

Lecture September 3

$$\hat{\beta}_{OLS} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{n} \|y - X\beta\|_2^2$$

$$\left\{ \|x\|_2^2 = \sum_{i=1}^n x_i^2 \right.$$

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$

$$\hat{\beta}_{OLS} = \underbrace{(X^T X)^{-1}}_{\text{}} X^T y$$

$$B = \lim_{\lambda \rightarrow 0} (X^T X + \lambda \underline{I}_p)^{-1} X^T$$

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

$$y \in \mathbb{R}^n$$

$$\beta \in \mathbb{R}^p$$

Moore-Penrose
pseudo inverse,

SVD

$$X = U \Sigma V'$$

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

$$U^T U = U U^T = I_m$$

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

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

$$V V^T = V^T V = I_p$$

$$m \geq p$$

example

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \quad \begin{matrix} \sigma_2 \\ \sigma_3 = 0 \end{matrix}$$

$$\sigma_1 > \sigma_2 > \sigma_3$$

$$\sigma_0 > \sigma_1 > \dots > \sigma_{p-1} \geq 0$$

$$\Sigma \Sigma^T \neq \Sigma^T \Sigma$$

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

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

$$\Sigma \Sigma^T = \begin{bmatrix} 4 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{aligned} \underline{X^T X} &= V \Sigma^T \underbrace{U^T U}_{I_n} \Sigma V^T \\ [A, D] &= 0 \\ &= V \boxed{\Sigma^T I_n \Sigma} V^T \\ &= \underline{\Sigma^T \Sigma} = \Sigma^2 \quad (p \times p) \\ V V^T &= V^T V = I_p \end{aligned}$$

$$\begin{aligned} \tilde{y}_{OLS} &= X \hat{\beta}_{OLS} \\ &= X (X^T X)^{-1} X^T y \end{aligned}$$

$$\begin{aligned} &= U \cancel{\Sigma} V^T \left(\frac{1}{\cancel{\Sigma \Sigma}} \right) V \cancel{\Sigma}^T U^T y \\ &= U \underline{U^T} y = \left(\sum_{i=0}^{n-1} u_i u_i^T \right) y \end{aligned}$$

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

$$u_i^T u_j = \delta_{ij}$$

$$\begin{aligned} X^T X &= \Sigma^T \Sigma = \Sigma^2 \in \mathbb{R}^{p \times p} \\ &= \begin{bmatrix} \sigma_0^2 & & & \\ & \sigma_1^2 & & \\ & & \ddots & \\ 0 & & & \sigma_{p-1}^2 \end{bmatrix} \end{aligned}$$

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

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

$$V = [v_0, v_1, \dots, v_{p-1}]$$

$$\begin{aligned} X X^T &= U \Sigma V^T V \Sigma^T U^T = \\ &= U (\Sigma \Sigma^T) U^T \end{aligned}$$

$$X U$$

$$(XX^T)u = u(\Sigma\Sigma^T)$$

u are the eigenvector
of XX^T

$X^T X$ has σ_i^2 has
eigenvalues.

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

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

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

eigenvalues = singular
values $\sigma_0 > \sigma_1 > \dots > \sigma_{p-1}$

> 0

(i) sample values

$$E[x] = \frac{1}{n} \sum_{i=0}^{n-1} x_i = \bar{\mu}_x$$

$$= \sum_{i=0}^{n-1} p(x_i) x_i = \mu_x$$

$$\text{var}(x) = \frac{1}{n} \sum_{i=0}^{n-1} (x_i - \bar{\mu}_x)^2$$

$$\underline{\text{cov}(x, y)} = \underline{\frac{1}{n} \sum_{i=0}^{n-1} (x_i - \bar{\mu}_x)(y_i - \bar{\mu}_y)}$$

$$p(x_1, x_2) = p(x_1) p(x_2)$$

$$\begin{aligned} \text{cov}(x, y) &= \iint dx dy \\ &\times (x - \mu_x)(y - \mu_y) \\ &p(x) p(y) \\ &= 0 \end{aligned}$$

$$\text{cov}(x, y) \sim \underline{x^T x}$$

(cov(x))

covariance matrix

$$X = [x_0 \ x_1 \ \dots \ x_{p-1}]$$

$$x_i = \begin{bmatrix} x_{0i} \\ x_{1i} \\ \vdots \\ x_{m-1i} \end{bmatrix}$$

scaling of matrix

$$\bar{x}_i = \begin{bmatrix} x_{0i} - \bar{\mu}_{x_i} \\ x_{1i} - \bar{\mu}_{x_i} \\ \vdots \\ x_{m-1i} - \bar{\mu}_{x_i} \end{bmatrix}$$

$$X^T X \rightarrow \bar{X}^T \bar{X} \quad (= X^T X)$$

$$\begin{aligned} & \bar{X}_n^T \bar{X}_j \\ &= \left[\sum_{l=0}^{n-1} (x_{li} - \bar{\mu}_{x_i})(x_{lj} - \bar{\mu}_{x_j}) \right] \\ &= n \cdot \text{cov}(\bar{x}_i, \bar{x}_j) \end{aligned}$$

$$X^T X = n \cdot \underline{\text{COV}(X)}$$

$$= \begin{bmatrix} \text{COV}(\bar{x}_0, \bar{x}_0) & \text{COV}(\bar{x}_0, \bar{x}_1) & \dots \\ \vdots & \ddots & \\ \text{COV}(\bar{x}_{p-1}, \bar{x}_0) & \dots & \end{bmatrix}$$

is this a symmetric matrix? = yes

Diagonal = $\text{var}(x_i)$

$$\text{corr}(x, y) = \frac{\text{COV}(x, y)}{\sqrt{\text{var}(x) \text{var}(y)}}$$

$$\in [-1, +1]$$

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

$$\begin{bmatrix} 1 & \text{cov}(\bar{x}_0, \bar{x}_1) & \dots & \dots \\ & \ddots & & \\ \text{cov}(\bar{x}_p, \bar{x}_0) & \dots & \dots & 1 \end{bmatrix}$$

$$\text{cov}(X) = \frac{1}{n} X^T X$$

Linear regression,
2nd derivative of
MSE = Hessian = $X^T X$

$$(iii) \quad \text{var}(\beta) \propto (X^T X)^{-1}$$

Ridge Regression

$$\hat{\beta}_{\text{Ridge}} = \arg \min_{\beta \in \mathbb{R}^p}$$

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

$$\begin{aligned}
 & \frac{1}{n} \text{MSE} \\
 & + \lambda \sum_{j=0}^{p-1} \beta_j^2 \\
 & \quad \underbrace{\lambda \|\beta\|_2^2}_{\lambda \beta^T \beta} \quad \text{--- Regularization penalty}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial C}{\partial \beta^T} &= - \frac{2}{n} x^T (y - x\beta) \\
 &+ \lambda 2\beta = 0
 \end{aligned}$$

$$x^T x \beta + \lambda \beta = x^T y \Rightarrow$$

$$\hat{\beta}_{\text{Ridge}} = \underline{\underline{\left(x^T x + \lambda \mathbf{I}_p \right)^{-1} x^T y}}$$

$\lambda > 0$ Hyperparameter

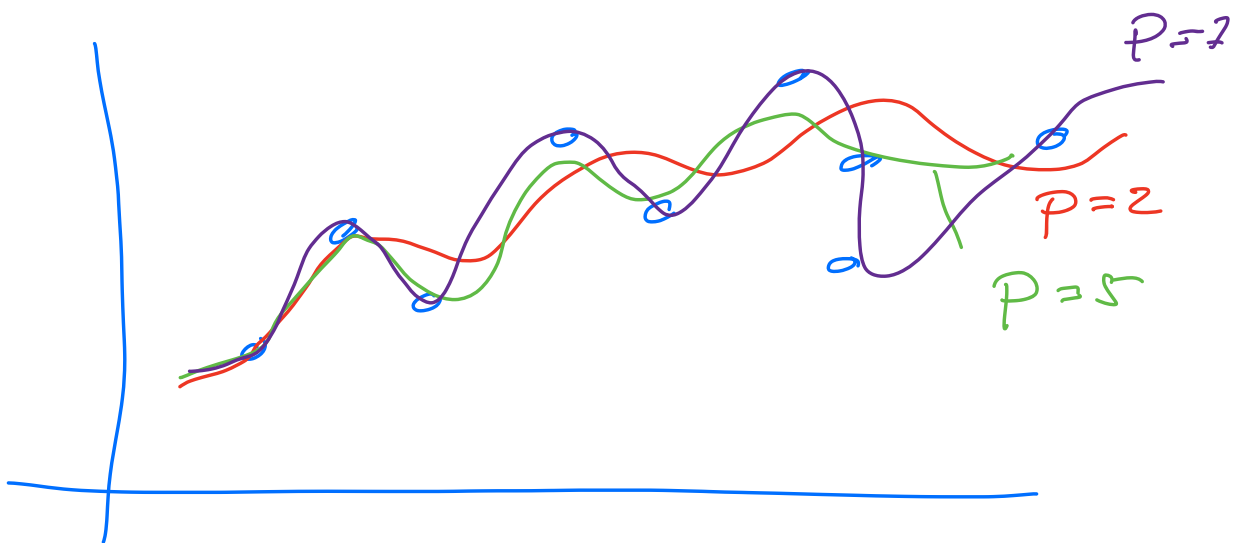
Lasso Regression

$$\hat{\beta}_{\text{Lasso}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \quad \begin{aligned} & \|y - X\beta\|_2^2 \\ & + \lambda \sum_{j=0}^{p-1} |\beta_j| \end{aligned} \quad (\lambda \|\beta\|_1)$$

$$\frac{d|\beta|}{d\beta} := \operatorname{sgn}(\beta) = \begin{cases} 1 & \beta > 0 \\ 0 & \beta = 0 \\ -1 & \beta < 0 \end{cases}$$

$$\frac{\partial \mathcal{L}}{\partial \beta} = 0 = -2X^T(y - X\beta) + \lambda \operatorname{sgn}(\beta)$$

why do we introduce λ as an additional parameter?



	$p =$	2	5	7
β_1		6.2	6.82	6.87
β_2		0.1	-0.7	-0.8
β_3			+30.0	+31.0
β_4			-10.0	+20.0
β_5			+2.0	+30.0
β_6				-10.0 ± 10
β_7				+5.0

$$f(x) = \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \dots + \beta_7 x^7$$

with a parameter γ ,
we can dampen the
fluctuations in β .