


# Improving Training

Geena Kim

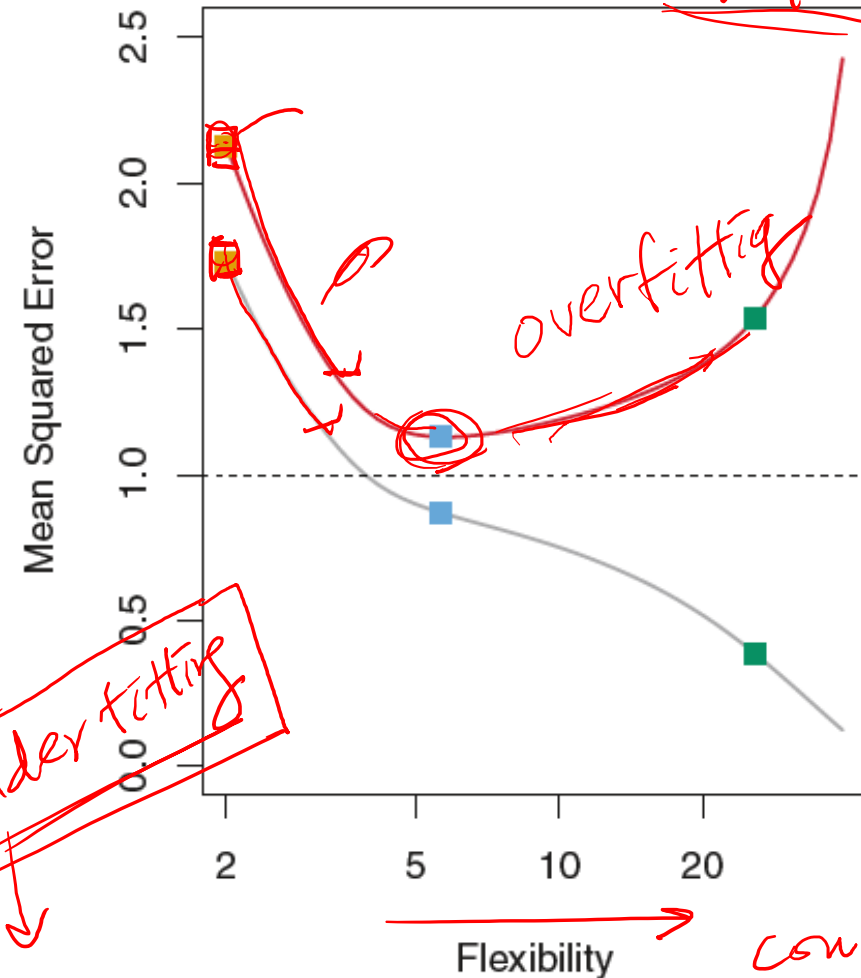
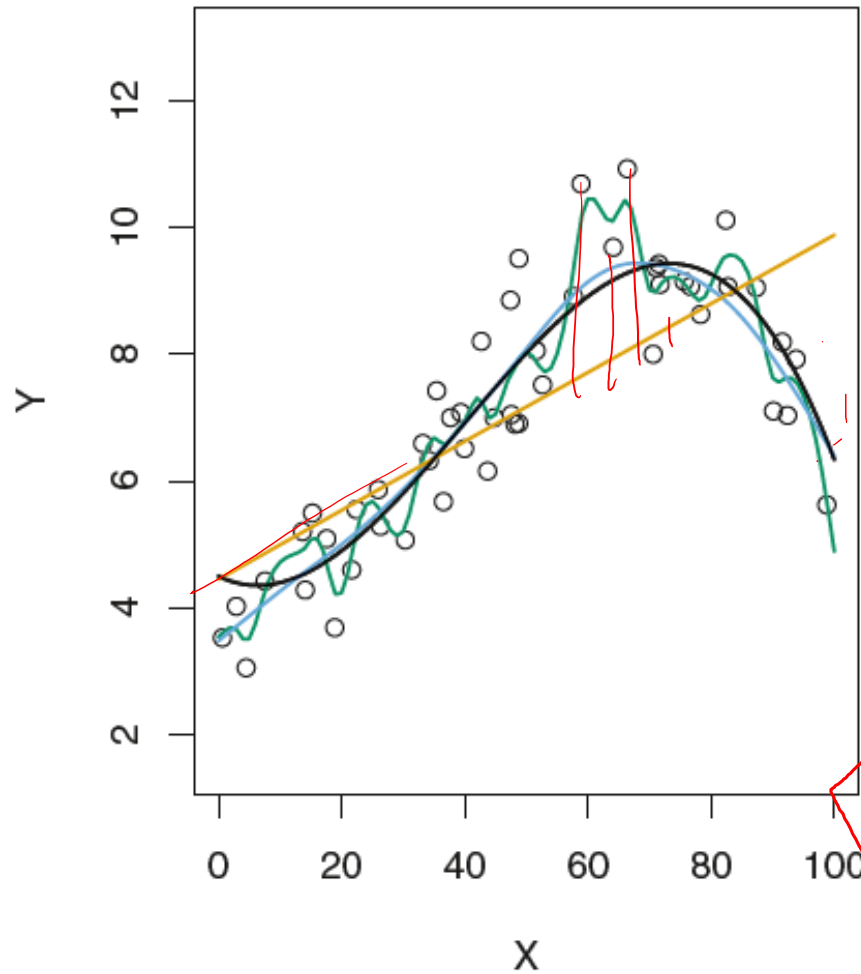


# Better training

## The Goals:

- Smallest generalization error 
- Better test performance score

# Generalization error



Underfitting

Generalization Error  
(Validation or test)  
test Model to Unseen data

Flexibility

complexity  
= capacity

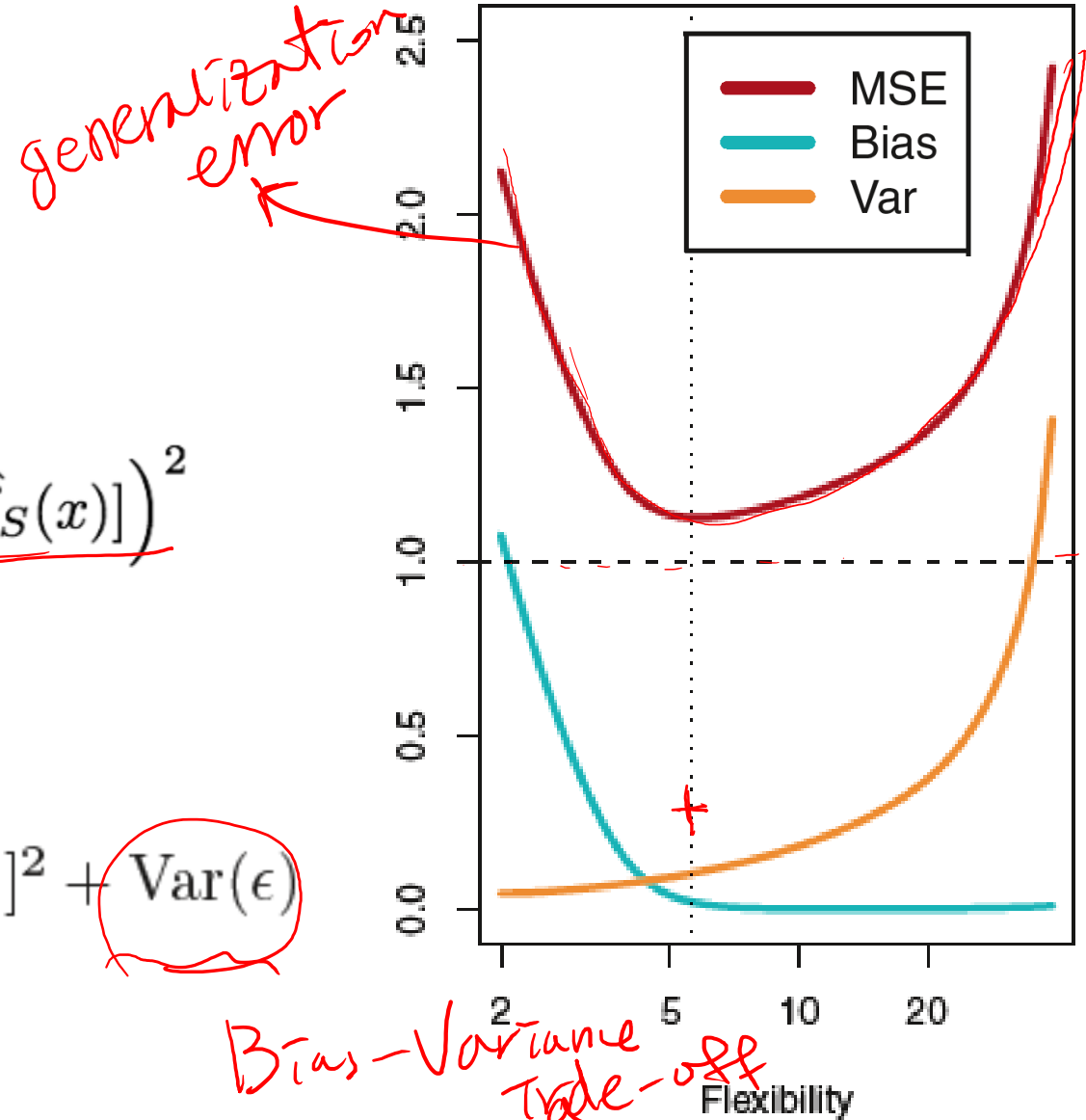
# Where is the error coming from?

E.g. In regression...

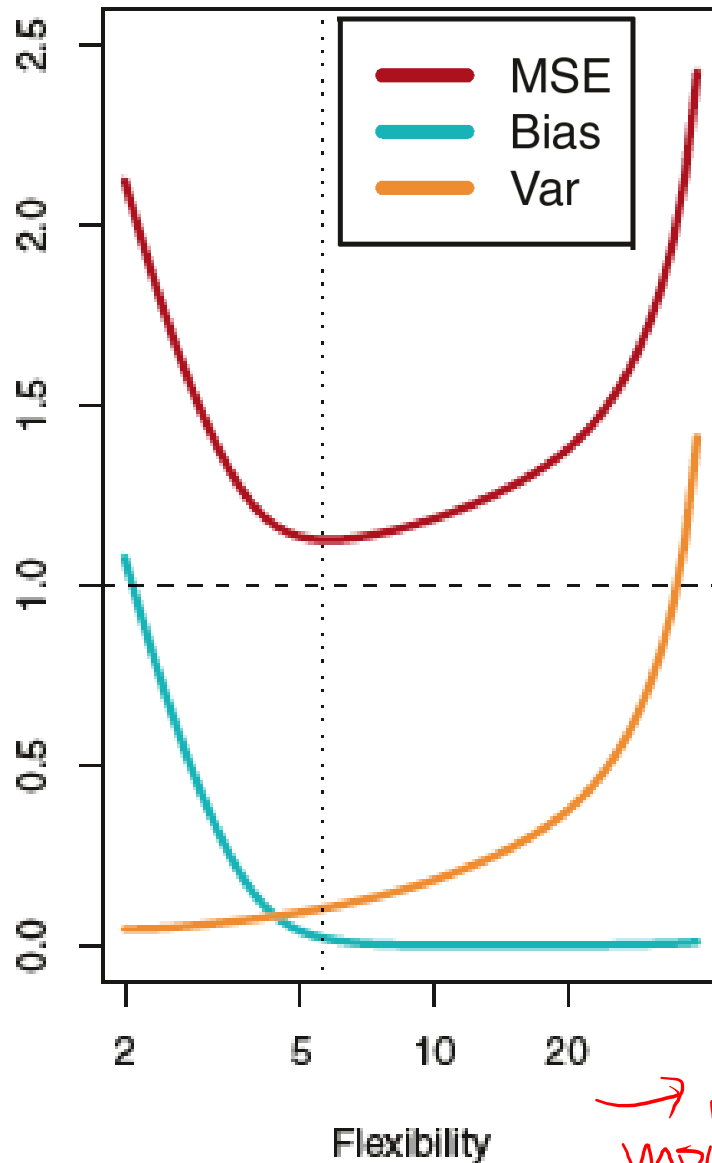
$$y = \underline{f(x)} + \underline{\epsilon}$$

$$\begin{aligned} \underline{MSE} &= \mathbb{E}[(y - \hat{f}_S(x))^2] \\ &= \underbrace{\text{Var}(f(x) - \hat{f}_S(x)) + \text{Var}(\epsilon)} + \underbrace{\left(\mathbb{E}[f(x)] - \mathbb{E}[\hat{f}_S(x)]\right)^2}_{\text{Bias}^2} \\ &\quad + \underbrace{\mathbb{E}^2[\epsilon] + 2\mathbb{E}[\epsilon]\mathbb{E}[f(x)] - 2\mathbb{E}[\epsilon]\mathbb{E}[\hat{f}_S(x)]}_{0} \end{aligned}$$

$$E(y_0 - \hat{f}(x_0))^2 = \underbrace{\text{Var}(\hat{f}(x_0))}_{\text{Variance}} + \underbrace{[\text{Bias}(\hat{f}(x_0))]^2}_{\text{Bias}^2} + \underbrace{\text{Var}(\epsilon)}_{\text{Noise}}$$



# How do we know which term to drop/include?



- Parameters
- Design parameters

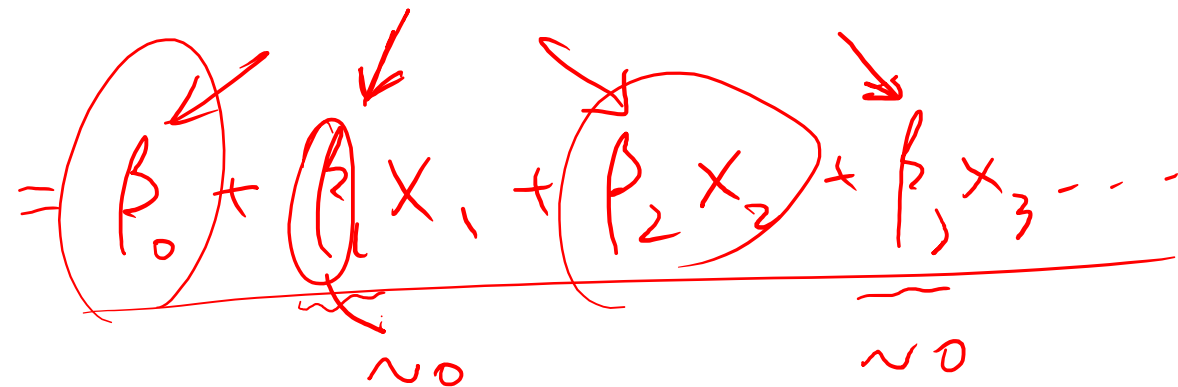
→ model complexity  
←

# What features to include?

## Method 1. Best subset method

- The idea: test all possible combinations
- Curse of dimensionality!

## Method 2. Regularization



A handwritten diagram illustrating the concept of regularization in linear regression. The equation shown is 
$$= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots$$
 Each coefficient  $\beta_0, \beta_1, \beta_2, \beta_3$  is circled in red. Red arrows point to each circle. Below the  $\beta_1$  and  $\beta_3$  terms, there are red wavy lines and the text "no", indicating that these terms are being penalized or "shrunk" towards zero by the regularization process.

# Regularization

Original loss function

$$\mathcal{L} = \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ij} \right)^2 + \boxed{\phantom{0000}}$$

Let's penalize some terms that are not necessary

With a L2 regularization

$$\mathcal{L} = \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \quad \lambda \geq 0$$

# L2 regularization (Ridge)

$$\mathcal{L} = \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

Also called Ridge regression

What does the lambda ( $\lambda$ ) do?



# L2 regularization

What does the lambda ( $\lambda$ ) do?

$$\mathcal{L} = \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

$\lambda$



$|\beta|$



Total Loss (L)

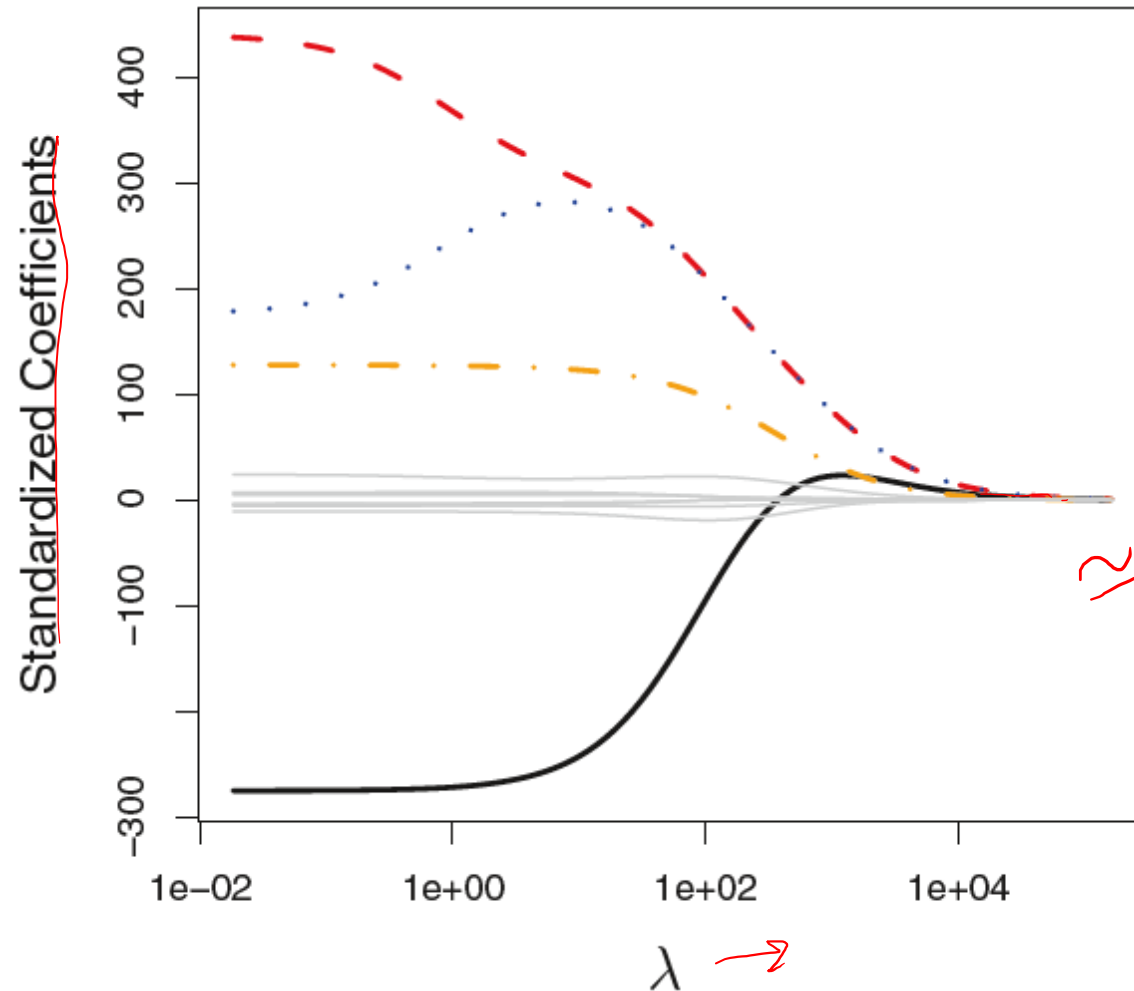


Original Loss ( $L_0$ )



# L2 regularization

What does the lambda ( $\lambda$ ) do?

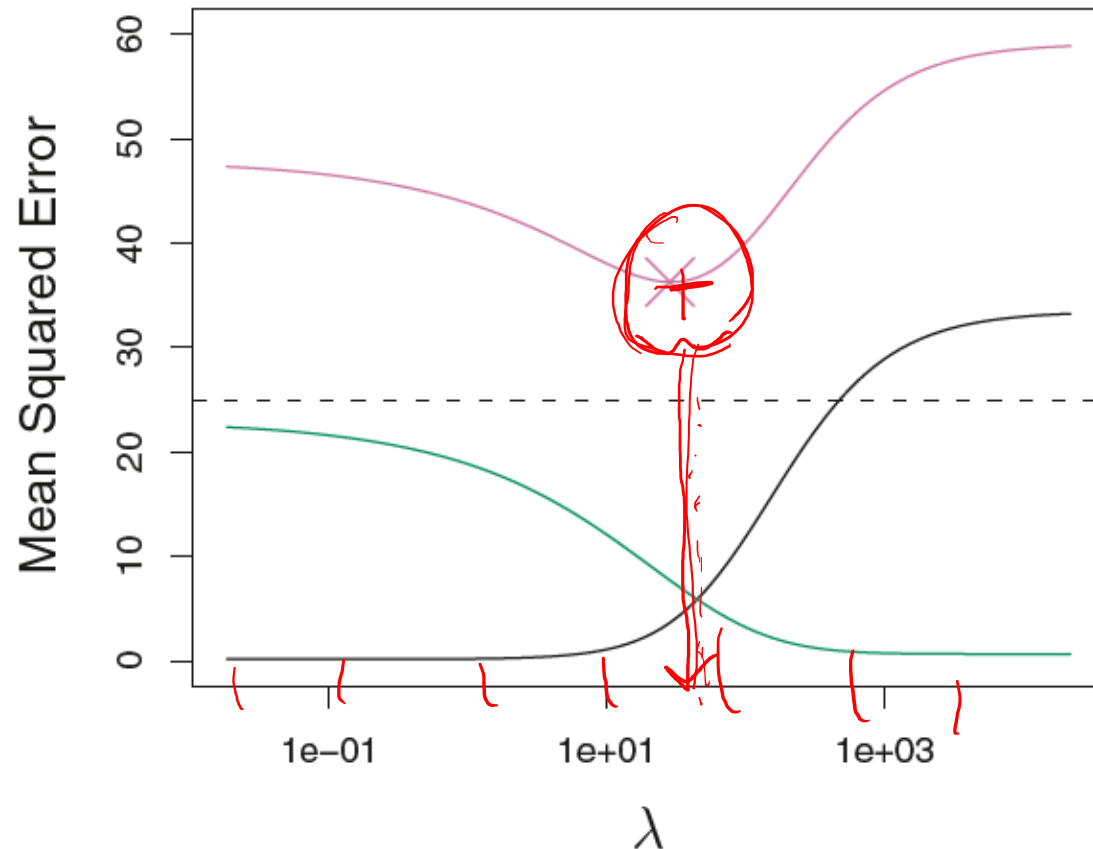


$\lambda$  vs.  $|\beta|_2$

$\lambda$  vs.  $\beta_j$

# L2 regularization

What does the lambda ( $\lambda$ ) do?



$\lambda$  vs. MSE ( $L_0$ )

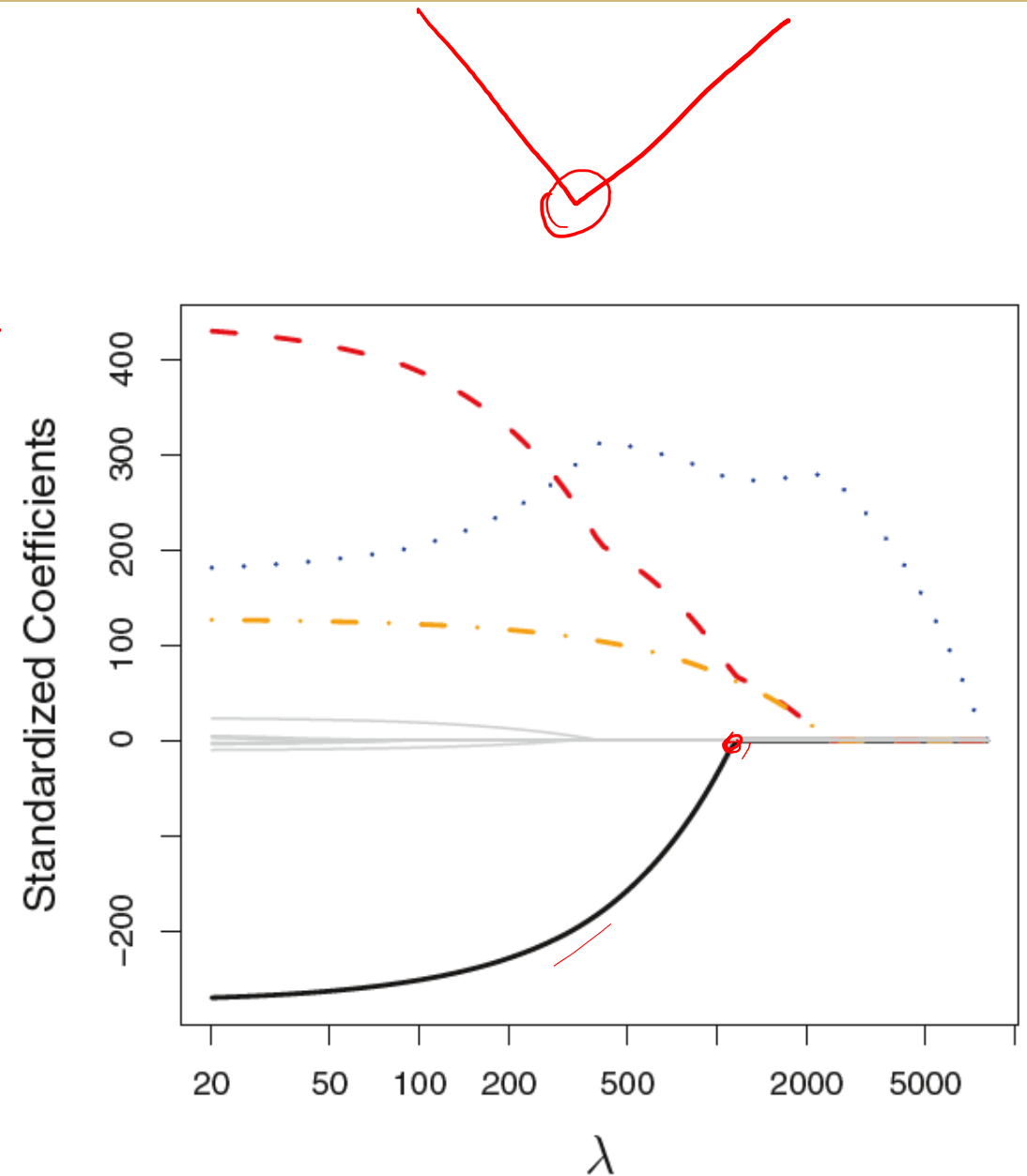
$\lambda$  vs. bias and variance

# L1 regularization (Lasso)

$$\mathcal{L} = \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

What does the lambda ( $\lambda$ ) do?

Lasso can make certain  $\beta$  0. Why?

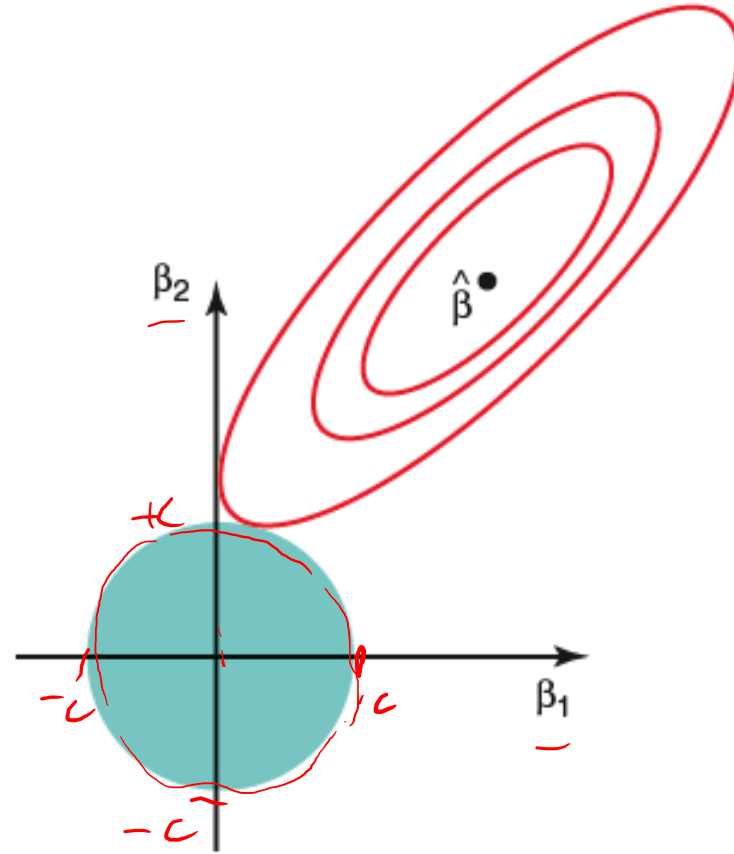
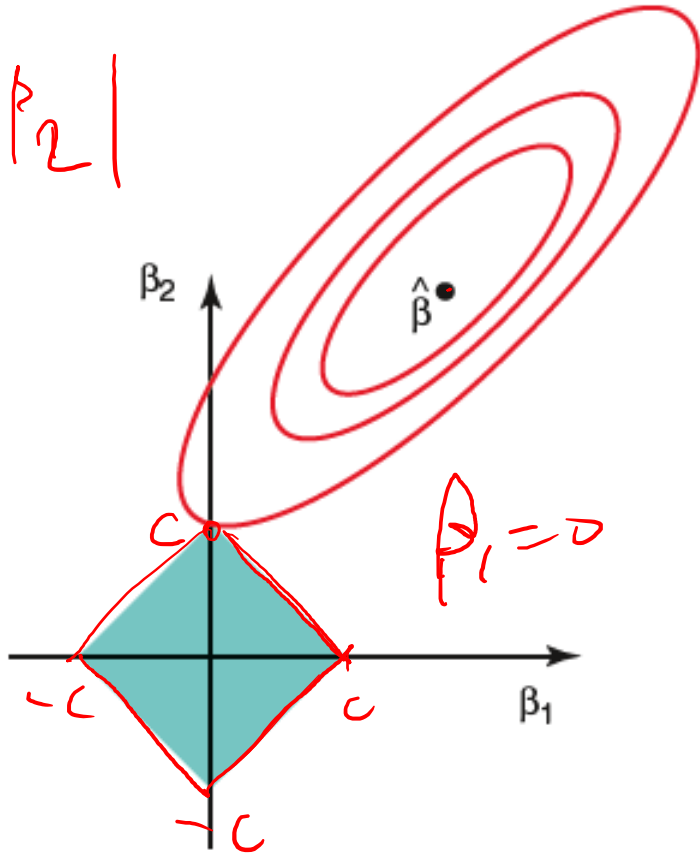


# Ridge and Lasso

$$\underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} + \lambda$$

$$\text{subject to } \sum_{j=1}^p \beta_j^2 \leq s$$

$$c = |\beta_1| + |\beta_2|$$



$$\beta_1^2 + \beta_2^2 = c$$

# Elastic Net

$$\mathcal{L} = \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ij} \right)^2 + \lambda \left( \underbrace{\alpha}_{\text{Lasso}} \sum_{j=1}^p |\beta_j| + \frac{1-\alpha}{2} \sum_{j=1}^p \beta_j^2 \right)$$

- Elastic Net is a convex combination of Ridge and Lasso
- Elastic Net > Ridge > Lasso

# What features to include?

## Method 1. Best subset method

- The idea: test all possible combinations
- Curse of dimensionality!



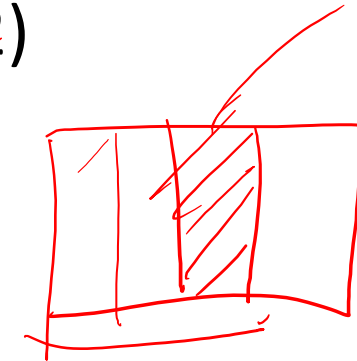
## Method 2. Regularization

- The idea: Penalize unnecessary complexity/features
- Hyperparameter lambda
- Ridge (L2), Lasso (L1), Elastic Net (L1+L2)

TIP: normalize the columns

$\sum W^2$

$\sum |w|$

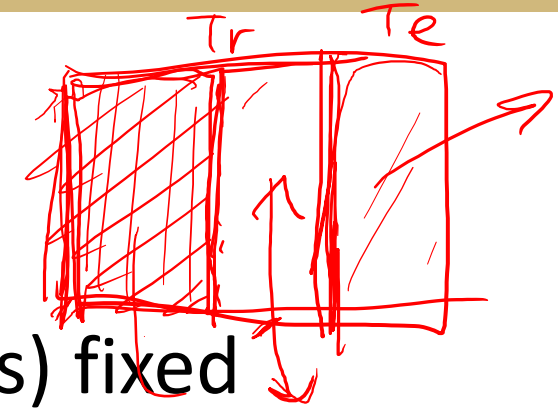


## Method 3. Cross-Validation

# Model validation during the training

The general idea:

- Split dataset into Train, Validation, Test
- Train using train data with a hyperparam(s) fixed
- Tune the hyperparameter(s) with validation
- When tuning is done, test with the test data



Validation

online(while training)

- How do I know my validation dataset was good or bad?



# Cross-Validation

*k-fold cross validation*

