

# Regularization & Cross Validation

Professor: Srikanth Krishnamurthy

Presented by: Chenlian Xu

Qianli Ma

Mar 3, 2018



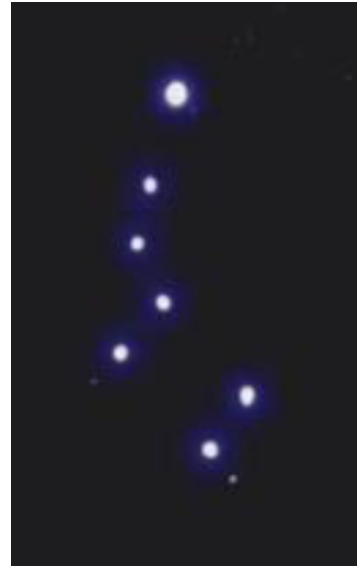
# Overfitting

*What is it?*

*How to prevent it?*

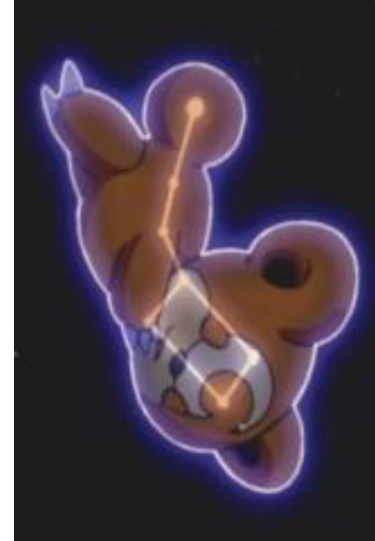


# Overfitting



Data

Normal  
Fitting



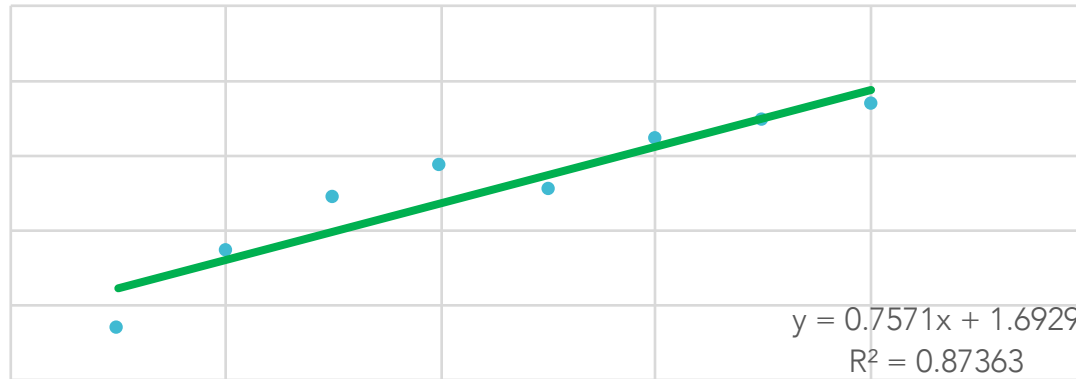
Serious Overfitting



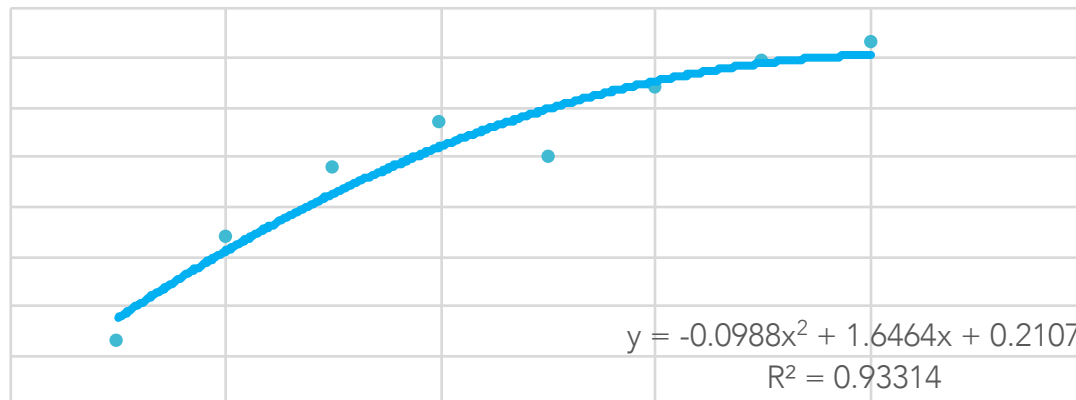
Overfitting



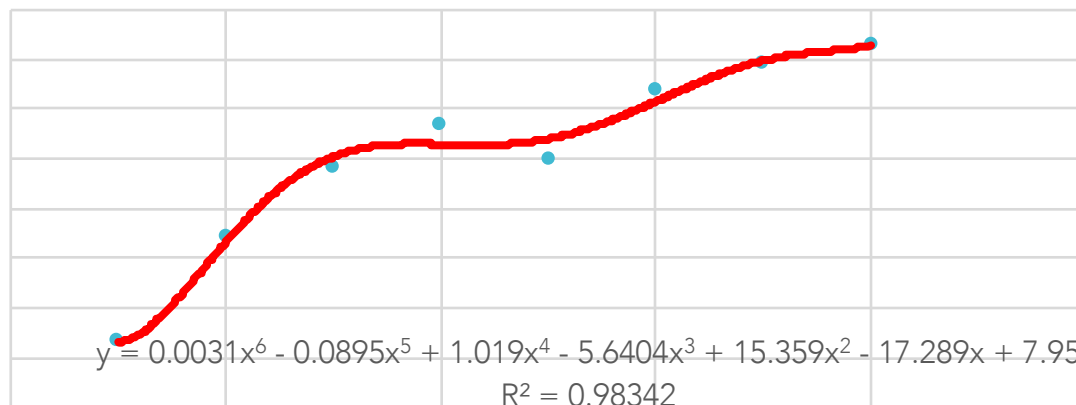
# Overfitting



Underfitting  
High Bias



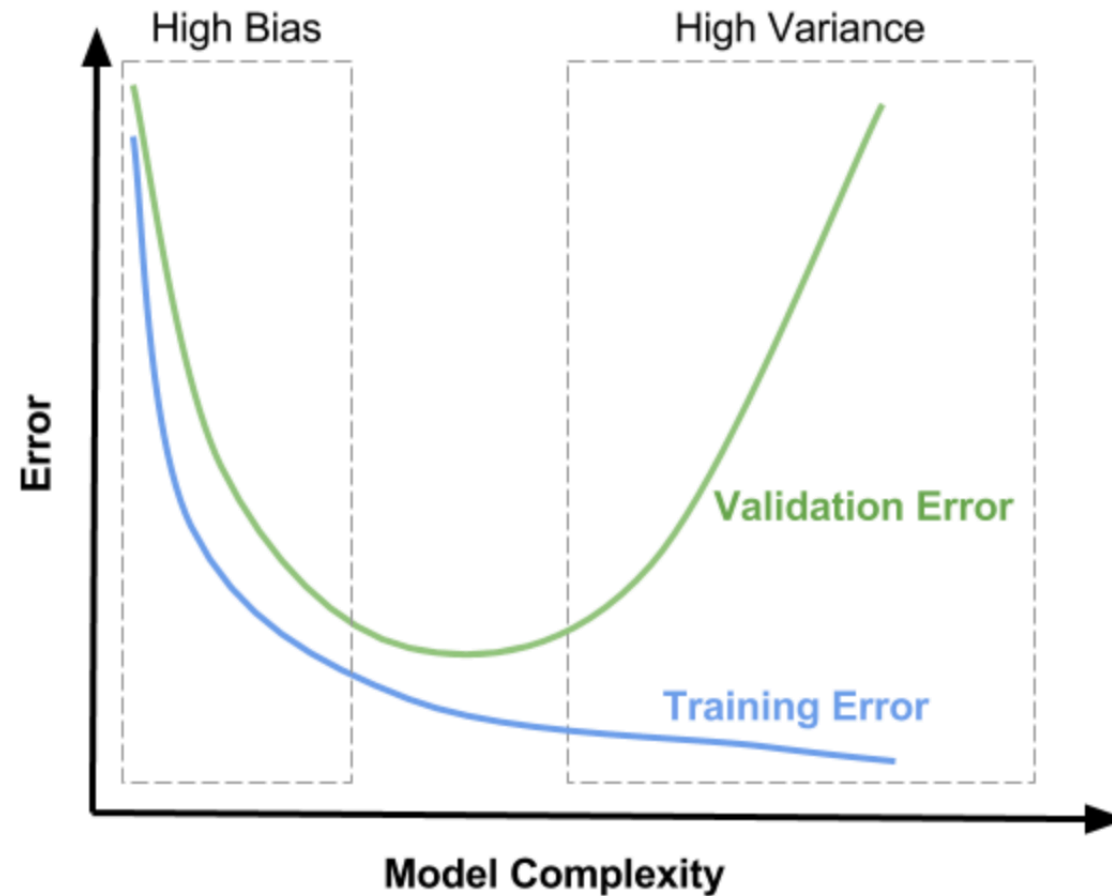
Normal Fitting  
Bias-Variance Tradeoff



Overfitting  
High Variance

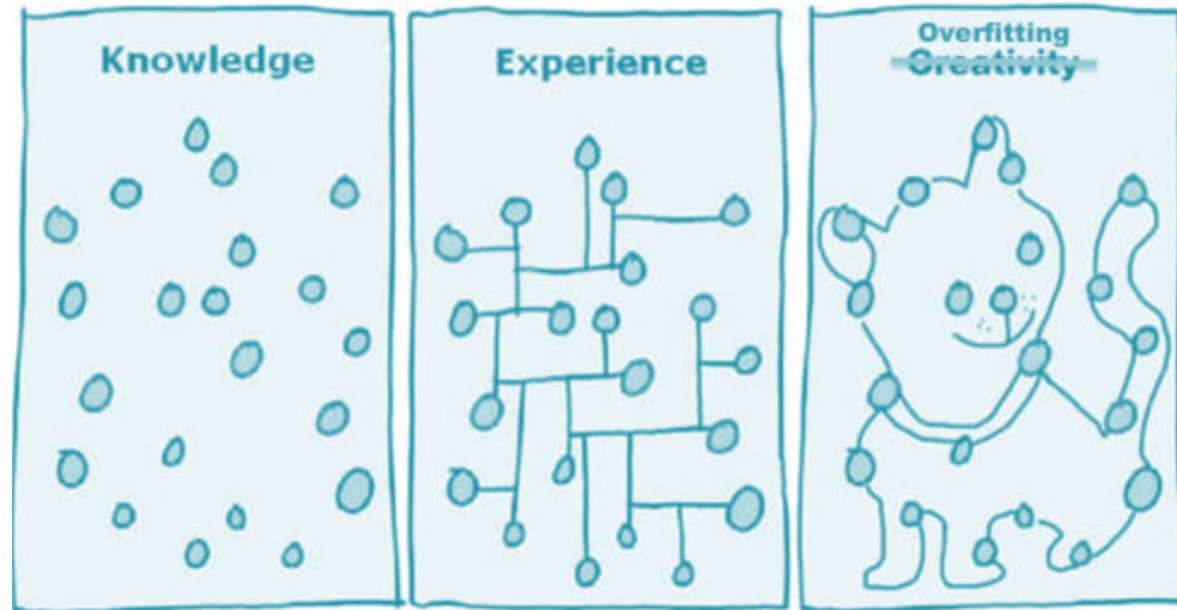
# Overfitting

- Bias-Variance Tradeoff



# Overfitting

- Too many features
- Fit the train set very well
- Fail to generalize to new examples



# How to prevent overfitting?

- 1 Reduce number of features
  - Manually select
  - Model selection algorithm
- 2 Regularization
  - sparsity
  - Reduce values of parameters



# Regularization

*L1, L2, Elastic Net*





# Regularization

- Reduces overfitting by adding a complexity **penalty** to the **loss function**

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left[ \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \textit{penalty} \right]$$

# L1 & L2 Regularization

- L1 LASSO (Least absolute shrinkage and selection operator)

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left[ \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^n |\theta_j| \right]$$


- L2 Ridge

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left[ \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

# Elastic Net

- Elastic Net
- linearly combines the L1 and L2 penalties of the lasso and ridge methods

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left[ \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^n |\theta_j| + \lambda \sum_{j=1}^n \theta_j^2 \right]$$



L1                  L2

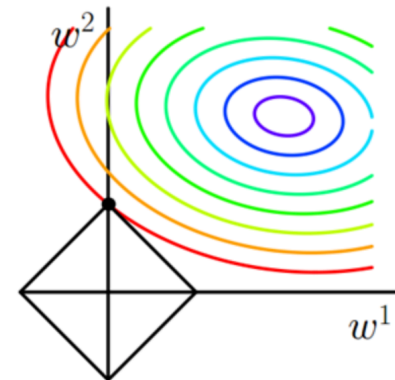
# Effect of L1 Regularization

- L1 LASSO (Least absolute shrinkage and selection operator)

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left[ \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^n |\theta_j| \right]$$

L1 Regularization encourages **sparsity**

$$\theta \rightarrow \theta' = \theta - \eta \lambda \cdot \text{sgn}(\theta) - \eta \frac{\partial C_0}{\partial \theta}$$



## Effect of L2 Regularization

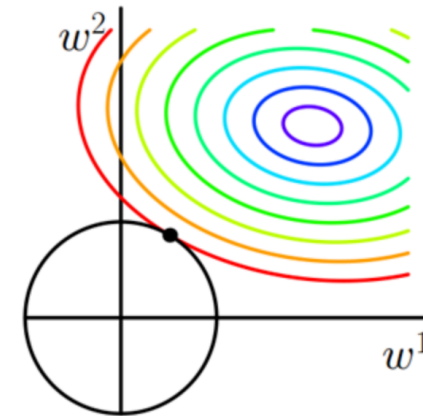
- L2 Ridge

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left[ \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

As  $\lambda$  increases, sum of squares decreases

Weight decay

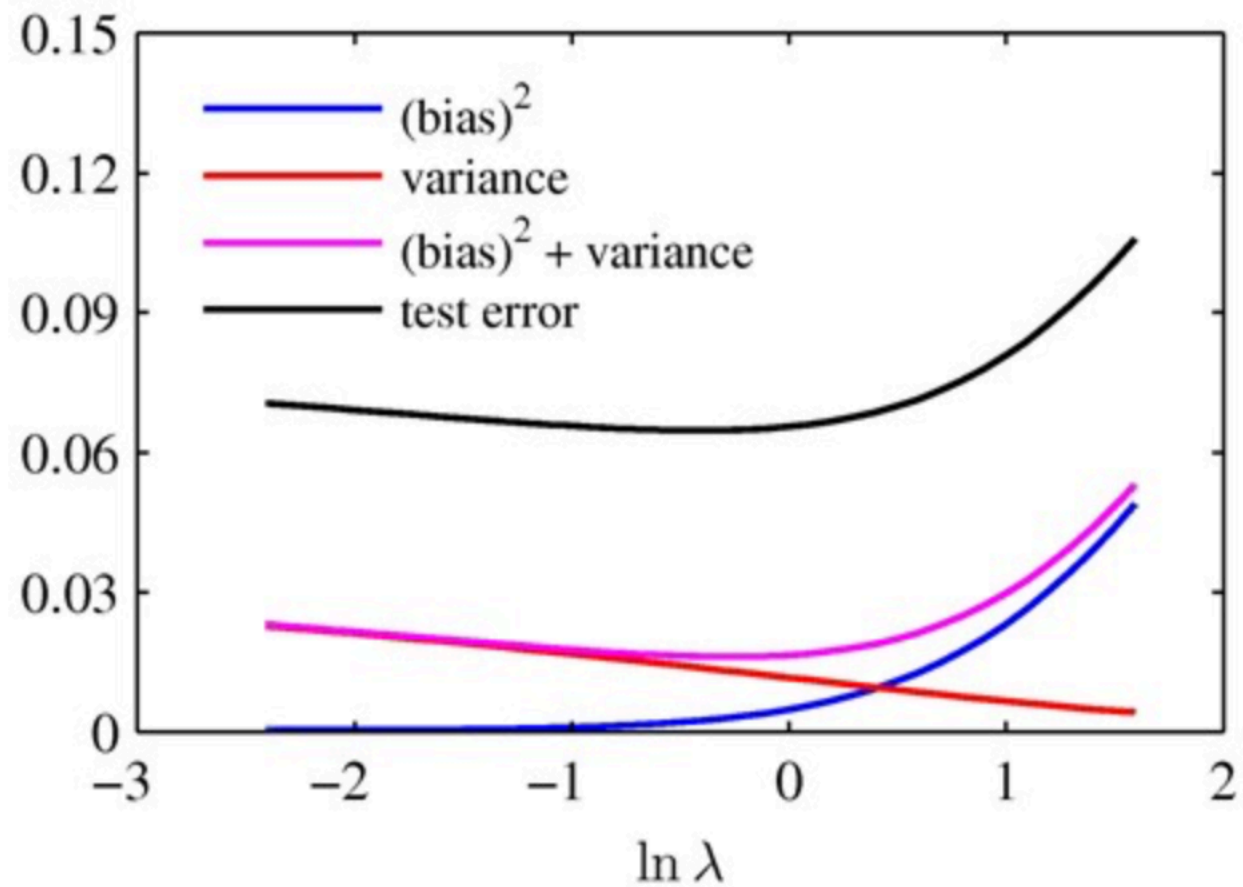
$$\theta \rightarrow \theta' = (1 - \eta\lambda)\theta - \eta \frac{\partial C_0}{\partial \theta}$$



# Regularization

- Reduces overfitting
- Reduces variance
- Minimizes the test-set error
- Minimizes the  $R_{in}$

# Regularization



## Only Regularization?

- “Optimal”  $\lambda$  ?
- We need a **validation** set

$\lambda_1, \lambda_2, \lambda_3 \dots$



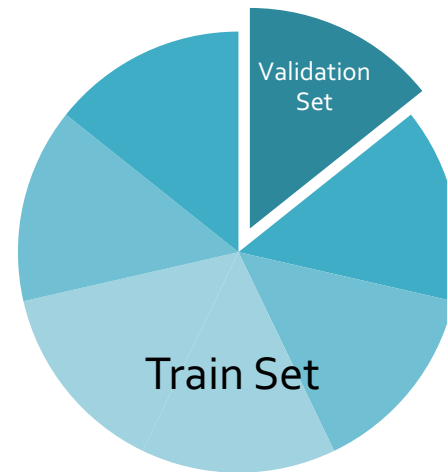
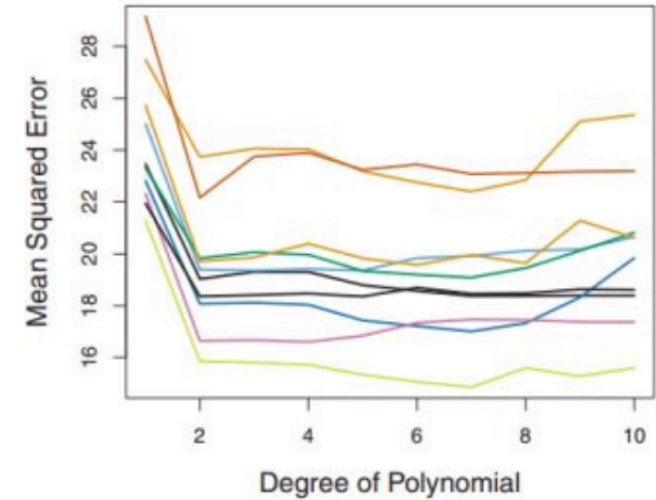
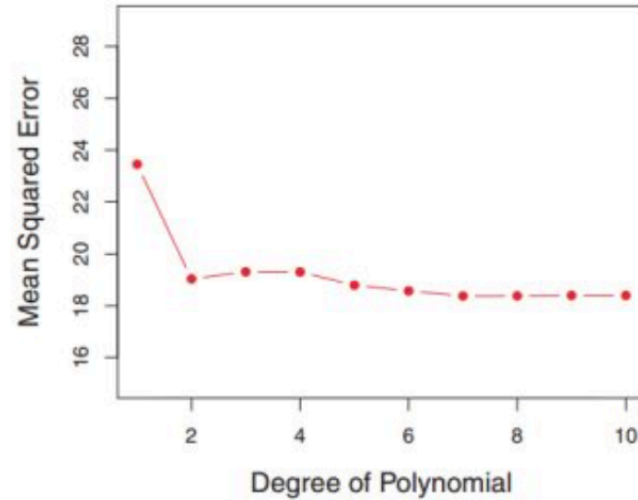
learn, test



$E_1, E_2, E_3 \dots$



# Why do we need Cross Validation?



You don't even use me for modeling!

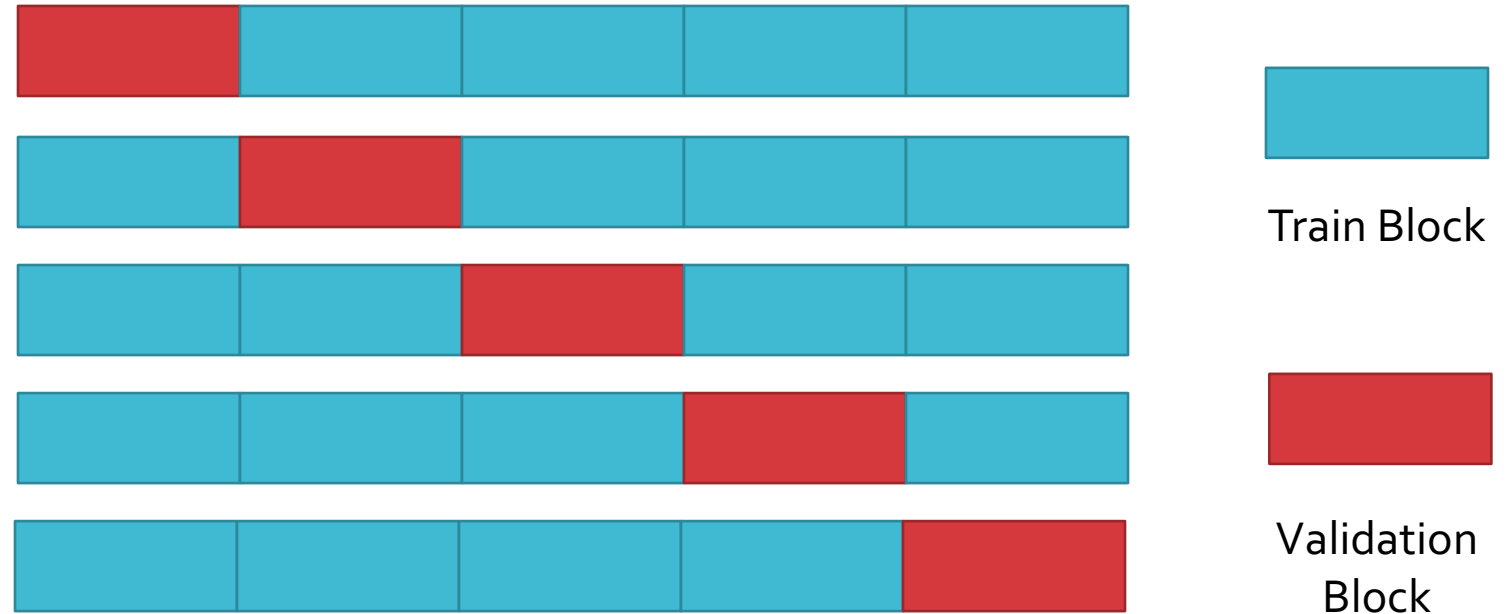
small validation set  $\Rightarrow$  large error in estimated loss  
large validation set  $\Rightarrow$  small training set  $\Rightarrow$  bad model

# Cross Validation

*Estimate the “optimal”  $\lambda$  by using it*

# K-fold cross validation

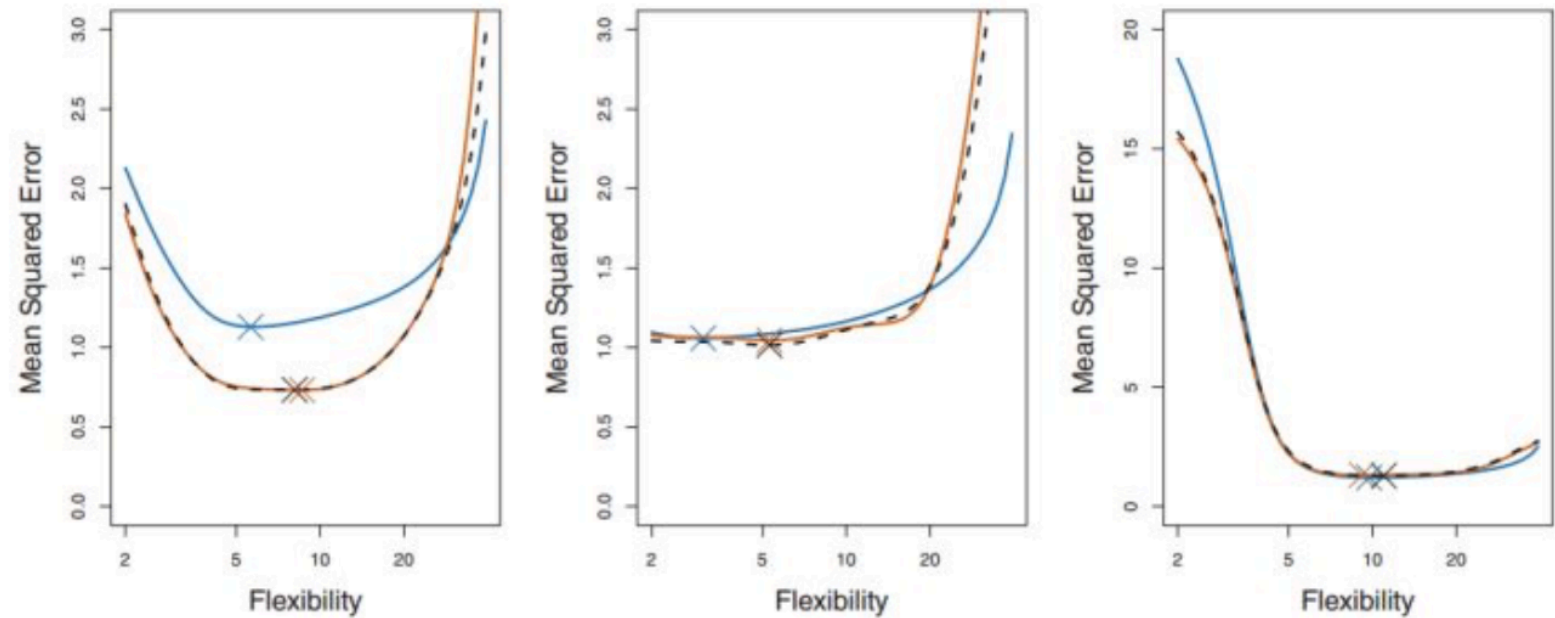
Divide the dataset into  $k$  blocks  
for  $k = 1$  to  $k$   
train on blocks except  $k$ th block, test on  $k$ th block  
average the results, choose best  $\lambda$ .



# K-fold cross validation

Common cases:  $K = 5$ ,  $10$  or  $K = N$  (LOOCV)

High computation cost:  $K$  folds  $\times$  many choices of model or  $\lambda$



----- LOOCV

———— 10-fold CV



# Summary

*Regularization  
&  
Cross Validation*

# Regularization & Cross Validation

- Trading off bias and variance is hard.
  - Degree of Polynomial  $\nearrow$  Bias  $\searrow$  Variance  $\nearrow$
- Regularization penalizes hypothesis complexity
  - L2 regularization leads to small weights
  - L1 regularization leads to many zero weights (sparsity)
- Cross-validation enables selection of regularization penalties by estimating test-set error on parts of the training set



# Demo

*Regularization  
&  
Cross Validation*

