

Lecture January 20

Shrinkage methods:

OLS

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

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

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

$$\hat{y}_{\text{OLS}} = x \beta^{\text{opt}}$$

Ridge

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

$$\| \beta \|_2^2 + \lambda \| \beta \|_2^2$$

$$\| \beta \|_2^2 = \sum_{i=0}^{p-1} \beta_i^2$$

$$\frac{\partial C}{\partial \beta} = 0 \Rightarrow \beta_{\text{Ridge}}^{\text{opt}} =$$

$$\frac{(x^T x + \lambda I)^{-1} x^T y}{\lambda \| \beta \|_2^2}$$

$$\lambda \| \beta \|_2^2 \leq C < \infty$$

norm-2 of vector x

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

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

Lasso

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

Example

- leave out intercept (β_0)

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

$$X = \begin{bmatrix} 1 & 0 & \dots \\ 0 & 1 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

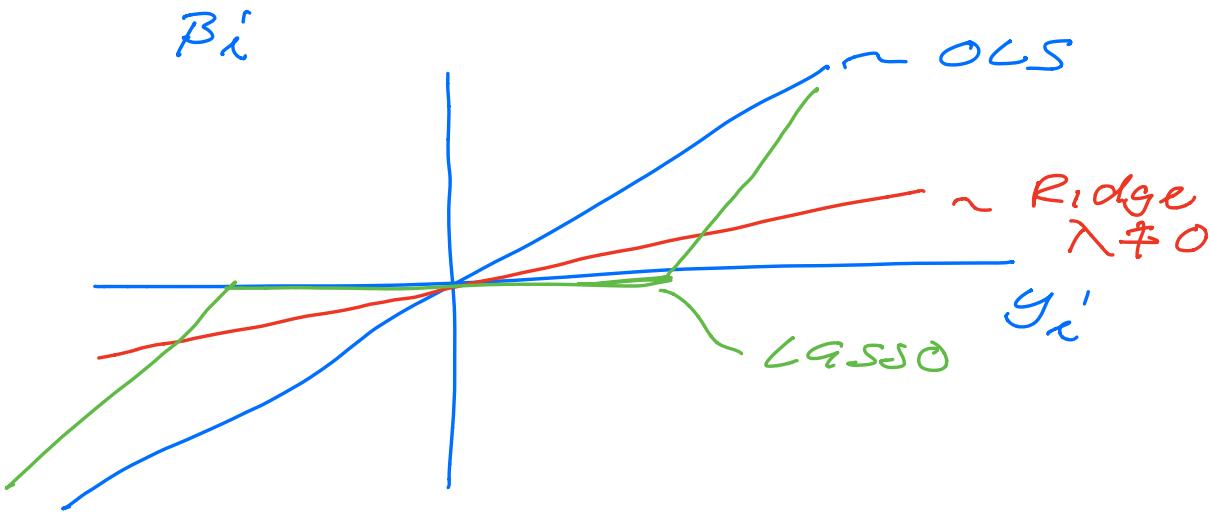
OLS (drop $\frac{1}{n}$)

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

$X_i \neq \beta$

$$\beta_i^{\text{opt}} = y_i$$

$$\hat{y} = X\beta^{\text{opt}}$$



Ridge

$$C(\beta) = \sum_{i=0}^{n-1} (y_i' - \beta_i)^2 + \lambda \sum_{j=0}^{n-1} \beta_j^2$$

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

$$\beta_i^{\text{Ridge}} = \frac{y_i'}{1 + \lambda}$$

$$\tilde{y}_i = \beta_0 + \underline{\beta_1} x_i + \underline{\beta_2} x_i^2 + \dots$$

LASSO

$$C(\beta) = \sum_i (y_i' - \beta_i)^2 + \lambda \sum_i |\beta_i|$$

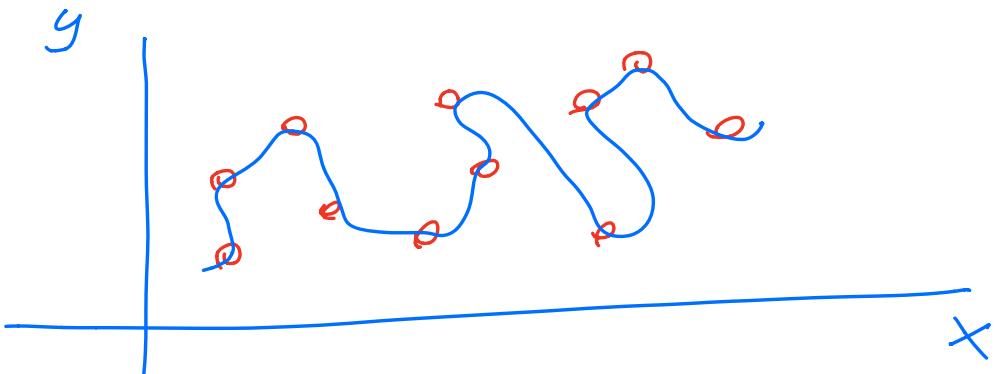
$$\lambda \sum_i \sqrt{\beta_i^2}$$

\$\sim\$

$$\frac{\partial \mathcal{L}}{\partial \beta} = 0 \Rightarrow$$

$$-2 \sum_i (y_i - \beta_i) + \lambda \sum |\beta_i|$$

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



$$\hat{y}_i = \beta_0 + \beta_1 x_i^1 + \beta_2 x_i^2 + \dots + \beta_{p-1} x_i^{p-1}$$

$$\text{OLS} \quad \hat{y} \approx \hat{g} + \epsilon$$

$$\epsilon \sim N(0, \sigma^2)$$

$$\text{var}(\beta^{\text{opt}}) = \frac{\sigma^2 (X^T X)^{-1}}{n}$$

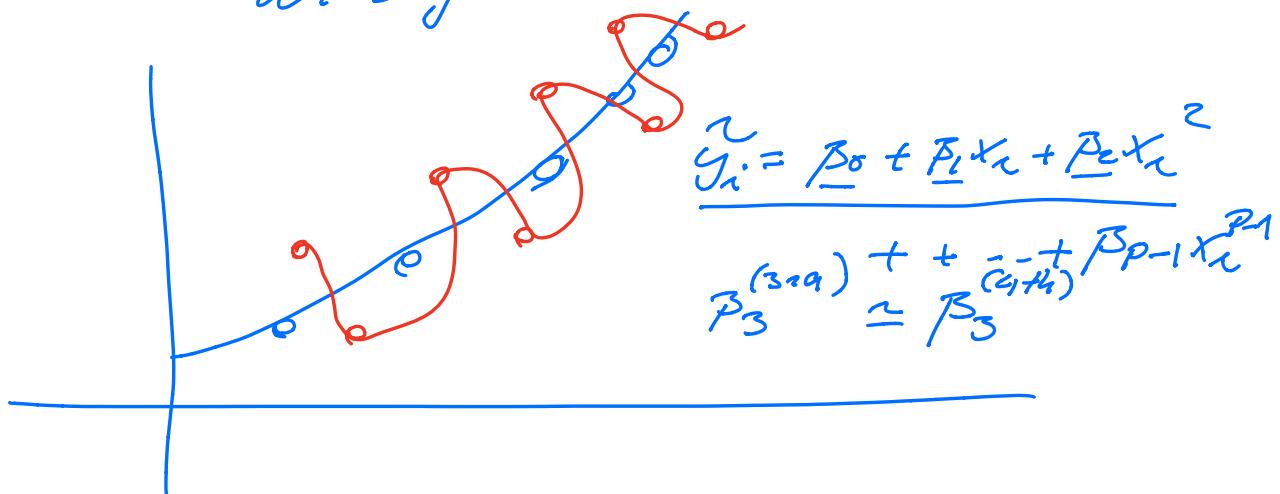
\hat{y}_i (5th-order)

$$\hat{y}_i = (\beta_0 \pm \text{STD} \beta_0) + (\beta_1 \pm \text{STD} \beta_1) x_i^1$$

+ ---

\hat{g}_i (7th-order)

$(\beta, \pm SSTD\beta)_{S7H}$ and
compare with $(\beta, \pm SSTD\beta)_{M7}$
with fluctuating velocity.



OLS and probability

$$D = \{(x_0, y_0), (x_1, y_1), \dots, (x_{m-1}, y_{m-1})\}$$

$$y_i = x_i * \beta + \varepsilon_i'$$

$$\downarrow \quad \begin{matrix} \uparrow \\ \text{iid } \sim N(0, \sigma^2) \end{matrix}$$
$$N(x_i * \beta, \sigma^2)$$

$$P(y_i | x_i; \beta, \sigma^2) \sim N(x_i; \beta, \sigma^2)$$

Maximum Likelihood:

$$P(y | x; \beta, \sigma^2) = \prod_{i=0}^{n-1} P(y_i | x_i; \beta, \sigma^2)$$

$$\beta^{\text{opt}} = \arg \max_{\beta \in \mathbb{R}^P} P(y | \beta)$$

$$= \arg \max_{\beta \in \mathbb{R}^P} \sum_{i=0}^{n-1} \log P(y_i | x_i; \beta)$$

$$C(\beta) = \sum_{i=0}^{n-1} \log P(y_i | x_i; \beta)$$

minimise $C(\beta)$

$$C(\beta) = - \sum_{i=0}^{n-1} \log \dots$$

- log is a convex function

$$y_i \sim N(x_i; \beta, \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y_i - x_i \beta)^2}{2\sigma^2} \right]$$

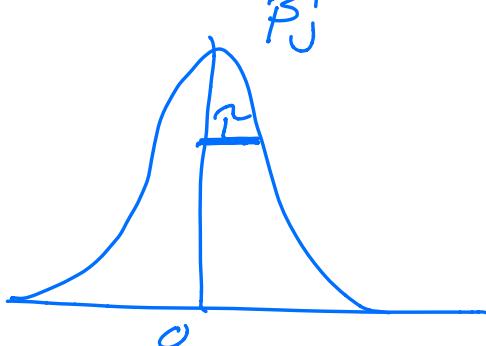
$$\begin{aligned}
 & = P(y_i | x_i, \beta, \sigma^2) \\
 \log P(y_i | x_i, \beta, \sigma^2) & = \\
 & - (y_i - x_i \cdot \beta)^2 / 2\sigma^2 - \frac{1}{2} \log(2\pi\sigma^2) \\
 & + \sum_{i=0}^{n-1} (y_i - x_i \cdot \beta)^2 / 2\sigma^2 \\
 & + \frac{n}{2} \log(2\pi\sigma^2)
 \end{aligned}$$

$$\frac{\partial C(\beta)}{\partial \beta} = 0 = \frac{\partial}{\partial \beta} \sum_{i=0}^{n-1} (y_i - x_i \cdot \beta)^2$$

$$\Rightarrow \beta^{\text{opt}} = (X^T X)^{-1} X^T y$$

Connection with Ridge?

assume $P(\beta) = \prod_{j=0}^{p-1} N(\beta_j | 0, \tau^2)$



Bayesian interpretation

$$P(Y | X \beta) = \prod_{i=0}^{n-1} N(y_i | x_i \beta, \sigma^2) \times \prod_{j=0}^{p-1} N(\beta_j | 0, \tau^2)$$

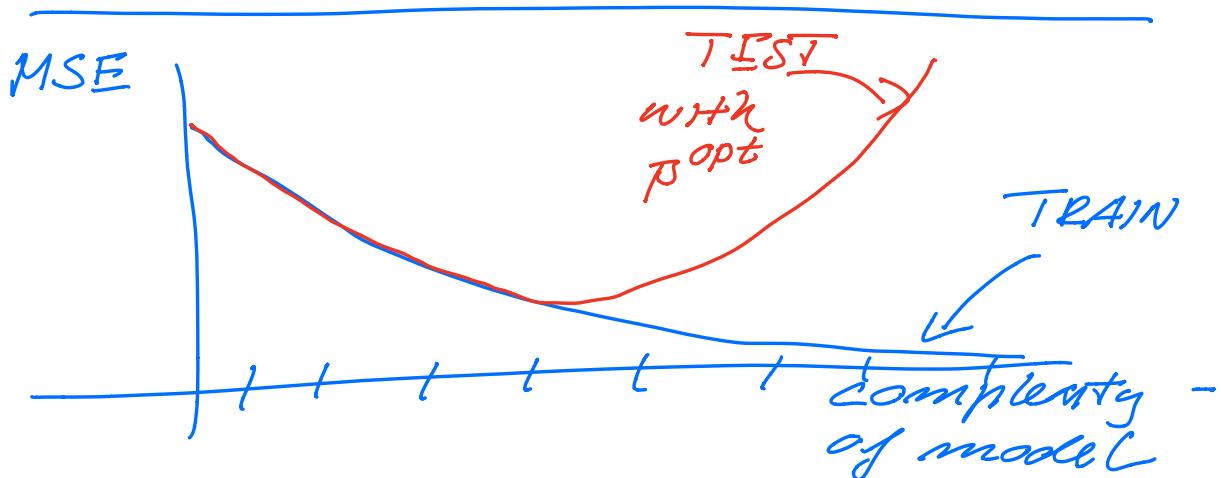
$$\hat{\beta}^{\text{opt}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{arg\,min}}$$

$$- \left(\sum_{i=0}^{n-1} \log(N(y_i | x_i \beta, \sigma^2)) \right) \\ + \sum_{j=0}^{p-1} \log N(\beta_j | 0, \tau^2)$$

$$C(\beta) \propto \sum_{i=0}^{n-1} (y_i - x_i \beta)^2 \\ + \lambda \|\beta\|_2^2$$
$$\lambda \approx 1/\sigma^2$$

Hastie et al chapter 3.4

Bias - Variance tradeoff



Hastie et al., figure 2.11

Bias-variance tradeoff:

Rewrite of MSE.

$$y = f(x) + \epsilon$$

$$f(x) \approx \tilde{g}(x) = X\beta$$

$$C(\beta) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{g}_i)^2$$

Expected value:

$$\mathbb{E}[x^i] = \int p(x) x^i dx$$

$$\mu = \int p(x) x dx$$

$$\left\{ \begin{array}{l} \text{Discrete version} \\ \sum_i p(x_i) x_i \end{array} \right.$$

In ML sample mean

$$\bar{\mu} = \frac{1}{n} \sum_i x_i' \neq \mu \\ = E[x] \quad (\langle x \rangle)$$

$$C(\beta) = E[(y - \hat{y})^2]$$

(add and subtract $E[y]$)

$$= \underbrace{E[(y - E[y])^2]}_{\text{bias}} + \underbrace{E[(\hat{y} - E[\hat{y}])^2]}_{\text{variance (of the model)}} + \sigma^2$$

How to calculate the

bias and variance

in a reliable?

- Resampling

* Bootstrapping

* cross validation.

Bootstrap

1) pick randomly n -
data from a sample
of m -data

$$\underline{D} = \{x_1, x_2, x_3, \dots, x_n\}$$

$$\overset{*}{D} = \underbrace{\{x_3, x_1, x_1, x_{10}, \dots, x_{m-1}, x_m\}}_{n \text{ in total}}$$

2) compute sample
expected value Θ_1^*

3) repeat as many as
 $\frac{1-2}{M-1}$ times you want

4) compute final
expected values

$$\Theta^* = \frac{1}{M} \sum_{j=1}^M \Theta_j^*$$

EFRON (1976)

used in small data
sets.