

High-Dimensional Statistics

I. Section Introduction

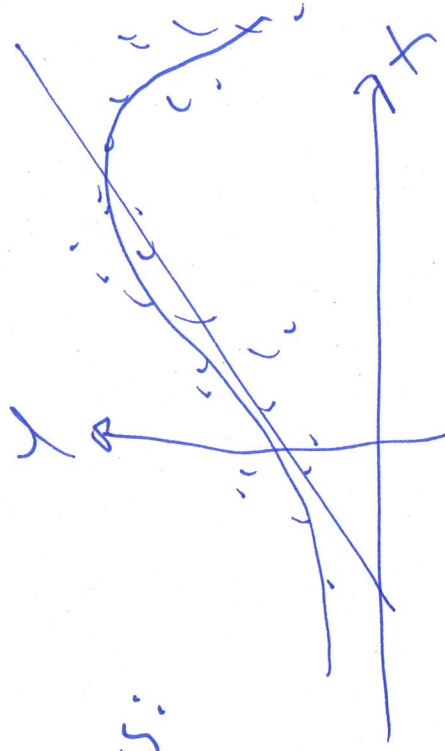
high-dimensional: more parameters (in a model) than sample size
no. parameters: p
 n

$$p \gg n$$

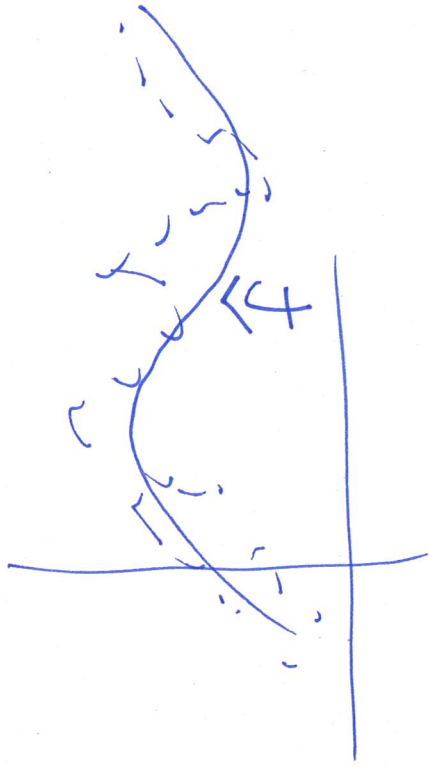
Very common in nonparametric statistics:

e.g. $Y = f(X) + \epsilon$

nonparametric smooth function $\mathbb{R} \rightarrow \mathbb{R}$: ∞ -dim. object



for smoothing spline method:



\hat{f} parameterized (with splines) with $p \approx n$ parameters
"order"

and smoothness of f is crucial to deal with
this problem

new view: sparsity (instead of smoothness)

p parameters but many of them are zero
(but unknown which of them are zero)

prime example:

$$y = X\beta + \epsilon$$

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} X_1^{(1)} & \dots & X_n^{(1)} \\ \vdots & \ddots & \vdots \\ X_1^{(p)} & \dots & X_n^{(p)} \end{pmatrix}$$

$$p \gg n$$

high-dimensional linear model

\uparrow \uparrow
 p different features, covariates

Sparsity: $S = \text{supp}(\beta) = \{j; \beta_j \neq 0\}$
active set of Covariates

Cardinality of S : $|S| \ll n$

$$y = X \beta + \varepsilon = \underset{n \times p \times 1}{X} \underset{n \times |S|}{\beta_S} + \varepsilon$$

Practical implication: "~~throw~~ in many potentially useless
Covariates — so ~~that~~ we reduce the chance to miss
an ~~important~~ one"

What is the price of collecting too many
covariates?

(besides the benefit to catch the important ones)

II. The Lasso for linear models (Ch. 2 in

Bühlmann & van de Geer)

II. 1. Introduction

data are realizations of

$$(X_1, Y_1), \dots, (X_n, Y_n)$$

fixed or random

Random

$$X_i \in \mathcal{X} \subseteq \mathbb{R}^p; \quad Y_i \in \mathcal{Y} \subseteq \mathbb{R}$$

useful linear model:

$$y_i = \sum_{j=1}^p \beta_j X_i^{(j)} + \varepsilon_i \quad (i=1, \dots, n)$$

$$X_i \rightarrow X_i^{(1)} \rightarrow X_i^{(p)T}, \quad X^{(j)} = (X_1^{(j)} \rightarrow X_n^{(j)})^T$$

$$\varepsilon_1, \dots, \varepsilon_n \text{ i.i.d., } E[\varepsilon_i] = 0, \quad \text{Var}(\varepsilon_i) = \sigma^2$$

independent of X_1, \dots, X_n

often in practice: centered and scaled

$$y_i \leftarrow y_i - \bar{y}; \quad X_i^{(j)} \leftarrow \frac{X_i^{(j)} - \bar{X}^{(j)}}{\frac{\sigma^{(j)}}{\sigma}}$$

$$\left(\frac{\sigma^{(j)}}{\sigma}\right)^2 = \frac{1}{n} \sum_{i=1}^n (X_i^{(j)} - \bar{X}^{(j)})^2$$

notation:
$$Y = X\beta + \varepsilon$$

$$n \times 1 \quad n \times p \quad p \times 1 \quad n \times 1$$

(if model is assumed correct: "true" parameters is β^0)

II. 2. The Lasso estimator

for $p > n$: OLS is not unique
will overfit the data

→ complexity regularization

$$\hat{\beta}(A) = \underset{\beta}{\operatorname{argmin}} \left(\|Y - X\beta\|_2^2 / n + A \|\beta\|_1 \right) \quad \underline{L_{\text{Lasso}}}$$

$A > 0$ regularization parameter

makes only sense if ~~$\hat{\beta}^{(i)}$~~ are (roughly) equal

different from Ridge Regression:

$$\hat{\beta}(\lambda) = \arg\min_{\beta} (\|y - X\beta\|_2^2 / n + \lambda \|\beta\|_2^2)$$

Fit properties of Lasso:

(1) sparse estimator

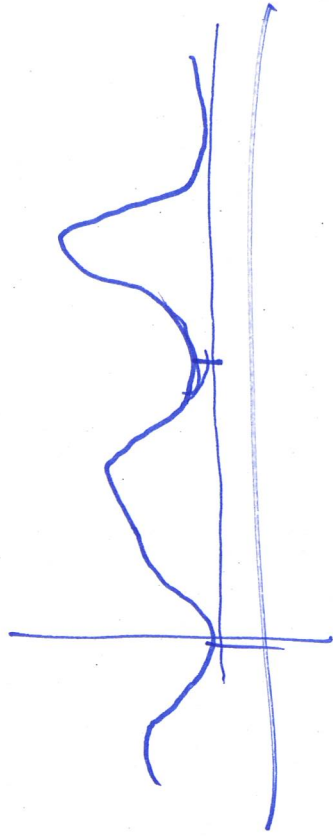
$\hat{\beta}_j(\lambda) = 0$ depending on j and λ

"Lasso is doing variable selection"

in contrast to Ridge regression

(2) Lasso involves convex optimization:
every local minimum is a global minimum

there are typically many of them
all of them being global



but under further assumptions: uniqueness

Explanation for (1):

$$\hat{\beta}(Y) = \arg \min_{\beta} \left(\|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right) \quad \text{Lagrangian form}$$

primal problem:

$$\hat{\beta}_{\text{primal}}(R) = \arg \min_{\beta} \underbrace{\|Y - X\beta\|_2^2 + \lambda \|\beta\|_1}_{\text{convex constraints}}$$

The two optimization problems are equivalent with a one-to-one correspondence between λ and R (depending on $(X, Y, \lambda)_{\lambda=1}^{\lambda=0}$)

for Ridge:

$$\hat{\beta}_{\text{ridge}}(R) = \arg \min_{\beta} \|Y - X\beta\|_2^2 + R \|\beta\|_2^2$$

Lawso: Tibshirani (1996)

Least Absolute Selection and Shrinkage Operator
sparsity

II.3 Orthogonal design

X with $n^{-1} X^T X = I_{p \times p}$ ($\rightarrow p \leq n$)

Then: Least

$$\hat{\beta}_j(A) = \text{sign}(z_j) (|z_j| - \lambda/2)_+ \quad j=1, \dots, p$$

$$z_j = (X^T Y)_j / n = \bar{y}_{0LS;j}$$