

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.

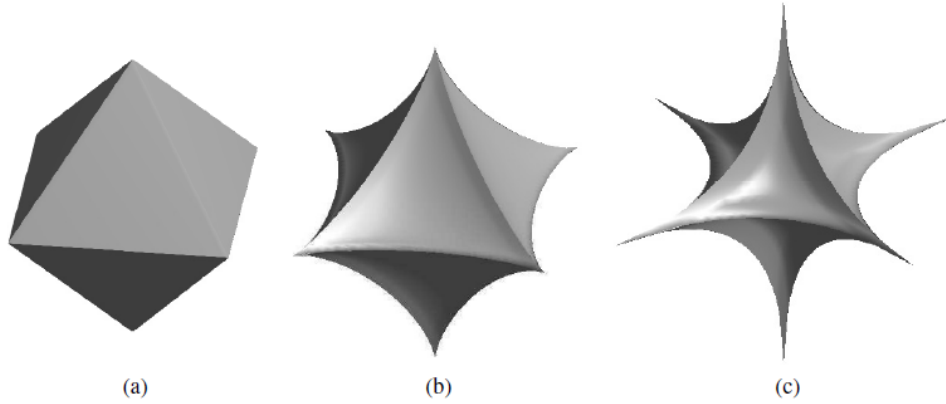


Figure 7.1: Illustrations of the l_q balls for different choices of parameter $q \in (0, 1]$. (a) $q = 1$, (b) $q = 0.75$ and (c) $q = 0.5$.

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).