

High-dimensional Statistical Methods

Skoltech

Q. Paris*

National Research University HSE
Faculty of Computer Science
Moscow, Russia

Chapter 1. High-dimensional regression

Lecture 1

Least squares and constrained least squares

1	Least squares	2
1.1	Basics	2
1.2	Performance bound	3
2	Constrained least-squares	4
2.1	ℓ_0 -constrained least-squares	4
2.2	ℓ_1 -constrained least-squares	5

*email: qparis@hse.ru, teaching material: <http://www.qparis-math.com/teaching>.

The material discussed below is inspired from [1] and [2]. The proofs of the results gathered below are provided during the lectures.

1 Least squares

1.1 Basics

Let us first recall the general problem considered as well as some notation from the introduction. Let $n \geq 1$ and z_1, \dots, z_n be deterministic input points, known (or given) to the statistician, in some input space \mathcal{Z} . To each of the z_i 's, corresponds an observation $Y_i \in \mathbf{R}$ of the form

$$Y_i = f^*(z_i) + \xi_i, \quad (1.1)$$

where ξ_1, \dots, ξ_n are real-valued and independent random variables. Here, $f^* : \mathcal{Z} \rightarrow \mathbf{R}$ denotes an unknown function and the goal is, based only on the observations Y_1, \dots, Y_n , to recover the true vector

$$\mu^* = \begin{bmatrix} \mu_1^* \\ \vdots \\ \mu_n^* \end{bmatrix} \in \mathbf{R}^n \quad \text{where} \quad \mu_i^* = f^*(z_i).$$

Let $\{\varphi_1, \dots, \varphi_p\}$ be a collection of known functions $\varphi_j : \mathcal{Z} \rightarrow \mathbf{R}$ (referred to as the dictionary). Suppose that the unknown function f^* can be expanded in the dictionary in the sense that, for some unknown coefficients $\beta_1^*, \dots, \beta_p^* \in \mathbf{R}$, we have

$$\forall i \in \{1, \dots, n\} : \quad f^*(z_i) = \sum_{j=1}^p \beta_j^* \varphi_j(z_i). \quad (1.2)$$

With the notation

$$\mathbf{x}_i = \begin{bmatrix} \varphi_1(z_i) \\ \vdots \\ \varphi_p(z_i) \end{bmatrix} \in \mathbf{R}^p \quad \text{and} \quad \boldsymbol{\beta}^* = \begin{bmatrix} \beta_1^* \\ \vdots \\ \beta_p^* \end{bmatrix} \in \mathbf{R}^p,$$

equation (1.1) therefore becomes $Y_i = \mathbf{x}_i^\top \boldsymbol{\beta}^* + \xi_i$. In matrix form, we obtain

$$\mathbf{Y} = \mathbf{X} \boldsymbol{\beta}^* + \boldsymbol{\xi},$$

where we recall that

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \in \mathbf{R}^n, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix} \in M_{n,p}(\mathbf{R}), \quad \text{and} \quad \boldsymbol{\xi} = \begin{bmatrix} \xi_1 \\ \vdots \\ \xi_n \end{bmatrix} \in \mathbf{R}^n.$$

In this context, the least squares estimator $\hat{\mu}^{\text{ls}}$ of the unknown vector $\mu^* \in \mathbf{R}^n$ is defined by

$$\hat{\mu}^{\text{ls}} = \mathbf{X}\hat{\beta}^{\text{ls}},$$

where

$$\hat{\beta}^{\text{ls}} \in \arg \min_{\beta \in \mathbf{R}^p} \mathcal{C}^{\text{ls}}(\beta) \quad \text{and} \quad \mathcal{C}^{\text{ls}}(\beta) = \frac{1}{n} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2. \quad (1.3)$$

We now review some basic facts about this estimation technique that should be familiar to the reader.

Theorem 1.1. *The following statements hold.*

- (1) *The function $\beta \mapsto \mathcal{C}^{\text{ls}}(\beta)$ is convex and $\nabla \mathcal{C}^{\text{ls}}(\beta) = 2\mathbf{X}^\top(\mathbf{X}\beta - \mathbf{Y})/n$.*
- (2) *The properties of convex functions guarantee that*

$$\begin{aligned} \hat{\beta} \in \arg \min_{\beta \in \mathbf{R}^p} \mathcal{C}^{\text{ls}}(\beta) &\Leftrightarrow \nabla \mathcal{C}^{\text{ls}}(\hat{\beta}) = 0 \\ &\Leftrightarrow \mathbf{X}^\top \mathbf{X} \hat{\beta} = \mathbf{X}^\top \mathbf{Y}. \end{aligned}$$

- (3) *If $\text{rk}(\mathbf{X}) = p$, then $\mathbf{X}^\top \mathbf{X} \in \mathbf{M}_p(\mathbf{R})$ is invertible and $\hat{\beta}^{\text{ls}}$ is uniquely defined by*

$$\hat{\beta}^{\text{ls}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}.$$

- (4) *If $\text{rk}(\mathbf{X}) < p$, then a solution (not unique) of (1.3) is defined by*

$$\hat{\beta}^{\text{ls}} = (\mathbf{X}^\top \mathbf{X})^+ \mathbf{X}^\top \mathbf{Y},$$

where, for any matrix A , we denote A^+ its pseudo inverse^a.

^aThe Moore-Penrose pseudo inverse of a matrix generalizes the notion of inverse for singular matrices. For any $A \in \mathbf{M}_{p,q}(\mathbf{R})$ its pseudo inverse A^+ is a matrix in $\mathbf{M}_{q,p}(\mathbf{R})$ such that $AA^+x = x$, $\forall x \in \text{Im}(A)$, and such that $A^+Ay = y$, $\forall y \in \text{Im}(A^\top)$. In particular $A^+ = A^{-1}$ when A is a square and invertible matrix.

1.2 Performance bound

In this paragraph, we study the performance of the estimator $\hat{\mu}^{\text{ls}}$ for any values of $n \geq 1$, the sample size, and $p \geq 1$, the dimension of the design points. In particular, the following results hold in the high-dimensional context, *i.e.* when p is much larger than $n \geq 1$.

Theorem 1.2. *Let $r = \text{rk}(\mathbf{X})$. Suppose that the noise vector $\boldsymbol{\xi} \in \mathbf{R}^n$ is sub-gaussian with variance proxy $\sigma^2 > 0$. Then the following statements hold.*

(1) *For all $n \geq 1$ and all $\delta \in (0, 1)$,*

$$\mathcal{E}(\hat{\mu}^{\text{ls}}) \leq \frac{16\sigma^2}{n} \left\{ 2r + \log \left(\frac{1}{\delta} \right) \right\},$$

with probability at least $1 - \delta$.

(2) *For all $n \geq 1$,*

$$\mathbf{E}\mathcal{E}(\hat{\mu}^{\text{ls}}) \leq 8e^{\frac{2}{e}} \frac{\sigma^2 r}{n}.$$

Note that $c = 8e^{\frac{2}{e}} \leq 17$.

The previous theorem, valid for any values of $n \geq 1$, $p \geq 1$ and $r \geq 1$, highlights an important drawback of the least squares method in the high-dimensional context. Indeed, if $p > n$ and the design matrix \mathbf{X} is of full rank n , then the previous upper bounds are of order σ^2 , which may be quite large.

2 Constrained least-squares

This section investigates the favorable case where some preliminary, or apriori, information on the unknown $\boldsymbol{\beta}^*$ can be implemented in the statistical procedure in the form of an explicit constraint, *i.e.* when one considers

$$\hat{\mu}_{\mathcal{K}}^{\text{ls}} = \mathbf{X}\hat{\boldsymbol{\beta}}_{\mathcal{K}}^{\text{ls}} \quad \text{where} \quad \hat{\boldsymbol{\beta}}_{\mathcal{K}}^{\text{ls}} \in \arg \min_{\boldsymbol{\beta} \in \mathcal{K}} \mathcal{C}^{\text{ls}}(\boldsymbol{\beta}),$$

for some explicit $\mathcal{K} \subset \mathbf{R}^p$. Below, we focus on two specific examples for the constraint \mathcal{K} , of specific importance in the high-dimensional setting.

2.1 ℓ_0 -constrained least-squares

This paragraph studies the case where the unknown $\boldsymbol{\beta}^*$ is apriori known to be s -sparse. Recall that $\boldsymbol{\beta} \in \mathbf{R}^p$ is said s -sparse if $\|\boldsymbol{\beta}\|_0 \leq s$ where

$$\|\boldsymbol{\beta}\|_0 = \sum_{j=1}^p \mathbf{1}\{\beta_j \neq 0\}.$$

This information corresponds obviously to the constraint

$$\mathcal{K} = \{\boldsymbol{\beta} \in \mathbf{R}^p : \|\boldsymbol{\beta}\|_0 \leq s\}.$$

It should be noted that, for this constraint, the computation of $\hat{\boldsymbol{\beta}}_{\mathcal{K}}^{\text{ls}}$ is unrealistic in practice for large values of p . Indeed, it requires to minimize $\binom{p}{s}$ least-squares criterions in dimension s . However, the results presented below stand as a benchmark for the next paragraph.

Theorem 2.1. *Let $s \geq 1$ be a integer smaller than $p/2$. Suppose that $\boldsymbol{\beta}^* \in \mathcal{K}$ and the noise vector $\boldsymbol{\xi} \in \mathbf{R}^n$ is sub-gaussian with variance proxy $\sigma^2 > 0$. Then the following statements hold.*

(1) *For all $n \geq 1$ and all $\delta \in (0, 1)$,*

$$\mathcal{E}(\hat{\mu}_{\mathcal{K}}^{\text{ls}}) \leq \frac{8\sigma^2}{n} \left\{ 2s \log \left(\frac{3ep}{s} \right) + \log \left(\frac{1}{\delta} \right) \right\},$$

with probability at least $1 - \delta$.

(2) *For all $n \geq 1$,*

$$\mathbf{E}\mathcal{E}(\hat{\mu}_{\mathcal{K}}^{\text{ls}}) \leq \frac{8\sigma^2}{n} \left\{ 1 + 2s \log \left(\frac{3ep}{s} \right) \right\}.$$

Contrary to the global LS estimator, the ℓ_0 -constrained LS exhibits remarkable properties. First, its (theoretical) performance is not affected by the rank of the design matrix \mathbf{X} and depends on the dimension p only through a log term. As mentioned above, the computation of this estimator is however computationally unrealistic.

2.2 ℓ_1 -constrained least-squares

This paragraph studies a computationally friendly alternative to the ℓ_0 -constrained LS. If the unknown $\boldsymbol{\beta}^*$ is supposed sparse and bounded, the inequality

$$\|\boldsymbol{\beta}\|_1 \leq \min \left\{ \|\boldsymbol{\beta}\|_2 \sqrt{\|\boldsymbol{\beta}\|_0}, \|\boldsymbol{\beta}\|_\infty \|\boldsymbol{\beta}\|_0 \right\},$$

suggests that $\boldsymbol{\beta}^*$ lies in an ℓ_1 -ball of small radius. Next, consider, therefore the constraint

$$\mathcal{K} = \{\boldsymbol{\beta} \in \mathbf{R}^p : \|\boldsymbol{\beta}\|_1 \leq \lambda\},$$

for some $\lambda > 0$. Contrary to the previous paragraph, the ℓ_1 -constrained LS can be computed very efficiently (see below). The next result, provides upper bounds on its theoretical performance.

Theorem 2.2. Suppose that $\beta^* \in \mathcal{K}$ and the noise vector $\xi \in \mathbf{R}^n$ is sub-gaussian with variance proxy $\sigma^2 > 0$. Denote

$$\varkappa := \sup_{1 \leq j \leq p} \|\mathbf{x}^j\|_2,$$

where $\mathbf{x}^1, \dots, \mathbf{x}^p \in \mathbf{R}^n$, denote the columns of the design matrix \mathbf{X} . Then the following statements hold.

(1) For all $n \geq 1$ and all $\delta \in (0, 1)$,

$$\mathcal{E}(\hat{\mu}_{\mathcal{K}}^{\text{ls}}) \leq \frac{4\sigma\varkappa}{n} \sqrt{2 \log \left(\frac{2p}{\delta} \right)},$$

with probability at least $1 - \delta$.

(2) For all $n \geq 1$,

$$\mathbf{E} \mathcal{E}(\hat{\mu}_{\mathcal{K}}^{\text{ls}}) \leq \frac{4\sigma\varkappa \sqrt{2 \log(2p)}}{n}.$$

In practice, normalize the design matrix such that $\varkappa \leq \sqrt{n}$.