

Lecture January 21

what have we done?

$$y = f(x) + \varepsilon \quad | \quad y_i = X_i^T \beta + \varepsilon_i$$

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

$$E[y_i] = X_i^T \beta$$

$$y_i \sim N(X_i^T \beta, \sigma^2)$$

OLS

$$\beta_{OLS}^{opt} = (X^T X)^{-1} X^T y$$

Ridge

$$\beta_{Ridge}^{opt} = (X^T X + \lambda I)^{-1} X^T y$$

Lasso

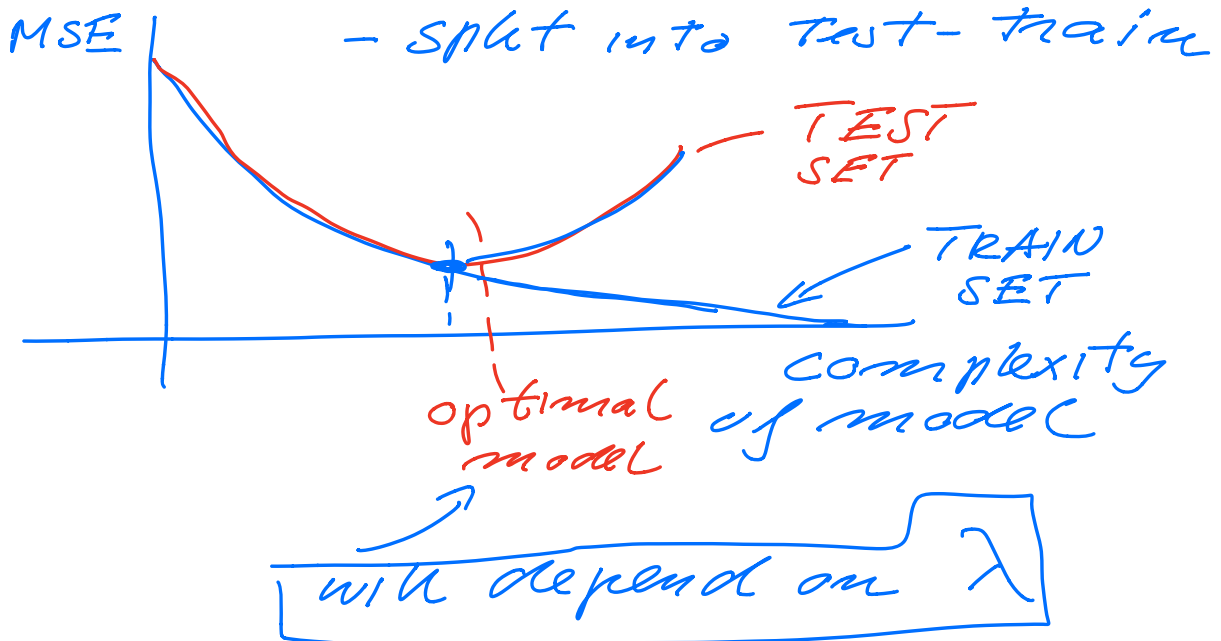
$$C(\beta) = \frac{1}{n} \| (y - X\beta) \|_2^2$$

$$+ \lambda \|\beta\|_1$$

hyperparameter.

Tuning λ ($\lambda > 0$) can
yield a MSE for the

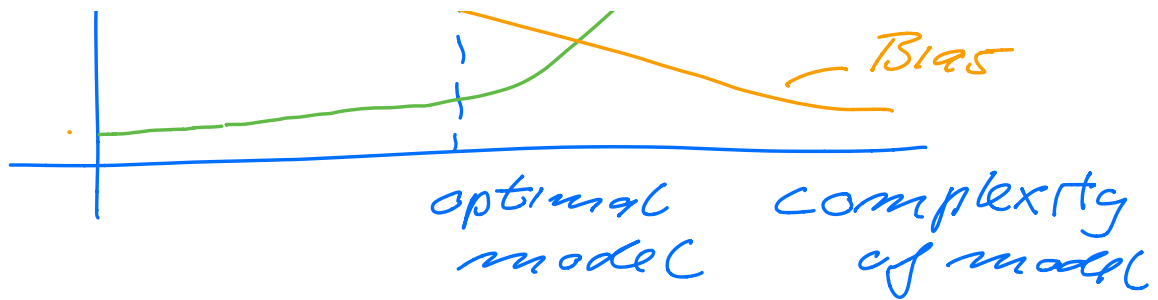
test data which is
smaller than $MSE(OLS)_{\lambda=0}$
when $\lambda > 0$



Bias - variance Tradeoff

$$MSE = \underbrace{E[(y - E[\tilde{y}])^2]}_{\text{BIAS}} + \underbrace{\text{var}[\tilde{y}]}_{\text{variance}} + \sigma^2$$





- Resampling (Reliable expected values)
- Bootstrap

- Cross-validation

$$\bar{\mu}(\text{sample}) = \frac{1}{n} \sum_{i=0}^{n-1} x_i' \neq \mu_{\text{true}}$$

cross-validation:

k-fold CV

$$k = 3$$

- 1)

TRAIN	TRAIN	TEST
-------	-------	------

 \rightarrow Predicted error on Test ϵ_1
- 2)

TRAIN	TEST	TRAIN
-------	------	-------

 $\rightarrow \epsilon_2$
- 3)

TEST	TRAIN	TRAIN
------	-------	-------

 $\rightarrow \epsilon_3$

$$\bar{\epsilon} = \frac{\epsilon_1 + \epsilon_2 + \epsilon_3}{K} \quad K=3$$

Typical choices of $K \approx 5-10$

LOOCV = leave one out CV

1)

TRAIN		n
-------	--	-----

2)

TRAIN		$n-1$		T
-------	--	-------	--	---

⋮

n)

1		TRAIN
---	--	-------

Bootstrap

TRAIN		TEST
-------	--	------

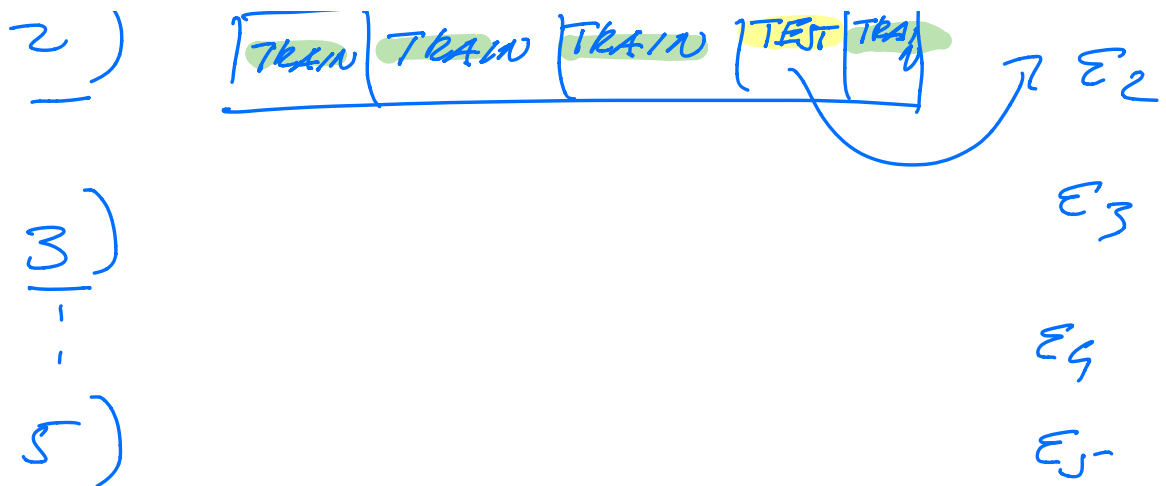
↑ Bootstrap unchanged

CV - $K=5$

1)

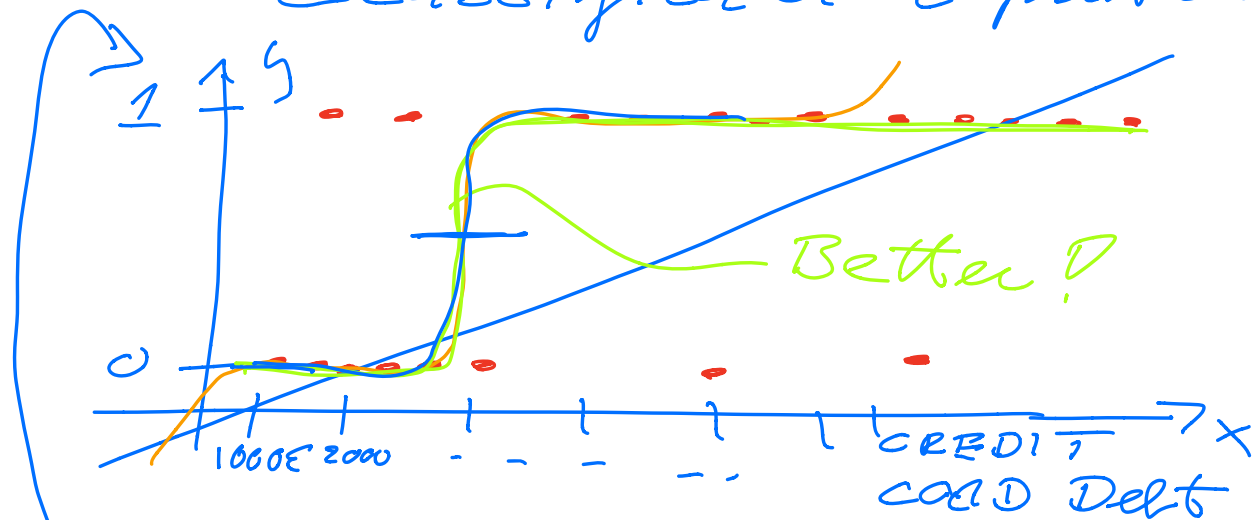
TRAIN		TRAIN		TRAIN		TRAIN		TEST
-------	--	-------	--	-------	--	-------	--	------

→ ϵ_1



$$\bar{\epsilon} = \frac{\epsilon_1 + \epsilon_2 + \epsilon_3 + \epsilon_4 + \epsilon_5}{5}$$

Logistic regression =
Classification problem



Binary case = $\begin{cases} 1 = \text{not pay} \\ 0 = \text{pay} \end{cases}$

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

$$y_i = \{0, 1\}$$

$$P(y_i | x_i, \beta) = ?$$

$$\boxed{\binom{n}{x} p^x q^{n-x}}$$

Binomial
Distribution

$$P(y_i | x_i, \beta) = \underbrace{p(y_i=1)}_{p=?}^{y_i} \underbrace{p(y_i=0)}_{q=?}^{1-y_i}$$

$$p(y_i=0) = 1 - p(y_i=1)$$

$$\boxed{p(y_i=0) + p(y_i=1)} = 1$$

Maximum Likelihood:

$$\underline{P(y | X, \beta)} = \prod_{i=0}^{n-1} \underbrace{p(y_i=1)}^{y_i} \underbrace{(1 - p(y_i=1))}^{1-y_i}$$

$P(y_i | x_i, \beta)$ | we need a model for this.

- tanh

- Sigmoid.

$$p(y_i | x_i, \beta) = \frac{e^{\beta_0 + \beta_1 x_i}}{e}$$

$$\frac{1}{1 + e^{\beta_0 + \beta_1 x_i}}$$

$$= \frac{e^{-t}}{1 + e^{-t}}$$

$$= \frac{1}{1 + e^{-t}}$$

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

$$= \sum_{i=0}^{n-1} \log(p(y_i=1)^{y_i} (1 - p(y_i=1))^{1-y_i})$$

$$\arg \min_{\beta \in \mathbb{R}^p} - \sum_{i=0}^{n-1} \log(p(y_i=1)^{y_i} (1 - p(y_i=1))^{1-y_i})$$

$$= \beta^{\text{opt}}$$

cost function
= $C(\beta)$

$$\frac{\partial C(\beta)}{\partial \beta} = 0 \Rightarrow$$

non-linear equations
in β , \Rightarrow
no analytical expression

for $\beta \Rightarrow$ compute
gradients and second
derivatives numerically
in order to minimize
 $C(\beta)$