

In this chapter, we discuss the most basic problem in high-dimensional statistics, which is the so-called high dimensional linear regression. In its low-dimensional setting, in which the number of predictors d is substantially less than the sample size n , the associated theory is classical. By contrast, the goal in this chapter is to develop theory applicable to the high-dimensional regime: (1) $n \asymp p$ and (2) $p \gg n$. As one might expect, if the model lacks any additional structure, then there is no hope of obtaining consistent estimator when p/n not goes to 0. This leads to the great interest in sparsity, we focus on different types of sparse models in this chapter.

7.1 Problem Settings

Let $\theta^* \in \mathbb{R}^d$ be the “true” parameter, $X \in \mathbb{R}^{n \times d}$ be the feature matrix and $Y \in \mathbb{R}^n$ be the response vector. The model we consider is in following form:

$$Y = X\theta^* + \epsilon \quad (7.1)$$

The focus of this chapter is settings in which the sample size n is smaller than the number of predictors d . Note when $n < d$, it is impossible to obtain any meaningful estimates of θ^* unless there is some additional assumption on the low-dimensional structure. One widely considered assumption is the *hard sparsity* assumption:

Assumption 7.1 (hard sparsity). Denote by $S(\theta^*) = \{j \in \{1, 2, \dots, d\} \mid \theta_j^* \neq 0\}$ the support of vector θ^* . We assume that the cardinality $|S(\theta^*)|$ is substantially smaller than d , i.e. $|S(\theta^*)| \leq d' \ll d$.

Note Assumption 7.1 is a bit strong since it regularizes the maximal number of non-zero terms to be some value substantially smaller than d . In real applications, we might prefer to use another notion of sparsity, named weak sparsity, where lots of terms are closed to zero. There are different ways to formalize such an idea, one is via l_q -norm, which is formalized in the following assumption.

Assumption 7.2 (weak sparsity). θ^* lies in the l_q -ball around 0, that is,

$$B_q(r_q) = \{\theta \in \mathbb{R}^d \mid \|\theta\|_q \leq r_q\} \quad (7.2)$$

Examples of Assumption 7.2 are given in Figure 7.1.

7.2 Applications of sparse linear models

Before studying the theoretical performance of estimators under sparse assumption, we first introduce some applications of sparse linear model.

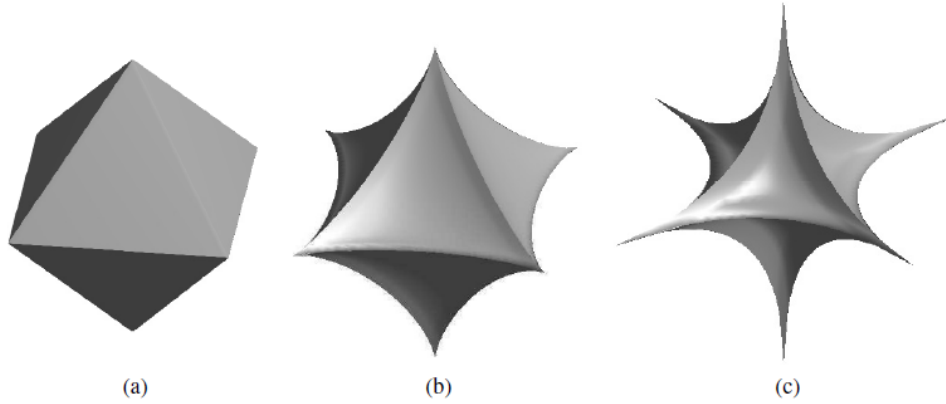


Figure 7.1: Illustrations of the l_q balls for different choices of parameter $q \in (0, 1]$. (a) $q = 1$, (b) $q = 0.75$ and (c) $q = 0.5$.

Example 7.1 (Gaussian sequence model). We assume the observations are generated from the following model sequentially:

$$Y_k = \sqrt{n}\theta_k^* + \epsilon_k \quad (7.3)$$

where $\epsilon_k \sim N(0, \sigma^2)$ are i.i.d. and $k = 1, \dots, n$. Note this model can be rewritten as

$$Y = \sqrt{n}I\theta^* + \epsilon \quad (7.4)$$

One important characteristic of the model above is, the number of parameters $p = n$, which means the number of parameters grows as the number of observations grows. Although it appears simple on the surface, it is a surprisingly rich model: indeed, many problems in nonparametric estimation, among them regression and density estimation, can be reduced to an “equivalent” instance of the Gaussian sequence model, in the sense that the optimal rates for estimation are the same under both models.

Example 7.2 (Selection of Gaussian graphical models). Let (X_1, X_2, \dots, X_n) be a zero-mean Gaussian random vector with a non-degenerate covariance matrix. The density function is represented as

$$p(x_1, \dots, x_n) = \frac{1}{\sqrt{(2\pi)^d \det(\Theta^{-1})}} \exp\left(-\frac{1}{2}x^\top \Theta x\right) \quad (7.5)$$

where Θ is the precision-matrix. By theory of Gaussian graphical model, we have that (i, j) -term in Θ is zero if and only if $x_i \perp x_j \mid x_{-\{i,j\}}$, which means there is no linkage between node i and j in graph representation. Hence it’s important to introduce some assumptions on the sparsity of precision matrix Θ .

7.3 Recovery in the noiseless setting

In this section, we consider the model that observations are perfect, i.e., we consider

$$Y = X\theta \quad (7.6)$$

Assume we are told that there is some vector θ^* with at most $s \ll d$ non-zero entries such that $y = X\theta^*$. Our goal is to recover the sparse solution to the linear system. That is, we aim to solve problem

$$\min_{\theta \in \mathbb{R}^d} \|\theta\|_0, \text{ such that } X\theta = Y \quad (7.7)$$

where $\|\theta\|_0 = \sum_{k=1}^d \mathbb{I}(\theta_k \neq 0)$ is the psedo-norm that represents the number of non-zero entries of vector θ . Note the object function is non-smooth and non-convex, which is computationally intractable. The most closed convex relaxation is to replace the l_0 norm by l_1 norm, that is,

$$\min_{\theta \in \mathbb{R}^d} \|\theta\|_1, \text{ such that } X\theta = Y \quad (7.8)$$

We refer this problem as basis pursuit linear program (see Chen, Donoho and Saunders (1998))

7.3.1 Exact recovery and restricted nullspace

We now turn to a theoretical question: when is solving the basis pursuit program (l_1 problem) (7.8) equivalent to solving the original l_0 -problem (7.7).

Assume that $\theta^* \in \mathbb{R}^d$ such that $Y = X\theta^*$, which means the linear restriction is feasible. Next, we denote by $S = \{j \in \{1, 2, \dots, d\} \mid \theta_j^* \neq 0\}$ the support vector and S^c the complement of S . Recall any solution to $X\theta = Y$ must be of form

$$\theta^* + \Delta, \text{ where } \Delta \in \mathbf{Null}(X) = \{\Delta \mid X\Delta = 0\} \quad (7.9)$$

To proceed, we define the tangent cone of the l_1 -ball at θ^* , that is,

$$T(\theta^*) = \{\Delta \in \mathbb{R}^d \mid \|\theta^* + t\Delta\|_1 \leq \|\theta^*\|_1, \text{ for some } t > 0\} \quad (7.10)$$

It's easy to verify that $T(\theta^*)$ is a cone. For a tangent cone, the most important information is the direction. Note $T(\theta^*) \cap \mathbf{Null}(X)$ is still a cone, meaning that if $T(\theta^*) \cap \mathbf{Null}(X) \neq \emptyset$, then θ^* is not a solution to basis pursuit problem (7.8). Here we present two possible cases between the solution of (7.8) and θ^* (see Figure 7.2).

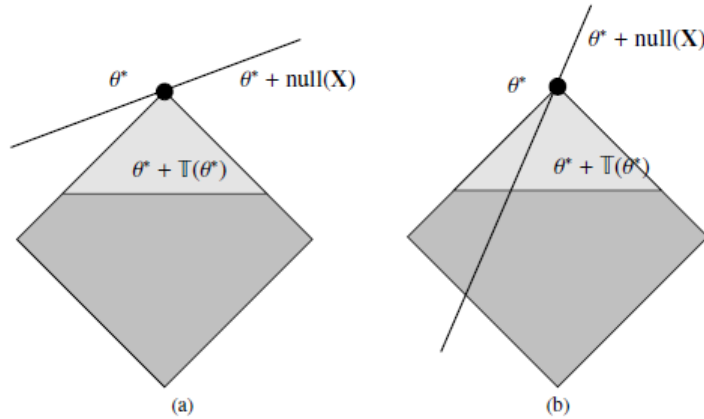


Figure 7.2: Two possible cases for whether θ^* is the unique solution.

Next, we investigate the condition for θ^* to be the unique solution of problem (7.8), which is known as the *restricted nullspace property*. More specifically, given support S , we define

$$C(S) = \{\Delta \in \mathbb{R}^d \mid \|\Delta_{S^c}\|_1 \leq \|\Delta_S\|_1\} \quad (7.11)$$

We have the following definition for restricted nullspace property.

Definition 7.1 (restricted nullspace property). Matrix X satisfies the restricted nullspace property with respect to S , if $C(S) \cap \mathbf{Null}(X) = \{0\}$.

One quick way to understand restricted nullspace property is: let $\theta^* = (\theta_1, \dots, \theta_s, 0, \dots, 0)$ be the underlying true vector. In order to make θ^* possible to be recovered exactly, the optimization problem (7.8) should have unique solution θ^* . This means for any nonzero $\Delta = (\Delta_1, \dots, \Delta_d) \in \mathbf{Null}(X)$, $\theta^* + \Delta = (\theta_1 + \Delta_1, \dots, \theta_s + \Delta_s, \Delta_{s+1}, \dots, \Delta_d)$ cannot have smaller l_1 norm than $\theta^* = (\theta_1, \dots, \theta_s, 0, \dots, 0)$. Hence

$$\sum_{i \in S^c} |\Delta_i| > \sum_{i \in S} |\theta_i| - \sum_{i \in S} |\theta_i + \Delta_i| \quad (7.12)$$

Since $X\Delta = X(-\Delta) = 0$, we have that

$$\sum_{i \in S^c} |\Delta_i| \geq \max\left\{\sum_{i \in S} (|\theta_i| - |\theta_i + \Delta_i|), \sum_{i \in S} (|\theta_i| - |\theta_i - \Delta_i|)\right\} \quad (7.13)$$

For a small enough Δ , (i.e., $|\Delta_i| < |\theta_i|$ for $i \in S$), this implies $\|\Delta_{S^c}\|_1 > \|\Delta_S\|_1$. Hence there is no interception between $C(S)$ and $\mathbf{Null}(X)$.

The following theorem justifies above arguments in more details.

Theorem 7.1. *The following two properties are equivalent:*

- a. *For any vector θ^* with support S , the basis pursuit program (7.8) applied with condition $X\theta^* = Y$ has unique solution $\hat{\theta} = \theta^*$.*
- b. *The matrix X satisfies the restricted nullspace property with respect to S .*

Proof. We first show that $b \Rightarrow a$. Denote by $\hat{\theta}$ a solution of problem (7.8). By optimality of $\hat{\theta}$, we have that $\|\hat{\theta}\|_1 \leq \|\theta^*\|_1$. Let $\Delta = \hat{\theta} - \theta^*$, note $X\Delta = 0$, by restricted nullspace property, we have either $\Delta = 0$ or $\|\Delta_{S^c}\|_1 > \|\Delta_S\|_1$ when $\Delta \neq 0$. For the second case, we have that $\|\theta_S^*\|_1 = \|\theta^*\|_1 \geq \|\hat{\theta}\|_1 \geq \|\Delta_{S^c}\|_1 + \|\theta_S^*\|_1 - \|\Delta_S\|_1$, which implies $\|\Delta_{S^c}\|_1 \leq \|\Delta_S\|_1$. Hence we get the contradiction. This means $\Delta = 0$, or equivalently, θ^* is the unique solution to the basis pursuit program (7.8), when condition $X\theta^* = Y$ is equipped.

Next, we show $a \Rightarrow b$. By condition (a.), for any $\Delta \in \mathbf{Null}(X)$ with $\Delta \neq 0$, $(\Delta_S, 0)^\top$ is the unique solution of the following problem:

$$\min_{\theta \in \mathbb{R}^d} \|\theta\|_1, \text{ such that } X\theta = X(\Delta_S, 0)^\top \quad (7.14)$$

Note $X(0, -\Delta_{S^c})^\top - X(\Delta_S, 0)^\top = -X\Delta = 0$, we must have that $(0, -\Delta_{S^c})^\top$ is a feasible solution. By optimality of $(\Delta_S, 0)^\top$, this means $\|\Delta_S\|_1 < \|\Delta_{S^c}\|_1$. Hence $\Delta \notin C(S)$, which means X satisfies the restricted nullspace property with respect to S .

This completes the proof.

Now we have know that, exact recovery is equivalent to the requirement that X satisfies the restricted nullspace property with respect to S . One question is: when will X satisfies the restricted nullspace property with respect to S ? Next, we give some sufficient conditions for restricted nullspace property to hold.

7.3.2 Two sufficient conditions for restricted nullspace

The earliest sufficient condition is on the pairwise *incoherence parameter*, which is defined as

$$\delta_{pw}(X) = \max_{i,j} |(\frac{1}{n}X^\top X - I)_{i,j}| \quad (7.15)$$

The smaller $\delta_{pw}(X)$ is, the more $\frac{1}{\sqrt{n}}X$ is closed to an orthonormal matrix. The next theorem guarantee a uniform version of the restricted nullspace property.

Theorem 7.2. *If the pairwise incoherence satisfies the bound $\delta_{pw}(X) \leq \frac{1}{3s}$, then the restricted nullspace property holds for all subsets S of cardinality at most s .*

The proof will be uploaded in a later version.

A more widely used sufficient condition is by the restricted isometry property (RIP), which is a natural generalization of the pairwise incoherence condition.

Definition 7.2 (Restricted isometry property). For a given integer $s \in \{1, 2, \dots, d\}$, we say that $X \in \mathbb{R}^{n \times d}$ satisfies a restricted isometry property of order s with constant $\delta_s(X) > 0$, if

$$\|\frac{1}{n}X_s^\top X_s - I_s\|_2 \leq \delta_s(X) \quad (7.16)$$

for all subsets S of size at most s . Here $\|A\|_2 = \sigma_{\max}(A)$ is the largest singular value.

Definition 7.2 implies that, for any subsets S of size at most s , $\frac{1}{n}X_s^\top X_s - I \preceq \delta_s^2(X)I$, or equivalently

$$\frac{1}{n}X^\top X - I \preceq \delta_s^2(X)I \quad (7.17)$$

Although RIP imposes constraints on much larger submatrices than pairwise incoherence, the magnitude of the constraints required to guarantee the uniform restricted nullspace property can be milder. Indeed, for any matrix X and sparsity level $s \in \{2, \dots, d\}$, we have the sandwich relation

$$\delta_{pw}(X) \leq \delta_s(X) \leq s\delta_{pw}(X) \quad (7.18)$$

where neither bound can be improved in general. The following theorem shows that suitable control on the RIP constants implies that the restricted nullspace property holds.

Theorem 7.3. *If the RIP constant of order $2s$ is bounded by $\delta_{2s}(X) < 1/3$, then the uniform restricted nullspace property holds for any subset S of cardinality $|S| \leq s$.*

The proof will be uploaded in a later version.

As a summary, both Theorem 7.2 and Theorem 7.3 can be used as sufficient conditions for restricted nullspace property. A major advantage of the RIP approach is that for various classes of random design matrices, of particular interest in compressed sensing, it can be used to guarantee exactness of basis pursuit using a sample size n , which is much smaller than that guaranteed by pairwise incoherence.

Another aspect that should be noted is, unlike the restricted nullspace property, neither the pairwise incoherence condition nor the RIP condition are necessary conditions.

7.4 Estimation in noisy settings

In the noiseless setting, we discuss the conditions that can guarantee an exact recovery by solving the l_1 relaxed optimization problem. In this section, however, we will discuss the following noisy setting:

$$Y = X\theta^* + \epsilon \quad (7.19)$$

Our estimation method is by *LASSO program*:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|Y - X\theta\|_2^2 + \lambda_n \|\theta\|_1 \quad (7.20)$$

In the noisy setting, we can no longer expect to achieve perfect recovery. Hence we will mainly focus on how to bound the l_2 estimation error $\|\hat{\theta} - \theta^*\|_2$. Recall in noiseless setting, the condition we require is the restricted nullspace property. However, in the presence of noise, we will require a slightly stronger condition: for any $\alpha \geq 1$, we define

$$C_\alpha(S) = \{\Delta \in \mathbb{R}^d \mid \|\Delta_{S^c}\|_1 \leq \alpha \|\Delta_S\|_1\} \quad (7.21)$$

It's easy to check that $C(S) \subseteq C_\alpha(S)$.

Definition 7.3 (restricted eigenvalue condition). The matrix X satisfies the restricted eigenvalue (RE) condition over subset S with parameters (κ, α) if

$$\frac{\frac{1}{n} \|X\Delta\|_2^2}{\|\Delta\|_2^2} \geq \kappa \quad (7.22)$$

for all $\Delta \in C_\alpha(S)$.

Indeed RE condition is a strengthening of the restricted nullspace property. Note if RE condition holds with parameters $(\kappa, 1)$, then for any $\Delta \in \text{Null}(X)$ with $\Delta \neq 0$, i.e. $X\Delta = 0$, we must have $\frac{1}{n} \|X\Delta\|_2^2 / \|\Delta\|_2^2 = 0 < \kappa$, which means $\Delta \notin C(S)$. Hence restricted nullspace property holds.

Another comments for the intuition of RE condition is on the curvature, which can be found on page 208 of the book.

7.4.1 Bound on l_2 -error for hard sparse models

In this section, we impose the following assumption.

- (A1). The vector θ^* is supported on $S \subseteq \{1, 2, \dots, d\}$ with $|S| = s$.
(A2). The design matrix satisfies the restricted eigenvalue condition over S with parameter $(\kappa, 3)$.

The following theorem justifies the bounds on the $\|\hat{\theta} - \theta^*\|_2$, where $\hat{\theta}$ is the LASSO estimator.

Theorem 7.4. *Under assumptions (A1) and (A2), any solution to (7.20) with $\lambda_n \geq 2\|\frac{1}{n}X^\top \epsilon\|_\infty$ must satisfy*

$$\|\hat{\theta} - \theta^*\|_2 \leq \frac{3}{\kappa} \sqrt{s} \lambda_n \quad (7.23)$$

Before talking about the proof of Theorem 7.4, we first investigate one example to get some intuitions.

Example 7.3 (Classical linear Gaussian model). Let $Y = X\theta^* + \epsilon$, where $X \in \mathbb{R}^{n \times d}$ is a fixed design matrix and ϵ has i.i.d entries from $N(0, \sigma^2 I)$. Suppose X satisfies the restricted eigenvalue property and is C -column normalized, i.e., $\max_{j=1, \dots, d} \frac{1}{\sqrt{n}} \|X_j\|_2 \leq C$, where X_j denotes the j -th column of X . Then random vector $\frac{1}{n}X^\top \epsilon$ satisfies $\text{Var}(\frac{1}{n}X^\top \epsilon) = \frac{\sigma^2}{n^2} X^\top X$ and each entry of $\frac{1}{n}X^\top \epsilon$ has variance no greater than $\frac{1}{n}\sigma^2 C^2$. Hence, we have that

$$\mathbb{P}(\|\frac{1}{n}X^\top \epsilon\|_\infty \geq C\sigma(\sqrt{\frac{2 \log d}{n}} + \delta)) \leq 2 \exp(-\frac{n\delta^2}{2}) \quad (7.24)$$

This means with probability at least $1 - 2 \exp(-\frac{n\delta^2}{2})$, by choosing $\lambda_n = 2C\sigma(\sqrt{\frac{2 \log d}{n}} + \delta)$ and applying Theorem 7.4, we have that

$$\|\hat{\theta} - \theta^*\|_2 \leq \frac{6C}{\kappa} \sqrt{s} \sigma (\sqrt{\frac{2 \log d}{n}} + \delta) \quad (7.25)$$

Now we are ready to show Theorem 7.4.

Proof. We first show that, when $\lambda_n \geq 2\|\frac{1}{n}X^\top \epsilon\|_\infty$, it holds that $\hat{\Delta} = \hat{\theta} - \theta^* \in C_3(S)$, i.e., $\|\hat{\Delta}_{S^c}\|_1 \leq 3\|\hat{\Delta}_S\|_1$. Toward this ends, let's define

$$L(\theta, \lambda_n) = \frac{1}{2n} \|Y - X\theta\|_2^2 + \lambda_n \|\theta\|_1 \quad (7.26)$$

By optimality of $\hat{\theta}$, we have that

$$\frac{1}{2n} \|Y - X\hat{\theta}\|_2^2 + \lambda_n \|\hat{\theta}\|_1 \leq \frac{1}{2n} \|\epsilon\|_2^2 + \lambda_n \|\theta^*\|_1 \quad (7.27)$$

or equivalently,

$$\frac{1}{2n} (\|\epsilon - X\hat{\Delta}\|_2^2 - \|\epsilon\|_2^2) = \frac{1}{2n} \|X\hat{\Delta}\|_2^2 - \frac{1}{n} \epsilon^\top X\hat{\Delta} \leq \lambda_n (\|\theta^*\|_1 - \|\hat{\theta}\|_1) \quad (7.28)$$

This further implies

$$\frac{1}{2n} \|X\hat{\Delta}\|_2^2 \leq \frac{1}{n} \epsilon^\top X\hat{\Delta} + \lambda_n (\|\theta^*\|_1 - \|\hat{\theta}\|_1) \quad (7.29)$$

By assumption that θ^* is S -sparse, we have that

$$\|\theta^*\|_1 - \|\hat{\theta}\|_1 = \|\theta_S^*\|_1 - \|\theta_S^* + \hat{\Delta}_S\|_1 - \|\hat{\Delta}_{S^c}\|_1 \quad (7.30)$$

which implies

$$\begin{aligned} 0 \leq \frac{1}{2n} \|X\hat{\Delta}\|_2^2 &\leq \frac{1}{n} \epsilon^\top X\hat{\Delta} + \lambda_n(\|\theta^*\|_1 - \|\hat{\theta}\|_1) \\ &\leq \left(\frac{1}{n} X^\top \epsilon\right)^\top \hat{\Delta} + \lambda_n(\|\theta_S^*\|_1 - \|\theta_S^* + \hat{\Delta}_S\|_1 - \|\hat{\Delta}_{S^c}\|_1) \\ &\leq \left\| \frac{1}{n} X^\top \epsilon \right\|_\infty \|\hat{\Delta}\|_1 + \lambda_n(\|\hat{\Delta}_S\|_1 - \|\hat{\Delta}_{S^c}\|_1) \\ &\leq \lambda_n \left(\frac{3}{2} \|\hat{\Delta}_S\|_1 - \frac{1}{2} \|\hat{\Delta}_{S^c}\|_1 \right) \end{aligned}$$

Hence we have that $\|\hat{\Delta}_{S^c}\|_1 \leq 3\|\hat{\Delta}_S\|_1$, which means $\hat{\Delta} \in C_3(S)$.

By using fact that X satisfies restricted eigenvalue condition with parameter $(\kappa, 3)$, we have that $\frac{1}{n} \|X\hat{\Delta}\|_2^2 \geq \kappa \|\hat{\Delta}\|_2^2$, which further implies

$$\kappa \|\hat{\Delta}\|_2^2 \leq 3\lambda_n \|\hat{\Delta}_S\|_1 \leq 3\lambda_n \sqrt{s} \|\hat{\Delta}\|_2 \quad (7.31)$$

Hence we obtain that

$$\|\hat{\theta} - \theta^*\|_2 \leq \frac{3}{\kappa} \sqrt{s} \lambda_n \quad (7.32)$$

This completes the proof.