

High-dimensional Statistical Methods

Skoltech

Q. Paris*

National Research University HSE
Faculty of Computer Science
Moscow, Russia

Introduction

1	Fixed design regression	3
2	Mean squared error	3
3	Expansion in a dictionary	3
4	The LS criterion and its variants	5
4.1	Classical LS	5
4.2	Constrained LS	6
4.3	Penalized LS	6
4.4	Sparsity, BIC and the Lasso	7
5	Notation	8

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

Foreword

These notes aim to introduce, in a simple framework, the foundations of some modern statistical methods in the context of high-dimensional data. The topics presented below are in no way an exhaustive overview of the subject and are biased by the personal choices of the author. The material gathered here relies on important reference books or lecture notes to which we refer the reader for additional material and details. As the main references for this course, we recommend the lecture notes by [Rigollet \(2015\)](#) or the book by [Giraud \(2015\)](#). The reader is also invited to use the additional and very valuable references [van Handel \(2016\)](#) or [Vershynin \(2016\)](#). Related material may also be found in [Buhlmann and van de Geer \(2011\)](#); [Koltchinskii \(2011\)](#); [Hastie et al. \(2015\)](#) and the references therein. Below is a tentative program for the course, that may be modified along the way.

1 High-dimensional regression

- 1.1 Least squares and constrained least squares
- 1.2 Bayes information criterion
- 1.3 Lasso
- 1.4 Additional topics

2 Matrix estimation

- 2.1 Matrix regression model
- 2.2 Thresholding estimator
- 2.3 Penalized approaches
- 2.4 Additional topics

3 Graphical models

- 3.1 Directed and non-directed acyclic graphs
- 3.2 Gaussian graphical models
- 3.3 Estimation of the precision matrix
- 3.4 Additional topics

A Appendices

- A.1 Sub-gaussian random variables
- A.2 Matrices and norms
- A.3 Convexity and optimization
- A.4 Algorithms

1 Fixed design regression

Let $n \geq 1$ and z_1, \dots, z_n be known and deterministic points (called design or input points) in some input space \mathcal{Z} fixed by the statistician or practitioner. Suppose that, to each of the z_i 's, corresponds an observation (or measurement) $Y_i \in \mathbf{R}$ of the form

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

where ξ_1, \dots, ξ_n denote real-valued, centered and independent random variables and f^* denotes an unknown function $f^* : \mathcal{Z} \rightarrow \mathbf{R}$. In this setting, Y_i stands for some physical measurement $f^*(z_i)$, relative to z_i , corrupted by some noise ξ_i . From a statistical point of view, the goal is to estimate (or recover) the true vector

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

based on the only knowledge of the observations Y_1, \dots, Y_n .

2 Mean squared error

Given an estimator $\hat{\mu} \in \mathbf{R}^n$ of μ^* based on the observations Y_1, \dots, Y_n , a natural measure of its performance is the mean squared error

$$\mathcal{E}(\hat{\mu}) = \frac{1}{n} \sum_{i=1}^n (\hat{\mu}_i - \mu_i^*)^2 = \frac{1}{n} \|\hat{\mu} - \mu^*\|_2^2,$$

where $\|\cdot\|_2$ denotes the Euclidean norm. We stress that, of course, since $\hat{\mu}$ is constructed from the random Y_i 's, the mean squared error $\mathcal{E}(\hat{\mu})$ is itself a (positive) random quantity.

3 Expansion in a dictionary

As will be made clear below, the nature of the input space \mathcal{Z} plays no role in the analysis insofar the input points z_1, \dots, z_n are considered as known or fixed by the statistician. In particular, \mathcal{Z} could be infinite-dimensional. For instance, the input points could be curves, images, signals, etc. The high-dimensionality of the problem, that will play a fundamental role in the sequel, will be related to a reformulation of the initial problem described below.

Suppose given a collection, also referred to as dictionary, $\varphi_1, \dots, \varphi_p$ of p known functions $\varphi_j : \mathcal{Z} \rightarrow \mathbf{R}$ and suppose for simplicity that the unknown function f^*

can be expanded in the dictionary in the following sense: for some unknown coefficients $\beta_1^*, \dots, \beta_p^* \in \mathbf{R}$, we have

$$\forall i \in \{1, \dots, n\} : \quad f^*(z_i) = \sum_{j=1}^p \beta_j^* \varphi_j(z_i). \quad (3.1)$$

Then, introducing the notation

$$\mathbf{x}_i = \begin{bmatrix} \varphi_1(z_i) \\ \vdots \\ \varphi_p(z_i) \end{bmatrix} \in \mathbf{R}^p \quad \text{and} \quad \boldsymbol{\beta}^* = \begin{bmatrix} \beta_1^* \\ \vdots \\ \beta_p^* \end{bmatrix} \in \mathbf{R}^p, \quad (3.2)$$

the reader may easily check that equation (1.1) becomes

$$Y_i = \mathbf{x}_i^\top \boldsymbol{\beta}^* + \xi_i. \quad (3.3)$$

The points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbf{R}^p$ will also be called the design points (by extension) and the unknown vector $\boldsymbol{\beta}^* \in \mathbf{R}^p$ the regression vector. We insist on the fact that the input points z_1, \dots, z_n and the dictionary $\varphi_1, \dots, \varphi_p$ being known and deterministic, the design points are quantities within the knowledge of the statistician during the analysis. The functions in the dictionary can be considered, typically, as some basis functions in a given functional space such as Splines, Fourier or Wavelet basis functions. Other examples of dictionary functions can be tailored for specific applications.

Example 3.1 (Gene expression level). Suppose $\mathcal{Z} = \mathbf{R}$ and f^* is a piece-wise constant function $f^* : \mathcal{Z} \rightarrow \mathbf{R}$. Suppose in addition that the positions $t_1 < \dots < t_p \in \mathbf{R}$ of potential jumps of f^* are known but that the value of f^* on each $[t_j, t_{j+1})$ is not. This situation arises, for instance, in biology (or genetics) when analysing the transcription level¹ of genes in a DNA sequence. In this context, the genes in the given sequence (represented by the t_j 's here) are known but their transcription level is to be discovered. Here, observe that the function f^* may be represented as

$$\forall z \in \mathbf{R} : \quad f^*(z) = \sum_{j=1}^p \beta_j^* \mathbf{1}\{z \geq t_j\},$$

for known dictionary functions $\varphi_j(z) = \mathbf{1}\{z \geq t_j\}$ and unknown coefficients β_j^* . Note therefore that identity (3.1) holds here for any choice of the input points $z_1, \dots, z_n \in \mathbf{R}$.

The interest of the previous reduction is obviously that, provided assumption (3.1) holds, the introduction of the dictionary $\varphi_1, \dots, \varphi_p$ reduces the statistical problematic of finding an unknown function f^* (of potentially infinite-dimensional nature), to that of estimating a p -dimensional vector $\boldsymbol{\beta}^*$. In particular, it should be expected that the estimation of μ^* , in this context, is

¹The transcription level of a gene in a cell at a given time is the quantity of ARNm associated to this gene present at this time in the cell

highly simplified. However, for a decomposition such as (3.1) to be realistic, the number p (corresponding to the dimension of both β^* and the \mathbf{x}_i 's) of basis functions is expected to be very large and, in particular, potentially much larger than the sample size n . This context conflicts severely with the scenario of classical statistics in which p is of reasonable magnitude and the sample size n is supposed very large. Here, statistical guarantees are required for fixed n and p with, possibly, $p \gg n$.

Another situation studied in the course corresponds to the case where the considered dictionary $\varphi_1, \dots, \varphi_p$ is miss-specified, *i.e.*

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

This situation arises naturally when, given too little information on the unknown function f^* , the statistician chooses a dictionary $\varphi_1, \dots, \varphi_p$ (*i.e.* formulates a model for the unknown f^*) in which one cannot expand the unknown function f^* as in (3.1). In this context, the goal of the statistician is to guarantee a form of robustness of the statistical methodology considered: given our fixed dictionary, our estimators should do almost as good as the best possible function \bar{f} in the model, *i.e.*

$$\bar{f}(z) = \sum_{j=1}^p \bar{\beta}_j \varphi_j(z),$$

where vector $\bar{\beta} = (\bar{\beta}_1, \dots, \bar{\beta}_p)^\top \in \mathbf{R}^p$ is defined as any

$$\bar{\beta} \in \arg \min_{\beta \in \mathbf{R}^p} \frac{1}{n} \sum_{i=1}^n (f^*(z_i) - \mathbf{x}_i^\top \beta)^2,$$

and \mathbf{x}_i is as in (3.2). From a technical point of view, the analysis will be carried out by deriving so called oracle inequalities, quantifying precisely the aforementioned robustness of the considered estimation techniques.

4 The LS criterion and its variants

4.1 Classical LS

In the context of representation (3.3), the reader should be familiar with the idea that, given the data points $(Y_1, \mathbf{x}_1), \dots, (Y_n, \mathbf{x}_n) \in \mathbf{R} \times \mathbf{R}^p$, a natural method to estimate β^* , and therefore μ^* , is to minimize the least-squares criterion, *i.e.* consider

$$\hat{\beta}^{\text{ls}} \in \arg \min_{\beta \in \mathbf{R}^p} \mathcal{C}^{\text{ls}}(\beta) \quad \text{where} \quad \mathcal{C}^{\text{ls}}(\beta) = \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \beta)^2, \quad (4.1)$$

and form the estimator $\hat{\mu}^{\text{ls}}$ of μ^\star with coordinates

$$\hat{\mu}_i^{\text{ls}} = \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}^{\text{ls}}.$$

As will be described in the first chapter, the least squares estimator performs quite badly in most of the cases, when $p \gg n$. In particular, we will prove that, up to universal constants, the expected mean squared error of the least squares estimator satisfies

$$\mathbf{E} \mathcal{E}(\hat{\mu}^{\text{ls}}) \approx \frac{\sigma^2 r}{n},$$

(where statement ' \approx ' will be precised) if the noise variables ξ_1, \dots, ξ_n are independent, centered and sub-gaussian variables with variance proxy σ^2 (see the appendix devoted to this notion) and if r denoted the dimension of the subspace spanned by the design points. This rate of convergence is especially bad if the design is of full rank $r = \min\{n, p\}$. Hence, new estimation techniques are called for to deal efficiently with the high-dimensional nature of the problem.

The bad performance of the least-squares estimator, even though proved formally in the course, can be understood intuitively as follows. The (convex) minimization problem (4.1) encodes absolutely no restriction on the potential candidates $\boldsymbol{\beta} \in \mathbf{R}^p$ leading to an optimal value of the criterion $\mathcal{C}^{\text{ls}}(\boldsymbol{\beta})$. To speak very roughly, the LS estimator results, in a sense, from an 'agnostic' search (no restriction on the candidates $\boldsymbol{\beta}$) of a pin (the vector $\boldsymbol{\beta}^\star$) in a haystack (the space \mathbf{R}^p) with only one eye (the sample size n is small compared to p).

4.2 Constrained LS

In some applications, the practitioner (or the expert) may have some intuition or knowledge on the location, or on some characteristics (usually of geometrical nature), of the unknown $\boldsymbol{\beta}^\star$. In this context, one is naturally brought to consider a constrained LS problem of the form

$$\hat{\boldsymbol{\beta}}_{\mathcal{K}}^{\text{ls}} \in \arg \min_{\boldsymbol{\beta} \in \mathcal{K}} \mathcal{C}^{\text{ls}}(\boldsymbol{\beta}), \quad (4.2)$$

where \mathcal{K} denotes a strict subset of \mathbf{R}^p . From a mathematical point of view, the analysis of such an estimator can be carried out (at least) for any totally bounded set \mathcal{K} . For practical and optimization purposes, one is however led to favor convex sets and a number of our lectures will be devoted to this problem.

4.3 Penalized LS

The framework of constrained LS can be, in some instances, considered as too restrictive since one has to specify exactly a subset \mathcal{K} on which to perform the minimization. A more general approach is to penalize the LS criterion.

Precisely, if the intuition on β^* can be formalized by saying that, for some given function $\Omega : \mathbf{R}^p \rightarrow [0, +\infty)$, the quantity $\Omega(\beta^*)$ is likely to be small, a natural alternative to constrained LS is to minimize a penalized version of the least-squares criterion

$$\mathcal{C}^{\text{pls}}(\beta) = \mathcal{C}^{\text{ls}}(\beta) + \lambda \Omega(\beta), \quad (4.3)$$

for some parameter $\lambda > 0$ (possibly depending on the sample size n) to be fixed by the statistician. A natural correspondence between the penalized and constrained approach, given by Lagrangian duality, will be established along the way. As for constrained LS, the statistician should (for practical purposes) favor a convex penalty function Ω as the resulting criterion $\beta \mapsto \mathcal{C}^{\text{pls}}(\beta)$, also convex, can therefore be optimized using the powerful tools from convex (and stochastic) optimization. This observation will be crucial in the subsequent analysis. On the subject of convex optimization and its recent developments, we refer the reader to [Bubeck \(2015\)](#).

4.4 Sparsity, BIC and the Lasso

In the high-dimensional framework, a possible geometric characteristic of the high-dimensional vector β^* , known to be of crucial impact on the performance of statistical methods, is sparsity. The vector β^* is said to be sparse if only a few of its components are non-zero. Formally, β^* is said to be k -sparse if its ℓ_0 -norm²

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

is smaller than k . When it is supposed that the unknown β^* is sparse, then the choice of the penalty term $\Omega(\beta) = \|\beta\|_0$ in (4.3) leads to the BIC³ estimator,

$$\hat{\beta}^{\text{bic}} \in \arg \min_{\beta \in \mathbf{R}^p} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \beta)^2 + \lambda^2 \|\beta\|_0 \right\}. \quad (4.4)$$

As we will see, the theoretical performance of the BIC estimator is remarkable in that it adapts to the unknown sparsity of β^* and, contrary to the classical least squares, is much less affected from the dimensionality p of the problem. Unfortunately, from a numerical point of view, computing the BIC estimator is usually unrealistic and no known computational method is known to be significantly better than an exhaustive search among the 2^p possible sparsity patterns of a p -dimensional vector.

²This quantity is abusively called a norm by convention.

³BIC stands for Bayes Information Criterion

A important alternative to the BIC is the Lasso⁴. The Lasso estimator,

$$\hat{\boldsymbol{\beta}}^{\text{lasso}} \in \arg \min_{\boldsymbol{\beta} \in \mathbf{R}^p} \left\{ \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 + 2\lambda \|\boldsymbol{\beta}\|_1 \right\}, \quad (4.5)$$

may be seen as a convex relaxation of the BIC. Both computationally realistic and performant, the Lasso is arguably one of the most popular method in high-dimensional statistics nowadays and a significant part of the lectures will be devoted to its study.

5 Notation

We list below the most important notation used along the course. For any integers $p, q \geq 1$, we denote $M_{p,q}(\mathbf{R})$ the set of $p \times q$ matrices with real coefficients and set $M_p(\mathbf{R}) = M_{p,p}(\mathbf{R})$. We define the vector and matrix notation

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \in \mathbf{R}^n, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix} \in M_{n,p}(\mathbf{R}), \quad \text{and} \quad \boldsymbol{\xi} = \begin{bmatrix} \xi_1 \\ \vdots \\ \xi_n \end{bmatrix} \in \mathbf{R}^n. \quad (5.1)$$

In matrix notation, the least-squares criterion therefore reads

$$\mathcal{C}^{\text{ls}}(\boldsymbol{\beta}) = \frac{1}{n} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2.$$

The columns of the design matrix \mathbf{X} will be denoted $\mathbf{x}^1, \dots, \mathbf{x}^p \in \mathbf{R}^n$ with index in exponent. We introduce

$$\hat{\Sigma} = \frac{1}{n} \mathbf{X}^\top \mathbf{X} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \in M_p(\mathbf{R}).$$

Note in particular that, for any $k, \ell \in \{1, \dots, p\}$, the entry (k, ℓ) of $\hat{\Sigma}$ is

$$\hat{\Sigma}_{k,\ell} = \frac{(\mathbf{x}^k)^\top \mathbf{x}^\ell}{n}.$$

In other words, the matrix $\hat{\Sigma}$ corresponds to the Gram matrix of the column vectors $\mathbf{x}^1, \dots, \mathbf{x}^p \in \mathbf{R}^n$ up to a scaling factor $1/n$. Finally, for any $k \geq 1$ and any vector $v = (v_1, \dots, v_k) \in \mathbf{R}^k$, we will use the classical notation $\|v\|_q$, $1 \leq q \leq +\infty$ to denote the norms

$$\|v\|_q = \left(\sum_{\ell=1}^k |v_\ell|^q \right)^{1/q}, \text{ if } 1 \leq q < +\infty, \quad \text{and} \quad \|v\|_\infty = \max_{1 \leq \ell \leq k} |v_\ell|.$$

⁴Lasso stands for Least Absolute Shrinkage and Selection Operator.

References

- S. Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends in Machine Learning*, 8(3-4):231–357, 2015.
- P. Buhlmann and S. van de Geer. *Statistics for High-Dimensional Data*. Springer, 2011.
- C. Giraud. *Introduction to High-dimensional Statistics*. CRC Press, 2015.
- T. Hastie, R. Tibshirani, and M. Wainwright. *Statistical Learning With Sparsity: The Lasso and Generalizations*. CRC Press, 2015.
- V. Koltchinskii. Oracle inequalities in empirical risk minimization and sparse recovery problems. In *Lectures from the 38th Summer School on Probability Theory held in Saint-Flour in July, 2008*. Springer, 2011.
- P. Rigollet. High-dimensional statistics. MIT lecture notes, 2015.
- R. van Handel. Probability in high dimension. Princeton lecture notes, 2016.
- R. Vershynin. High-dimensional probability: An introduction with applications in data science. Book in preparation, 2016.