

# High-Dimensional Statistics

## I. Introduction

high-dimensional (model):  $\underbrace{\text{more parameters than sample size}}_{=n}$

no. par. =  $p$

$p \gg n$

Notation: asymptotic viewpoint

$p = p_n$  ,  $p_n \rightarrow \infty$  ( $n \rightarrow \infty$ ), e.g.  $p_n = n^\alpha$   $\alpha > 0$

$p \gg n \iff \frac{p_n}{n} = \frac{p_n}{n} \rightarrow \infty$  ( $n \rightarrow \infty$ )

$$p = O(a_n) \iff \frac{p}{a_n} = \frac{p_n}{a_n} \in M < \infty \quad \forall n$$

$$p = o(a_n) \iff \frac{p}{a_n} = \frac{p_n}{a_n} \rightarrow 0 \quad (n \rightarrow \infty)$$

$$p \asymp a_n \iff p_n = O(a_n), \quad a_n = O(p_n)$$

$$p \sim a_n \iff p_n/a_n \rightarrow 1 \quad (n \rightarrow \infty)$$

$p \gg n$  is common in nonparametric statistics

example:  $Y_i = f(X_i) + \varepsilon_i$ ,  $X_i \in \mathbb{R}^1$ ,  $\varepsilon_i \in \mathbb{R}^1$   
 $f(\cdot)$  smooth,  $E[\varepsilon_i] = 0$

$f(\cdot)$  is an infinite-dimensional parameter

e.g.  $f(x) = \sum_{j=1}^{\infty} b_j \phi_j(x)$   $\rightarrow$  basis functions:  $\text{par} = b_1, b_2, \dots$

Assuming smoothness of  $f(\cdot)$  is typically required to get a good estimator  $\hat{f}(\cdot)$  for  $f(\cdot)$

for example: smoothing splines uses  $p \leq n$  which need to be estimated

new view: sparsity assumption (in statistics: Donoho, Johnstone  $\approx 1990$ )

that is:  $p$  parameters but many of them are zero (approx. zero)

this is another assumption which guarantees good estimators in the  $p \gg n$  setting

$$Y = X\beta + \varepsilon \quad \text{linear model}$$

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

$$n = 115$$
$$p = 4088$$

→ ordinary least squares  
"does not work"

# I. 1. High-dimensional linear model

a) a prime example

$$Y = X\beta^0 + \varepsilon, \quad X \text{ fixed (or random and conditioning on } X)$$

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

$$E[\varepsilon] = 0$$

$$S^0 = \text{supp}(\beta^0) = \{j; \beta_j^0 \neq 0\}$$

sparsity index  $s_0 = |S^0|$ , assume  $s_0 = |S^0| \ll n$

in many practical applications: assuming sparsity is "reasonable"

example: biomarker discovery

Riboflavin production data

~ "hope" that only a few genes are relevant

for  $Y = \log\text{-production rate}$

$$Y_i = \sum_{j=1}^p \beta_j X_i^{(j)} + \epsilon_i$$

$\uparrow$  expr. of gene  $j$

$\log\text{-prod. rate}$

## II. Lasso for linear models

### II. 1. Introduction

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

in practice: centered and scaled variables

$$y_i \leftarrow y_i - \bar{y}$$

$$X_i^{(j)} \leftarrow \frac{X_i^{(j)} - \bar{X}^{(j)}}{\hat{\sigma}_j}$$

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

## II. 2. The Lasso estimator

most popular estimator for high-dim. lin. model

(Tibshirani, 1996)

for  $p > n$  :  $\text{rank}(X) < p$

→ ordinary least squares is not unique :

if overfits and produces residual sum of squares = 0

$$y \in \mathbb{R}^n = \hat{y}_{OLS}$$

$\text{span}(X^{(2)}, \dots, X^{(n)})$  ;  $\dim = n$



→ need complexity regularization

$$\sum_{j=1}^p |\beta_j|$$

Lasso:

$$\hat{\beta}(\lambda) = \underset{\beta}{\operatorname{argmin}} \left( \|Y - X\beta\|_2^2 / n + \lambda \|\beta\|_1 \right)$$

$\geq 0$  penalty parameter

differs from Ridge regression (Tikhonov regularization)

$$\hat{\beta}_{\text{Ridge}}(\lambda) = \underset{\beta}{\operatorname{argmin}} \left( \|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right)$$

$\geq 0$

## First properties:

(1) Lasso is sparse estimator

$\hat{\beta}_j(\lambda) = 0$  for "many"  $j$ 's (depending on  $\lambda$ )

"Lasso is doing variable selection"

LASSO = Least Absolute Shrinkage and Selection Operator

not true for Ridge estimator

(2) Lasso involves convex optimization: "easy"

every local minimum is a global one

but in general not unique

(under additional conditions  $\rightarrow$  uniqueness, see later)

Explanation for (1): due to convexity

$$\hat{\beta}(\lambda) = \arg \min_{\beta} (\|Y - X\beta\|_2^2 / n + \lambda \|\beta\|_1) \quad \text{"Lagrange form"}$$

$\Leftrightarrow$

primal problem:

$$\hat{\beta}_{\text{primal}}(\lambda) = \arg \min_{\beta: \|\beta\|_2 \leq R} \|Y - X\beta\|_2^2 / n$$

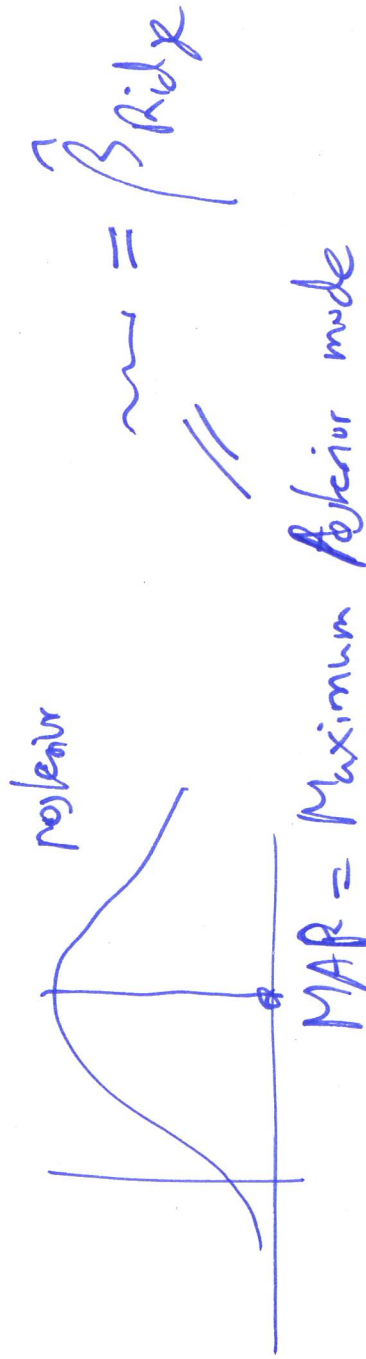
equivalent optimizations with a one-to-one correspondence between  $\lambda$  and  $R$  (not explicit; depends on data  $(Y_i, X_i)_{i=1 \rightarrow n}$ )

software: glmnet in R  
including CV for choosing  $\lambda$

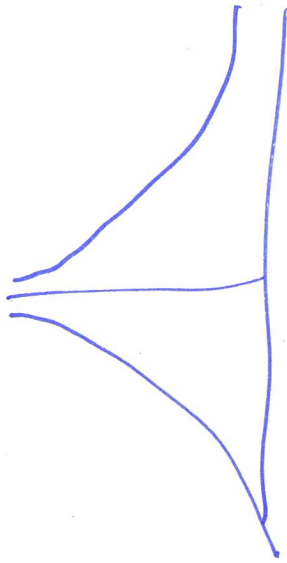
## Bayesian:

$\beta_1 \rightarrow \beta$  i.i.d.  $\sim \mathcal{N}(0, \tau^2)$  prior

$\rightarrow$  posterior accuracy  $\propto \sim$  Gaussian



$\beta_1 \rightarrow \beta_0$  i.i.d.  $\sim$  Double-Exp. Laplace      mean 0      variance  $\tau^2$



$\xi \sim \text{Cauchy}$   
 $\rightarrow \text{MAR} = \text{Lasso}$