

LASSO REGRESSION

MA8701, 03.02.2023

• = added after class

$$Y = X\beta + \epsilon$$

$N \times 1 \quad N \times p \quad p \times 1 \quad N \times 1$

$$E(\epsilon) = 0$$
$$\text{Var}(\epsilon) = \sigma^2 I$$

use centered Y and X to avoid β_0 in the model

$$\text{minimize}_{\beta} \sum_{i=1}^N (y_i - \sum_{j=1}^p x_{ij} \beta_j)^2 \quad \text{subject to} \quad \sum_{j=1}^p |\beta_j| \leq t$$

$$\hat{\beta}_{\text{lasso}} = \underset{\beta}{\text{argmin}} (Y - X\beta)^T (Y - X\beta) + \lambda \sum_{j=1}^p |\beta_j|$$

$$\text{For } t \geq t_0 = \sum_{j=1}^p |\hat{\beta}_{\text{OLS},j}| \Rightarrow \hat{\beta}_{\text{lasso}} = \hat{\beta}_{\text{OLS}}$$

$$t \leq t_{\text{max}} = \max_j X_j^T Y \Rightarrow \hat{\beta}_{\text{lasso}} = 0$$

j \downarrow j th col of X

Group discussion: log PSA and 8 covariates

Discuss what do see

Ridge $X(X^T X + \lambda I)^{-1} X^T$

• X -axis: $df(\lambda) = k(H(\lambda))$
 $[0, 8]$

- $df \rightarrow 0 \Rightarrow \hat{\beta}_{\text{ridge}} = 0$
- $df = 8 \Rightarrow \hat{\beta}_{\text{ridge}} = \hat{\beta}_{\text{OLS}}$
- reentry of coeffs change not monotone

Lasso

• X -axis is $\frac{t}{\sum_{j=1}^p |\hat{\beta}_{\text{OLS},j}|}$ $[0, 1]$

- $t=1: \hat{\beta}_{\text{OLS}} = \hat{\beta}_{\text{lasso}}$
- when $\hat{\beta}$ shrink to 0, do not jump back to value > 0 .

- linear between $\ln w_k$ = change in active set = not-zero coeffs
- the nature of shrinkage is complex

PARAMETER ESTIMATION

In general no analytic solution to the β zeros - except for 3 special cases: one covariate, two covariates, orthogonal design matrix.

ONE COVARIATE

$$\min_{\beta} \left(\underbrace{(Y - X\beta)^T (Y - X\beta)}_{\substack{N \times 1 \quad N \times 1 \quad 1 \times 1}} + \lambda |\beta| \right) \quad \beta \text{ scalar}$$

$1; \hat{w}, \hat{z}_0$

$$\min_{\beta} \left(Y^T Y - 2\beta^T X^T Y + \beta^T X^T X \beta + \lambda |\beta| \right)$$

scalar

$$\hat{\beta}_{OLS} = (X^T X)^{-1} X^T Y = \frac{1}{N} X^T Y$$

standardised cov

$$x_i = \frac{x_i - \bar{x}}{s}$$

$(\sum x_i^2)^{-1} = N$ or $N-1$

$$\sum_{i=1}^N x_i = \frac{1}{s} \left(\sum_{i=1}^N x_i - N\bar{x} \right) = 0$$

$$\sum_{i=1}^N x_i^2 = \frac{1}{s^2} \left(\sum_{i=1}^N (x_i - \bar{x})^2 \right)$$

$$s^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

$$s^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$$

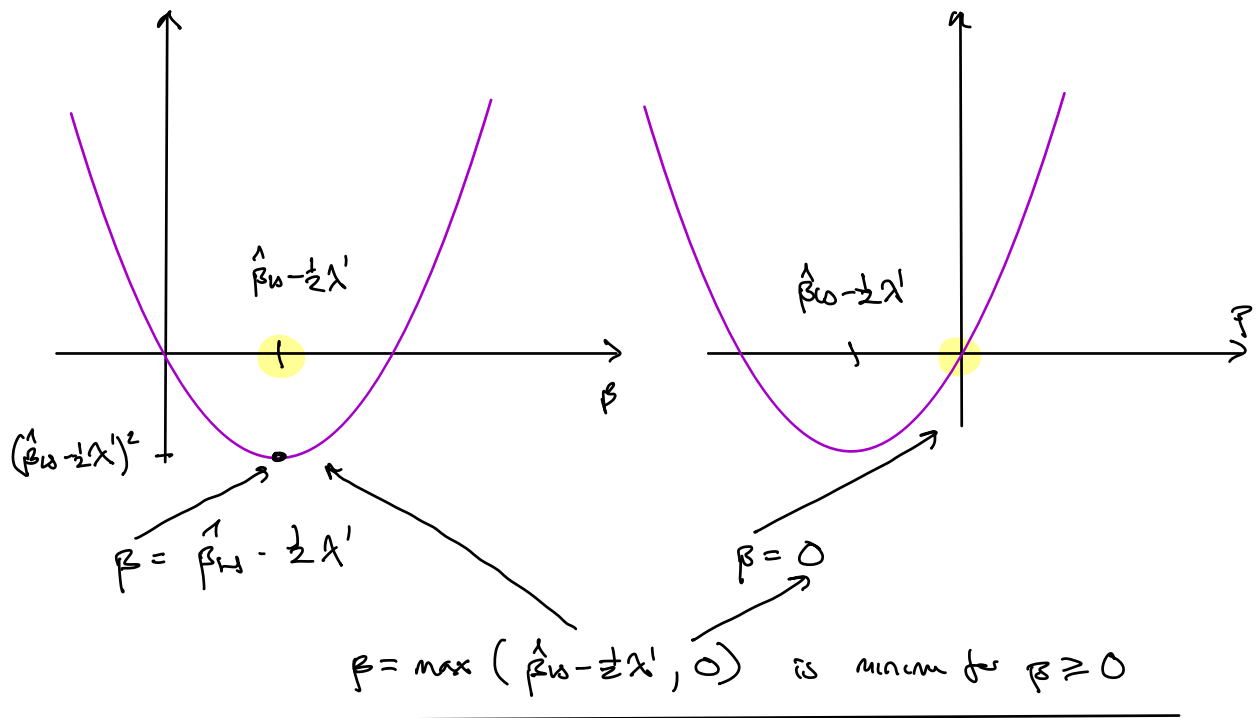
$$\min_{\beta} \left(-2\beta N\hat{\beta}_{OLS} + N\beta^2 + \lambda |\beta| \right) = \min_{\beta} \left(-2\beta \hat{\beta}_{OLS} + \beta^2 + \frac{\lambda}{N} |\beta| \right)$$

Due to the $|\beta|$ term we look separately at $\beta \geq 0$, $\beta \leq 0$

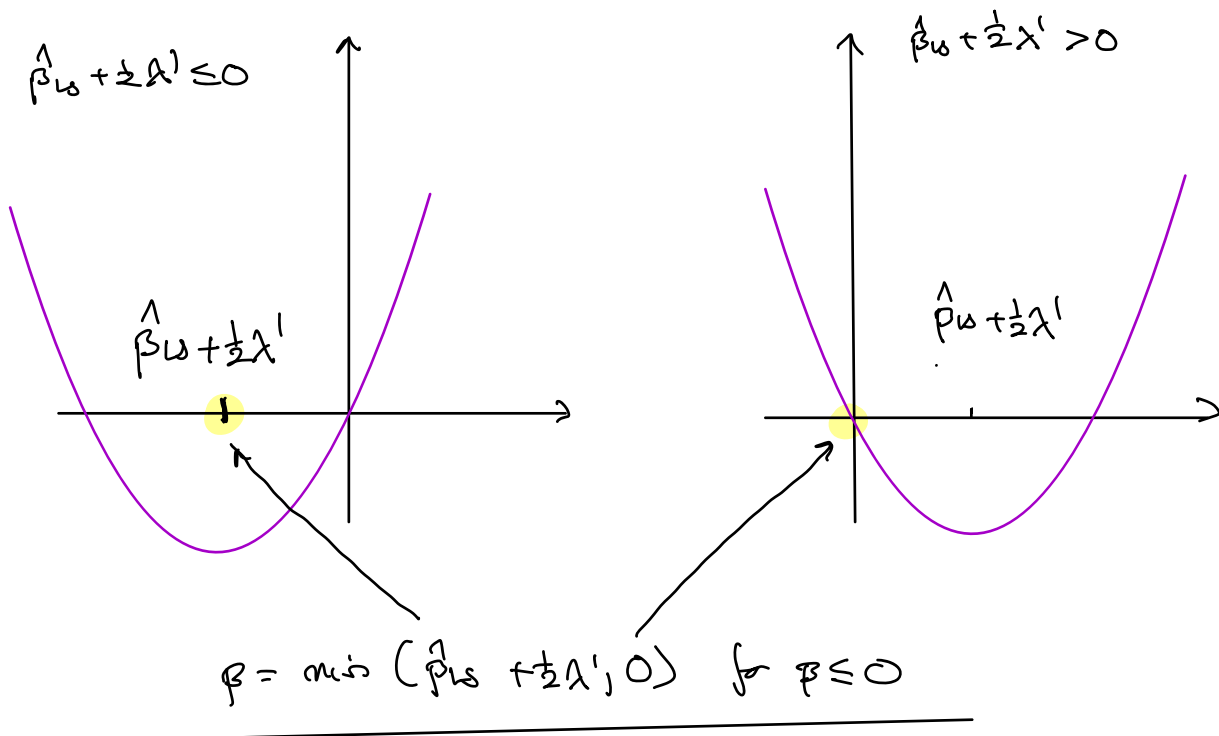
$\beta \geq 0$: loss $-2\beta \hat{\beta}_{OLS} + \beta^2 + \lambda \beta = \beta^2 - 2(\hat{\beta}_{OLS} - \frac{\lambda}{2}) \beta$

$$= (\beta - (\hat{\beta}_{OLS} - \frac{\lambda}{2}))^2 - (\hat{\beta}_{OLS} - \frac{\lambda}{2})^2$$

$$\beta^2 - 2a\beta = (\beta - a)^2 - a^2$$

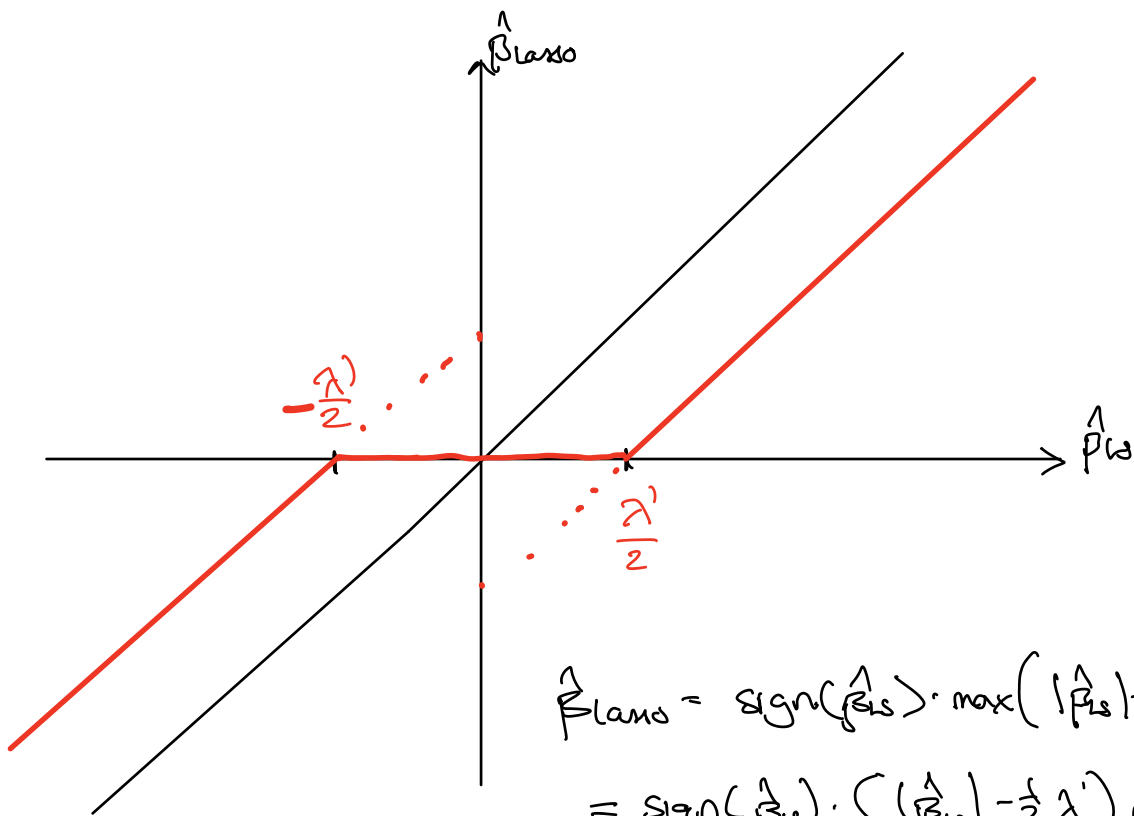
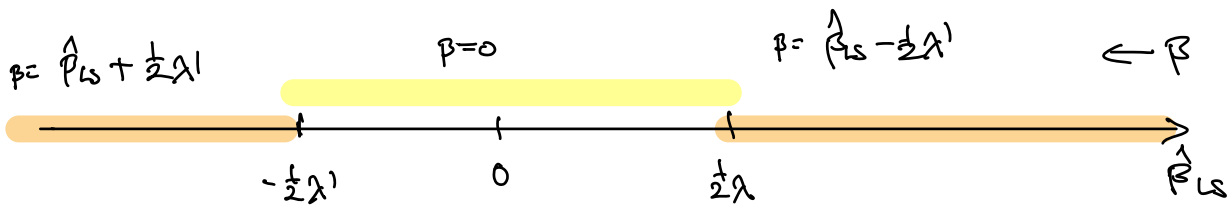


$\beta \leq 0$: loss $-2\beta\hat{\rho}_0 + \beta^2 - \lambda'\beta = \beta^2 - 2(\hat{\rho}_0 + \frac{1}{2}\lambda')\beta$
 $= (\beta - (\hat{\rho}_0 + \frac{1}{2}\lambda'))^2 - (\hat{\rho}_0 + \frac{1}{2}\lambda')^2$



Next: combine and do this conditional on $\hat{\beta}_{LS}$ not β *as a function*

$$\left. \begin{aligned} \beta &= \max(\hat{\beta}_{LS} - \frac{1}{2}\lambda', 0) & \beta \geq 0 \\ \beta &= \min(\hat{\beta}_{LS} + \frac{1}{2}\lambda', 0) & \beta \leq 0 \end{aligned} \right\}$$



$$\begin{aligned} \hat{\beta}_{Lasso} &= \text{sign}(\hat{\beta}_{LS}) \cdot \max(|\hat{\beta}_{LS}| - \frac{1}{2}\lambda', 0) \\ &= \text{sign}(\hat{\beta}_{LS}) \cdot (|\hat{\beta}_{LS}| - \frac{1}{2}\lambda')_+ \end{aligned}$$

Soft threshold operator: $S_\lambda(x) = \text{sign}(x) (|x| - \lambda)_+$

Orthogonal design matrix

$$X^T X = I = (X^T X)^T$$

$$\hat{\beta}_{LS} = (X^T X)^{-1} X^T Y = X^T Y$$

$$\min_{\beta} \left((Y - X\beta)^T (Y - X\beta) + \lambda \sum_{j=1}^p |\beta_j| \right)$$

$$= \min_{\beta} \left(Y^T Y - 2\beta^T \underbrace{X^T Y}_{\hat{\beta}_{LS}} + \beta^T \underbrace{X^T X}_{I} \beta + \lambda \sum_{j=1}^p |\beta_j| \right)$$

$$= \min_{\beta} \left(-2\beta^T \hat{\beta}_{LS} + \beta^T \beta + \lambda \sum_{j=1}^p |\beta_j| \right)$$

$\beta_1 \hat{\beta}_{LS,1} + \beta_2 \hat{\beta}_{LS,2} + \dots \quad \beta_1^2 + \beta_2^2 + \dots$

$$= \min_{\beta} \left(\sum_{j=1}^p (-2\beta_j \hat{\beta}_{LS,j} + \beta_j^2 + \lambda |\beta_j|) \right)$$

$$= \sum_{j=1}^p \min_{\beta_j} \left(-2\beta_j \hat{\beta}_{LS,j} + \beta_j^2 + \lambda |\beta_j| \right)$$

\Rightarrow can handle each β_j separately

$$\hat{\beta}_{lasso}: \text{ we find } \text{sign}(\hat{\beta}_{LS,j}) \left(|\hat{\beta}_{LS,j}| - \frac{1}{2} \lambda \right)$$

Pseudocode for cyclic coordinate descent

λ given

$t=0$: initialize $\beta_1^t, \dots, \beta_p^t$

define the order of the $J = \{k_1, k_2, \dots, k_p\}$

for (t in 1 :convergence) [or while]

↓

for (j in J)

⌈

$\tilde{y} = y - X_{-j} \beta_{-j}^t$ partial residuals

minimize $\beta_j \quad (\tilde{y} - X_j \beta_j)^T (\tilde{y} - X_j \beta_j) + \lambda |\beta_j|$

↓ } → soft thresholding $\beta_j^{(t+1)}$