

**Erasmus+ machine
Learning lecture,
October 9, 2023**

SVD analysis of OLS
Design matrix $X \in \mathbb{R}^{n \times p}$

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

$$\cancel{X} = \begin{bmatrix} x_0 & x_1 & x_2 & \dots & x_{p-1} \end{bmatrix}$$

$$X = U \Sigma V^T$$
$$U \in \mathbb{R}^{n \times n}$$
$$V \in \mathbb{R}^{p \times p}$$
$$U U^T = U^T U = \mathbf{I}$$
$$V V^T = V^T V = \mathbf{I}$$
$$n \geq p$$

$$\Sigma = \begin{bmatrix} \tau_0 & & & \\ & \ddots & & 0 \\ & & -\tau_{p-1} & \\ & & & 0 \end{bmatrix}$$

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

$$\tau_0 > \tau_1 > \tau_2 > \dots > \tau_{p-1} > 0$$

Example

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

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

\Downarrow

$$\begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}$$

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

$$\frac{\partial^2 C}{\partial \beta \partial \beta^T} = \frac{2}{n} X^T X$$

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

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

$$= V \Sigma^T \Sigma V^T \quad . \quad \checkmark$$

$$[X^T X] \cdot V = V \Sigma^T \Sigma V^T V = V \underbrace{\Sigma^T \Sigma}_{\mathbb{R}^{P \times P}} V^T$$

$$V = \begin{bmatrix} V_0 & V_1 & V_2 & \dots & V_{p-1} \end{bmatrix}$$

$$V_i^T V_j = S_{ij}$$

$$[X^T X] V = V \begin{bmatrix} -\sigma_0^2 & & & \\ & -\sigma_1^2 & & 0 \\ & & \ddots & \\ & & & -\sigma_{p-1}^2 \end{bmatrix}$$

$$[X^T X] v_i = v_i \sigma_i^2$$

eigenvalues of $X^T X$ are the singular values of X . σ_i^2 with eigenvectors v_i . $\sigma_i^2 > 0 \Rightarrow X^T X$ is

a positive definite matrix \Rightarrow
convex optimization problem.

$\frac{1}{n} X^T X$ is the covariance matrix

For completeness

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

$$X X^T = u \underbrace{\Sigma V^T V \cdot \Sigma^T u^T}_{\Sigma}$$

$$= u \Sigma \Sigma^T u^T + u$$

$$[X X^T] u = u \Sigma \Sigma^T = u \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix}$$
$$u = [u_0 \ u_1 \ u_2 \dots \ u_n]$$

$$\Rightarrow [xx^T] u_i = \tau_i^2 u_i'$$

$$\tau_0^2, \tau_1^2, \dots, \tau_{p-1}^2, 0, 0, \dots, 0$$

OLS

$$\hat{\beta} = (x^T x)^{-1} x^T y$$

$$\tilde{y} = x \hat{\beta} = x (x^T x)^{-1} x^T y$$

SVD

$$= u \Sigma v^T [v \Sigma^T v^T]^{-1} v \Sigma^T u^T y$$

assume A and B , are square
and both are invertible

$$V; V^T = V^{-1}$$
$$\underbrace{\Sigma^T \Sigma}_{\text{is also}} = \begin{bmatrix} \sigma_0^2 & & \\ & \ddots & \\ & & \sigma_{p-1}^2 \end{bmatrix}$$

invertible $|R|^{P \times P}$

$$\underline{\frac{1}{AB}} = (AB)^{-1} = B^{-1} \cdot A^{-1}$$

$$VV^T = V^T V = \underline{\underline{1}}$$

$$(V \Sigma^T \Sigma V^T)^{-\frac{1}{2}} = (V^T)^{-\frac{1}{2}} (\Sigma^T \Sigma)^{-\frac{1}{2}} V^{-\frac{1}{2}}$$

$$= V (\Sigma^T \Sigma)^{-\frac{1}{2}} V^T \quad (\in \mathbb{R}^{P \times P})$$

$$\hat{y} = u \sum \underbrace{V^T V}_{\equiv} (\Sigma^T \Sigma)^{-\frac{1}{2}} \underbrace{V^T V \Sigma^T u g}_{\equiv}$$

$$= u \sum \left(\frac{1}{\Sigma^T \Sigma} \right) \Sigma^T u^T y$$

$$\hat{y} = \left[\sum_{j=0}^{P-1} \left(\begin{matrix} u \\ u^T \end{matrix} \right) \right] y$$

$$\hat{y} = u \Sigma (\Sigma^T \Sigma)^{-1} \Sigma^T u^T y$$

$$\Sigma = \begin{bmatrix} \sigma_0^2 & \sigma_1 & \sigma_2 & \dots & \sigma_{P-1} \end{bmatrix} \in \mathbb{R}^{m \times p}$$

$$u = [u_0, u_1, u_2, \dots, u_{n-1}]$$

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

$$\Sigma^T u^T = \begin{bmatrix} \tau_0 u_0 & \tau_1 u_1 & \dots & \tau_{q-1} u_{q-1} \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} z & 0 \\ c & 1 \\ c & 0 \end{bmatrix}$$

$$\Sigma^T = \begin{bmatrix} z & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

$$u = \begin{bmatrix} u_{00} & u_{01} & u_{02} \\ u_{10} & u_{11} & u_{12} \\ u_{20} & u_{21} & u_{22} \end{bmatrix} \quad n = 3$$

$$= \begin{bmatrix} u_0 & u_1 & u_2 \end{bmatrix}$$

$$\tilde{y} = \left[\sum_{j=0}^{p-1} u_j u_j^T \right] y$$

$$p \leq n \left[\sum_{j=0}^{p-1} u_i u_j^T \right] y = u u^T y \\ = y$$

projection of y in terms
of the orthogonal vectors u_i
but including only the
singular values different
from zero.

$X^T X$ when it is singular
cannot invert. Cheap trick
is to simply a positive
small constant λ to
the diagonals.

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

$$X^T X + \lambda I_{P \times P} \Rightarrow$$

$$\tilde{y} = \underbrace{x}_{\in \mathbb{R}^{n \times p}} \left(x^T x + \lambda I \right)^{-1} x^T y$$

Ridge Regression

Cost function (objective)

$$C(\beta) = \frac{1}{n} \sum_{i=0}^{n-1} \left(y_i - \sum_{j=0}^{p-1} x_{ij} \beta_j \right)^2$$

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

— Regularization
or shrinkage

$$= \frac{1}{n} (y - X\beta)^T (y - X\beta)$$

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

$$= \frac{1}{n} \| (y - X\beta) \|_2^2 + \lambda \| \beta \|_2^2$$

$\| x \|_2 = \sqrt{\sum_{i=1}^n x_i^2}$

$$\frac{\partial C}{\partial \beta} = -\frac{1}{m} X^T (y - X\beta) + 2\lambda \beta = 0$$

$$X \geq 0$$

$$\lambda \cdot a = \lambda \Rightarrow \lambda$$

$$0 = -X^T y + X^T X \beta + \lambda \beta \Rightarrow$$

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

Ridge Regression

λ is a hyperparameter
(shrinkage parameter)

can show that (exercises
this week)

$$\tilde{y} = \left[\sum_{j=0}^{P-1} \frac{u_j u_j^T}{\tau_j^2 + \lambda} \right] y$$

when λ is large, we
"shrink" the contribution
from a given component
- j .

Example

$$X = \mathbb{1}$$

$$\tilde{y} = \beta ; n = p$$

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

$$\tilde{y} = X\beta = \beta$$

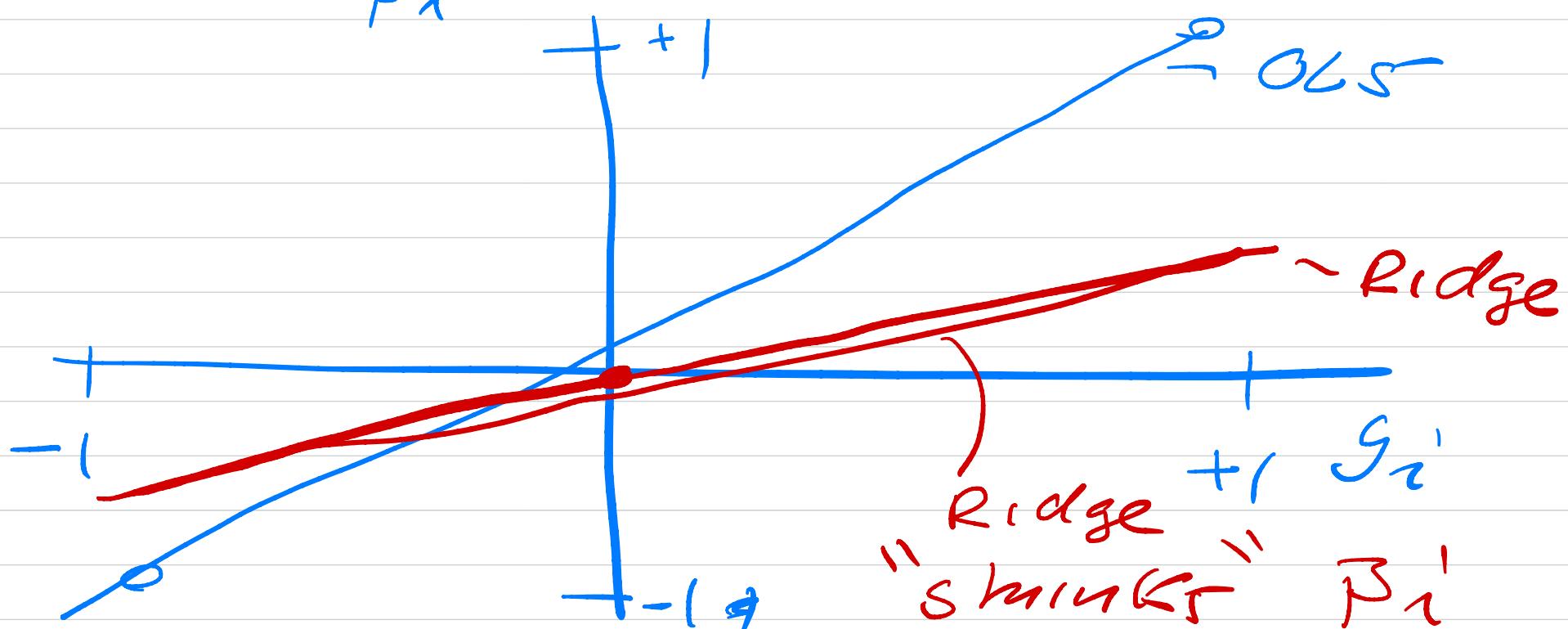
minimize wrt so $\beta \Rightarrow$

$$y_i = \beta_i$$

Ridge : $C(\beta) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \beta_i)^2 + \lambda \sum_{i=0}^{n-1} \beta_i^2$

$$\frac{\partial C}{\partial \beta_i} = 0 = -\frac{2}{n} (y_i - \hat{y}_i) + 2\lambda \beta_i$$

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



LASSO Regression

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

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

$$= \frac{1}{n} \| (y - X\beta) \|_2^2 + \lambda \|\beta\|_1$$

we cannot find an analytical expression for
 $\hat{\beta}_{\text{LASSO}}$ due to

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

Example

$$\tilde{y} = \beta$$

$$C(\beta) = \frac{1}{n} \sum_i (y_i - \beta_i)^2$$

$$+ \lambda \sum_i |P_i|$$

$$\frac{\partial C}{\partial \beta_i} = -\frac{2}{n} (y_i - \beta_i) + \lambda \frac{P_i}{|P_i|}$$

$$\beta_{\text{LASSO}}^{\text{LASSO}} = \begin{cases} y_i - \lambda/2 & \text{if } y_i > \lambda/2 \\ y_i + \lambda/2 & \text{if } y_i < -\lambda/2 \\ 0 & \text{if } |y_i| \leq \lambda/2 \end{cases}$$

