenvector that does not correspond to  $\lambda_1$ :  $\pm \mathbf{u}_1$  are the only local minimizers of  $\varphi$ . All other critical points have

<sup>2</sup>Curves of this form are *geodesics* on  $\mathbb{S}^{n-1}$ .

a direction of strict negative curvature. This benign geometry implies that a simple projected gradient method converges to a global optimizer from almost any initialization. Exercise 4.6 gives a more algorithmic analysis, using the power method. We will return more formally to questions about global optimization of nonconvex functions in Chapter 9 in the context of learning sparse models for data. For now, we take these observations as an intuitive indication of why the SVD is amenable to efficient computation.

*Implications and history.*

Whichever rationale we adopt, the fact that the SVD is both optimal (in a precisely defined and often quite relevant sense) and efficient (at least for moderate problems) makes it a very useful element in the numerical computing toolbox. The canonical example application of the SVD is *Principal Component Analysis* (PCA). Outlined in 1901 and 1933 papers by Pearson and Hotelling [Pearson, 1901, Hotelling, 1933], respectively, PCA finds a best-fitting low-dimensional subspace, which can be computed via the SVD, as suggested by Theorem 4.2.4 below. Remarkably, Pearson’s 1901 paper asserts that PCA is “well-suited to numerical computation” – meaning hand calculations!

#### 4.2.2 Best Low-Rank Matrix Approximation

We are interested in recovering a low-rank matrix that is consistent with certain linear observations. Because the rank has a similar characteristic to the  $\ell^0$  norm, one should expect that these problems would be computationally difficult in general, as in the case with recovering a sparse solution (see Theorem 2.2.8).

Remarkably, there *are* however a few special instances of rank minimization that we *can* solve efficiently, with virtually no assumptions on the input. The most important is the *best rank- $r$  approximation* problem, in which we try to approximate an arbitrary input matrix  $\mathbf{Y}$  with a matrix  $\mathbf{X}$  of rank at most  $r$  such that the approximation error  $\|\mathbf{X} - \mathbf{Y}\|_F$  is as small as possible. The optimal solution to this problem can be obtained by simply retaining the first  $r$  leading singular values/vectors of  $\mathbf{Y}$ :

**Theorem 4.2.4.** *Let  $\mathbf{Y} \in \mathbb{R}^{n_1 \times n_2}$ , and consider the following optimization problem*

$$\begin{aligned} & \text{minimize} && \|\mathbf{X} - \mathbf{Y}\|_F, \\ & \text{subject to} && \text{rank}(\mathbf{X}) \leq r. \end{aligned} \tag{4.2.10}$$

*Every optimal solution  $\hat{\mathbf{X}}$  to the above problem has the form  $\hat{\mathbf{X}} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^*$ , where  $\mathbf{Y} = \sum_{i=1}^{\min(n_1, n_2)} \sigma_i \mathbf{u}_i \mathbf{v}_i^*$  is a (full) singular value decomposition of  $\mathbf{Y}$ .*

In fact, the same solution (truncating the SVD) also solves the low-rank approximation problem when the error is measured in the operator norm,

or any other orthogonal-invariant matrix norm. Please see Exercise 4.3 for guidance on how to prove Theorem 4.2.4.

The problem (4.2.10) can be turned around and cast as one of minimizing the rank of the unknown matrix, subject to a data fidelity constraint:

$$\begin{aligned} & \text{minimize} && \text{rank}(\mathbf{X}), \\ & \text{subject to} && \|\mathbf{X} - \mathbf{Y}\|_F \leq \varepsilon. \end{aligned} \quad (4.2.11)$$

This is an example of a *matrix rank minimization* problem – we seek a matrix of minimum rank that is consistent with some given observations. Because of its very special nature, this particular rank minimization can be solved optimally via the SVD. We leave the solution to this problem as an exercise to the reader (see Exercise 4.4).<sup>3</sup>

## 4.3 Recovering a Low-Rank Matrix

### 4.3.1 General Rank Minimization Problems

In the previous section, we saw that for certain very specific rank minimization problems, globally optimal solutions could be obtained using efficient algorithms based on the singular value decomposition. However, all of the applications discussed above (and many others!) force us to attempt to minimize the rank of  $\mathbf{X}$  over much more complicated sets. One model example problem is the *affine rank minimization* [Fazel et al., 2004] problem

$$\begin{aligned} & \text{minimize} && \text{rank}(\mathbf{X}), \\ & \text{subject to} && \mathcal{A}[\mathbf{X}] = \mathbf{y}. \end{aligned} \quad (4.3.1)$$

Here  $\mathbf{y} \in \mathbb{R}^m$  is an observation, and  $\mathcal{A} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^m$  is a linear map. When  $m \ll n_1 n_2$ , the linear system of equations  $\mathcal{A}[\mathbf{X}] = \mathbf{y}$  is underdetermined. The notion of a linear map  $\mathcal{A}$  from  $n_1 \times n_2$  matrices to  $m$ -dimensional vectors may seem somewhat abstract. Any linear map of this form can be represented using the matrix inner product<sup>4</sup>:

$$\mathcal{A}[\mathbf{X}] = (\langle \mathbf{A}_1, \mathbf{X} \rangle, \dots, \langle \mathbf{A}_m, \mathbf{X} \rangle). \quad (4.3.2)$$

Here, the set of matrices  $\mathbf{A}_1, \dots, \mathbf{A}_m \in \mathbb{R}^{n_1 \times n_2}$  define our “measurements”  $\mathbf{y}$ , through their inner products with the unknown matrix  $\mathbf{X}$ .<sup>5</sup>

A mathematically simple and natural assumption on these “measurement” matrices is that they are i.i.d. Gaussian matrices. Such an

<sup>3</sup>Hint: you may first try to guess what the optimal solution is and then show its optimality.

<sup>4</sup>Recall that the standard inner product between matrices  $\mathbf{P}, \mathbf{Q} \in \mathbb{R}^{n_1 \times n_2}$  is defined by  $\langle \mathbf{P}, \mathbf{Q} \rangle = \sum_{ij} \mathbf{P}_{ij} \mathbf{Q}_{ij} = \text{trace}[\mathbf{Q}^* \mathbf{P}]$ .

<sup>5</sup>You can think of the measurements  $\mathbf{A}_i$  as analogous to the *rows* of the matrix  $\mathbf{A}$  in the equation  $\mathbf{y} = \mathbf{A}\mathbf{x}$  studied in Chapters 2-3.

assumption will allow us to understand the conditions under which one could expect to recover a low-rank matrix with generic measurements. Hence, our first attempt to understand the low-rank recovery problem will rely on such a simplifying assumption. However, in many practical problems of interest, the operator  $\mathcal{A}$  has particular structures that make it behave differently. As a concrete example, in the matrix completion problems discussed above, we would have  $p = |\Omega|$ , and  $\mathbf{A}_l = \mathbf{e}_{i_l} \mathbf{e}_{j_l}^*$ , with  $\Omega = \{(i_1, j_1), \dots, (i_m, j_m)\}$ . We will also thoroughly analyze this important special case and provide conditions under which the recovery can be successful.

*Connection to  $\ell^0$ , NP-hardness.*

To make the connection to sparse recovery explicit, using the observation that  $\text{rank}(\mathbf{X}) = \|\boldsymbol{\sigma}(\mathbf{X})\|_0$ , we can rewrite the affine rank minimization problem as

$$\begin{aligned} & \text{minimize} && \|\boldsymbol{\sigma}(\mathbf{X})\|_0 \\ & \text{subject to} && \mathcal{A}[\mathbf{X}] = \mathbf{y}. \end{aligned} \quad (4.3.3)$$

Moreover, if  $\mathbf{X}$  is a diagonal matrix, then  $\text{rank}(\mathbf{X}) = \|\mathbf{X}\|_0$ . So, every  $\ell^0$  minimization problem can be converted into a rank minimization problem with a diagonal constraint. This means that in the worst case, the rank minimization problem is at least as hard as the  $\ell^0$  minimization problem: it is NP-hard (as shown in Theorem 2.2.8).

As was the case for  $\ell^0$  minimization, we could simply give up here in searching for tractable algorithms. However, given the close analogy between rank minimization and  $\ell^0$  minimization, we might hope that there could be some fairly broad class of “nice enough” problems that we *can* solve efficiently.

#### 4.3.2 Convex Relaxation of Rank Minimization

The close analogy to  $\ell^0$  minimization suggests a natural strategy: replace the rank, which is the  $\ell^0$  norm  $\boldsymbol{\sigma}(\mathbf{X})$  with the  $\ell^1$  norm of  $\boldsymbol{\sigma}(\mathbf{X})$ :

$$\|\boldsymbol{\sigma}(\mathbf{X})\|_1 = \sum_i \sigma_i(\mathbf{X}). \quad (4.3.4)$$

We call this function the *nuclear norm* of  $\mathbf{X}$ , and reserve the special notation

$$\|\mathbf{X}\|_* = \sum_i \sigma_i(\mathbf{X}). \quad (4.3.5)$$

When  $\mathbf{X}$  is a symmetric positive semidefinite matrix,  $\mathbf{X}$  has real nonnegative eigenvalues, and  $\sigma_i(\mathbf{X}) = \lambda_i(\mathbf{X})$ . Since  $\sum_i \lambda_i(\mathbf{X}) = \text{trace}(\mathbf{X})$ , in the special case when  $\mathbf{X}$  is semidefinite,  $\|\mathbf{X}\|_* = \text{trace}[\mathbf{X}]$ . For this reason, the nuclear norm is sometimes also referred to as the *trace norm*. Other names

in various literatures include the *Schatten 1-norm* and *Ky-Fan  $m$ -norm*.<sup>6</sup> When  $\mathbf{X}$  is not a semidefinite matrix, the function  $\|\mathbf{X}\|_*$  depends on the entries in a very complicated way. It may not be at all obvious that this function is a norm, or even a convex function.

To allay any suspicion, we give a quick proof that  $\|\cdot\|_*$  is indeed a norm:

**Theorem 4.3.1.** *For  $\mathbf{M} \in \mathbb{R}^{n_1 \times n_2}$ , let  $\|\mathbf{M}\|_* = \sum_{i=1}^{\min\{n_1, n_2\}} \sigma_i(\mathbf{M})$ . Then  $\|\cdot\|_*$  is a norm. Moreover, the nuclear norm and the  $\ell^2$  operator norm are dual norms:*

$$\|\mathbf{M}\|_* = \sup_{\|\mathbf{N}\| \leq 1} \langle \mathbf{M}, \mathbf{N} \rangle, \quad \text{and} \quad \|\mathbf{M}\| = \sup_{\|\mathbf{N}\|_* \leq 1} \langle \mathbf{M}, \mathbf{N} \rangle. \quad (4.3.6)$$

*Proof.* We begin by proving the first equality in (4.3.6). Let

$$\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^* \quad (4.3.7)$$

be a full singular value decomposition of  $\mathbf{M}$ , with  $\mathbf{U} \in O(n_1)$ ,  $\mathbf{V} \in O(n_2)$ , and  $\mathbf{\Sigma} \in \mathbb{R}^{n_1 \times n_2}$ , and note that

$$\begin{aligned} \sup_{\|\mathbf{N}\| \leq 1} \langle \mathbf{N}, \mathbf{M} \rangle &= \sup_{\|\mathbf{N}\| \leq 1} \langle \mathbf{N}, \mathbf{U}\mathbf{\Sigma}\mathbf{V}^* \rangle \\ &= \sup_{\|\mathbf{N}\| \leq 1} \left\langle \mathbf{U}^* \mathbf{N} \mathbf{V}, \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_{n_2} & \\ 0 & 0 & 0 & \\ & \vdots & & \end{bmatrix} \right\rangle \\ &\geq \sum_{i=1}^{n_2} \sigma_i, \end{aligned} \quad (4.3.8)$$

where the last line follows by making a particular choice

$$\mathbf{N} = \mathbf{U} \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ 0 & 0 & 0 & \\ & \vdots & & \end{bmatrix} \mathbf{V}^*. \quad (4.3.9)$$

So,  $\sup_{\mathbf{N}} \langle \mathbf{N}, \mathbf{M} \rangle \geq \|\mathbf{M}\|_*$ .

For the opposite direction, notice if matrix  $\mathbf{N} \in \mathbb{R}^{n_1 \times n_2}$  satisfies  $\|\mathbf{N}\| \leq 1$ , then  $\bar{\mathbf{N}} \doteq \mathbf{U}^* \mathbf{N} \mathbf{V}$  has columns of  $\ell^2$  norm at most one. Thus, for each

---

<sup>6</sup>For  $p \in [1, \infty]$ , the Schatten  $p$ -norm of a matrix is  $\|\mathbf{X}\|_{S_p} = \|\boldsymbol{\sigma}(\mathbf{X})\|_p$ . The Ky-Fan  $k$ -norm is  $\|\mathbf{X}\|_{KF_k} = \sum_{i=1}^k \sigma_i(\mathbf{X})$ . Both of these functions are examples of *orthogonal invariant matrix norms*.

$i$ ,  $\bar{N}_{ii} \leq 1$ , and

$$\langle \mathbf{N}, \mathbf{M} \rangle = \langle \bar{\mathbf{N}}, \boldsymbol{\Sigma} \rangle = \sum_{i=1}^{n_2} \bar{N}_{ii} \sigma_i \leq \sum_i \sigma_i = \|\mathbf{M}\|_*. \quad (4.3.10)$$

This establishes the result.

For the second equality in (4.3.6), notice that for any nonzero  $\mathbf{M}$ ,

$$\langle \mathbf{M}, \mathbf{N} \rangle = \|\mathbf{M}\| \left\langle \frac{\mathbf{M}}{\|\mathbf{M}\|}, \mathbf{N} \right\rangle \leq \|\mathbf{M}\| \|\mathbf{N}\|_*. \quad (4.3.11)$$

Hence,  $\sup_{\|\mathbf{N}\|_* \leq 1} \langle \mathbf{M}, \mathbf{N} \rangle \leq \|\mathbf{M}\|$ . To show that this inequality is actually an equality, let take  $\mathbf{N} = \mathbf{u}_1 \mathbf{v}_1^*$ , and notice that  $\|\mathbf{N}\|_* = 1$  and  $\langle \mathbf{M}, \mathbf{N} \rangle = \mathbf{u}_1^* \mathbf{M} \mathbf{v}_1 = \sigma_1(\mathbf{M}) = \|\mathbf{M}\|$ . This completes the proof of (4.3.6).

To see that  $\|\cdot\|_*$  is indeed a norm, we just use (4.3.6) to verify that the three axioms of a norm are satisfied. Since the singular values are nonnegative, and  $\sigma_1(\mathbf{M}) = 0$  if and only if  $\mathbf{M} = \mathbf{0}$ , it is immediate that  $\|\mathbf{M}\|_* \geq 0$  with equality iff  $\mathbf{M} = \mathbf{0}$ . For nonnegative homogeneity, notice that for  $t \in \mathbb{R}_+$ ,

$$\|t\mathbf{M}\|_* = \sup_{\|\mathbf{N}\|_* \leq 1} \langle t\mathbf{M}, \mathbf{N} \rangle = t \sup_{\|\mathbf{N}\|_* \leq 1} \langle \mathbf{M}, \mathbf{N} \rangle = t \|\mathbf{M}\|_*. \quad (4.3.12)$$

Finally, for the triangle inequality, consider two matrices  $\mathbf{M}$  and  $\mathbf{M}'$ , and notice that

$$\begin{aligned} \|\mathbf{M} + \mathbf{M}'\|_* &= \sup_{\|\tilde{\mathbf{N}}\|_* \leq 1} \langle \mathbf{M} + \mathbf{M}', \tilde{\mathbf{N}} \rangle \\ &\leq \sup_{\|\mathbf{N}\|_* \leq 1} \langle \mathbf{M}, \mathbf{N} \rangle + \sup_{\|\mathbf{N}'\|_* \leq 1} \langle \mathbf{M}', \mathbf{N}' \rangle \\ &= \|\mathbf{M}\|_* + \|\mathbf{M}'\|_*, \end{aligned} \quad (4.3.13)$$

verifying the triangle inequality. This shows that  $\|\cdot\|_*$  is indeed a norm.  $\square$

The above proof highlights a useful fact about  $\|\cdot\|_*$ : it is the dual norm of the operator norm  $\|\mathbf{X}\| = \sigma_1(\mathbf{X})$ . The fact that  $\|\cdot\|_*$  is the dual of  $\|\cdot\|$  explains the  $*$  notation – this symbol is often used for duality.

Because  $\|\cdot\|_*$  is a norm, it is convex. Hence, a natural convex replacement for the rank minimization problem is the *nuclear norm minimization* problem

$$\begin{aligned} &\text{minimize} && \|\mathbf{X}\|_*, \\ &\text{subject to} && \mathcal{A}[\mathbf{X}] = \mathbf{y}. \end{aligned} \quad (4.3.14)$$

This problem is convex, and moreover is efficiently solvable. In Chapter 8, we will see how to use the special structure of this problem to give practical, efficient algorithms which work well at moderate scales.

### 4.3.3 Nuclear Norm as A Convex Envelope of Rank

From the analogy to  $\ell^0/\ell^1$  minimization, we might guess that the nuclear norm is a good convex surrogate for the rank, over some appropriate set. Recall that we have proved in Theorem 2.3.3 that the  $\ell^1$  norm was the convex envelope of the  $\ell^0$  norm over the  $\ell^\infty$  ball. Since for a matrix  $\mathbf{X}$ ,  $\|\boldsymbol{\sigma}(\mathbf{X})\|_\infty = \sigma_1(\mathbf{X}) = \|\mathbf{X}\|$ , you might guess the following relationship:

**Theorem 4.3.2.**  $\|\mathbf{M}\|_*$  is the convex envelope of  $f(\mathbf{M}) = \text{rank}(\mathbf{M})$  over

$$\mathcal{B}_{op} = \{\mathbf{M} \mid \|\mathbf{M}\| \leq 1\}. \quad (4.3.15)$$

*Proof.* We prove that any convex function  $f$  which satisfies

$$f(\mathbf{M}) \leq \text{rank}(\mathbf{M}) \quad (4.3.16)$$

for all  $\mathbf{M} \in \mathcal{B}_{op}$ , is majorized by the nuclear norm:  $f(\mathbf{M}) \leq \|\mathbf{M}\|_*$ .

Write the SVD  $\mathbf{M} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^*$ . Notice that

$$\boldsymbol{\Sigma} \in \text{conv} \left\{ \text{diag}(\mathbf{w}) \mid \mathbf{w} \in \{0, 1\}^{\min\{n_1, n_2\}} \right\}, \quad (4.3.17)$$

and for any  $\mathbf{w} \in \{0, 1\}^{\min\{n_1, n_2\}}$ ,

$$\|\mathbf{U} \text{diag}(\mathbf{w}) \mathbf{V}^*\|_* = \sum_i w_i = \text{rank}(\mathbf{U} \text{diag}(\mathbf{w}) \mathbf{V}^*). \quad (4.3.18)$$

Writing

$$\boldsymbol{\Sigma} = \sum_i \lambda_i \text{diag}(\mathbf{w}_i) \quad (4.3.19)$$

with  $\mathbf{w}_i \in \{0, 1\}^{\min\{n_1, n_2\}}$  with  $\lambda_i \geq 0$  and  $\sum_i \lambda_i = 1$ , and applying Jensen's inequality, we obtain

$$f(\mathbf{M}) = f\left(\mathbf{U} \sum_i \lambda_i \text{diag}(\mathbf{w}_i) \mathbf{V}^*\right) \quad (4.3.20)$$

$$\leq \sum_i \lambda_i f(\mathbf{U} \text{diag}(\mathbf{w}_i) \mathbf{V}^*) \quad (4.3.21)$$

$$\leq \sum_i \lambda_i \text{rank}(\mathbf{U} \text{diag}(\mathbf{w}_i) \mathbf{V}^*) \quad (4.3.22)$$

$$= \sum_i \lambda_i \|\mathbf{w}_i\|_1 \quad (4.3.23)$$

$$= \left\| \mathbf{U} \sum_i \lambda_i \text{diag}(\mathbf{w}_i) \mathbf{V}^* \right\|_* \quad (4.3.24)$$

$$= \|\mathbf{M}\|_* \quad (4.3.25)$$

as desired.  $\square$

Note that this proof essentially mirrored our argument for  $\ell^1$  and  $\ell^\infty$ . This is not a coincidence!

#### 4.3.4 Success of Nuclear Norm under Rank-RIP

For now, content that we can solve nuclear norm minimization problems efficiently, we turn our attention to whether nuclear norm minimization actually gives the correct answers. Namely, if we know that  $\mathbf{y} = \mathcal{A}[\mathbf{X}_{\text{true}}]$ , with  $r = \text{rank}(\mathbf{X}_{\text{true}}) \ll n$ , is it true that  $\mathbf{X}_{\text{true}}$  is the unique optimal solution to the nuclear norm minimization problem (4.3.14)? What we can say depends strongly on what we know about the operator  $\mathcal{A}$ .

By analogy to the *sparse* recovery problem, we can ask if it is enough for  $\mathcal{A}$  to preserve the geometry of a small set of structured objects – here, the low-rank matrices. Formally, we can define a *rank-restricted isometry property*, under which for every rank- $r$   $\mathbf{X}$ ,  $\|\mathcal{A}[\mathbf{X}]\|_2 \approx \|\mathbf{X}\|_F$ .

**Definition 4.3.3** (Rank-Restricted Isometry Property [Recht et al., 2010]). *The operator  $\mathcal{A}$  has the rank-restricted isometry property of rank  $r$  with constant  $\delta$ , if  $\forall \mathbf{X}$  s.t.  $\text{rank}(\mathbf{X}) \leq r$ , we have*

$$(1 - \delta)\|\mathbf{X}\|_F^2 \leq \|\mathcal{A}[\mathbf{X}]\|_2^2 \leq (1 + \delta)\|\mathbf{X}\|_F^2. \quad (4.3.26)$$

*The rank- $r$  restricted isometry constant  $\delta_r(\mathcal{A})$  is the smallest  $\delta$  such that the above property holds.*

As with the RIP for sparse vectors, the rank-RIP implies uniqueness of structured (low-rank) solutions:

**Theorem 4.3.4.** *If  $\mathbf{y} = \mathcal{A}[\mathbf{X}_{\text{true}}]$ , with  $r = \text{rank}(\mathbf{X}_{\text{true}})$  and  $\delta_{2r}(\mathcal{A}) < 1$ , then  $\mathbf{X}_{\text{true}}$  is the unique optimal solution to the rank minimization problem*

$$\begin{aligned} &\text{minimize} && \text{rank}(\mathbf{X}) \\ &\text{subject to} && \mathcal{A}[\mathbf{X}] = \mathbf{y}. \end{aligned} \quad (4.3.27)$$

We leave the proof of this claim as an exercise to the reader (see Exercise 4.14). The key property is the subadditivity of the matrix rank, namely,

$$\text{rank}(\mathbf{X} + \mathbf{X}') \leq \text{rank}(\mathbf{X}) + \text{rank}(\mathbf{X}'). \quad (4.3.28)$$

Moreover, like the RIP for sparse vectors, when the rank-RIP holds with sufficiently small constant  $\delta$ , we can conclude that nuclear norm minimization will recover the desired low-rank solution:

**Theorem 4.3.5** (Recht, Fazel, Parillo [Recht et al., 2010]). *Suppose that  $\mathbf{y} = \mathcal{A}[\mathbf{X}_{\text{true}}]$  with  $\text{rank}(\mathbf{X}_{\text{true}}) \leq r$ , and that  $\delta_{4r}(\mathcal{A}) \leq \sqrt{2} - 1$ . Then  $\mathbf{X}_{\text{true}}$  is the unique optimal solution to the nuclear norm minimization problem*

$$\begin{aligned} &\text{minimize} && \|\mathbf{X}\|_* \\ &\text{subject to} && \mathcal{A}[\mathbf{X}] = \mathbf{y}. \end{aligned} \quad (4.3.29)$$

There is nothing so special here about the numbers  $4r$  and  $\sqrt{2} - 1$ . The interesting part is the qualitative statement: if  $\mathcal{A}$  respects the geometry of low-rank matrices in a sufficiently strong sense, then nuclear norm minimization succeeds. The proof is analogous to the proof we have given in the



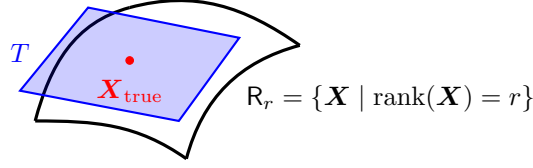


Figure 4.5. **“Support” of a Low-rank Matrix  $\mathbf{X}_{\text{true}}$ .** Consider a rank- $r$  matrix  $\mathbf{X}_0$  with compact singular value decomposition  $\mathbf{X}_{\text{true}} = \mathbf{U}\Sigma\mathbf{V}^*$ . The subspace  $T = \{\mathbf{U}\mathbf{W}^* + \mathbf{Q}\mathbf{V}^*\}$  can be interpreted as the *tangent space* to the collection  $\mathbf{R}_r$  of rank- $r$  matrices at  $\mathbf{X}_0$ .

previous chapter for the success of  $\ell^1$  minimization for recovering sparse signals. However, to extend the proof techniques from  $\ell^1$  to nuclear norm, we need to generalize a few concepts from vectors to matrices.

**“Support” and “Signs” of a Low-rank Matrix  $\mathbf{X}_0$ .**

Let  $\mathbf{X}_{\text{true}} = \mathbf{U}\Sigma\mathbf{V}^*$  denote the compact SVD of the true solution  $\mathbf{X}_{\text{true}}$ . Let

$$T \doteq \{\mathbf{U}\mathbf{W}^* + \mathbf{Q}\mathbf{V}^* \mid \mathbf{W} \in \mathbb{R}^{n_2 \times r}, \mathbf{Q} \in \mathbb{R}^{n_1 \times r}\} \subseteq \mathbb{R}^{n_1 \times n_2}. \quad (4.3.30)$$

Notice that  $T$  is a linear subspace. In the analogy between  $\ell^1$  minimization and nuclear norm minimization, the subspace  $T$  plays the role of the “support” of  $\mathbf{X}_{\text{true}}$ . Geometrically,  $T$  represents the *tangent space* to the set of rank- $r$  matrices at  $\mathbf{X}_{\text{true}}$  – see Figure 4.5 and Exercise 4.9. The subspace  $T$  is generated by matrices  $\mathbf{U}\mathbf{W}^*$  whose column space is contained in  $\text{col}(\mathbf{X}_{\text{true}})$  and matrices  $\mathbf{Q}\mathbf{V}^*$  whose row space is contained in  $\text{row}(\mathbf{X}_{\text{true}})$ . Notice that elements in  $T$  have rank no more than  $2r$ . Meanwhile the matrix  $\mathbf{U}\mathbf{V}^*$  plays the role of the “signs” of  $\mathbf{X}_{\text{true}}$  since  $\mathbf{U}\mathbf{V}^* \in T$  and

$$\langle \mathbf{X}_{\text{true}}, \mathbf{U}\mathbf{V}^* \rangle = \|\mathbf{X}_{\text{true}}\|_*. \quad (4.3.31)$$

The orthogonal complement of  $T$  is

$$T^\perp \doteq \{\mathbf{M} \mid \text{col}(\mathbf{M}) \perp \text{col}(\mathbf{X}), \text{row}(\mathbf{M}) \perp \text{row}(\mathbf{X})\}. \quad (4.3.32)$$

Let  $\mathbf{P}_U = \mathbf{U}\mathbf{U}^*$  and  $\mathbf{P}_V = \mathbf{V}\mathbf{V}^*$  be the orthogonal projections onto the column space and row space of  $\mathbf{X}_{\text{true}}$ , respectively. Then the orthogonal projections onto these subspaces are given by<sup>7</sup>

$$\mathcal{P}_T[\mathbf{M}] = \mathbf{P}_U\mathbf{M} + \mathbf{M}\mathbf{P}_V - \mathbf{P}_U\mathbf{M}\mathbf{P}_V, \quad (4.3.33)$$

and

$$\mathcal{P}_{T^\perp}[\mathbf{M}] = (\mathbf{I} - \mathbf{P}_U)\mathbf{M}(\mathbf{I} - \mathbf{P}_V). \quad (4.3.34)$$

<sup>7</sup>Equations (4.3.33) and (4.3.34) can be derived from the condition that at  $\mathcal{P}_T[\mathbf{M}]$ , the error  $\mathbf{M} - \mathcal{P}_T[\mathbf{M}]$  is orthogonal to  $T$ .

Notice that because the orthogonal projections  $\mathbf{P}_{U^\perp} = \mathbf{I} - \mathbf{P}_U$  and  $\mathbf{P}_{V^\perp} = \mathbf{I} - \mathbf{P}_V$  have norm at most one,  $\mathcal{P}_{T^\perp}$  does not increase the operator norm:

$$\|\mathcal{P}_{T^\perp}[\mathbf{M}]\| \leq \|\mathbf{M}\|. \quad (4.3.35)$$

**Feasible Cone Restriction.**

Note that any matrix  $\mathbf{M} \in T^\perp$  has columns that are orthogonal to the columns of  $\mathbf{U}$  and rows that are orthogonal to the rows of  $\mathbf{V}^*$ . This implies that

$$\|\mathbf{M} + \mathbf{UV}^*\| = \max\{\|\mathbf{M}\|, \|\mathbf{UV}^*\|\} = \max\{\|\mathbf{M}\|, 1\}. \quad (4.3.36)$$

So, for any matrix  $\mathbf{X}$ ,

$$\|\mathbf{X}\|_* = \sup_{\|\mathbf{Q}\| \leq 1} \langle \mathbf{X}, \mathbf{Q} \rangle \quad (4.3.37)$$

$$\geq \sup_{\|\mathbf{M}\| \leq 1} \langle \mathbf{X}, \mathbf{UV}^* + \mathcal{P}_{T^\perp}[\mathbf{M}] \rangle \quad (4.3.38)$$

$$= \langle \mathbf{X}, \mathbf{UV}^* \rangle + \sup_{\|\mathbf{M}\| \leq 1} \langle \mathcal{P}_{T^\perp}[\mathbf{X}], \mathbf{M} \rangle \quad (4.3.39)$$

$$= \langle \mathbf{X}, \mathbf{UV}^* \rangle + \|\mathcal{P}_{T^\perp}[\mathbf{X}]\|_*. \quad (4.3.40)$$

Let  $\hat{\mathbf{X}}$  be any optimal solution to our problem (4.3.29). It can be written as  $\hat{\mathbf{X}} = \mathbf{X}_{\text{true}} + \mathbf{H}$ , with  $\mathbf{H} = \hat{\mathbf{X}} - \mathbf{X}_{\text{true}} \in \text{null}(\mathcal{A})$ . From the above calculation, we have

$$\|\mathbf{X}_{\text{true}} + \mathbf{H}\|_* \geq \langle \mathbf{X}_{\text{true}} + \mathbf{H}, \mathbf{UV}^* \rangle + \|\mathcal{P}_{T^\perp}[\hat{\mathbf{X}}]\|_* \quad (4.3.41)$$

$$= \|\mathbf{X}_{\text{true}}\|_* + \langle \mathbf{H}, \mathbf{UV}^* \rangle + \|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_* \quad (4.3.42)$$

$$\geq \|\mathbf{X}_{\text{true}}\|_* - \|\mathcal{P}_T[\mathbf{H}]\|_* + \|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_*. \quad (4.3.43)$$

So, if a better solution than  $\mathbf{X}_{\text{true}}$  exists, the feasible perturbation  $\mathbf{H}$  must satisfy the following cone restriction:

$$\|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_* \leq \|\mathcal{P}_T[\mathbf{H}]\|_*. \quad (4.3.44)$$

**Matrix Restricted Strong Convexity Property.**

As in the proof of  $\ell^1$  success, we want to show that feasible perturbations  $\mathbf{H} \in \text{null}(\mathcal{A})$  must have  $\|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_* \gg \|\mathcal{P}_T[\mathbf{H}]\|_*$ . This is true if the operator  $\mathcal{A}$  satisfies the following (uniform) *matrix restricted strong convexity* property (RSC) property:

**Definition 4.3.6** (Matrix Restricted Strong Convexity). *The linear operator  $\mathcal{A}$  satisfies the matrix restricted strong convexity (RSC) condition of rank  $r$  with constant  $C$  if for every support  $T$  of a matrix of rank  $k$  and for all nonzero  $\mathbf{H}$  satisfying*

$$\|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_* \leq C \cdot \|\mathcal{P}_T[\mathbf{H}]\|_*. \quad (4.3.45)$$

with some constant  $C \geq 1$ , we have

$$\|\mathcal{A}[\mathbf{H}]\|_2^2 > \mu \cdot \|\mathbf{H}\|_F^2 \quad (4.3.46)$$

for some constant  $\mu > 0$ .

The following theorem says that if  $\mathcal{A}$  satisfies the rank-RIP, then it satisfies the matrix RSC:

**Theorem 4.3.7** (Rank-RIP Implies Matrix RSC). *If a linear operator  $\mathcal{A}$  satisfies rank-RIP with  $\delta_{4r} < \frac{1}{1+C\sqrt{2}}$ , then  $\mathcal{A}$  satisfies the matrix-RSC condition of rank  $r$  with constant  $C$ .*

Both the statement and proof of Theorem 4.3.7 parallel Theorem 3.3.10 for the  $\ell^1$  norm. Theorem 4.3.7 involves  $\delta_{4r}$ , as opposed to  $\delta_{2k}$  for  $k$ -sparse vectors. The bigger constant in  $4r = r + 3r$  reflects the need to account for all three components of the singular value decomposition in the proof:

*Proof.* Using the parallelogram identity, similar to Lemma 3.3.9, it is not difficult to show that for any  $\mathbf{Z}, \mathbf{Z}'$  such that  $\mathbf{Z} \perp \mathbf{Z}'$ , and  $\text{rank}(\mathbf{Z}) + \text{rank}(\mathbf{Z}') \leq 4r$ ,

$$|\langle \mathcal{A}[\mathbf{Z}], \mathcal{A}[\mathbf{Z}'] \rangle| \leq \delta_{4r}(\mathcal{A}) \|\mathbf{Z}\|_F \|\mathbf{Z}'\|_F. \quad (4.3.47)$$

Let  $T$  denote the support subspace for some matrix of rank  $r$ . Take any  $\mathbf{H}$  that satisfies the cone restriction  $\|\mathcal{P}_{T^\perp}[\mathbf{Z}]\|_* \leq C \cdot \|\mathcal{P}_T \mathbf{Z}\|_*$ , and write

$$\mathbf{H} = \mathcal{P}_T[\mathbf{H}] + \mathcal{P}_{T^\perp}[\mathbf{H}]. \quad (4.3.48)$$

Let  $\mathbf{H}_T$  denote  $\mathcal{P}_T[\mathbf{H}]$ . For the second term,  $\mathcal{P}_{T^\perp}[\mathbf{H}]$ , write its compact singular value decomposition

$$\mathcal{P}_{T^\perp}[\mathbf{H}] = \sum_i \eta_i \phi_i \zeta_i^*, \quad (4.3.49)$$

where  $\phi_1, \phi_2, \dots$  are the left singular vectors,  $\zeta_1, \zeta_2, \dots$  the right singular vectors, and  $\eta_1 \geq \eta_2 \geq \dots > 0$  the singular values. From the variational characterization of the singular vectors, each  $\phi_i$  is orthogonal to the columns of  $\mathbf{U}$ , and each  $\zeta_i$  is orthogonal to the columns of  $\mathbf{V}$ . So, if we partition  $\mathcal{P}_{T^\perp}[\mathbf{H}]$  as

$$\mathcal{P}_{T^\perp}[\mathbf{H}] = \underbrace{\sum_{i=1}^r \eta_i \phi_i \zeta_i^*}_{\doteq \Phi_1} + \underbrace{\sum_{i=r+1}^{2r} \eta_i \phi_i \zeta_i^*}_{\doteq \Phi_2} + \dots, \quad (4.3.50)$$

we have  $\Phi_i \perp \Phi_j$  for  $i \neq j$ ,  $\Phi_i \perp \mathbf{H}_T$  for every  $T$ .

Since the singular values  $\eta_i$  are non-increasing, the largest singular value of the  $(i+1)$ -th block is bounded by the average of the singular values in the  $i$ -th block.

$$\forall i \geq 1, \quad \|\Phi_{i+1}\| \leq \frac{\|\Phi_i\|_*}{r}. \quad (4.3.51)$$

So, noting that as an element in  $T$ , we have  $\text{rank}(\mathbf{H}_T) \leq 2r$  and so  $\text{rank}(\mathbf{H}_T + \Phi_1) \leq 3r$ . Notice that

$$\mathcal{A}[\mathbf{H}_T] + \mathcal{A}[\Phi_1] = \mathcal{A}[\mathbf{H}] - \mathcal{A}[\Phi_2] - \mathcal{A}[\Phi_3] - \dots \quad (4.3.52)$$

Then, very similar to the derivation of inequalities (3.3.32) in Theorem 3.3.10, and by applying the rank-RIP to matrices of rank bounded by at most  $4r$ , we have

$$\begin{aligned} & (1 - \delta_{4r}) \|\mathbf{H}_T + \Phi_1\|_F^2 \\ & \leq \langle \mathcal{A}[\mathbf{H}_T + \Phi_1], \mathcal{A}[\mathbf{H}_T + \Phi_1] \rangle \\ & = \langle \mathcal{A}[\mathbf{H}_T + \Phi_1], \mathcal{A}[\mathbf{H}] - \mathcal{A}[\Phi_2] - \mathcal{A}[\Phi_3] - \dots \rangle \\ & \leq \sum_{j \geq 2} |\langle \mathcal{A}[\mathbf{H}_T], \mathcal{A}[\Phi_j] \rangle| + |\langle \mathcal{A}[\Phi_1], \mathcal{A}[\Phi_j] \rangle| + \langle \mathcal{A}[\mathbf{H}_T + \Phi_1], \mathcal{A}[\mathbf{H}] \rangle \\ & \leq \delta_{4r} (\|\mathbf{H}_T\|_F + \|\Phi_1\|_F) \sum_{j \geq 2} \|\Phi_j\|_F + \|\mathcal{A}[\mathbf{H}_T + \Phi_1]\|_2 \|\mathcal{A}[\mathbf{H}]\|_2 \\ & \leq \delta_{4r} \sqrt{2} \|\mathbf{H}_T + \Phi_1\|_F \sum_{j \geq 2} \|\Phi_j\|_F + (1 + \delta_{4r}) \|\mathbf{H}_T + \Phi_1\|_F \|\mathcal{A}[\mathbf{H}]\|_2 \\ & \leq \delta_{4r} \sqrt{2} \|\mathbf{H}_T + \Phi_1\|_F \frac{\|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_*}{\sqrt{r}} + (1 + \delta_{4r}) \|\mathbf{H}_T + \Phi_1\|_F \|\mathcal{A}[\mathbf{H}]\|_2. \end{aligned}$$

Note that  $\mathbf{H}$  is restricted by the cone condition (4.3.45), which leads to:

$$\|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_* \leq C\sqrt{r} \|\mathbf{H}_T + \Phi_1\|_F. \quad (4.3.53)$$

Combining this with the previous inequality, we obtain:

$$\|\mathcal{A}[\mathbf{H}]\|_2 \geq \frac{1 - \delta_{4r}(1 + C\sqrt{2})}{1 + \delta_{4r}} \|\mathbf{H}_T + \Phi_1\|_F. \quad (4.3.54)$$

Since the singular values  $\eta_i$  are non-increasing, the  $i$ th singular value in  $\Phi_2 + \Phi_3 + \dots$  is no larger than the mean of the first  $i$  singular values in  $\mathcal{P}_{T^\perp}[\mathbf{H}]$ . So we have

$$\forall i \geq r + 1, \eta_i \leq \|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_* / i. \quad (4.3.55)$$

This leads to

$$\|\Phi_2 + \Phi_3 + \dots\|_F^2 = \sum_{i=r+1}^{\infty} \eta_i^2 \quad (4.3.56)$$

$$\leq \|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_*^2 \sum_{i=r+1}^{\infty} \frac{1}{i^2} \quad (4.3.57)$$

$$\leq \frac{\|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_*^2}{r} \leq \frac{C^2 \|\mathbf{H}_T\|_*^2}{r} \quad (4.3.58)$$

$$\leq C^2 \|\mathbf{H}_T\|_F^2 \leq C^2 \|\mathbf{H}_T + \Phi_1\|_F^2. \quad (4.3.59)$$

Since  $\Phi_i$  with  $i \geq 2$  are orthogonal to  $\mathbf{H}_T + \Phi_1$ , this gives us

$$\|\mathbf{H}\|_F^2 \leq (1 + C^2) \|\mathbf{H}_T + \Phi_1\|_F^2. \quad (4.3.60)$$

Combining this with the previous bound (4.3.54) on  $\|\mathcal{A}[\mathbf{H}]\|_2$ , we obtain

$$\|\mathcal{A}[\mathbf{H}]\|_2 \geq \frac{1 - \delta_{4r}(1 + C\sqrt{2})}{(1 + \delta_{4r})\sqrt{1 + C^2}} \|\mathbf{H}\|_F. \quad (4.3.61)$$

This concludes the proof.  $\square$

Note that for the nuclear norm minimization problem, the feasible perturbation  $\mathbf{H}$  satisfies the cone restriction (4.3.44). Thus Theorem 4.3.5 is essentially a corollary to Theorem 4.3.7 with constant  $C = 1$  for the cone restriction.

#### 4.3.5 Rank-RIP of Random Measurements

Theorem 4.3.5 indicates that the rank-RIP implies a very strong conclusion: nuclear norm minimization exactly recovers low-rank matrices. Moreover the recovery is *uniform* in the sense that a single set of measurements  $\mathcal{A}$  suffices to recover any sufficiently low-rank matrix  $\mathbf{X}_{\text{true}}$ . The remaining question is what measurement operators satisfy the rank-RIP?

##### Random Gaussian Measurements

A simple and natural choice is to consider the random Gaussian measurements:

$$\mathcal{A}[\mathbf{X}] = (\langle \mathbf{A}_1, \mathbf{X} \rangle, \dots, \langle \mathbf{A}_m, \mathbf{X} \rangle), \quad (4.3.62)$$

where the entries of the matrices  $\mathbf{A}_1, \dots, \mathbf{A}_m \in \mathbb{R}^{n_1 \times n_2}$  are all i.i.d. Gaussian  $\mathcal{N}(0, \frac{1}{m})$ . This is equivalent to viewing  $\mathcal{A}$  as an  $m \times n_1 n_2$  matrix with entries  $\mathcal{A}_{ij}$  sampled i.i.d.  $\mathcal{N}(0, \frac{1}{m})$ . We demonstrate that such random maps satisfy the rank-RIP with high probability, using ideas and techniques similar to the proof of the (regular) RIP of random Gaussian matrices in Section 3.4.2:

**Theorem 4.3.8** (Rank-RIP of Gaussian Measurements). *If the measurement operator  $\mathcal{A}$  is a random Gaussian map with entries i.i.d.  $\mathcal{N}(0, \frac{1}{m})$ , then  $\mathcal{A}$  satisfies the rank-RIP with constant  $\delta_r(\mathcal{A}) \leq \delta$  with high probability, provided  $m \geq Cr(n_1 + n_2) \times \delta^{-2} \log \delta^{-1}$ , where  $C > 0$  is a numerical constant.*

*Proof.* Let

$$S_r = \{\mathbf{X} \mid \text{rank}(\mathbf{X}) \leq r, \|\mathbf{X}\|_F = 1\}.$$

Notice that  $\delta_r(\mathcal{A}) \leq \delta$  if and only if

$$\sup_{\mathbf{X} \in S_r} |\langle \mathcal{A}(\mathbf{X}), \mathcal{A}(\mathbf{X}) \rangle - 1| \leq \delta. \quad (4.3.63)$$

We complete the rest of the proof in three steps.

1. *Constructing a Covering  $\epsilon$ -Net for  $S_r$ .*

Notice that for any rank- $r$  matrix  $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$ , it can be represented by its SVD;  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ . So to construct a covering of all rank- $r$  matrices, we can try to construct a covering for each of the terms  $\mathbf{U}$ ,  $\mathbf{V}$  and  $\mathbf{\Sigma}$ , respectively.

**Lemma 4.3.9.** *There is a covering  $\varepsilon$ -net  $N_U$  for the  $H = \{\mathbf{U} \in \mathbb{R}^{n_1 \times r} | \mathbf{U}^* \mathbf{U} = \mathbf{I}\}$ , of size  $|N_U| \leq (6/\varepsilon)^{n_1 r}$ .*

*Proof.* Let  $N'$  be an  $\varepsilon/2$ -net for  $\{\mathbf{U} \in \mathbb{R}^{n_1 \times r} | \|\mathbf{U}\|_F^2 = 1\}$  of size  $|N'| \leq (6/\varepsilon)^{n_1 r}$  (from Lemma 3.4.5). Let

$$Q = \{\mathbf{U}' \in N' | \exists \mathbf{U} \in H \text{ with } \|\mathbf{U} - \mathbf{U}'\|_F \leq \varepsilon/2\}.$$

For each  $\mathbf{U}' \in Q$ , let  $\hat{\mathbf{U}}(\mathbf{U}')$  be the nearest element of  $H$ . Set  $N_U = \{\hat{\mathbf{U}}(\mathbf{U}') | \mathbf{U}' \in Q\} \subseteq H$ . By the triangle inequality,  $N_U$  is an  $\varepsilon$ -net for  $H$ .  $\square$

Similarly, one can construct an  $\varepsilon$ -net  $N_V$  for  $H' = \{\mathbf{V} \in \mathbb{R}^{n_2 \times r} | \mathbf{V}^* \mathbf{V} = \mathbf{I}\}$  of size  $|N_V| \leq (6/\varepsilon)^{n_2 r}$ . With this lemma, we have the following result.

**Lemma 4.3.10.** *There is a covering  $\varepsilon$ -net  $N_r$  for the set  $S_r$ , of size  $|N_r| \leq \exp((n_1 + n_2)r \log(18/\varepsilon) + r \log(9/\varepsilon))$ .*

*Proof.* Choose covering  $\varepsilon/3$ -nets  $N_U$  and  $N_V$  for  $H$  and  $H'$ , respectively. According to the above lemma, the sizes of the nets can be less than  $(18/\varepsilon)^{n_1 r}$  and  $(18/\varepsilon)^{n_2 r}$ , respectively. Form a covering  $\varepsilon/3$ -net  $N_\Sigma$  for

$$D = \{\mathbf{\Sigma} \in \mathbb{R}^{r \times r} | \mathbf{\Sigma} \text{ diagonal, } \|\mathbf{\Sigma}\|_F = 1\}.$$

According to Lemma 3.4.5, the size of the net can be less than  $|N_\Sigma| \leq (9/\varepsilon)^r$ .

Now consider the following net for the whole set  $S_r$ :

$$N_r = \{\mathbf{U}\mathbf{\Sigma}\mathbf{V}^* | \mathbf{U} \in N_U, \mathbf{\Sigma} \in N_\Sigma, \mathbf{V} \in N_V\}.$$

Its size is bounded by the product of all three nets, hence the expression in the Lemma. Now we only have to show that this is indeed a covering  $\varepsilon$ -net for  $S_r$ . For any given  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ , we can find  $\hat{\mathbf{X}} = \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}\hat{\mathbf{V}}^* \in N_r$  with  $\|\mathbf{U} - \hat{\mathbf{U}}\|_F \leq \varepsilon/3$ ,  $\|\mathbf{V} - \hat{\mathbf{V}}\|_F \leq \varepsilon/3$ , and  $\|\mathbf{\Sigma} - \hat{\mathbf{\Sigma}}\|_F \leq \varepsilon/3$ .

The triangle inequality gives

$$\begin{aligned} & \|\mathbf{X} - \hat{\mathbf{X}}\|_F \\ & \leq \|\mathbf{U} - \hat{\mathbf{U}}\|_F \|\mathbf{\Sigma}\mathbf{V}^*\| + \|\hat{\mathbf{U}}\| \|\mathbf{\Sigma} - \hat{\mathbf{\Sigma}}\|_F \|\mathbf{V}^*\| + \|\hat{\mathbf{U}}\hat{\mathbf{\Sigma}}\| \|\mathbf{V}^* - \hat{\mathbf{V}}^*\|_F \\ & \leq \varepsilon, \end{aligned}$$

where we have used that each of the approximation errors is bounded by  $\varepsilon/3$ , and each of the operator norms in the above expression are bounded by 1.  $\square$

### 2. Discretization.

As in the  $\ell^1$  case for sparse signals in Section 3.4.2, the goal of discretization is trying to show that if  $\mathcal{A}$  is restricted isometric on the finite set of (discrete) points in the covering net  $N_r$  with a constant  $\delta_{N_r}$ , so is  $\mathcal{A}$  on the whole set  $S_r$ , with a constant  $\delta_r$  possibly slightly larger than  $\delta_{N_r}$ .

Now consider a point  $\mathbf{X}$  in  $S_r$  and its closest point  $\hat{\mathbf{X}}$  in  $N_r$ . Thus, we have  $\|\mathbf{X} - \hat{\mathbf{X}}\|_F \leq \varepsilon$ . Also we have<sup>8</sup>

$$\begin{aligned} & |\langle \mathcal{A}[\mathbf{X}], \mathcal{A}[\mathbf{X}] \rangle - \langle \mathcal{A}[\hat{\mathbf{X}}], \mathcal{A}[\hat{\mathbf{X}}] \rangle| \\ = & |\langle \mathcal{A}[\mathbf{X}], \mathcal{A}[\mathbf{X} - \hat{\mathbf{X}}\mathbf{P}_V] \rangle + \langle \mathcal{A}[\mathbf{X} - \hat{\mathbf{X}}\mathbf{P}_V], \mathcal{A}[\hat{\mathbf{X}}\mathbf{P}_V] \rangle \\ & + \langle \mathcal{A}[\hat{\mathbf{X}}\mathbf{P}_V - \hat{\mathbf{X}}], \mathcal{A}[\hat{\mathbf{X}}\mathbf{P}_V] \rangle + \langle \mathcal{A}[\hat{\mathbf{X}}], \mathcal{A}[\hat{\mathbf{X}}\mathbf{P}_V - \hat{\mathbf{X}}] \rangle|. \end{aligned}$$

To bound the first term in the above expression, notice that

$$\|\mathbf{X} - \hat{\mathbf{X}}\mathbf{P}_V\|_F = \|(\mathbf{X} - \hat{\mathbf{X}})\mathbf{P}_V\|_F \leq \|\mathbf{X} - \hat{\mathbf{X}}\|_F \leq \varepsilon.$$

Also,  $\mathbf{X} - \hat{\mathbf{X}}\mathbf{P}_V$  is of rank  $r$ . So we have

$$|\langle \mathcal{A}[\mathbf{X}], \mathcal{A}[\mathbf{X} - \hat{\mathbf{X}}\mathbf{P}_V] \rangle| \leq (1 + \delta_r(\mathcal{A}))\varepsilon.$$

For the second term, since  $\mathbf{P}_{\hat{U}}$  is an orthogonal projection onto the space of matrices whose columns are the same as  $\hat{\mathbf{X}}$ , we have

$$\|\mathbf{X} - \mathbf{P}_{\hat{U}}\mathbf{X}\|_F \leq \|\mathbf{X} - \hat{\mathbf{X}}\|_F \leq \varepsilon.$$

Also, since  $\mathbf{X}$  and  $\mathbf{P}_{\hat{U}}\mathbf{X}$  have the same row space, so  $\mathbf{X} - \mathbf{P}_{\hat{U}}\mathbf{X}$  is of rank  $r$  or less. Therefore, we also have

$$|\langle \mathcal{A}[\mathbf{X} - \mathbf{P}_{\hat{U}}\mathbf{X}], \mathcal{A}[\hat{\mathbf{X}}\mathbf{P}_V] \rangle| \leq (1 + \delta_r(\mathcal{A}))\varepsilon.$$

Similarly for the third and fourth terms, each is bounded by the same bound. Therefore, we get

$$|\langle \mathcal{A}[\mathbf{X}], \mathcal{A}[\mathbf{X}] \rangle - \langle \mathcal{A}[\hat{\mathbf{X}}], \mathcal{A}[\hat{\mathbf{X}}] \rangle| \leq 4(1 + \delta_r(\mathcal{A}))\varepsilon.$$

From this we have

$$\delta_r(\mathcal{A}) - \delta_{N_r} \leq 4(1 + \delta_r(\mathcal{A}))\varepsilon. \quad (4.3.64)$$

This gives

$$\delta_r(\mathcal{A}) \leq \frac{4\varepsilon + \delta_{N_r}}{1 - 4\varepsilon}. \quad (4.3.65)$$

### 3. Union Bound.

For each  $\mathbf{X} \in N_r$ ,  $\mathcal{A}[\mathbf{X}] \in \mathbb{R}^m$  is a random vector with entries independent  $\mathcal{N}(0, 1/m)$ . We have

$$\mathbb{P} \left[ \left| \|\mathcal{A}[\mathbf{X}]\|_2^2 - 1 \right| > t \right] \leq 2 \exp(-mt^2/8). \quad (4.3.66)$$

---

<sup>8</sup>Notice that here the derivation is more subtle than the  $\ell^1$  case because  $\mathbf{X} - \hat{\mathbf{X}}$  is not necessarily of rank  $r$ .

Hence, summing the probabilities over all elements of  $N_r$ , we have

$$\begin{aligned}\mathbb{P}[\delta_{N_r} > t] &\leq 2|N_r| \exp(-mt^2/8) \\ &= 2 \exp\left(-\frac{mt^2}{8} + (n_1 + n_2)r \log(18/\varepsilon) + r \log(9/\varepsilon)\right).\end{aligned}$$

If we choose  $\varepsilon = c \cdot \delta$  and  $t = c \cdot \delta$  for some small constant  $c$  and ensure  $m \geq Cr(n_1 + n_2)\delta^{-2} \log \delta^{-1}$  for some large enough  $C$ , the above failure probability is bounded by  $2 \exp(-c'm\delta^2)$ . On the complement of this “failure” event,  $\delta_{N_r} \leq c \cdot \delta$ , and due to (4.3.65) we have  $\delta_r(\mathcal{A}) \leq \delta$ . This concludes the proof of the Theorem 4.3.8.  $\square$

The number of measurements  $m = O(r(n_1 + n_2))$  required is nearly optimal, since an  $n_1 \times n_2$  rank- $r$  matrix has  $r(n_1 + n_2 - r)$  degrees of freedom. Of course, the big  $O$  notation hides a numerical constant. Like the  $\ell^1$  minimization for sparse recovery, when the dimension is high, nuclear norm minimization exhibits a phase transition between success and failure. Identifying this transition yields more precise estimates of the number  $m$  of measurements required to reconstruct a low rank matrix. We discuss this issue in more detail below.

### *Random Submatrix of a Unitary Basis*

Although random Gaussian measurements have very nice properties such as (rank) RIP, the lack of structure in such measurements makes it rather expensive to generate, store and apply such operators in practice. Hence it is natural to ask if there exist other more structured measurements that have similarly good RIP properties. In Section 3.4.3, we have seen that given any unitary matrix that is incoherent from sparse signals, then a randomly selected subset of its rows will satisfy the RIP with high probability. An important special case that has been widely used in practice for compressive sensing is a randomly chosen submatrix of the discrete Fourier transform basis. It is then natural to ask what are the Fourier-type bases for matrices.

In the case of sparse recovery, we start with a unitary basis  $\mathbf{U} \in \mathbb{C}^{n \times n}$  and show that if the rows  $\{\mathbf{u}_i\}_{i=1}^n$  of the basis is incoherent with sparse signals:

$$\forall i \quad \|\mathbf{u}_i\|_\infty = \sup_{\mathbf{x}: \|\mathbf{x}\|_0=1} \langle \mathbf{u}_i, \mathbf{x} \rangle \leq \zeta/\sqrt{n}$$

for some constant  $\zeta$ , then a randomly selected (sufficient) number of rows of  $\mathbf{U}$  will satisfy RIP.

To simplify the discussion of matrices, we will assume  $n_1 = n_2 = n$  for the rest of this subsection; a similar approach applies when  $n_1 \neq n_2$ . Let us assume  $\{\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_{n^2}\} \subset \mathbb{C}^{n \times n}$  form a unitary basis for the matrix space  $\mathbb{C}^{n \times n}$ . Similarly we want each of the matrix  $\mathbf{U}_i$  to be incoherent with



low-rank matrices. Note that for any  $\mathbf{X} \in \mathbb{C}^{n \times n}$ ,

$$\|\mathbf{U}_i\| = \sup_{\mathbf{X}: \|\mathbf{X}\|_2=1, \text{rank}(\mathbf{X})=1} \langle \mathbf{U}_i, \mathbf{X} \rangle. \quad (4.3.67)$$

Hence in order for each  $\mathbf{U}_i$  to be incoherent with low-rank matrices, we could require:

$$\forall i \quad \|\mathbf{U}_i\| \leq \zeta/\sqrt{n}. \quad (4.3.68)$$

Then to construct the measurement operator  $\mathcal{A}$ , we randomly select a subset of  $m$  bases from  $\{\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_{n^2}\}$  and properly scale them as:<sup>9</sup>

$$\mathcal{A}: \quad \mathbf{A}_i = \frac{n}{\sqrt{m}} \mathbf{U}_i, \quad i = 1, \dots, m. \quad (4.3.69)$$

Then one should expect that when  $m$  is large enough, with high probability, the so-defined  $\mathcal{A}$  satisfies the rank-RIP. The following theorem makes this precise:

**Theorem 4.3.11.** *Let us assume  $\{\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_{n^2}\} \subset \mathbb{C}^{n \times n}$  be a unitary basis for the matrix space  $\mathbb{C}^{n \times n}$  and with  $\|\mathbf{U}_i\| \leq \zeta/\sqrt{n}$  for some constant  $\zeta$ . Let  $\mathcal{A}$  to be defined as per (4.3.69). Then if*

$$m \geq C\zeta^2 \cdot rn \log^6 n, \quad (4.3.70)$$

*then with high probability,  $\mathcal{A}$  satisfies the rank-RIP over the set of all rank- $r$  matrices.*

The proof of this theorem is out of the scope of this book and interested readers may refer to the work of [Liu, 2011].

According to this statement, from an incoherent unitary basis, with high probability we could find a (compressive) sensing operator  $\mathcal{A}$  such that it is rank-RIP. Hence with this operator, one can recover all rank- $r$  matrices via the nuclear norm minimization. The remaining question is what type of structured bases (of the matrix space) are rank-incoherent as per (4.3.68)? To this end, one should seek a matrix analogue to the Fourier basis.

In the case of MRI imaging, we have seen that measurements that one can physically take are essentially the Fourier coefficients of the brain image. As it turns out, the matrix analogue to Fourier basis also has a natural origin from physics. In quantum-state tomography, a system of  $k$  qubits is of dimension  $n = 2^k$ . The quantum state of such a system is described by a density matrix  $\mathbf{X}_0 \in \mathbb{C}^{n \times n}$  which is positive semidefinite with trace 1. When the state is early pure,  $\mathbf{X}_0$  is a very low-rank matrix with  $\text{rank}(\mathbf{X}_0) = r \ll n$ .

One problem in quantum physics is how to recover the quantum state  $\mathbf{X}_0$  of a system from linear measurements. As it turns out, a set of experimentally feasible measurements are given by the so-called *Pauli observables*.

---

<sup>9</sup>The scaling is to ensure that the “column” of  $\mathcal{A}$  to be of unit norm.

Each Pauli measurement is given by the inner product of  $\mathbf{X}_0$  with matrices of the form  $\mathbf{P}_1 \otimes \cdots \otimes \mathbf{P}_k$  where  $\otimes$  is the tensor (Kronecker) product and each  $\mathbf{P}_i = \frac{1}{\sqrt{2}}\boldsymbol{\sigma}$  where  $\boldsymbol{\sigma}$  is a  $2 \times 2$  matrix chosen from the following four possibilities:

$$\boldsymbol{\sigma}_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \boldsymbol{\sigma}_2 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \boldsymbol{\sigma}_3 = \begin{bmatrix} 0 & -i \\ i & 0 \end{bmatrix}, \quad \boldsymbol{\sigma}_4 = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}.$$

It is easy to see that there are a total of  $4^k$  possible choices for the tensor product, denoted as  $\{\mathbf{U}_i\}_{i=1}^{4^k}$  and they together form an orthonormal basis for the matrix space  $\mathbb{C}^{n \times n}$  where  $n = 2^k$ .

One can show that for each basis  $\mathbf{U}_i = \mathbf{P}_1 \otimes \cdots \otimes \mathbf{P}_k$ , its operator norm is bounded as  $\|\mathbf{U}_i\| \leq 1/\sqrt{n}$  hence incoherent with low-rank matrices. Then according to Theorem 4.3.11, a randomly selected  $m \geq Crn \log^6 n$  rows of the Pauli bases will satisfy the rank-RIP property with high probability. Hence, such a sensing operator will be able to uniformly recover all pure quantum states less than rank  $r$ .

#### 4.3.6 Noise, Inexact Low Rank, and Phase Transition

Above, we established that, under fairly broad conditions, nuclear norm minimization correctly recovers a low-rank matrix  $\mathbf{X}_{\text{true}}$  from ideal measurements  $\mathbf{y} = \mathcal{A}[\mathbf{X}_{\text{true}}]$ . In practice, the measurements can be corrupted by noise or measurement errors. In some cases,  $\mathbf{X}_{\text{true}}$  may not be exactly low rank. It is desirable to understand whether nuclear norm minimization still gives reasonably good estimates of  $\mathbf{X}_{\text{true}}$  in these situations.

In Section 3.5, we established that  $\ell^1$  minimization accurately estimates sparse signals under deterministic noise, random noise, and even inexact sparsity. As we will see in this section, essentially the same analysis and results generalize to the case of nuclear norm minimization for recovering low-rank matrices.

##### *Deterministic Noise.*

Here we still assume the matrix  $\mathbf{X}_{\text{true}}$  is perfectly low-rank, but the measurements  $\mathbf{y}$  is corrupted by small additive noise:

$$\mathbf{y} = \mathcal{A}[\mathbf{X}_{\text{true}}] + \mathbf{z}, \quad \|\mathbf{z}\|_2 \leq \varepsilon. \quad (4.3.71)$$

Similar to Theorem 3.5.1, for recovering low-rank matrices with (deterministic) noise, we have the following result.

**Theorem 4.3.12** (Stable Low-rank Recovery via BPDN). *Suppose that  $\mathbf{y} = \mathcal{A}[\mathbf{X}_{\text{true}}] + \mathbf{z}$ , with  $\|\mathbf{z}\|_2 \leq \varepsilon$ , and let  $\text{rank}(\mathbf{X}_{\text{true}}) = r$ . If  $\delta_{4r}(\mathcal{A}) < \sqrt{2} - 1$ , then any solution  $\hat{\mathbf{X}}$  to the optimization problem*

$$\begin{aligned} & \text{minimize} && \|\mathbf{X}\|_* \\ & \text{subject to} && \|\mathcal{A}[\mathbf{X}] - \mathbf{y}\|_2 \leq \varepsilon. \end{aligned} \quad (4.3.72)$$

satisfies

$$\left\| \hat{\mathbf{X}} - \mathbf{X}_{\text{true}} \right\|_F \leq C\varepsilon. \quad (4.3.73)$$

Here,  $C$  is a numerical constant.

*Proof.* The proof of this theorem parallels that for Theorem 3.5.1 and we leave the details for the reader as an exercise (see Exercise 4.15).  $\square$

### Random Noise.

Now let us consider the case when the noise in the above measurement model (4.3.71) is random (Gaussian):

$$\mathbf{y} = \mathcal{A}[\mathbf{X}_{\text{true}}] + \mathbf{z}, \quad (4.3.74)$$

where entries of  $\mathbf{z}$  are random i.i.d. Gaussian  $\mathcal{N}(0, \frac{\sigma^2}{m})$ . Then we have the following theorem that parallels Theorem 3.5.3 for the  $\ell^1$  case.

**Theorem 4.3.13** (Stable Sparse Recovery via Lasso). *Suppose that  $\mathcal{A} \sim_{\text{iid}} \mathcal{N}(0, \frac{1}{m})$ , and  $\mathbf{y} = \mathcal{A}[\mathbf{X}_{\text{true}}] + \mathbf{z}$ , with  $\mathbf{X}_{\text{true}}$  of rank  $r$  and  $\mathbf{z} \sim_{\text{iid}} \mathcal{N}(0, \frac{\sigma^2}{m})$ . Solve the matrix Lasso*

$$\text{minimize } \frac{1}{2} \|\mathbf{y} - \mathcal{A}[\mathbf{X}]\|_2^2 + \lambda_m \|\mathbf{X}\|_*, \quad (4.3.75)$$

with regularization parameter  $\lambda = c \cdot 2\sigma \sqrt{\frac{r(n_1+n_2)}{m}}$  for a large enough  $c$ . Then with high probability,

$$\left\| \hat{\mathbf{X}} - \mathbf{X}_{\text{true}} \right\|_F \leq C' \sigma \sqrt{\frac{r(n_1+n_2)}{m}}. \quad (4.3.76)$$

Notice in contrast to deterministic noise, random noise leads to a much favorable scaling  $\sqrt{\frac{r(n_1+n_2)}{m}}$  in the estimation error. In a typical compressive sensing setting (as suggested by Theorem 4.3.8), the sampling dimension  $m$  needs to be at least  $C \cdot r(n_1 + n_2)$  for some large constant  $C$ . Hence the scaling factor is proportional to  $1/\sqrt{C}$  and it becomes small when  $C$  is large.

*Proof.* The overall proof strategy is quite similar to that of Theorem 3.5.3 for the stability of Lasso estimate. We will lay out the key places that are different from the  $\ell^1$  case and leave the details to the reader an exercise.

In the proof of Lemma 3.5.2, we see that in order to establish the cone condition for the Lasso type minimization, one of the key steps is to bound  $|\langle \mathbf{A}^T \mathbf{z}, \mathbf{h} \rangle|$  via

$$|\langle \mathbf{A}^T \mathbf{z}, \mathbf{h} \rangle| \leq \left\| \mathbf{A}^T \mathbf{z} \right\|_\infty \|\mathbf{h}\|_1.$$

Following similar arguments, in the matrix Lasso case here, we need to bound  $|\langle \mathcal{A}^* \mathbf{z}, \mathbf{H} \rangle|$  instead as

$$|\langle \mathcal{A}^* \mathbf{z}, \mathbf{H} \rangle| \leq \|\mathcal{A}^* \mathbf{z}\| \|\mathbf{H}\|_*,$$

where  $\|\mathcal{A}^* \mathbf{z}\|$  is the operator norm (largest singular value) of the matrix  $\mathcal{A}^* \mathbf{z} = \sum_{i=1}^m z_i \mathbf{A}_i$ . To this end, we need to provide a tight bound for the operator norm of  $\mathcal{A}^* \mathbf{z}$ .

Notice that

$$M \doteq \left\| \sum_{i=1}^m z_i \mathbf{A}_i \right\| = \sup_{\mathbf{u} \in \mathbb{S}^{n_1-1}, \mathbf{v} \in \mathbb{S}^{n_2-1}} \mathbf{u}^T \sum_{i=1}^m z_i \mathbf{A}_i \mathbf{v} \quad (4.3.77)$$

$$= \sup_{\mathbf{u} \in \mathbb{S}^{n_1-1}, \mathbf{v} \in \mathbb{S}^{n_2-1}} \langle \mathbf{z}, \mathcal{A}[\mathbf{u}\mathbf{v}^T] \rangle. \quad (4.3.78)$$

The  $\mathbf{u}_*$  and  $\mathbf{v}_*$  that achieve the maximum value in (4.3.78) depend on  $\mathbf{z}$  and  $\mathcal{A}$ . So in order to eliminate this dependency and provide a bound for  $\|\sum_{i=1}^m z_i \mathbf{A}_i\|$ , we cover the two spheres  $\mathbb{S}^{n_1-1}$  and  $\mathbb{S}^{n_2-1}$  with two  $\varepsilon$ -nets  $N_1$  and  $N_2$  respectively. According to Lemma 3.4.5, the sizes of the nets can be less than  $(3/\varepsilon)^{n_1}$  and  $(3/\varepsilon)^{n_2}$  respectively.

Let us denote

$$M_N \doteq \sup_{\mathbf{u} \in N_1, \mathbf{v} \in N_2} \mathbf{u}^T \sum_{i=1}^m z_i \mathbf{A}_i \mathbf{v},$$

and then it is easy to show that<sup>10</sup>

$$M \leq \frac{M_N}{1 - 2\varepsilon}. \quad (4.3.79)$$

Notice that given any  $\mathbf{u} \in N_1, \mathbf{v} \in N_2$ ,  $\langle \mathbf{z}, \mathcal{A}[\mathbf{u}\mathbf{v}^T] \rangle$  is a Gaussian variable of distribution  $\mathcal{N}(0, \|\mathcal{A}[\mathbf{u}\mathbf{v}^T]\|_2^2 (\sigma^2/m))$ . Since  $\mathcal{A}$  is rank-RIP and  $\mathbf{u}\mathbf{v}^T$  is a rank-1 matrix of unit Frobenius norm, we have

$$\|\mathcal{A}[\mathbf{u}\mathbf{v}^T]\|_2^2 \leq (1 + \delta) \leq 2. \quad (4.3.80)$$

Thus, we have

$$\mathbb{P} \left[ \left| \mathbf{u}^T \sum_{i=1}^m z_i \mathbf{A}_i \mathbf{v} \right| > t \right] \leq 2 \exp \left( -\frac{mt^2}{4\sigma^2} \right). \quad (4.3.81)$$

Apply the union bound on all possible pairs of  $(\mathbf{u}, \mathbf{v})$  from the two nets and choose  $t = \alpha\sigma\sqrt{\frac{n_1+n_2}{m}}$  for some large enough  $\alpha$ , then we have  $M_N > t$  with diminishing probability as  $n_1$  or  $n_2$  becomes large. Therefore, we have

$$M = \left\| \sum_{i=1}^m z_i \mathbf{A}_i \right\| \leq \beta\sigma\sqrt{\frac{n_1+n_2}{m}} \quad (4.3.82)$$

for some constant  $\beta$  with high probability.

Now, similar to the proof of Lemma 3.5.2, if we choose  $\lambda_m$  to be in the order of  $O(\sigma\sqrt{\frac{n_1+n_2}{m}})$ , then the feasible perturbation  $\mathbf{H}$  satisfies the

<sup>10</sup>We leave the details of proving this inequality to the reader as an exercise.

cone restriction. Since  $\mathcal{A}$  is rank-RIP, it implies that  $\mathcal{A}$  satisfies the matrix restricted strong convexity (RSC) property. That leads to the bound on the estimation error:

$$\|\mathbf{H}\|_F = \|\hat{\mathbf{X}} - \mathbf{X}_{\text{true}}\|_F \leq C' \sigma \sqrt{\frac{r(n_1 + n_2)}{m}}. \quad (4.3.83)$$

The details of the proof for this follows essentially the same steps as those in the proof of Theorem 3.5.3 for the  $\ell^1$  case. We leave those to the reader as an exercise (see Exercise 4.17.)  $\square$

The error bound given in the above theorem can actually be shown to be nearly optimal as it is close to the best error that can be achieved by any estimator over all rank- $r$  matrices. The following theorem, due to [Candes and Plan, 2011] makes this precise:

**Theorem 4.3.14.** *Suppose that  $\mathcal{A} \sim_{\text{iid}} \mathcal{N}(0, \frac{1}{m})$  and we observe  $\mathbf{y} = \mathcal{A}[\mathbf{X}_{\text{true}}] + \mathbf{z}$  where entries of  $\mathbf{z}$  are i.i.d.  $\mathcal{N}(0, \frac{\sigma^2}{m})$  random variables. Set*

$$M^*(\mathcal{A}) = \inf_{\hat{\mathbf{X}}(\mathbf{y})} \sup_{\text{rank}(\mathbf{X}) \leq r} \mathbb{E} \left\| \hat{\mathbf{X}}(\mathbf{y}) - \mathbf{X} \right\|_F^2. \quad (4.3.84)$$

Then we have

$$M^*(\mathcal{A}) \geq c \sigma^2 \frac{rn}{m}, \quad (4.3.85)$$

for  $n = \max\{n_1, n_2\}$ , where  $c > 0$  is a numerical constant.

The proof of this theorem is beyond the scope of this book; we refer interested readers to [Candes and Plan, 2011] for a proof. According to Theorem 4.3.13, the worst error of the matrix Lasso matches the best achievable by any estimator, up to constants.

#### ***Inexact Low-rank Matrices.***

In the case when  $\mathbf{X}_{\text{true}}$  is not exactly low-rank, let  $[\mathbf{X}_{\text{true}}]_r$  be the best rank- $r$  approximation of  $\mathbf{X}_{\text{true}}$ . We can rewrite the observation model

$$\mathbf{y} = \mathcal{A}[\mathbf{X}_{\text{true}}] + \mathbf{z}, \quad \|\mathbf{z}\|_2 \leq \varepsilon. \quad (4.3.86)$$

as:

$$\mathbf{y} = \mathcal{A}[[\mathbf{X}_{\text{true}}]_r] + \mathcal{A}[\mathbf{X}_{\text{true}} - [\mathbf{X}_{\text{true}}]_r] + \mathbf{z}, \quad \|\mathbf{z}\|_2 \leq \varepsilon.$$

**Theorem 4.3.15.** *Let  $\mathbf{y} = \mathcal{A}[\mathbf{X}_{\text{true}}] + \mathbf{z}$ , with  $\|\mathbf{z}\|_2 \leq \varepsilon$ . Let  $\hat{\mathbf{X}}$  solve the denoising problem*

$$\begin{aligned} & \text{minimize} && \|\mathbf{X}\|_* \\ & \text{subject to} && \|\mathbf{y} - \mathcal{A}[\mathbf{X}]\|_2 \leq \varepsilon. \end{aligned} \quad (4.3.87)$$

Then for any  $r$  such that  $\delta_{4r}(\mathcal{A}) < \sqrt{2} - 1$ ,

$$\left\| \hat{\mathbf{X}} - \mathbf{X}_{\text{true}} \right\|_2 \leq C \frac{\|\mathbf{X}_{\text{true}} - [\mathbf{X}_{\text{true}}]_r\|_*}{\sqrt{r}} + C' \varepsilon \quad (4.3.88)$$

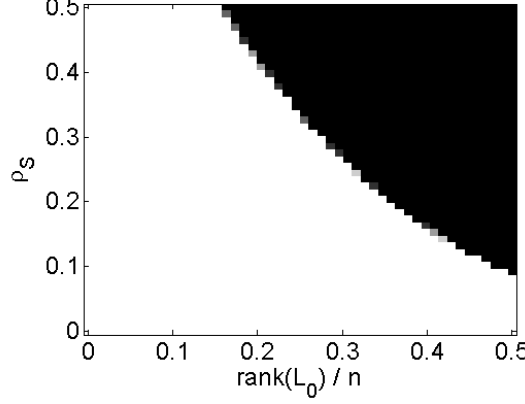


Figure 4.6. **Phase transitions in Low Rank Matrix Recovery.** We plot the probability of successfully recovering an  $n \times n$  matrix from Gaussian measurements. Vertical axis: undersampling factor  $1 - m/n^2$ . Horizontal axis: rank  $r/n$ . The success of nuclear norm minimization exhibits a very sharp transition from success to failure.

for some constants  $C$  and  $C'$ .

*Proof.* The proof of this theorem parallels that for Theorem 3.5.5 for the inexact sparse recovery problem. We here only setup some analogous concepts and key ideas that allow us to extend that proof to the matrix case here. But we leave details of the proof as an exercise to the reader.

Let  $\mathbf{X}_{\text{true}} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$  denote the compact SVD of the true solution  $\mathbf{X}_{\text{true}}$ . Then its best rank- $r$  approximation is  $[\mathbf{X}_{\text{true}}]_r = \mathbf{U}_r\mathbf{\Sigma}_r\mathbf{V}_r^*$ . Now let

$$T \doteq \{\mathbf{U}_r\mathbf{W}^* + \mathbf{Q}\mathbf{V}_r^* \mid \mathbf{W} \in \mathbb{R}^{n_2 \times r}, \mathbf{Q} \in \mathbb{R}^{n_1 \times r}\} \subseteq \mathbb{R}^{n_1 \times n_2}. \quad (4.3.89)$$

Show that in the inexact low-rank case, instead of the cone restriction (4.3.44), we have the following restriction for the feasible perturbation  $\mathbf{H} = \hat{\mathbf{X}} - \mathbf{X}_{\text{true}}$ :

$$\|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_* \leq \|\mathcal{P}_T[\mathbf{H}]\|_* + 2\|\mathcal{P}_{T^\perp}[\mathbf{X}_{\text{true}}]\|_*. \quad (4.3.90)$$

Notice that  $\mathcal{P}_{T^\perp}[\mathbf{X}_{\text{true}}] = \mathbf{X}_{\text{true}} - [\mathbf{X}_{\text{true}}]_r$ . Then, similar to the proof of Theorem 3.5.5, simply carry the extra term  $2\|\mathcal{P}_{T^\perp}[\mathbf{X}_{\text{true}}]\|_*$  at the places in the proof of Theorem 4.3.7 where the cone restriction is applied. One can reach the conclusion of the theorem. We leave details of the proof as an exercise to the reader (see Exercise 4.16).  $\square$

#### **Phase Transition for Low Rank Matrix Recovery.**

Thus far, we have seen strong parallels between sparse vector recovery using  $\ell^1$  norm minimization and low-rank matrix recovery using nuclear norm minimization. In both cases, we saw how an appropriate notion of

restricted isometry property could be used to guarantee exact recovery from a near-minimal number of random measurements – about  $k \log(n/k)$  for  $k$ -sparse vectors, and about  $nr$  for rank- $r$  matrices. However, just like in the sparse vector case, this tool does not yield sharp constants.

In fact, there is a phase transition phenomenon for low-rank recovery, which mirrors that for sparse recovery: as the dimension grows, the transition between success and failure in low-rank recovery becomes increasingly sharp. Figure 4.6 illustrates this.

Just as we did for sparse recovery, we can use the “coefficient space” geometry of the low-rank recovery problem to derive very sharp estimates of this transition. This geometry is phrased in terms of the descent cone  $\mathcal{D}$  of the nuclear norm at the target solution  $\mathbf{X}_{\text{true}}$ :

$$\mathcal{D} = \{\mathbf{H} \mid \|\mathbf{X}_{\text{true}} + \mathbf{H}\|_* \leq \|\mathbf{X}_{\text{true}}\|_*\}. \quad (4.3.91)$$

As for sparse recovery,  $\mathbf{X}_{\text{true}}$  is the unique optimal solution to the nuclear norm minimization problem if and only if  $\mathcal{D} \cap \text{null}(\mathcal{A}) = \{\mathbf{0}\}$ . Hence, quantifying the probability of success under a random linear projection becomes equivalent to quantifying the probability that the two convex cones  $\mathcal{D}$  and  $\text{null}(\mathcal{A})$  have only trivial intersection. Deploying Theorem 3.6.4, we find that there is a sharp transition between success and failure around

$$m^* \sim \delta(\mathcal{D}), \quad (4.3.92)$$

the statistical dimension of the descent cone. Moreover, the theorem tells us that the width of the transition region is roughly  $O(\sqrt{n_1 n_2})$ . The *location* of the transition region can be characterized using the same machinery that we deployed in Section 3.6 to estimate the statistical dimension of the descent cone of the  $\ell^1$  norm. This machinery involves estimating the expected squared distance of a random vector (here, random matrix) to the polar cone, which is spanned by the subdifferential of the nuclear norm. For convenience, for a matrix  $\mathbf{M}$  with singular value decomposition  $\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ , let us define the *singular value thresholding operator* as

$$\mathcal{D}_\tau[\mathbf{M}] \doteq \mathbf{U}\mathcal{S}_\tau[\mathbf{\Sigma}]\mathbf{V}^*, \quad (4.3.93)$$

where  $\mathcal{S}_\tau[\cdot]$  is the entry-wise *soft thresholding* operator:

$$\forall \mathbf{X}, \quad \mathcal{S}_\tau[\mathbf{X}] = \text{sign}(\mathbf{X}) \circ (|\mathbf{X}| - \tau)_+,$$

where  $\circ$  is the entry-wise (Hadamard) product of two matrices. An intermediate result produced by these calculations is as follows:

**Theorem 4.3.16** (Phase Transition in Low-rank Recovery). *Let  $\mathcal{D}$  denote the descent cone of the nuclear norm at any matrix  $\mathbf{X}_{\text{true}} \in \mathbb{R}^{n_1 \times n_2}$  of rank  $r$ . Let  $\mathbf{G}$  be an  $(n_1 - r) \times (n_2 - r)$  matrix with entries i.i.d.  $\mathcal{N}(0, 1)$ . Set*

$$\psi(n_1, n_2, r) = \inf_{\tau \geq 0} \left\{ r(n_1 + n_2 - r + \tau^2) + \mathbb{E}_{\mathbf{G}} \left[ \|\mathcal{D}_\tau[\mathbf{G}]\|_F^2 \right] \right\}. \quad (4.3.94)$$

Then

$$\psi(n_1, n_2, r) - 2\sqrt{n_2/r} \leq \delta(\mathbf{D}) \leq \psi(n_1, n_2, r). \quad (4.3.95)$$

This theorem identifies a sharp transition in low-rank recovery. It is possible to use asymptotic results on the limiting distribution of the singular values of a random matrix to give a formula for  $\psi(n_1, n_2, r)/(n_1 n_2)$ , which is valid when  $n_1 \rightarrow \infty$ ,  $n_1/n_2 \rightarrow \alpha \in (0, \infty)$  and  $r/n_1 \rightarrow \rho \in (0, 1)$ . In the exercises, we guide the interested reader through this derivation. Here, we merely display the result of this calculation in Figure 4.6, and note the excellent agreement between this theoretical prediction and numerical experiment: *for the idealized setting of “generic” measurements, we have a very precise prediction of the phase transition!*

## 4.4 Low Rank Matrix Completion

We have seen how concepts from sparse recovery transpose directly to the low-rank recovery problem. The concept of sparsity had a natural analogue in the concept of rank-deficiency. The  $\ell^1$  minimization problem for sparse recovery had a natural analogue in the nuclear norm minimization problem for low-rank recovery. Moreover, these convex relaxations succeed under analogous conditions involving restricted isometry properties of the observation operator.

However, in many of the most interesting applications of nuclear norm minimization, the RIP does not hold! In the introduction to this chapter, we sketched applications to recommender systems, in which we had access to *a subset* of the entries of a low-rank user-item matrix. We also sketched problems in reconstructing 3D shape, in which we observed *a subset* of the pixels of the rank-3 matrix  $\mathbf{NL}$ . Finally, we sketched a problem in Euclidean embedding, in which we observe *a subset* of the distances between some objects of interest. In all of these problems, the object of interest is a low-rank matrix  $\mathbf{X}_{\text{true}} \in \mathbb{R}^{n \times n}$ ; the observation selects a subset  $\Omega \subset [n] \times [n]$  of the entries of  $\mathbf{X}_{\text{true}}$ . The *matrix completion* asks us to fill in the missing entries:

**Problem 4.4.1** (Matrix Completion). *Let  $\mathbf{X}_{\text{true}} \in \mathbb{R}^{n \times n}$  be a low-rank matrix. Suppose we are given  $\mathbf{y} = \mathcal{P}_\Omega[\mathbf{X}_{\text{true}}]$ , where  $\Omega \subseteq [n] \times [n]$ . Fill in the missing entries of  $\mathbf{X}_{\text{true}}$ .*

In matrix completion  $\mathcal{A} = \mathcal{P}_\Omega$  is the restriction onto some small subset  $\Omega \subseteq [n] \times [n]$  of the entries. In this situation, if  $(i, j) \notin \Omega$ ,  $\mathcal{P}_\Omega[\mathbf{E}_{ij}] = \mathbf{0}$ . That is to say, if  $\Omega$  is a strict subset of  $[n] \times [n]$ , then  $\mathcal{P}_\Omega$  has matrices of rank one in its null space! So, the rank-RIP cannot hold for any positive rank  $r$  with any nontrivial  $\delta < 1$ .

At a more basic level, the example of  $\mathbf{X}_{\text{true}} = \mathbf{E}_{ij}$  suggests that there are some (very sparse) matrices that are impossible to complete from only a few



entries. This is in contrast to our discussion of low-rank matrix recovery thus far, in which the only factor that dictates the ease or difficulty of recovering a target  $\mathbf{X}_{\text{true}}$  is the complexity  $\text{rank}(\mathbf{X}_{\text{true}})$ . Nevertheless, our development thus far suggests that even for the more challenging problem of matrix completion, there may be some class of *well-structured* matrices  $\mathbf{X}_{\text{true}}$  of interest for applications, which *can* be efficiently completed from just a few entries. In this section, we will see that this is indeed the case.

#### 4.4.1 Nuclear Norm Minimization for Matrix Completion

In light of our previous study of matrix recovery, a natural approach to completing a low-rank matrix from a small subset  $\mathbf{Y} = \mathcal{P}_{\Omega}[\mathbf{X}_{\text{true}}]$  of its entries is to look for the matrix  $\mathbf{X}$  of minimum nuclear norm that agrees with the observation:

$$\begin{aligned} & \text{minimize} && \|\mathbf{X}\|_* \\ & \text{subject to} && \mathcal{P}_{\Omega}[\mathbf{X}] = \mathbf{Y}. \end{aligned} \quad (4.4.1)$$

This is a special instance of the general nuclear norm minimization problem (4.3.14), with observation operator  $\mathcal{A} = \mathcal{P}_{\Omega}$ . As such, it is a semidefinite program, and can be solved to high accuracy in polynomial time. In practice, though, it is more important to have methods that scale to large problem instances. In the next section, we sketch one approach to achieving this, using Lagrange multiplier techniques. This approach has pedagogical value: it introduces several objects that will be used for analyzing when we can solve matrix completion problems efficiently. It also yields reasonably scalable algorithms. For practical matrix completion at the scale of  $n \sim 10^6$  and beyond, even more scalable methods are needed; we discuss these issues in Chapters 8-9.

#### 4.4.2 Algorithm via Augmented Lagrangian Method

There are two basic challenges in solving problem (4.4.1) at large scale. The first arises from the nonsmoothness of the nuclear norm  $\|\cdot\|_*$ ; the second is due to the need to satisfy the constraint  $\mathcal{P}_{\Omega}[\mathbf{X}] = \mathbf{Y}$  exactly.<sup>11</sup> The fundamental technology for handling constraints in optimization is Lagrange duality.

The basic object is the *Lagrangian*, which introduces a matrix  $\mathbf{\Lambda}$  of Lagrange multipliers for the constraint  $\mathcal{P}_{\Omega}[\mathbf{X}] = \mathbf{Y}$ . The Lagrangian for (4.4.1) is

$$\mathcal{L}(\mathbf{X}, \mathbf{\Lambda}) = \|\mathbf{X}\|_* + \langle \mathbf{\Lambda}, \mathcal{P}_{\Omega}[\mathbf{Y}] - \mathcal{P}_{\Omega}[\mathbf{X}] \rangle. \quad (4.4.2)$$

<sup>11</sup>In practice, when observations are noisy, exactly satisfying  $\mathcal{P}_{\Omega}[\mathbf{X}] = \mathbf{Y}$  is neither necessary nor desirable. We study the noisy matrix completion in Section 4.4.5, and develop dedicated algorithms for it in Chapter 8.

Optimal  $\mathbf{X}$  are characterized by *saddle points* at which the Lagrangian is *minimized* with respect to  $\mathbf{X}$ , and maximized with respect to  $\mathbf{\Lambda}$ . A basic approach to solving a constrained problem such as (4.4.1) is to seek such a saddle point. In practice, more robustly convergent algorithms can be derived by instead working with the *augmented* Lagrangian

$$\mathcal{L}_\mu(\mathbf{X}, \mathbf{\Lambda}) = \|\mathbf{X}\|_* + \langle \mathbf{\Lambda}, \mathcal{P}_\Omega[\mathbf{Y}] - \mathcal{P}_\Omega[\mathbf{X}] \rangle + \frac{\mu}{2} \|\mathcal{P}_\Omega[\mathbf{Y}] - \mathcal{P}_\Omega[\mathbf{X}]\|_F^2, \quad (4.4.3)$$

which encourages satisfaction of the constraint by adding an additional quadratic penalty term  $\frac{\mu}{2} \|\mathcal{P}_\Omega[\mathbf{Y}] - \mathcal{P}_\Omega[\mathbf{X}]\|_F^2$ . A more general introduction to augmented Lagrangian method (ALM) is given in section sec:convex:ALM of Chapter 8.

The augmented Lagrangian method seeks a saddle point of  $\mathcal{L}_\mu$  by alternating between minimizing with respect to the “primal variables”  $\mathbf{X}$  and taking one step of gradient ascent to increase  $\mathcal{L}_\mu$  using the “dual variables”  $\mathbf{\Lambda}$ :

$$\mathbf{X}^{(k+1)} \in \arg \min_{\mathbf{X}} \mathcal{L}_\mu(\mathbf{X}, \mathbf{\Lambda}^{(k)}), \quad (4.4.4)$$

$$\mathbf{\Lambda}^{(k+1)} = \mathbf{\Lambda}^{(k)} + \mu \mathcal{P}_\Omega[\mathbf{Y} - \mathbf{X}^{(k+1)}]. \quad (4.4.5)$$

Here,  $\mathcal{P}_\Omega[\mathbf{Y} - \mathbf{X}^{(k+1)}] = \nabla_{\mathbf{\Lambda}} \mathcal{L}_\mu(\mathbf{X}^{(k+1)}, \mathbf{\Lambda})$ . The ALM algorithm makes a very special choice of the step size ( $\mu$ ) for updating  $\mathbf{\Lambda}$ . This choice is important in general: it ensures that  $\mathbf{\Lambda}$  stays dual feasible, an issue that we will explain in more depth in section 8.4 of Chapter 8.

Under very general conditions, the iteration (4.4.4)-(4.4.5) converges to a primal dual optimal pair  $(\mathbf{X}_*, \mathbf{\Lambda}_*)$ , and hence yields a solution to (4.4.1). While this algorithm appears simple, some caution is necessary: the first step is itself a nontrivial optimization problem! This subproblem has a characteristic form, which we encountered in our study of sparse recovery in noise: the objective function is a sum of a smooth convex term  $f(\mathbf{X})$ , and a nonsmooth convex function  $g(\mathbf{X}) = \|\mathbf{X}\|_*$ :

$$\min_{\mathbf{X}} \underbrace{\|\mathbf{X}\|_*}_{g(\mathbf{X}) \text{ convex}} + \underbrace{\langle \mathbf{\Lambda}, \mathcal{P}_\Omega[\mathbf{Y}] - \mathcal{P}_\Omega[\mathbf{X}] \rangle + \frac{\mu}{2} \|\mathcal{P}_\Omega[\mathbf{Y}] - \mathcal{P}_\Omega[\mathbf{X}]\|_F^2}_{f(\mathbf{X}) \text{ smooth, convex}}. \quad (4.4.6)$$

Here,

$$\nabla f(\mathbf{X}) = -\mathcal{P}_\Omega[\mathbf{\Lambda}] + \mu \mathcal{P}_\Omega[\mathbf{X} - \mathbf{Y}]. \quad (4.4.7)$$

This is  $\mu$ -Lipschitz, in the sense that for any pair of matrices  $\mathbf{X}$  and  $\mathbf{X}'$ ,

$$\|\nabla f(\mathbf{X}) - \nabla f(\mathbf{X}')\|_F \leq \mu \|\mathbf{X} - \mathbf{X}'\|_F. \quad (4.4.8)$$

This class of problem is amenable to the *proximal gradient method*.

The general proximal gradient iteration applies to objectives of the form  $F(\mathbf{X}) = g(\mathbf{X}) + f(\mathbf{X})$ , where  $g$  is convex, and  $f$  is convex, smooth, and has  $L$ -Lipschitz gradient. See section 8.2 of Chapter 8. Here we have the

---

**Algorithm 4.1 (Matrix Completion by Augmented Lagrange Multiplier (ALM))**


---

```

1: initialize:  $\mathbf{X}^{(0)} = \mathbf{\Lambda}^{(0)} = 0, \mu > 0$ .
2: while not converged do
3:   compute  $\mathbf{X}^{(k+1)} \in \arg \min_{\mathbf{X}} \mathcal{L}_{\mu}(\mathbf{X}, \mathbf{\Lambda}^{(k)})$ ;
4:   compute  $\mathbf{\Lambda}^{(k+1)} = \mathbf{\Lambda}^{(k)} + \mu(\mathcal{P}_{\Omega}[\mathbf{Y}] - \mathcal{P}_{\Omega}[\mathbf{X}^{(k+1)}])$ ;
5: end while

```

---

Lipschitz constant  $L = \mu$ . So the iteration takes the form

$$\mathbf{X}^{(k+1)} = \arg \min_{\mathbf{X}} \left\{ g(\mathbf{X}) + \frac{\mu}{2} \left\| \mathbf{X} - \left( \mathbf{X}^{(k)} - \frac{1}{\mu} \nabla f(\mathbf{X}^{(k)}) \right) \right\|_F^2 \right\}. \quad (4.4.9)$$

In particular, it requires us to solve a sequence of “proximal problems”

$$\min_{\mathbf{X}} \left\{ g(\mathbf{X}) + \frac{\mu}{2} \|\mathbf{X} - \mathbf{M}\|_F^2 \right\}, \quad (4.4.10)$$

for particular choices of the matrix  $\mathbf{M}$ . When  $g$  is the nuclear norm, this problem can be solved in closed form from the SVD of  $\mathbf{M}$ . Recall from (4.3.93), for a matrix  $\mathbf{M}$  with the singular value decomposition  $\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ , its singular value thresholding operator is defined to be

$$\mathcal{D}_{\tau}[\mathbf{M}] = \mathbf{U}\mathcal{S}_{\tau}[\mathbf{\Sigma}]\mathbf{V}^*,$$

where  $\mathcal{S}_{\tau}[\mathbf{X}] = \text{sign}(\mathbf{X}) \circ (|\mathbf{X}| - \tau)_+$  is the soft thresholding operator.

**Theorem 4.4.2.** *The unique solution  $\mathbf{X}_{\star}$  to the program:*

$$\min_{\mathbf{X}} \left\{ \|\mathbf{X}\|_* + \frac{\mu}{2} \|\mathbf{X} - \mathbf{M}\|_F^2 \right\}, \quad (4.4.11)$$

*is given by*

$$\mathbf{X}_{\star} = \mathcal{D}_{\mu^{-1}}[\mathbf{M}]. \quad (4.4.12)$$

The proof of this result follows from Exercise 4.11. The resulting procedures are stated as Algorithms 4.1-4.2. We have neglected important issues such as the choice of stopping conditions, and the effect of inexact solution to subproblem (4.4.4) on the convergence of the basic ALM iteration in Algorithm 4.1.

To understand when the convex program (4.4.1) and the above algorithm correctly recover a matrix  $\mathbf{X}$  from a part of its entries, we vary the rank of matrix  $\mathbf{X}$  as a fraction of the dimension  $n$  and a fraction  $\rho_s \in (0, 1)$  of (randomly chosen) unobserved entries. In other words,  $\rho_s$  is the probability that an entry is omitted from the observation. Figure 4.7 shows the simulation results of using the above algorithm to recover a matrix  $\mathbf{X} = \mathbf{L}_0$  under different settings.

**Algorithm 4.2 (Minimizing the Augmented Lagrangian via Proximal Gradient)**

- 
- 1: **initialize:**  $\mathbf{X}^{(0)}$  starts with the  $\mathbf{X}^{(k)}$  from the outer loop.
  - 2: **while** not converged **do**
  - 3:   compute

$$\begin{aligned}\mathbf{X}^{(\ell+1)} &= \text{prox}_{g/\mu}(\mathbf{X}^{(\ell)} - \mu^{-1} \nabla f(\mathbf{X}^{(\ell)})) \\ &= \mathcal{D}_{\mu^{-1}} \left[ \mathcal{P}_{\Omega^c}[\mathbf{X}^{(\ell)}] + \mathcal{P}_{\Omega}[\mathbf{Y}] + \mu^{-1} \mathcal{P}_{\Omega}[\mathbf{\Lambda}^{(k)}] \right].\end{aligned}$$

- 4: **end while**
- 

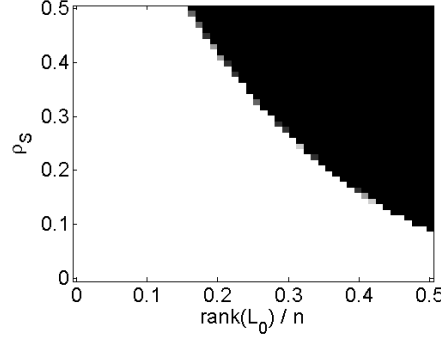


Figure 4.7. **Correct matrix completion for varying rank and sparsity.** Fraction of correct recoveries across 10 trials, as a function of  $\text{rank}(\mathbf{L}_0)$  (x-axis) and fraction  $\rho_s$  of unobserved entries (y-axis). Here,  $n = 400$ . In all cases,  $\mathbf{L}_0 = \mathbf{A}\mathbf{B}^*$  is a product of independent  $n \times r$  i.i.d.  $\mathcal{N}(0, 1/n)$  matrices. Trials are considered successful if  $\|\hat{\mathbf{X}} - \mathbf{L}_0\|_F / \|\mathbf{L}_0\|_F < 10^{-3}$ .

We may draw a few observations from the above simulations: 1. the convex program (4.4.1) and the above algorithm indeed succeed under a surprisingly wide range of conditions, as long as the rank of the matrix is relatively low and a fraction of the entries are observed. 2. the success and failure of the convex program (4.4.1) exhibit a sharp phase-transition phenomenon.

#### 4.4.3 When Nuclear Norm Minimization Succeed?

The above simulations encourage us to understand the conditions under which nuclear norm minimization (4.4.1) is guaranteed to succeed for matrix completion?<sup>12</sup> It may be easier to first think about when it fails. It may

---

<sup>12</sup>Or ultimately, if possible, to precisely characterize the phase transition behavior we have observed through experiments.

fail if (i)  $\mathbf{X}_0$  is *sparse* (as in the example of  $\mathbf{E}_{ij}$ ), or (ii) if the sampling pattern  $\Omega$  is chosen adversarially (e.g., if we miss an entire row or column of  $\mathbf{X}_0$ ). Below, we will state a theorem that makes this intuition precise – namely, if  $\mathbf{X}_0$  is low-rank, and not too “spiky”, and  $\Omega$  is chosen at random, then nuclear norm minimization succeeds with high probability. Below we make these assumptions precise.

***Incoherent Low-rank Matrices.***

Although our intuition is that  $\mathbf{X}_0$  itself should not be too “sparse”, for technical reasons it will be necessary to enforce this condition on the singular vectors of  $\mathbf{X}_0$ , rather than on  $\mathbf{X}_0$  itself. Let  $\mathbf{X}_0 = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$  be the (reduced) singular value decomposition of  $\mathbf{X}_0$ . We say that  $\mathbf{X}_0$  is  $\nu$ -*incoherent* if the following hold:

$$\forall i \in [n], \quad \|\mathbf{e}_i^* \mathbf{U}\|_2^2 \leq \nu r/n, \quad (4.4.13)$$

$$\forall j \in [n], \quad \|\mathbf{e}_j^* \mathbf{V}\|_2^2 \leq \nu r/n. \quad (4.4.14)$$

These two conditions control the “spikiness” of the singular vectors of  $\mathbf{X}_0$ . To understand them better, note that  $\mathbf{U}$  is an  $n \times r$  matrix whose columns have unit  $\ell^2$  norm. Hence,  $\sum_i \|\mathbf{e}_i^* \mathbf{U}\|_2^2 = \|\mathbf{U}\|_F^2 = r$ . There are  $n$  rows, and so at least one of them must have  $\ell^2$  norm at least as large as the average,  $r/n$ . Hence, for any matrix  $\mathbf{U}$  with unit norm columns,  $\max_i \|\mathbf{e}_i^* \mathbf{U}\|_2^2 \geq r/n$ . The *incoherence parameter*  $\nu$  quantifies how much we lose with respect to this optimal bound. So, if  $\nu$  is small, the singular vectors are, in a sense, spread around. To give a sense of scale, notice that it is always true that

$$1 \leq \nu \leq n/r. \quad (4.4.15)$$

If  $\mathbf{U}$  and  $\mathbf{V}$  are chosen uniformly at random (say by orthogonalizing the columns of a Gaussian matrix), then with high probability  $\nu$  is bounded by  $C \log(n)$ . However, the definition does not require  $\mathbf{U}$  and  $\mathbf{V}$  to be random.

One important implication of this definition for matrix completion is that *when  $\nu$  is small, there are no sparse matrices close to the tangent space  $T$* . Indeed, let  $\mathbf{E}_{ij} = \mathbf{e}_i \mathbf{e}_j^*$  denote the one-sparse matrix whose nonzero element occurs in entry  $(i, j)$ . Then, using the expression (4.3.33) for the projection operator  $\mathcal{P}_T$  onto the tangent space  $T$ , we have that

$$\begin{aligned} \|\mathcal{P}_T[\mathbf{E}_{ij}]\|_F^2 &= \|\mathbf{U}\mathbf{U}^* \mathbf{E}_{ij}\|_F^2 + \|(\mathbf{I} - \mathbf{U}\mathbf{U}^*) \mathbf{E}_{ij} \mathbf{V}\mathbf{V}^*\|_F^2 \\ &\leq \|\mathbf{U}^* \mathbf{e}_i\|_2^2 + \|\mathbf{e}_j^* \mathbf{V}\|_2^2 \\ &\leq \frac{2\nu r}{n}. \end{aligned} \quad (4.4.16)$$

This indicates that no standard basis matrix  $\mathbf{E}_{ij}$  is too close to the subspace  $T$ . Strangely enough, this implies the standard basis  $\{\mathbf{E}_{ij}\}$  is a good choice for reconstructing elements from  $T$ . This is similar in spirit to our observations on incoherent operator bases: if no  $\mathbf{E}_{ij}$  is too close to  $T$ , information about any particular element  $\mathbf{X}_0 \in T$  must be spread across many different

$\mathbf{E}_{ij}$ . It will only take a few of these projections to be able to reconstruct  $\mathbf{X}_0$ . Note, however, a crucial difference between this notion of incoherence and our previous notions for matrix and vector recovery: here, the subspace  $T$  depends on  $\mathbf{X}_0$  itself. The discussion in this section suggests that random sampling will be effective for reconstructing the particular matrix  $\mathbf{X}_0$ . We make this intuition formal below.

***Exact Matrix Completion from Random Samples.***

We assume that each entry  $(i, j)$  belongs to the set  $\Omega$  independently with probability  $p$ . We call this a *Bernoulli* sampling model, since the indicators  $\mathbb{1}_{(i,j) \in \Omega}$  are independent  $\text{Ber}(p)$  random variables. Under this model, the expected number of observed entries is

$$m = \mathbb{E}[|\Omega|] = pn^2. \quad (4.4.17)$$

Under this model, nuclear norm minimization succeeds even when the number  $m$  of observations is close to the number of intrinsic degrees of freedom in the rank- $r$  matrix  $\mathbf{X}_0$ . The following theorem makes this precise:

**Theorem 4.4.3** (Matrix Completion via Nuclear Norm Minimization). *Let  $\mathbf{X}_0 \in \mathbb{R}^{n \times n}$  be a rank- $r$  matrix with coherence  $\nu$ . Suppose that we observe  $\mathbf{y} = \mathcal{P}_\Omega[\mathbf{X}_0]$ , with  $\Omega$  sampled according to the Bernoulli model with probability*

$$p \geq C_1 \frac{\nu r \log^2(n)}{n}. \quad (4.4.18)$$

*Then with probability at least  $1 - C_2 n^{-c_3}$ ,  $\mathbf{X}_0$  is the unique optimal solution to*

$$\text{minimize } \|\mathbf{X}\|_* \quad \text{subject to } \mathcal{P}_\Omega[\mathbf{X}] = \mathbf{y}. \quad (4.4.19)$$

There are several things to notice about the above theorem. First, the expected number of measurements is

$$m = pn^2 = C_1 \nu n r \log^2(n). \quad (4.4.20)$$

Since a rank- $r$  matrix has  $O(nr)$  degrees of freedom, the oversampling factor is only about  $C\nu \log^2(n)$  – the number of samples we must see is nearly minimal.<sup>13</sup> Second, the number of samples required scales with the coherence of the matrix  $\mathbf{X}_0$ . So, if we want to recover a very coherent (think, “nearly sparse”)  $\mathbf{X}_0$ , we will simply need more observations. Finally, the probability of success is in all the possible choices of the observed subset but

---

<sup>13</sup>According to Theorem 1.7 of [Candes and Tao, 2009], if the sampling probability  $p < \frac{\nu r \log(2n)}{2n}$ , there will be infinitely many matrices of rank at most  $r$  that satisfy the incoherence condition and all have the same entries on  $\Omega$ .

is only for a given low-rank matrix  $\mathbf{X}_0$ .<sup>14</sup> Of course, the precise conditions of the above theorem can only be interpreted as an idealized mathematical abstraction of real matrix completion or collaborative filtering problems. In particular, in real problems there may be noise in the observation, and, more importantly the observations may not be uniformly distributed. This has been a very intense area of recent investigation.

#### 4.4.4 Proving Correctness of Nuclear Norm Minimization

In this section, we prove Theorem 4.4.3. Our approach is analogous to our proof that  $\ell^1$  recovers sparse vectors under incoherence (in Section 3.2.2) – we simply write down the optimality conditions and try to show that they are satisfied! Carrying this program through will be trickier, though.

##### *Subdifferential of the Nuclear Norm.*

To get started, we need an optimality condition for the nuclear norm minimization problem (4.4.19). As for the  $\ell^1$  norm, this means we need an expression for the subdifferential. The following lemma provides one:

**Lemma 4.4.4.** *Let  $\mathbf{X} \in \mathbb{R}^{n \times n}$  have compact singular value decomposition  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ . The subdifferential of the nuclear norm at  $\mathbf{X}$  is given by*

$$\partial \|\cdot\|_*(\mathbf{X}) = \{\mathbf{Z} \mid \mathcal{P}_T[\mathbf{Z}] = \mathbf{U}\mathbf{V}^*, \|\mathcal{P}_{T^\perp}[\mathbf{Z}]\| \leq 1\}. \quad (4.4.21)$$

*Proof.* Consider any  $\mathbf{Z}$  satisfying  $\mathcal{P}_T[\mathbf{Z}] = \mathbf{U}\mathbf{V}^*$ , and  $\|\mathcal{P}_{T^\perp}[\mathbf{Z}]\| \leq 1$ . Notice that  $\|\mathbf{Z}\| = 1$ . Since  $\mathbf{X} \in T$ ,

$$\langle \mathbf{X}, \mathbf{Z} \rangle = \langle \mathbf{X}, \mathbf{U}\mathbf{V}^* \rangle = \langle \mathbf{U}^* \mathbf{X} \mathbf{V}, \mathbf{I} \rangle = \langle \mathbf{\Sigma}, \mathbf{I} \rangle = \|\mathbf{X}\|_*. \quad (4.4.22)$$

For every  $\mathbf{X}'$ ,

$$\|\mathbf{X}\|_* + \langle \mathbf{Z}, \mathbf{X}' - \mathbf{X} \rangle = \langle \mathbf{Z}, \mathbf{X}' \rangle \leq \|\mathbf{Z}\| \|\mathbf{X}'\|_* = \|\mathbf{X}'\|_*. \quad (4.4.23)$$

Thus  $\mathbf{Z}$  is a subgradient of the nuclear norm at  $\mathbf{X}$ :  $\mathbf{Z} \in \partial \|\cdot\|_*(\mathbf{X})$ . To complete the proof, we need to show that every element  $\mathbf{Z} \in \partial \|\cdot\|_*(\mathbf{X})$  satisfies  $\mathcal{P}_T[\mathbf{Z}] = \mathbf{U}\mathbf{V}^*$  and  $\|\mathcal{P}_{T^\perp}[\mathbf{Z}]\| \leq 1$ . We leave the converse as an exercise (see Exercise 4.18).  $\square$

If we compare to the expression for the subdifferential of the  $\ell^1$  norm, here, the subspace  $T$  plays the role of the *support* of the matrix, while the matrix  $\mathbf{U}\mathbf{V}^*$  is playing the role of the *signs*. Indeed, in this language,  $\partial \|\cdot\|_*$  consists of those  $\mathbf{Z}$  that are equal to the “sign”  $\mathbf{U}\mathbf{V}^*$  on the support  $T$ , and whose dual norm  $\|\cdot\|$  is bounded by one on the orthogonal complement  $T^\perp$  of the support.

<sup>14</sup>This is to contrast with the probability of success in the generic case, where an incoherent sampling operator is good for recovering the set of all matrices of rank less than  $r$ .

**Optimality Conditions.**

Once we have the subdifferential in hand, we can fairly immediately write down an optimality condition for the convex program of interest. Indeed, consider the optimization problem

$$\begin{aligned} & \text{minimize} && \|\mathbf{X}\|_* \\ & \text{subject to} && \mathcal{P}_\Omega[\mathbf{X}] = \mathcal{P}_\Omega[\mathbf{X}_0]. \end{aligned} \quad (4.4.24)$$

Any feasible  $\mathbf{X}$  can be written as  $\mathbf{X}_0 + \mathbf{H}$ , where  $\mathbf{H} \in \text{null}(\mathcal{P}_\Omega)$ , i.e.,  $\mathbf{H}$  is supported on the set  $\Omega^c$  of entries that we do not observe. Similar to the  $\ell^1$  case in Section 3.2.2, if we can find some  $\mathbf{\Lambda}$ , called a dual certificate, such that it satisfies (the KKT condition):

- (i)  $\mathbf{\Lambda}$  is supported on  $\Omega$  and
- (ii)  $\mathbf{\Lambda} \in \partial \|\cdot\|_*(\mathbf{X}_0)$  – i.e.,  $\mathcal{P}_T[\mathbf{\Lambda}] = \mathbf{U}\mathbf{V}^*$  and  $\|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}]\| \leq 1$ ,

then we have

$$\|\mathbf{X}_0 + \mathbf{H}\|_* \geq \|\mathbf{X}_0\|_* + \langle \mathbf{\Lambda}, \mathbf{H} \rangle = \|\mathbf{X}_0\|_*, \quad (4.4.25)$$

where the final equality holds because  $\mathbf{\Lambda}$  is supported on  $\Omega$  and  $\mathbf{H}$  is supported on  $\Omega^c$ . In addition, if we further have  $\|\mathcal{P}_{\Omega^c}\mathcal{P}_T\| < 1$  and  $\|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}]\| < 1$ , then one can show that  $\mathbf{X}_0$  is the *unique* optimal solution. The proof is similar to that in the  $\ell^1$  case (see the proof of Theorem 3.2.3) and we leave to the reader as an exercise (see Exercise 4.13).

A natural idea for constructing  $\mathbf{\Lambda}$  might be to simply follow the program that has worked before (in the  $\ell^1$  minimization case) and look for a matrix  $\mathbf{\Lambda}$  of smallest 2-norm that satisfies the equality constraints

$$\mathcal{P}_{\Omega^c}[\mathbf{\Lambda}] = \mathbf{0}, \quad \mathcal{P}_T[\mathbf{\Lambda}] = \mathbf{U}\mathbf{V}^*, \quad (4.4.26)$$

and then hope to check that it satisfies the inequality constraints  $\|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}]\| \leq 1$ . For example, we could take  $\mathbf{\Lambda} = \mathcal{P}_\Omega[\mathbf{G}]$ , with  $\mathbf{G} = (\mathcal{P}_T\mathcal{P}_\Omega)^\dagger[\mathbf{U}\mathbf{V}^*]$ , where  $(\cdot)^\dagger$  denotes the pseudo inverse. We are then left to check that

$$\|\mathcal{P}_{T^\perp}\mathcal{P}_\Omega(\mathcal{P}_T\mathcal{P}_\Omega)^\dagger[\mathbf{U}\mathbf{V}^*]\| \quad (4.4.27)$$

is small. This is a random matrix, but it is an exceedingly complicated one. It actually *is* possible to analyze its norm, but the analysis is quite intricate. The challenge arises because the thing that is random here is the support  $\Omega$ . It is repeated in several places, creating probabilistic dependencies, which complicates the analysis.

**Relaxed Optimality Conditions.**

As it is difficult to directly find a dual certificate satisfying the KKT conditions exactly, we might want to relax these conditions and see if we could still find another certificate for the optimality. The following proposition



suggests that we can ensure the optimality of  $\mathbf{X}_0$  with an alternative set of (relaxed) conditions:

**Proposition 4.4.5** (KKT Conditions – Approximate Version). *The matrix  $\mathbf{X}_0$  is the unique optimal solution to the nuclear minimization problem (4.4.19) if the following set of conditions hold*

1. *The operator norm of the operator  $p^{-1}\mathcal{P}_T\mathcal{P}_\Omega\mathcal{P}_T - \mathcal{P}_T$  is small:*  

$$\|p^{-1}\mathcal{P}_T\mathcal{P}_\Omega\mathcal{P}_T - \mathcal{P}_T\| \leq \frac{1}{2}.$$
2. *There exists a dual certificate  $\mathbf{\Lambda} \in \mathbb{R}^{n \times n}$  that satisfies  $\mathcal{P}_\Omega(\mathbf{\Lambda}) = \mathbf{\Lambda}$  and*
  - (a)  $\|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}]\| \leq \frac{1}{2}.$
  - (b)  $\|\mathcal{P}_T[\mathbf{\Lambda}] - \mathbf{UV}^*\|_F \leq \frac{1}{4n}.$

Conditions 2(a) and 2(b) above trade off between the degree of satisfaction of the equality constraint  $\mathcal{P}_T[\mathbf{\Lambda}] = \mathbf{UV}^*$  and the inequality constraint for the dual norm  $\|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}]\| \leq 1$  in the original KKT conditions. This is possible under the additional assumption that  $\|p^{-1}\mathcal{P}_T\mathcal{P}_\Omega\mathcal{P}_T - \mathcal{P}_T\|$  is not too large. This assumption is satisfied whenever the sampling map  $p^{-1}\mathcal{P}_\Omega$  nearly preserves the lengths of all elements  $\mathbf{X} \in T$ . It can be considered a strengthening of the condition that  $T \cap \Omega^\perp = \{\mathbf{0}\}$ , which was needed for unique optimality.

To prove Proposition 4.4.5, we will need another lemma. This says that provided  $\mathcal{P}_\Omega$  acts nicely on matrices from  $T$ , every feasible perturbation  $\mathbf{H}$  (i.e.,  $\mathbf{H}$  such that  $\mathcal{P}_\Omega[\mathbf{H}] = \mathbf{0}$ ) must have a nonnegligible component along  $T^\perp$ :

**Lemma 4.4.6.** *Suppose that the operator  $\mathcal{P}_\Omega$  satisfies*

$$\|\mathcal{P}_T - p^{-1}\mathcal{P}_T\mathcal{P}_\Omega\mathcal{P}_T\| \leq \frac{1}{2}. \quad (4.4.28)$$

*Then for any  $\mathbf{H}$  satisfying  $\mathcal{P}_\Omega[\mathbf{H}] = \mathbf{0}$ , we have*

$$\|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_F \geq \sqrt{\frac{p}{2}} \|\mathcal{P}_T[\mathbf{H}]\|_F. \quad (4.4.29)$$

*Proof.* We have

$$\begin{aligned} \langle \mathcal{P}_\Omega\mathcal{P}_T[\mathbf{H}], \mathcal{P}_\Omega\mathcal{P}_T[\mathbf{H}] \rangle &= \langle \mathcal{P}_T[\mathbf{H}], \mathcal{P}_\Omega\mathcal{P}_T[\mathbf{H}] \rangle \\ &= p \langle \mathcal{P}_T[\mathbf{H}], p^{-1}\mathcal{P}_\Omega\mathcal{P}_T[\mathbf{H}] \rangle \\ &= p \langle \mathcal{P}_T[\mathbf{H}], \mathcal{P}_T p^{-1}\mathcal{P}_\Omega\mathcal{P}_T[\mathbf{H}] \rangle \\ &\geq p \left( 1 - \|\mathcal{P}_T - \mathcal{P}_T p^{-1}\mathcal{P}_\Omega\mathcal{P}_T\| \right) \|\mathcal{P}_T[\mathbf{H}]\|_F^2 \\ &\geq \frac{p}{2} \|\mathcal{P}_T[\mathbf{H}]\|_F^2, \end{aligned} \quad (4.4.30)$$

Then from  $\mathcal{P}_\Omega \mathcal{P}_T[\mathbf{H}] + \mathcal{P}_\Omega \mathcal{P}_{T^\perp}[\mathbf{H}] = \mathcal{P}_\Omega[\mathbf{H}] = \mathbf{0}$ , we have

$$0 = \|\mathcal{P}_\Omega \mathcal{P}_T[\mathbf{H}] + \mathcal{P}_\Omega \mathcal{P}_{T^\perp}[\mathbf{H}]\|_F \quad (4.4.31)$$

$$\geq \|\mathcal{P}_\Omega \mathcal{P}_T[\mathbf{H}]\|_F - \|\mathcal{P}_\Omega \mathcal{P}_{T^\perp}[\mathbf{H}]\|_F \quad (4.4.32)$$

$$\geq \sqrt{\frac{p}{2}} \|\mathcal{P}_T[\mathbf{H}]\|_F - \|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_F, \quad (4.4.33)$$

giving the conclusion.  $\square$

We are now ready to prove the optimality of  $\mathbf{X}_0$  under the conditions given by Proposition 4.4.5.

*Proof.* We want to show that under the above conditions, for any feasible perturbation  $\mathbf{H} \neq \mathbf{0}$  and  $\mathbf{X} = \mathbf{X}_0 + \mathbf{H}$ , we have  $\|\mathbf{X}\|_* > \|\mathbf{X}_0\|_*$ . Let  $\mathcal{P}_{T^\perp}[\mathbf{H}] = \bar{\mathbf{U}}\bar{\Sigma}\bar{\mathbf{V}}^*$ . Then we have  $\bar{\mathbf{U}}\bar{\mathbf{V}}^* \in T^\perp$  and  $\|\bar{\mathbf{U}}\bar{\mathbf{V}}^*\| \leq 1$ . Therefore, we have  $\mathbf{U}\mathbf{V}^* + \bar{\mathbf{U}}\bar{\mathbf{V}}^* \in \partial\|\cdot\|_*(\mathbf{X}_0)$  is a subgradient of the nuclear norm at  $\mathbf{X}_0$ .

Also, we have  $\langle \bar{\mathbf{U}}\bar{\mathbf{V}}^*, \mathcal{P}_{T^\perp}[\mathbf{H}] \rangle = \|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_*$  and  $\langle \mathbf{\Lambda}, \mathbf{H} \rangle = 0$  and apply them to the following inequalities:

$$\begin{aligned} \|\mathbf{X}_0 + \mathbf{H}\|_* &\geq \|\mathbf{X}_0\|_* + \langle \mathbf{U}\mathbf{V}^* + \bar{\mathbf{U}}\bar{\mathbf{V}}^*, \mathbf{H} \rangle, \\ &= \|\mathbf{X}_0\|_* + \langle \mathbf{U}\mathbf{V}^* + \bar{\mathbf{U}}\bar{\mathbf{V}}^* - \mathbf{\Lambda}, \mathbf{H} \rangle, \\ &= \|\mathbf{X}_0\|_* + \langle \mathbf{U}\mathbf{V}^* - \mathcal{P}_T[\mathbf{\Lambda}], \mathbf{H} \rangle + \langle \bar{\mathbf{U}}\bar{\mathbf{V}}^* - \mathcal{P}_{T^\perp}[\mathbf{\Lambda}], \mathbf{H} \rangle, \\ &\geq \|\mathbf{X}_0\|_* - \frac{1}{4n} \|\mathcal{P}_T[\mathbf{H}]\|_F + \frac{1}{2} \|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_*, \\ &\geq \|\mathbf{X}_0\|_* + \underbrace{\left( \frac{1}{2} - \frac{1}{4n} \sqrt{\frac{2}{p}} \right)}_{> 0, \text{ since } p > n^{-2}} \|\mathcal{P}_{T^\perp}[\mathbf{H}]\|_F. \end{aligned} \quad (4.4.34)$$

In the final inequality, we have invoked Lemma 4.4.6.

Hence, for feasible perturbations  $\mathbf{H}$ ,  $\|\mathbf{X}_0 + \mathbf{H}\|_* \geq \|\mathbf{X}_0\|_*$ , with equality if and only if  $\mathcal{P}_{T^\perp}[\mathbf{H}] = \mathbf{0}$ . But via Lemma 4.4.6,  $\mathcal{P}_{T^\perp}[\mathbf{H}] = \mathbf{0} \implies \mathbf{H} = \mathbf{0}$ . Thus, for any nonzero feasible perturbation  $\mathbf{H}$ ,  $\|\mathbf{X}_0 + \mathbf{H}\|_* > \|\mathbf{X}_0\|_*$ , establishing the desired condition.  $\square$

### ***The Optimality Condition is Satisfied with High Probability.***

To complete the proof, we simply need to show that the optimality condition can be satisfied with high probability. To do this, we need to verify two claims: first, that with high probability the sampling operator  $\Omega$  acts nicely on  $T$ , in the sense that  $\|p^{-1}\mathcal{P}_T\mathcal{P}_\Omega\mathcal{P}_T - \mathcal{P}_T\|$  is small. We then need to show that with high probability we can construct the desired dual certificate  $\mathbf{\Lambda}$ .

#### *The sampling operator acts nicely on $T$*

We next prove that the sampling operator  $\mathcal{P}_\Omega$  preserves some part of every element of  $T$ , in the sense that  $\|p^{-1}\mathcal{P}_T\mathcal{P}_\Omega\mathcal{P}_T - \mathcal{P}_T\|$  is small. This phenomenon is a consequence of the incoherence of the matrix  $\mathbf{X}_0$  and the

uniform random model on  $\Omega$ . The proof of the following lemma uses the matrix (operator) Bernstein inequality to show this rigorously.

**Lemma 4.4.7.** *Let  $\mathcal{P}_\Omega : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$  denote the operator*

$$\mathcal{P}_\Omega[\mathbf{X}] = \sum_{ij} \mathbf{X}_{ij} \mathbb{1}_{(i,j) \in \Omega} \mathbf{E}_{ij} \quad (4.4.35)$$

with  $\mathbb{1}_{(i,j) \in \Omega}$  independent Bernoulli random variables with probability  $p$ . Fix any  $\varepsilon$  with  $c \frac{\sqrt{\log n}}{n} \leq \varepsilon \leq 1$ . There is a numerical constant  $C$  such that if  $p > C \frac{\nu r \log n}{\varepsilon^2 n}$ , then with high probability,

$$\|\mathcal{P}_T - p^{-1} \mathcal{P}_T \mathcal{P}_\Omega \mathcal{P}_T\| \leq \varepsilon. \quad (4.4.36)$$

*Proof.* We apply the matrix Bernstein inequality in Theorem E.4.4 to bound the norm of

$$\mathcal{P}_T - p^{-1} \mathcal{P}_T \mathcal{P}_\Omega \mathcal{P}_T = \sum_{ij} \underbrace{\mathcal{P}_T \left( \frac{\mathcal{I}}{n^2} - p^{-1} \mathbb{1}_{(i,j) \in \Omega} \mathbf{E}_{ij} \langle \mathbf{E}_{ij}, \cdot \rangle \right) \mathcal{P}_T}_{\doteq \mathcal{W}_{ij}}.$$

Here,  $\mathcal{W}_{ij} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$  are independent random linear maps, and  $\mathbb{E} \left[ \sum_{ij} \mathcal{W}_{ij} \right] = 0$ . The matrix Bernstein inequality requires (i) an almost sure bound  $R$  on  $\max_{ij} \|\mathcal{W}_{ij}\|$ , and (ii) control of the “variance”

$$\sum_{ij} \mathbb{E} [\mathcal{W}_{ij}^* \mathcal{W}_{ij}]. \quad (4.4.37)$$

We provide these as follows.

(i) Almost sure control of the summands:

$$\begin{aligned} \|\mathcal{W}_{ij}\| &\leq \max \left\{ \|n^{-2} \mathcal{P}_T\|, \|p^{-1} \mathcal{P}_T[\mathbf{E}_{ij}] \langle \mathcal{P}_T[\mathbf{E}_{ij}], \cdot \rangle\| \right\}, \quad \text{almost surely} \\ &= \max \left\{ n^{-2}, p^{-1} \|\mathcal{P}_T[\mathbf{E}_{ij}]\|_F^2 \right\}, \\ &\leq \max \left\{ n^{-2}, \frac{2\nu r}{np} \right\}, \\ &\leq \max \left\{ \frac{1}{n^2}, \frac{2\varepsilon^2}{C \log n} \right\}, \\ &= \frac{2\varepsilon^2}{C \log n}. \end{aligned} \quad (4.4.38)$$

We may take  $R = \frac{2\varepsilon^2}{C \log n}$ .

(ii) Control of the “operator variance”. Note that

$$\begin{aligned}
& \sum_{ij} \mathbb{E} [\mathcal{W}_{ij}^* \mathcal{W}_{ij}] \\
&= \sum_{ij} \mathbb{E} \left[ \frac{1}{n^4} \mathcal{P}_T - \frac{2p^{-1}}{n^2} \mathbb{1}_{(i,j) \in \Omega} \mathcal{P}_T \mathbf{E}_{ij} \langle \mathbf{E}_{ij}, \cdot \rangle \mathcal{P}_T \right. \\
&\quad \left. + \mathbb{1}_{(i,j) \in \Omega} p^{-2} \mathcal{P}_T \mathbf{E}_{ij} \|\mathcal{P}_T \mathbf{E}_{ij}\|_F^2 \langle \mathbf{E}_{ij}, \cdot \rangle \mathcal{P}_T \right] \\
&\preceq p^{-1} \sum_{ij} \mathcal{P}_T \mathbf{E}_{ij} \|\mathcal{P}_T \mathbf{E}_{ij}\|_F^2 \langle \mathbf{E}_{ij}, \cdot \rangle \mathcal{P}_T \\
&\preceq p^{-1} \frac{2\nu r}{n} \sum_{ij} \mathcal{P}_T \mathbf{E}_{ij} \langle \mathbf{E}_{ij}, \cdot \rangle \mathcal{P}_T \\
&\preceq \frac{2\varepsilon^2}{C \log n} \mathcal{P}_T. \tag{4.4.39}
\end{aligned}$$

The operator  $\sum_{ij} \mathbb{E} [\mathcal{W}_{ij}^* \mathcal{W}_{ij}]$  is self-adjoint and positive semidefinite. The above calculation therefore implies that

$$\begin{aligned}
\sigma^2 &= \max \left\{ \left\| \sum_{ij} \mathbb{E} [\mathcal{W}_{ij}^* \mathcal{W}_{ij}] \right\|, \left\| \sum_{ij} \mathbb{E} [\mathcal{W}_{ij} \mathcal{W}_{ij}^*] \right\| \right\} \\
&\leq \frac{2\varepsilon^2}{C \log n}. \tag{4.4.40}
\end{aligned}$$

Using these calculations, we obtain a bound

$$\mathbb{P} \left[ \left\| \sum_{ij} \mathcal{W}_{ij} \right\| > t \right] \leq 2n \exp \left( \frac{-t^2/2}{\frac{2\varepsilon^2}{C \log n} + t \frac{2\varepsilon^2}{3C \log n}} \right). \tag{4.4.41}$$

The probability of failure for  $t = \varepsilon$  is bounded by  $n^{-\rho}$ ; the exponent  $\rho$  can be made as large as desired by choosing  $C$  appropriately.  $\square$

Choosing  $\varepsilon = 1/2$  in the statement of the above lemma, we obtain the desired condition needed for Lemma 4.4.6.

*Construction of a dual certificate by the golfing scheme.*

From the above discussion, in order to prove Theorem 4.4.3, we only have to show that under the conditions of the theorem, we can find a dual certificate that satisfies two conditions 2 (a) and 2 (b) of Proposition 4.4.5. In this section, we show how to construct such a dual certificate,  $\mathbf{\Lambda}$ . In the next chapter, we will reuse this construction to analyze the related problem of *robust matrix recovery*, in which a fraction of the entries of a low-rank matrix have been corrupted. For this purpose, we give a complete summary of the properties of our construction in the following proposition.

Here, properties (i) and (ii) are essential for matrix completion; property (iii) will be used in the following chapters for matrix recovery.

**Proposition 4.4.8** (Dual certificate for low-rank recovery). *Let  $\mathbf{L}_0 \in \mathbb{R}^{n \times n}$  be a rank- $r$  matrix, with coherence  $\nu$ . Let  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n \times r}$  be matrices whose columns are leading left- and right singular vectors of  $\mathbf{L}_0$ . Let*

$$T = \{\mathbf{U}\mathbf{X}^* + \mathbf{Y}\mathbf{V}^* \mid \mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}\}. \quad (4.4.42)$$

Then if  $\Omega \sim \text{Ber}(p)$ , with

$$p > C_0 \frac{\nu r \log^2 n}{n}, \quad (4.4.43)$$

there exists a matrix  $\mathbf{\Lambda}$  supported on  $\Omega$ , satisfying

- (i)  $\|\mathcal{P}_T \mathbf{\Lambda} - \mathbf{U}\mathbf{V}^*\|_F \leq \frac{1}{4n}$ ,
- (ii)  $\|\mathcal{P}_{T^\perp} \mathbf{\Lambda}\| \leq \frac{1}{4}$ ,
- (iii)  $\|\mathbf{\Lambda}\|_\infty < \frac{C_1 \log n}{p} \times \|\mathbf{U}\mathbf{V}^*\|_\infty$ ,

with high probability. Here,  $C_1$  is a positive numerical constant.

We prove this proposition using an iterative construction. Let

$$\Omega_1, \dots, \Omega_k \quad (4.4.44)$$

be *independent* random subsets, chosen according to the Bernoulli model with parameter  $q$ . Set

$$\Omega = \bigcup_{i=1}^k \Omega_i. \quad (4.4.45)$$

Then  $\Omega$  is *also* a Bernoulli subset, with parameter

$$p = 1 - (1 - q)^k. \quad (4.4.46)$$

The parameter  $p$  is the probability that a given entry is in *at least one* of the subsets  $\Omega_i$ . Hence,  $p \leq kq$ . The argument that we develop below will lead us to choose  $k = C_g \log(n)$ , with  $C_g$  a constant. Because  $k$  is not too large, this implies that the parameter  $q$  is also not too small:

$$q \geq \frac{p}{k} = \frac{C_0}{C_g} \frac{\nu r \log n}{n}. \quad (4.4.47)$$

Provided  $C_0$  is large enough compared to  $C_g$ , the subsets  $\Omega_i$  *all* satisfy the conditions of Lemma 4.4.7, and so with high probability

$$\|\mathcal{P}_T - q^{-1} \mathcal{P}_T \mathcal{P}_{\Omega_j} \mathcal{P}_T\| \leq \frac{1}{2}, \quad j = 1, \dots, k. \quad (4.4.48)$$

We will construct a sequence of matrices  $\mathbf{\Lambda}_0, \mathbf{\Lambda}_1, \dots, \mathbf{\Lambda}_k$ , in which each  $\mathbf{\Lambda}_j$  depends only on  $\Omega_1, \dots, \Omega_j$ . We let  $\mathbf{\Lambda}_0 = \mathbf{0}$ . And let

$$\mathbf{E}_j = \mathcal{P}_T[\mathbf{\Lambda}_j] - \mathbf{U}\mathbf{V}^*. \quad (4.4.49)$$

Since our goal is to obtain  $\mathbf{\Lambda}$  such that  $\mathcal{P}_T[\mathbf{\Lambda}] \approx \mathbf{UV}^*$ ,  $\mathbf{E}_j$  should be considered the *error* at iteration  $j$ . To get our next  $\mathbf{\Lambda}$ , we simply try to correct the error:

$$\mathbf{\Lambda}_j = \mathbf{\Lambda}_{j-1} - (q^{-1}\mathcal{P}_{\Omega_j})[\mathbf{E}_{j-1}]. \quad (4.4.50)$$

This construction is known as the *golfing scheme*, as it tries to reach the goal by reducing error step by step.

There are several things worth noting about this construction. First, it produces  $\mathbf{\Lambda}_j$  supported only on  $\Omega_1 \cup \dots \cup \Omega_j$ . Thus, as desired,  $\mathbf{\Lambda}_k$  is supported on  $\Omega$ . Second, because  $\mathbf{UV}^* \in T$ ,  $\mathbf{E}_j \in T$  for each  $j$ . This means that

$$\begin{aligned} \mathbf{E}_j &= \mathcal{P}_T[\mathbf{\Lambda}_j] - \mathbf{UV}^* \\ &= \mathcal{P}_T[\mathbf{\Lambda}_{j-1}] - \mathbf{UV}^* - q^{-1}\mathcal{P}_T\mathcal{P}_{\Omega_j}[\mathbf{E}_{j-1}] \\ &= \mathbf{E}_j - q^{-1}\mathcal{P}_T\mathcal{P}_{\Omega_j}[\mathbf{E}_{j-1}] \\ &= (\mathcal{P}_T - q^{-1}\mathcal{P}_T\mathcal{P}_{\Omega_j}\mathcal{P}_T)[\mathbf{E}_{j-1}]. \end{aligned}$$

Since  $\mathbb{E}[q^{-1}\mathcal{P}_{\Omega_j}] = \mathcal{I}$ , in expectation, this iterative process drives the error to zero:  $\mathbb{E}[\mathbf{E}_j] = \mathbf{0}$ .

As it turns out, due to the fact that  $\|\mathcal{P}_T - q^{-1}\mathcal{P}_T\mathcal{P}_{\Omega_j}\mathcal{P}_T\| \leq \frac{1}{2}$ , after  $k$  steps, the error reduces to

$$\|\mathcal{P}_T[\mathbf{\Lambda}_k] - \mathbf{UV}^*\|_F = \|\mathbf{E}_k\|_F \leq 2^{-k} \|\mathbf{E}_0\|_F \quad (4.4.51)$$

with high probability.

So, based on the golfing scheme, to achieve the desired accuracy as suggested by the above lemma, we want  $2^{-k} \|\mathbf{E}_0\|_F = 2^{-k} \sqrt{r} \leq \frac{1}{4n}$ . Since  $r < n$ , we only need to have  $2^{-k} \sim O(1/n^2)$ , that is to choose  $k = C_g \log(n)$  for some large enough constant  $C_g$ , say  $C_g = 20$ . Therefore, under these conditions, the dual certificate constructed after  $k$  iterations  $\mathbf{\Lambda}_k$  satisfies condition 2 (b) of Proposition 4.4.5:

$$\|\mathcal{P}_T[\mathbf{\Lambda}_k] - \mathbf{UV}^*\|_F \leq \frac{1}{4n}. \quad (4.4.52)$$

Finally, to satisfy Condition 2(a) of Proposition 4.4.5, we need to show that the operator norm of the random matrix  $\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k]$  is bounded as

$$\|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k]\| \leq 1/4.$$

Notice that from the construction of  $\mathbf{\Lambda}_k$ , we have

$$\begin{aligned} \mathbf{\Lambda}_k &= \sum_{j=1}^k -q^{-1}\mathcal{P}_{\Omega_j}[\mathbf{E}_{j-1}], \\ \mathbf{E}_j &= (\mathcal{P}_T - \mathcal{P}_T q^{-1}\mathcal{P}_{\Omega_j}\mathcal{P}_T)[\mathbf{E}_{j-1}], \quad \text{with } \mathbf{E}_0 = -\mathbf{UV}^*. \end{aligned}$$

The matrix of interest can be expressed as

$$\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k] = \sum_{j=1}^k -q^{-1} \mathcal{P}_{T^\perp} \mathcal{P}_{\Omega_j}[\mathbf{E}_{j-1}] = \sum_{j=1}^k \mathcal{P}_{T^\perp} (\mathcal{P}_T - q^{-1} \mathcal{P}_{\Omega_j} \mathcal{P}_T)[\mathbf{E}_{j-1}], \quad (4.4.53)$$

where the second identity is due to  $\mathcal{P}_{T^\perp} \mathcal{P}_T = 0$  and  $\mathcal{P}_T[\mathbf{E}_j] = \mathbf{E}_j$ .

Since we are interested in bounding the norm of  $\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k]$ , it would help if we know good bounds on various norms of  $\mathcal{P}_{\Omega_j}$  and its interaction with the operator  $\mathcal{P}_T$  or  $\mathcal{P}_{T^\perp}$ . Notice each  $\mathcal{P}_{\Omega_j}$  is a summation of independent random operators. A very powerful tool we can use to bound the norm of summation of random matrices (or operators) is the so-called matrix Bernstein inequality introduced in the Appendix E.4, which we have used once before in Lemma 4.4.7.

To bound the norm of  $\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k]$ , we need good bounds on three additional operators similar to that in Lemma 4.4.7. The proofs of these bounds<sup>15</sup> are all similar to that of Lemma 4.4.7 by utilizing the matrix Bernstein inequality. We hence leave their derivations as exercises to the reader to get familiar with the matrix Bernstein inequality.

We phrase these bounds in terms of

$$\|\mathbf{Z}\|_\infty = \max_{ij} |\mathbf{Z}_{ij}|, \quad (4.4.54)$$

and the maximum of the largest  $\ell^2$  norm of a row and the largest  $\ell^2$  norm of a column, which we denote by  $\|\cdot\|_{rc}$ :

$$\|\mathbf{Z}\|_{rc} = \max \left\{ \max_i \|\mathbf{e}_i^* \mathbf{Z}\|_2, \max_j \|\mathbf{Z} \mathbf{e}_j\|_2 \right\}. \quad (4.4.55)$$

**Lemma 4.4.9.** *Let  $\mathbf{Z}$  be any fixed  $n \times n$  matrix, and  $\Omega$  a  $\text{Ber}(q)$  subset, with*

$$q > C_0 \frac{\nu r \log n}{n}. \quad (4.4.56)$$

*Then with high probability*

$$\|(q^{-1} \mathcal{P}_\Omega - \mathcal{I}) \mathbf{Z}\| \leq C \left( \frac{n}{C_0 \nu r} \|\mathbf{Z}\|_\infty + \sqrt{\frac{n}{C_0 \nu r}} \|\mathbf{Z}\|_{rc} \right), \quad (4.4.57)$$

*where  $C$  is a numerical constant.*

*Proof.* Exercise 4.21. □

**Lemma 4.4.10.** *Let  $\mathbf{Z}$  be any fixed  $n \times n$  matrix. There exists a numerical constant  $C_0$  such that if  $\Omega$  is a  $\text{Ber}(q)$  subset with*

$$q > C_0 \frac{\nu r \log n}{n}, \quad (4.4.58)$$

---

<sup>15</sup>following the work of [Chen et al., 2013].

then with high probability

$$\|(q^{-1}\mathcal{P}_T\mathcal{P}_\Omega - \mathcal{P}_T)\mathbf{Z}\|_{rc} \leq \frac{1}{2} \left( \sqrt{\frac{n}{\nu r}} \|\mathbf{Z}\|_\infty + \|\mathbf{Z}\|_{rc} \right). \quad (4.4.59)$$

*Proof.* Exercise 4.22. □

**Lemma 4.4.11.** Suppose  $\mathbf{Z}$  is a fixed  $n \times n$  matrix in  $T$ . There exists a constant  $C_0$  such that if  $\Omega$  is a Bernoulli( $q$ ) subset with

$$q > C_0 \frac{\nu r \log n}{n}. \quad (4.4.60)$$

Then with high probability we have

$$\|(\mathcal{P}_T - q^{-1}\mathcal{P}_T\mathcal{P}_\Omega\mathcal{P}_T)\mathbf{Z}\|_\infty \leq \frac{1}{2} \|\mathbf{Z}\|_\infty. \quad (4.4.61)$$

*Proof.* Exercise 4.23. □

With these three lemmas in hand, we are now ready to show that the spectral norm of  $\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k]$  is very small, in particular can be bounded as  $\|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k]\| \leq 1/4$ :

*Proof.* From the golfing construction,  $\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k]$  can be expressed as the series given in (4.4.53). Hence we have

$$\begin{aligned} \|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k]\| &\leq \sum_{j=1}^k \|\mathcal{P}_{T^\perp}(\mathcal{P}_T - q^{-1}\mathcal{P}_{\Omega_j}\mathcal{P}_T)[\mathbf{E}_{j-1}]\| \\ &\leq \sum_{j=1}^k \|(\mathcal{P}_T - q^{-1}\mathcal{P}_{\Omega_j}\mathcal{P}_T)[\mathbf{E}_{j-1}]\| \\ &= \sum_{j=1}^k \|(\mathcal{I} - q^{-1}\mathcal{P}_{\Omega_j})[\mathbf{E}_{j-1}]\|. \end{aligned} \quad (4.4.62)$$

Notice that in the construction of the golfing scheme, we have ensured that each subset  $\Omega_j$  is sampled according to the Bernoulli model, with parameter  $q > C_0 \frac{\nu r \log n}{n}$  for some large enough  $C_0$ . This means each of the  $k$  subsets  $\Omega_j$  satisfies the conditions of the above lemmas. We first apply Lemma 4.4.9 to the right hand side of the last inequality and obtain (assuming  $C_0 > 1$ ):

$$\|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k]\| \leq \frac{C}{\sqrt{C_0}} \sum_{j=1}^k \left( \frac{n}{\nu r} \|\mathbf{E}_{j-1}\|_\infty + \sqrt{\frac{n}{\nu r}} \|\mathbf{E}_{j-1}\|_{rc} \right). \quad (4.4.63)$$



To bound  $\|\mathbf{E}_{j-1}\|_\infty$  we apply Lemma 4.4.11 and obtain

$$\begin{aligned}\|\mathbf{E}_{j-1}\|_\infty &= \left\| (\mathcal{P}_T - \tfrac{1}{q} \mathcal{P}_T \mathcal{P}_{\Omega_{j-1}} \mathcal{P}_T) \cdots (\mathcal{P}_T - \tfrac{1}{q} \mathcal{P}_T \mathcal{P}_{\Omega_1} \mathcal{P}_T) \mathbf{E}_0 \right\|_\infty \\ &\leq \left(\tfrac{1}{2}\right)^{j-1} \|\mathbf{U}\mathbf{V}^*\|_\infty.\end{aligned}\quad (4.4.64)$$

Using this together with the fact that  $\mathbf{\Lambda}_k = -\sum_j q^{-1} \mathcal{P}_{\Omega_j} \mathbf{E}_{j-1}$ , we obtain

$$\|\mathbf{\Lambda}_k\|_\infty \leq q^{-1} \sum_j \|\mathbf{E}_{j-1}\|_\infty \quad (4.4.65)$$

$$\leq 2q^{-1} \|\mathbf{U}\mathbf{V}^*\|_\infty. \quad (4.4.66)$$

Since  $q > p/C_q \log n$ , this establishes property (iii) of Proposition 4.4.8 for  $\mathbf{\Lambda}_k$ .

To bound  $\|\mathbf{E}_{j-1}\|_{rc}$  we apply Lemma 4.4.10 and obtain

$$\begin{aligned}\|\mathbf{E}_{j-1}\|_{rc} &= \left\| (\mathcal{P}_T - \tfrac{1}{q} \mathcal{P}_T \mathcal{P}_{\Omega_{j-1}} \mathcal{P}_T) \mathbf{E}_{j-2} \right\|_{rc} \\ &\leq \tfrac{1}{2} \sqrt{\tfrac{n}{\nu r}} \|\mathbf{E}_{j-2}\|_\infty + \tfrac{1}{2} \|\mathbf{E}_{j-1}\|_{rc}.\end{aligned}\quad (4.4.67)$$

Combine the above two inequalities and apply them recursively to  $j-1, j-2, \dots, 0$  and we obtain

$$\|\mathbf{E}_{j-1}\|_{rc} \leq j \left(\tfrac{1}{2}\right)^{j-1} \sqrt{\tfrac{n}{\nu r}} \|\mathbf{U}\mathbf{V}^*\|_\infty + \left(\tfrac{1}{2}\right)^{j-1} \|\mathbf{U}\mathbf{V}^*\|_{rc}. \quad (4.4.68)$$

Substitute the bounds (4.4.64) and (4.4.68) to the right and side of (4.4.63) and we obtain

$$\begin{aligned}\|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k]\| &\leq \frac{C}{\sqrt{C_0}} \frac{n}{\nu r} \|\mathbf{U}\mathbf{V}^*\|_\infty \sum_{j=1}^k (j+1) \left(\tfrac{1}{2}\right)^{j-1} \\ &\quad + \frac{C}{\sqrt{C_0}} \sqrt{\tfrac{n}{\nu r}} \|\mathbf{U}\mathbf{V}^*\|_{rc} \sum_{j=1}^k \left(\tfrac{1}{2}\right)^{j-1} \\ &\leq \frac{6C}{\sqrt{C_0}} \frac{n}{\nu r} \|\mathbf{U}\mathbf{V}^*\|_\infty + \frac{2C}{\sqrt{C_0}} \sqrt{\tfrac{n}{\nu r}} \|\mathbf{U}\mathbf{V}^*\|_{rc}.\end{aligned}\quad (4.4.69)$$

As the matrix  $\mathbf{X}_0$  satisfies the incoherence conditions (4.4.13) and (4.4.14), we have

$$\begin{aligned}\|\mathbf{U}\mathbf{V}^*\|_\infty &\leq \max_{i,j} \left\{ \|\mathbf{U}^* \mathbf{e}_i\|_2 \times \|\mathbf{V}^* \mathbf{e}_j\|_2 \right\} \leq \frac{\nu r}{n}, \\ \|\mathbf{U}\mathbf{V}^*\|_{rc} &\leq \max \left\{ \max_i \|\mathbf{e}_i^* \mathbf{U}\mathbf{V}^*\|_2, \max_j \|\mathbf{U}\mathbf{V}^* \mathbf{e}_j\|_2 \right\} \leq \sqrt{\tfrac{\nu r}{n}}.\end{aligned}$$

Therefore,

$$\|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k]\| \leq \frac{6C}{\sqrt{C_0}} + \frac{2C}{\sqrt{C_0}} \leq \frac{1}{4} \quad (4.4.70)$$

for large enough  $C_0$ . This establishes property (ii) of Proposition 4.4.8 for  $\mathcal{P}_{T^\perp}[\mathbf{\Lambda}_k]$ .  $\square$

The above derivations and results show that the relaxed KKT conditions in Proposition 4.4.5 can be satisfied with high probability, proving Theorem 4.4.3.

#### 4.4.5 Stable Matrix Completion with Noise

So far in the matrix completion problem, we have assumed the observed entries are precise. In real world matrix completion problems, the observed entries are often corrupted with some noise:

$$\mathbf{Y}_{ij} = [\mathbf{X}_0]_{ij} + \mathbf{Z}_{ij}, \quad (i, j) \in \Omega, \quad (4.4.71)$$

where  $\mathbf{Z}_{ij}$  can be some small noise. Or equivalently, we can write

$$\mathcal{P}_\Omega[\mathbf{Y}] = \mathcal{P}_\Omega[\mathbf{X}_0] + \mathcal{P}_\Omega[\mathbf{Z}], \quad (4.4.72)$$

where  $\mathbf{Z}$  is an  $n \times n$  matrix of noises. We may assume the overall noise level is small  $\|\mathcal{P}_\Omega[\mathbf{Z}]\|_F < \delta$ . As in the stable matrix recovery case, we could expect to recover a low rank matrix close to  $\mathbf{X}_0$  via solving the following convex program:

$$\begin{aligned} & \text{minimize} && \|\mathbf{X}\|_* \\ & \text{subject to} && \|\mathcal{P}_\Omega[\mathbf{X}] - \mathcal{P}_\Omega[\mathbf{Y}]\|_F < \delta. \end{aligned} \quad (4.4.73)$$

The following theorem states that under the same conditions of Theorem 4.4.3 when the nuclear norm minimization recovers the correct low rank matrix from noiseless measurements, the above program gives a stable estimate  $\hat{\mathbf{X}}$  of the true low-rank matrix  $\mathbf{X}_0$ :

**Theorem 4.4.12** (Stable Matrix Completion). *Let  $\mathbf{X}_0 \in \mathbb{R}^{n \times n}$  be a rank- $r$ ,  $\nu$ -incoherent matrix. Suppose that we observe  $\mathcal{P}_\Omega[\mathbf{Y}] = \mathcal{P}_\Omega[\mathbf{X}_0] + \mathcal{P}_\Omega[\mathbf{Z}]$ , where  $\Omega$  is a subset of  $[n] \times [n]$ . If  $\Omega$  is uniformly sampled from subsets of size*

$$m \geq C_1 \nu n r \log^2(n), \quad (4.4.74)$$

*then with high probability, the optimal solution  $\hat{\mathbf{X}}$  to the convex program (4.4.73) satisfies*

$$\|\hat{\mathbf{X}} - \mathbf{X}_0\|_F \leq c \frac{n\sqrt{n}\log(n)}{\sqrt{m}} \delta \leq c' \frac{n}{\sqrt{r}} \delta \quad (4.4.75)$$

for some constant  $c > 0$ .

*Proof.* Similar to the proof of Theorem 4.4.3 in the noiseless case which has the same incoherence condition on  $\mathbf{X}_0$  and the sampling condition, we know the sampling operator  $\mathcal{P}_\Omega$  and the dual certificate  $\mathbf{\Lambda}_k$  constructed via the golfing scheme satisfies the properties in Proposition 4.4.5. All we

need to show here is that these properties also imply the conclusion of this theorem for the case with noisy measurements.

Let  $\mathbf{H} = \hat{\mathbf{X}} - \mathbf{X}_0$ . Notice that we can split  $\mathbf{H}$  into two parts  $\mathbf{H} = \mathcal{P}_\Omega[\mathbf{H}] + \mathcal{P}_{\Omega^c}[\mathbf{H}]$ . For the first part, we have

$$\begin{aligned} \|\mathcal{P}_\Omega[\mathbf{H}]\|_F &= \|\mathcal{P}_\Omega[\hat{\mathbf{X}} - \mathbf{X}_0]\|_F \\ &\leq \|\mathcal{P}_\Omega[\hat{\mathbf{X}} - \mathbf{Y}]\|_F + \|\mathcal{P}_\Omega[\mathbf{Y} - \mathbf{X}_0]\|_F \\ &\leq 2\delta. \end{aligned} \quad (4.4.76)$$

Notice that the second part  $\mathcal{P}_{\Omega^c}[\mathbf{H}]$  is a feasible perturbation to the noiseless matrix completion problem. From the proof of Proposition 4.4.5 and in particular (4.4.34), we have

$$\|\mathbf{X}_0 + \mathcal{P}_{\Omega^c}[\mathbf{H}]\|_* \geq \|\mathbf{X}_0\|_* + \left(\frac{1}{2} - \frac{1}{4C_2\sqrt{nr}}\right) \|\mathcal{P}_{T^\perp}[\mathcal{P}_{\Omega^c}[\mathbf{H}]]\|_F \quad (4.4.77)$$

and based on triangle inequality, we also have

$$\|\hat{\mathbf{X}}\|_* = \|\mathbf{X}_0 + \mathbf{H}\|_* \geq \|\mathbf{X}_0 + \mathcal{P}_{\Omega^c}[\mathbf{H}]\|_* - \|\mathcal{P}_\Omega[\mathbf{H}]\|_*. \quad (4.4.78)$$

Since  $\|\hat{\mathbf{X}}\|_* \leq \|\mathbf{X}_0\|_*$ , we have

$$\|\mathcal{P}_\Omega[\mathbf{H}]\|_* \geq \left(\frac{1}{2} - \frac{1}{4C_2\sqrt{nr}}\right) \|\mathcal{P}_{T^\perp}[\mathcal{P}_{\Omega^c}[\mathbf{H}]]\|_F. \quad (4.4.79)$$

This leads to

$$\|\mathcal{P}_{T^\perp}[\mathcal{P}_{\Omega^c}[\mathbf{H}]]\|_F \leq 4\|\mathcal{P}_\Omega[\mathbf{H}]\|_* \leq 4\sqrt{n}\|\mathcal{P}_\Omega[\mathbf{H}]\|_F \leq 4\sqrt{n}\delta. \quad (4.4.80)$$

Since  $\mathcal{P}_{\Omega^c}[\mathbf{H}] = \mathcal{P}_{T^\perp}[\mathcal{P}_{\Omega^c}[\mathbf{H}]] + \mathcal{P}_T[\mathcal{P}_{\Omega^c}[\mathbf{H}]]$ , we remain to bound the term  $\mathcal{P}_T[\mathcal{P}_{\Omega^c}[\mathbf{H}]]$ . Applying the proof of Lemma 4.4.6 to  $\mathcal{P}_{\Omega^c}[\mathbf{H}]$ , we have  $\|\mathcal{P}_{T^\perp}[\mathcal{P}_{\Omega^c}[\mathbf{H}]]\|_F \geq C_1 \frac{\sqrt{m}}{n \log(n)} \|\mathcal{P}_T[\mathcal{P}_{\Omega^c}[\mathbf{H}]]\|_F$  for some large enough  $C_1$ . Therefore, we have

$$\|\mathcal{P}_T[\mathcal{P}_{\Omega^c}[\mathbf{H}]]\|_F \leq \frac{n \log(n)}{C_1 \sqrt{m}} \|\mathcal{P}_{T^\perp}[\mathcal{P}_{\Omega^c}[\mathbf{H}]]\|_F \leq c \frac{n \sqrt{n} \log(n)}{\sqrt{m}} \delta. \quad (4.4.81)$$

This bound dominates the bounds of all the other terms, leading to the conclusion of the theorem.  $\square$

## 4.5 Summary

In this chapter, we have studied the problem of recovering a low rank matrix from incomplete observations. This problem arises in a range of applications. It generalizes in a natural way the problem of recovering a sparse vector. We described a convex relaxation of the low rank recovery problem, in which we minimize the nuclear norm, which is the sum ( $\ell^1$  norm) of the singular values of a matrix. We proved that roughly  $(n_1 + n_2)r$  generic measurements suffice to recover all rank- $r$  matrices. We also studied a specific,

more structured measurement model, in which we observe a subset of the entries of a matrix. This matrix completion problem captures the structure of many of the most important practical low rank recovery problems. It is mathematically more challenging, because certain very sparse low rank matrices cannot be completed without seeing almost all of their entries. Nevertheless, we saw that for low rank matrices whose singular vectors are not too concentrated on any coordinate, nuclear norm minimization indeed succeeds. Paralleling our development for sparse vectors, we showed that these algorithmic and theoretical results can be extended to cope with non-idealities, such as measurement noise. Moreover, in the next chapter we will see how these ideas combine naturally with ideas from sparse recovery to generate even richer classes of models and more robust algorithms.

## 4.6 Exercises

**4.1** (Proof of Schoenberg's Theorem). *In this exercise, we invite the interested reader to prove Schoenberg's Euclidean embedding theorem (Theorem 4.1.1). Let  $\mathbf{D}$  be a Euclidean distance matrix for some point set  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ , i.e.,  $\mathbf{D}_{ij} = \|\mathbf{x}_i\|_2^2 + \|\mathbf{x}_j\|_2^2 - 2\langle \mathbf{x}_i, \mathbf{x}_j \rangle$ . Let  $\mathbf{1} \in \mathbb{R}^n$  denote the vector of all ones, and  $\Phi = \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^*$ . Using that  $\Phi\mathbf{1} = \mathbf{0}$ , argue that  $\Phi\mathbf{D}\Phi^*$  satisfies the conditions of Schoenberg's theorem, i.e., it is negative semidefinite and has rank at most  $d$ .*

*For the converse, let  $\mathbf{D}$  be a symmetric matrix with zero diagonal, and suppose that  $\Phi\mathbf{D}\Phi^*$  is negative semidefinite and has rank at most  $d$ . Argue that there exists some matrix  $\mathbf{X} \in \mathbb{R}^{d \times n}$  for which  $\mathbf{D}_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$ .*

**4.2** (Derivation of the SVD). *Let  $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$  be a matrix of rank  $r$ . Argue that there exists matrices  $\mathbf{U} \in \mathbb{R}^{n_1 \times r}$ ,  $\mathbf{V} \in \mathbb{R}^{n_2 \times r}$ , with orthonormal columns and a diagonal matrix  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{r \times r}$ , with  $\sigma_1 \geq \dots \geq \sigma_r > 0$ , such that*

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^*. \quad (4.6.1)$$

*Hint: what is the relationship between the singular values  $\sigma_i$  and singular vectors  $\mathbf{v}_i$  and the eigenvalues / eigenvectors of the matrix  $\mathbf{X}^*\mathbf{X}$ ?*

**4.3** (Best Rank- $r$  Approximation). *We prove Theorem 4.2.4. First, consider the special case in which  $\mathbf{Y} = \Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$  with  $\sigma_1 > \sigma_2 > \dots > \sigma_n$ . An arbitrary rank- $r$  matrix  $\mathbf{X}$  can be expressed as  $\mathbf{X} = \mathbf{F}\mathbf{G}^*$  with  $\mathbf{F} \in \mathbb{R}^{n_1 \times r}$ ,  $\mathbf{F}^*\mathbf{F} = \mathbf{I}$  and  $\mathbf{G} \in \mathbb{R}^{n_2 \times r}$ .*

1. *Argue that for any fixed  $\mathbf{F}$ , the solution to the optimization problem*

$$\min_{\mathbf{G} \in \mathbb{R}^{n_2 \times r}} \|\mathbf{F}\mathbf{G}^* - \Sigma\|_F^2 \quad (4.6.2)$$

is given by  $\hat{\mathbf{G}} = \mathbf{\Sigma}^* \mathbf{F}$ , and the optimal cost is

$$\|(\mathbf{I} - \mathbf{F}\mathbf{F}^*)\mathbf{\Sigma}\|_F^2. \quad (4.6.3)$$

2. Let  $\mathbf{P} = \mathbf{I} - \mathbf{F}\mathbf{F}^*$ , and write  $\nu_i = \|\mathbf{P}\mathbf{e}_i\|_2^2$ . Argue that  $\sum_{i=1}^n \nu_i = n_1 - r$  and  $\nu_i \in [0, 1]$ . Conclude that

$$\|\mathbf{P}\mathbf{\Sigma}\|_F^2 = \sum_{i=1}^{n_1} \sigma_i^2 \nu_i \geq \sum_{i=r+1}^{n_1} \sigma_i^2, \quad (4.6.4)$$

with equality if and only if  $\nu_1 = \nu_2 = \dots = \nu_r = 0$  and  $\nu_{r+1} = \dots = \nu_n$ . Conclude that Theorem 4.2.4 holds in the special case  $\mathbf{Y} = \mathbf{\Sigma}$ .

3. Extend your argument to the situation in which the  $\sigma_i$  are not distinct (i.e.,  $\sigma_i = \sigma_{i+1}$  for some  $i$ ).
4. Extend your argument to any  $\mathbf{Y} \in \mathbb{R}^{n \times n}$ . Hint: use the fact that the Frobenius norm  $\|\mathbf{M}\|_F$  is unchanged by orthogonal transformations of the rows and columns:  $\|\mathbf{M}\|_F = \|\mathbf{R}\mathbf{M}\mathbf{S}\|_F$  for any orthogonal matrices  $\mathbf{R}, \mathbf{S}$ .

**4.4** (Minimal Rank Approximation). We consider a variant of Theorem 4.2.4 in which we are given a data matrix  $\mathbf{Y}$  and we want to find a matrix  $\mathbf{X}$  of minimum rank that approximates  $\mathbf{Y}$  up to some given fidelity:

$$\begin{aligned} & \text{minimize} && \text{rank}(\mathbf{X}), \\ & \text{subject to} && \|\mathbf{X} - \mathbf{Y}\|_F \leq \varepsilon. \end{aligned} \quad (4.6.5)$$

Give an expression for the optimal solution(s) to this problem, in terms of the SVD of  $\mathbf{Y}$ . Prove that your expression is correct.

**4.5** (Multiple and repeated eigenvalues). Consider the eigenvector problem

$$\text{minimize} \quad -\frac{1}{2} \mathbf{q}^* \mathbf{\Gamma} \mathbf{q} \quad \text{subject to} \quad \|\mathbf{q}\|_2^2 = 1, \quad (4.6.6)$$

where  $\mathbf{\Gamma}$  is a symmetric matrix. In the text, we argued that when the eigenvalues of  $\mathbf{\Gamma}$  are distinct, every local minimizer of this problem is global.

(i) Argue that even when  $\mathbf{\Gamma}$  has repeated eigenvalues, every local minimum of this problem is global. (ii) Now suppose we wish to find multiple eigenvector/eigenvalue pairs. Consider the optimization problem

$$\begin{aligned} & \text{minimize} && -\frac{1}{2} \mathbf{Q}^* \mathbf{\Gamma} \mathbf{Q} \\ & \text{subject to} && \mathbf{Q} \in \text{St}(n, p) = \{\mathbf{Q} \in \mathbb{R}^{n \times p} \mid \mathbf{Q}^* \mathbf{Q} = \mathbf{I}\}. \end{aligned} \quad (4.6.7)$$

Argue that every local minimizer of this problem has the form

$$\mathbf{Q} = [\mathbf{u}_1, \dots, \mathbf{u}_p] \mathbf{\Pi}, \quad (4.6.8)$$

where  $\mathbf{u}_1, \dots, \mathbf{u}_p$  are eigenvectors of  $\mathbf{\Gamma}$  associated with the  $p$  largest eigenvalues, and  $\mathbf{\Pi}$  is a permutation matrix.

**4.6** (The Power Method). We show how to compute eigenvectors (and hence singular vectors) using the power method. Let  $\mathbf{\Gamma} \in \mathbb{R}^{n \times n}$  be a symmetric positive semidefinite matrix. Let  $\mathbf{q}^{(0)}$  be a random vector that is

uniformly distributed on the sphere  $\mathbb{S}^{n-1}$  (we can generate such a random vector by taking an  $n$ -dimensional iid  $\mathcal{N}(0, 1)$  vector and then normalizing it to have unit  $\ell^2$  norm). Generate a sequence of vectors  $\mathbf{q}^{(1)}, \mathbf{q}^{(2)}, \dots$  via the iteration

$$\mathbf{q}^{(k+1)} = \frac{\mathbf{\Gamma} \mathbf{q}^{(k)}}{\|\mathbf{\Gamma} \mathbf{q}^{(k)}\|_2}. \quad (4.6.9)$$

This iteration is called the power method.

Suppose that there is a gap between the first and second eigenvalues of  $\mathbf{\Gamma}$ :  $\lambda_1(\mathbf{\Gamma}) > \lambda_2(\mathbf{\Gamma})$ .

1. What does  $\mathbf{q}^{(k)}$  converge to? Hint: write  $\mathbf{\Gamma} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^*$  in terms of its eigenvectors/values. How does  $\mathbf{V}^* \mathbf{q}^{(k)}$  evolve?
2. Obtain a bound on the error  $\|\mathbf{q}^{(k)} - \mathbf{q}^{(\infty)}\|_2$  in terms of the spectral gap  $\frac{\lambda_1 - \lambda_2}{\lambda_1}$ .
3. Your bound in (ii) should suggest that as long as there is a gap between  $\lambda_1$  and  $\lambda_2$ , the power method converges rapidly. How does the method behave if  $\lambda_1 = \lambda_2$ ?
4. How can we use the power method to compute the singular values of a matrix  $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$ ?

**4.7** (Convex Envelope Property via the Bidual). In Theorem 4.3.2, we proved that the nuclear norm  $\|\mathbf{X}\|_*$  is the convex envelope of  $\text{rank}(\mathbf{X})$  over the operator norm ball  $\mathbf{B}_{\text{op}} = \{\mathbf{X} \mid \|\mathbf{X}\| \leq 1\}$ . Here, we give an alternative derivation of this result, using the fact that the biconjugate of a function over a set  $\mathbf{B}$  is the convex envelope. Let  $f(\mathbf{X}) = \text{rank}(\mathbf{X})$  denote the rank function.

1. Prove that the Fenchel dual

$$f^*(\mathbf{Y}) = \sup_{\mathbf{X} \in \mathbf{B}} \{\langle \mathbf{X}, \mathbf{Y} \rangle - f(\mathbf{X})\}$$

can be expressed as

$$f^*(\mathbf{Y}) = \|\mathcal{D}_1[\mathbf{Y}]\|_*,$$

where  $\mathcal{D}_\tau[\mathbf{M}]$  is the singular value thresholding operator, given by  $\mathcal{D}_\tau[\mathbf{M}] = \mathbf{U} \mathcal{S}_\tau[\mathbf{S}] \mathbf{V}^*$  for any singular value decomposition  $\mathbf{M} = \mathbf{U} \mathbf{S} \mathbf{V}^*$  of  $\mathbf{M}$ .

2. Prove that the dual of  $f^*$ ,

$$f^{**}(\mathbf{X}) = \sup_{\mathbf{Y}} \langle \mathbf{X}, \mathbf{Y} \rangle - f^*(\mathbf{Y})$$

can satisfies

$$f^{**}(\mathbf{X}) = \|\mathbf{X}\|_*.$$

3. Use Proposition B.2.10 of Appendix B to conclude that  $\|\cdot\|_*$  is the convex envelope of  $\text{rank}(\cdot)$  over  $\mathbf{B}$ .

**4.8** (Convexifying Low-rank Approximation). Consider the following optimization problem:

$$\begin{aligned} & \text{minimize} && \|\Pi Y\|_F^2 \\ & \text{subject to} && \mathbf{0} \preceq \Pi \preceq \mathbf{I}, \text{ trace}[\Pi] = m - r. \end{aligned} \quad (4.6.10)$$

Prove that if  $\sigma_r(Y) > \sigma_{r+1}(Y)$ , this problem has a unique optimal solution  $\Pi_*$ , which is the orthoprojector onto the linear span of the  $n_1 - r$  trailing singular vectors  $\mathbf{u}_{r+1}, \mathbf{u}_{r+2}, \dots, \mathbf{u}_{n_1}$ . The matrix  $(\mathbf{I} - \Pi_*)Y$  is the best rank- $r$  approximation to  $Y$ .

**4.9** (Tangent Space to the Rank- $r$  Matrices). Consider a matrix  $\mathbf{X}_0$  of rank  $r$  with compact singular value decomposition  $\mathbf{X}_0 = \mathbf{U}\Sigma\mathbf{V}^*$ . Argue that the tangent space to the collection  $\mathbf{R}_r = \{\mathbf{X} \mid \text{rank}(\mathbf{X}) = r\}$  at  $\mathbf{X}_0$  is given by  $T = \{\mathbf{U}\mathbf{W}^* + \mathbf{H}\mathbf{V}^*\}$ . Hint: consider generating a nearby low-rank matrix by writing  $\mathbf{X}' = (\mathbf{U} + \Delta\mathbf{U})(\Sigma + \Delta\Sigma)(\mathbf{V} + \Delta\mathbf{V})^*$ .

**4.10** (Quadratic Measurements). Consider a target vector  $\mathbf{x}_0 \in \mathbb{R}^{n \times n}$ . In many applications, the observation can be modeled as a quadratic function of the vector  $\mathbf{x}_0$ . In notation, we see the squares  $y_1 = \langle \mathbf{a}_1, \mathbf{x}_0 \rangle^2$ ,  $y_2 = \langle \mathbf{a}_2, \mathbf{x}_0 \rangle^2$ ,  $\dots$ ,  $y_m = \langle \mathbf{a}_m, \mathbf{x}_0 \rangle^2$  of the projections of  $\mathbf{x}_0$  onto vectors  $\mathbf{a}_1 \dots \mathbf{a}_m$ . Notice that from this observation, it is only possible to reconstruct  $\mathbf{x}_0$  up to a sign ambiguity:  $-\mathbf{x}_0$  produces exactly the same observation.

1. Consider the quadratic problem

$$\min_{\mathbf{x}} \sum_{i=1}^n \left( y_i - \langle \mathbf{a}_i, \mathbf{x} \rangle^2 \right)^2. \quad (4.6.11)$$

Is this problem convex in  $\mathbf{x}$ ?

2. Convert this to a convex problem, by replacing the vector valued variable  $\mathbf{x}$  with a matrix valued variable  $\mathbf{X} = \mathbf{x}\mathbf{x}^*$ : convert the problem to

$$\min_{\mathbf{X}} \sum_{i=1}^n \left( y_i - \langle \mathbf{A}_i, \mathbf{X} \rangle \right)^2. \quad (4.6.12)$$

How should we choose the matrices  $\mathbf{A}_1, \dots, \mathbf{A}_m$ ? Show that if  $m < n^2$ ,  $\mathbf{X}_0 = \mathbf{x}_0\mathbf{x}_0^*$  is not the unique optimal solution to this problem. How can we use the fact that  $\text{rank}(\mathbf{X}_0) = 1$  to improve this?

3. In the absence of noise, we can attempt to solve for  $\mathbf{X}_0$  by solving the convex program

$$\min \|\mathbf{X}\|_* \quad \text{such that} \quad \mathcal{A}[\mathbf{X}] = \mathbf{y}. \quad (4.6.13)$$

Implement this optimization using a custom algorithm or CVX. Does it typically recover  $\mathbf{X}_0$ ?

4. Does the operator  $\mathcal{A}$  satisfy the rank RIP?

**4.11** (Proof of Theorem 4.4.2). We prove Theorem 4.4.2. The goal here is to show that the solution to

$$\min_{\mathbf{X}} \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X} - \mathbf{M}\|_F^2 \quad (4.6.14)$$

is given by  $\mathcal{D}_1[\mathbf{M}]$

1. Argue that Problem (4.6.14) is strongly convex, and hence has a unique optimal solution.
2. Show that a solution  $\mathbf{X}_*$  is optimal if and only if  $\mathbf{X}_* \in \mathbf{M} - \partial \|\cdot\|_*(\mathbf{X}_*)$ .
3. Using the condition from part (ii), show that if  $\mathbf{M}$  is diagonal, i.e.,  $\mathbf{M}_{ij} = 0$  for  $i \neq j$ , then  $\mathcal{S}_1[\mathbf{M}]$  is the unique optimal solution to (4.6.14).
4. Use the SVD to argue that in general,  $\mathcal{D}_1[\mathbf{M}]$  is the unique optimal solution to (4.6.14).

**4.12** (Uniform Matrix Completion?). Let  $\Omega$  be a strict subset of  $[n] \times [n]$ . Show that there exist two matrix  $\mathbf{X}_0$  and  $\mathbf{X}'_0$  of rank one such that  $\mathcal{P}_\Omega[\mathbf{X}_0] = \mathcal{P}_\Omega[\mathbf{X}'_0]$ . The implication of this is that it is not possible to reconstruct all low-rank matrices from the same observation  $\Omega$ .

**4.13** (Unique Optimality for Matrix Completion). Consider the optimization problem

$$\begin{aligned} & \text{minimize} && \|\mathbf{X}\|_* \\ & \text{subject to} && \mathcal{P}_\Omega[\mathbf{X}] = \mathbf{P}_\Omega[\mathbf{X}_0]. \end{aligned} \quad (4.6.15)$$

Suppose that  $\|\mathcal{P}_\Omega \mathcal{P}_T\| < 1$ . Assume we can find some  $\mathbf{\Lambda}$  such that

1.  $\mathbf{\Lambda}$  is supported on  $\Omega$  and
2.  $\mathbf{\Lambda} \in \partial \|\cdot\|_*(\mathbf{X}_0)$  – i.e.,  $\mathcal{P}_T[\mathbf{\Lambda}] = \mathbf{U}\mathbf{V}^*$  and  $\|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}]\| < 1$ .

Show that  $\mathbf{X}_0$  is the unique optimal solution to the optimization problem.

**4.14.** Prove Theorem 4.3.5.

**4.15.** Prove Theorem 4.3.12

**4.16.** Fill in the detailed steps of proof for Theorem 4.3.15.

**4.17.** Derive detailed steps that prove the error bound (4.3.83) in the proof of Theorem 4.3.13.

**4.18.** Show that in Lemma 4.4.4, any subdifferential of nuclear norm must be of the form given in (4.4.21).



**4.19.** Let  $\mathcal{R}_\Omega[\mathbf{X}_0] = \sum_{\ell=1}^q [\mathbf{X}_0]_{i_\ell, j_\ell} \mathbf{e}_{i_\ell} \mathbf{e}_{j_\ell}^*$  with each  $(i_\ell, j_\ell)$  chosen iid at random from the uniform distribution on  $[n] \times [n]$ . Use the matrix Bernstein inequality to show that if  $q > C\nu nr \log n$  for sufficiently large  $C$ , we have

$$\left\| \mathcal{P}_{T^\perp} \frac{n^2}{q} \mathcal{R}_\Omega \mathcal{P}_T \right\| \leq t. \quad (4.6.16)$$

for any arbitrarily small constant  $t$  with high probability. [Hint: similar to the proof of Lemma 4.4.7.]

**4.20.** For the dual certificate  $\mathbf{\Lambda}$  constructed from the golfing scheme, use the fact in Exercise 4.19 and the fact that  $\left\| \frac{n^2}{q} \mathcal{P}_{T^\perp} \mathcal{R}_{\Omega_j} [\mathbf{E}_j] \right\|_F \leq \left\| \frac{n^2}{q} \mathcal{P}_{T^\perp} \mathcal{R}_{\Omega_j} \mathcal{P}_T \right\| \|\mathbf{E}_j\|_F$ , show that if

$$m > C\nu nr^2 \log^2 n$$

for a large enough constant  $C$ , we have  $\|\mathcal{P}_{T^\perp}[\mathbf{\Lambda}]\| \leq 1/2$  with high probability.

**4.21.** Prove Lemma 4.4.9. Hint: write:

$$(q^{-1} \mathcal{P}_\Omega - \mathcal{I}) \mathbf{Z} = \sum_{ij} \underbrace{Z_{ij} (q^{-1} \mathbb{1}_{ij \in \Omega} - 1)}_{\doteq \mathbf{W}_{ij}} \mathbf{E}_{ij},$$

and apply the operator Bernstein inequality, controlling the operator norm of  $\mathbf{W}_{ij}$  in terms of  $\|\mathbf{Z}\|_\infty$  and controlling the matrix variance in terms of  $\|\mathbf{Z}\|_{rc}$ .

**4.22.** Prove Lemma 4.4.10. Use the matrix Bernstein inequality to obtain a bound on the probability that the  $\ell$ -th row  $\|\mathbf{e}_\ell^* (q^{-1} \mathcal{P}_T \mathcal{P}_\Omega - \mathcal{P}_T) \mathbf{Z}\|$  is large, repeat for each column, and then sum the failure probabilities over all rows and columns to obtain a bound on the probability that the  $\|\cdot\|_{rc}$  is large. Hint: apply the matrix Bernstein inequality to the random vector:

$$\mathbf{e}_\ell^* (q^{-1} \mathcal{P}_T \mathcal{P}_\Omega - \mathcal{P}_T) \mathbf{Z} = \sum_{ij} \underbrace{Z_{ij} (q^{-1} \mathbb{1}_{ij \in \Omega} - 1) \mathbf{e}_\ell^* \mathcal{P}_T [\mathbf{E}_{ij}]}_{\doteq \mathbf{W}_{ij}}.$$

**4.23.** Prove Lemma 4.4.11. Apply the standard Bernstein inequality to bound the probability that the  $k, l$  entry of  $(\mathcal{P}_T - q^{-1} \mathcal{P}_T \mathcal{P}_\Omega \mathcal{P}_T) \mathbf{Z}$  is large, and then sum this probability over all entries  $k, l$  to bound the probability that the  $\ell^\infty$  norm is large. For the  $k, l$  entry work with the sum of independent random variables

$$\begin{aligned} [(\mathcal{P}_T - q^{-1} \mathcal{P}_T \mathcal{P}_\Omega \mathcal{P}_T) \mathbf{Z}]_{kl} &= Z_{kl} - [q^{-1} \mathcal{P}_T \mathcal{P}_\Omega \mathbf{Z}]_{kl} \\ &= \sum_{ij} \underbrace{n^{-2} Z_{kl} - q^{-1} \mathbb{1}_{ij \in \Omega} \langle \mathcal{P}_T \mathbf{E}_{kl}, \mathcal{P}_T \mathbf{E}_{ij} \rangle}_{\doteq \mathbf{W}_{ij}} Z_{ij}. \end{aligned}$$