# Lecture October 25

OLS : $C(\tilde{\beta}) = \frac{1}{m} \sum_{i=0}^{m-1} (y_i - \tilde{y}_i)^2$

$\tilde{y}_i = \sum_{j=0}^{p-1} x_{ij} \, \tilde{\beta}_j$

$\quad = x_{i*} \beta$

$\beta = [\beta_0, \beta_1 \, -- \, \beta_{p-1}]^T$

$X \in \mathbb{R}^{m \times p}$

$y, \tilde{y} \in \mathbb{R}^m \qquad \Rightarrow \hat{\beta} = (X^T X)^{-1} X^T y$

## Ridge

$C(\beta) = \frac{1}{m} \sum_{i=0}^{m-1} (y_i - x_{i*}\beta)^2$

$\qquad + \lambda \sum_{j=0}^{p-1} \beta_j^2$

$\qquad \lambda > 0 \wedge \sum_{j=0}^{p-1} \beta_j^2 < t$

$\hat{\beta}_{Ridge} = (X^T X + \lambda I_p)^{-1} X^T y$

$I_p = \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \cdot_{\cdot 1} \end{bmatrix} \in \mathbb{R}^{p \times p}$

$$X^T X \in \mathbb{R}^{P \times P}$$

## Lasso

$$C(\beta) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - x_{i*}\beta)^2$$

$$+ \lambda \sum_{j=0}^{P-1} |\beta_j|$$

Scaling / normalisation of data.

Scikit-learn:

Ridge and Lasso do by default not include $\beta_0$ in the fitting

$$\lambda \sum_{j=1}^{P-1} \beta_j^2 \quad \text{Ridge}$$

$$\lambda \sum_{j=1}^{P-1} |\beta_j| \quad \text{Lasso}.$$

_____ * _____

$$\tilde{y}_i = \sum_{j=0}^{P-1} \beta_j x_i^j \quad \Rightarrow$$

$$X = \begin{bmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^{\,} \\ 1 & x_1 & x_1^2 & \cdots & x_1^{P-1} \\ \vdots & & & & \\ & & & & \\ & & & & \\ 1 & & & & x_{m-1}^{P-1} \end{bmatrix}$$

1st example includes
the intercept.

2nd example

$$\tilde{X} = \begin{bmatrix} x_0 & x_0^1 & \cdots & x_0^{P-1} \\ x_1 & & & \\ x_2 & & & \\ \vdots & & & P-1 \\ x_{m-1} & \cdots & & x_{m-1}^{P-1} \end{bmatrix}$$

Singular value decomp
( SVD )

$$X \in \mathbb{R}^{m \times P}$$

$$X = U \Sigma V^T$$

$$U^T U = U U^T = \mathbb{I}_m$$

$$U \in \mathbb{R}^{m \times m}$$

$$
U = \begin{bmatrix} | & | & & | \\ u_0 & u_1 & \cdots & u_{m-1} \\ | & | & & | \end{bmatrix}
$$

$$
\Sigma = \begin{bmatrix} \sigma_0 & & & \\ & \sigma_1 & & O \\ & & \ddots & \\ O & & & \sigma_p \\ & & & \vdots \\ & & & 0 \end{bmatrix}
$$

$$
\Sigma \in \mathbb{R}^{m \times p}
$$

$$
\sigma_0 > \sigma_1 > \cdots > \sigma_p > 0
$$

then for $\sigma_{p-1} \cdots \sigma_{m-1} = 0$

$$
V^T V = V V^T = \mathbb{I}_p
$$

$$
V \in \mathbb{R}^{p \times p}
$$

$$
V = \begin{bmatrix} | & | & & | \\ v_0 & v_1 & \cdots & v_{p-1} \\ | & 1 & & | \end{bmatrix}
$$

$$
\Sigma = \begin{bmatrix} \tilde{\Sigma} \\ 0 \end{bmatrix} \qquad \tilde{\Sigma} = \begin{bmatrix} \sigma_0 & & O \\ & \ddots & \\ O & & \sigma_{p-1} \end{bmatrix}
$$

$$
= \begin{bmatrix} \sigma_0 & \sigma_1 & & & & \mathcal{O} \\ & & \ddots & & & 0 \\ \mathcal{O} & & & & \sigma_{P-1} & \\ & & & & & 0 \\ & & & & & 0 \end{bmatrix}
$$

Simple example to illustrate Lasso, Ridge and OLS

$$ X \in \mathbb{R}^{n \times P} $$

$$ X = I = \begin{bmatrix} 1 & 0 \\ 0 & 1_{-n} \end{bmatrix} $$

## OLS (skip $1/n$)

$$ C(\bar{\beta}) = \sum_{i=0}^{n-1} (y_i^i - \beta_i)^2 $$

$$ = \sum_{i=0}^{P-1} (y_i^i - \beta_i)^2 $$

$$ \hat{\beta}^{OLS} = y_i^i $$

## Ridge

$$ \cdots \quad \sum_{}^{P-1} (\cdots)^2 \; ) \; \sum_{}^{P-1} \beta^2 $$

$$C(\beta) = \sum_{i \geq 0} (y_i - \beta_i) + \lambda \sum_{i \geq 0} \beta_i$$

$$\frac{\partial C}{\partial \beta} = 0 \qquad \Rightarrow$$

$$\hat{\beta}_i^{Ridge} = \frac{y_i}{1 + \lambda} \qquad \lambda > 0$$

## Lasso

$$C(\beta) = \sum_{i=0}^{P-1} (y_i - \beta_i)^2 \qquad \left| \quad \frac{\partial C}{\partial \beta_i} = 0 \right.$$
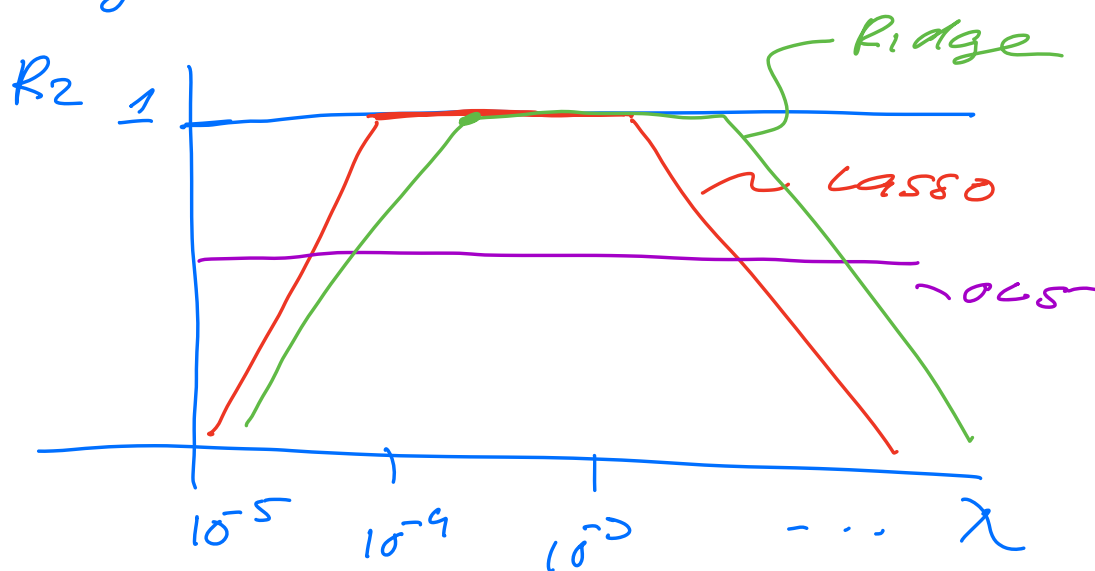$$+ \lambda \sum_{i=0}^{P-1} |\beta_i|$$

$$\frac{d|x|}{dx} = \begin{cases} +1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$$

$$\hat{\beta}_i^{Lasso} = \begin{cases} y_i - \lambda/2 & \text{if } y_i > \lambda/2 \\ y_i + \lambda/2 & \text{if } y_i < -\lambda/2 \\ 0 & \text{if } |y_i| \leq \lambda/2 \end{cases}$$



$\hat{\beta}_i^{OLS}$

$\beta_i$

$\hat{\beta}_i^{Ridge} = \frac{y_i}{1+\lambda}$

$-\lambda/2$

$y_i$

$\hat{\beta}_\lambda^{12}$ $\Big\{ \begin{matrix} \hat{\beta}_\lambda^{Lasso} \end{matrix}$

# Mathematics of SVD and Ridge.



$$X^T X = \quad \quad X = U \Sigma V^T$$

$$= V \Sigma^T \underbrace{U^T U}_{I_m} \Sigma V^T$$

$$= V \Sigma^T \Sigma V^T$$

$$= V \tilde{\Sigma}^2 V^T$$

$$\Sigma^T u^T u \Sigma = \Sigma^T \underbrace{I_m} \Sigma \in \mathbb{R}^{p \times p}$$

$$\Sigma \in \mathbb{R}^{n \times p}$$

$$\tilde{\Sigma}^2 = \begin{bmatrix} \sigma_0^2 & & & \\ & \sigma_1^2 & & \\ & & \ddots & \\ & & & \sigma_{p-1}^2 \end{bmatrix}$$

$$\tilde{y}_{OLS} = X \hat{\beta}_{OLS}$$

$$= X (X^T X)^{-1} X^T y$$

$$= u \Sigma' v^T (v \tilde{\Sigma}^2 v^T)^{-1} v \Sigma^T u^T y$$

$$= \underline{u u^T} y = \sum_{i=0}^{p-1} u_i u_i^T y$$

$$u = \begin{bmatrix} | & | & | & & | \\ u_0 & u_1 & u_2 & -- & u_{m-1} \\ | & | & | & & | \end{bmatrix}$$

## Ridge

$$\tilde{y}_{Ridge} = \left[ \sum_{i=0}^{p-1} u_i u_i^T \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \right] y$$

# Further properties

$$X^T X = V \tilde{\Sigma}^2 V^T = V \Sigma^T \Sigma V^T$$

multiply from the right
with $V$ $\quad$ $(V V^T = V^T V = I)$

$$(X^T X) V = V \tilde{\Sigma}^2$$

$$V = \begin{bmatrix} | & | & & \\ v_0 & v_1 & -- & v_{p-1} \\ | & | & & \end{bmatrix}$$

$$(X^T X) v_i = v_i \sigma_i^2$$

The eigenvalues and eigen-
vectors of $X^T X$ are the
singular values $\sigma_i^2$ and
the orthogonal vectors $v_i$.

## OLS

$$\frac{\partial C}{\partial \beta} = 0 = -X^T (X\beta - y) \frac{2}{n}$$

$$\frac{\partial^2 c}{\partial \beta \partial \beta^T} = \frac{2}{n} X^T X \quad (\text{curvature})$$

Hessian matrix

$$\text{cov}(x, y) = \frac{1}{n} \sum_{\ell=0}^{m-1} (x_\ell - \mu_x)(y_\ell - \mu_y)$$

$$X = [x_0, x_1 - - x_{m-1}]$$

$$y = [y_0, y_1 - - y_{m-1}]$$

$$X = \begin{bmatrix} x_{00} & x_{01} & \cdots & & x_{0\,p-1} \\ x_{10} & & & & \\ & & & & \\ & & & & \\ x_{m-1\,0} & x_{m-1\,1} & - & - & x_{m-1\,p-1} \end{bmatrix}$$

$$x_0 \quad x_1$$

$$\text{cov}(x_0, x_1) = \frac{1}{n} \sum_{\ell=0} (x_{\ell 0} - \mu_{x_0})(x_{\ell 1} - \mu_{x_1})$$

$$\Rightarrow \text{cov}(X)$$

$$= \frac{1}{n} X^T X$$