# Regularization & Cross Validation

Professor: Srikanth Krishnamurthy

Presented by: Chenlian Xu

Qianli Ma

Mar 3, 2018
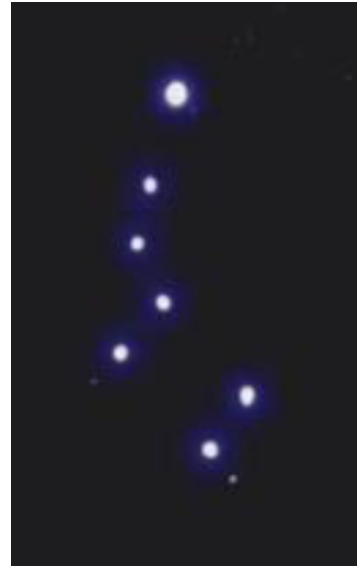
# Overfitting

*What is it?*
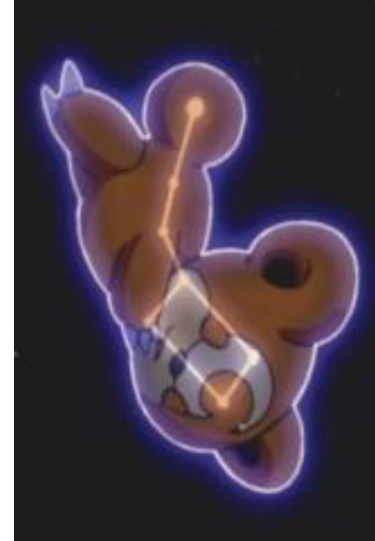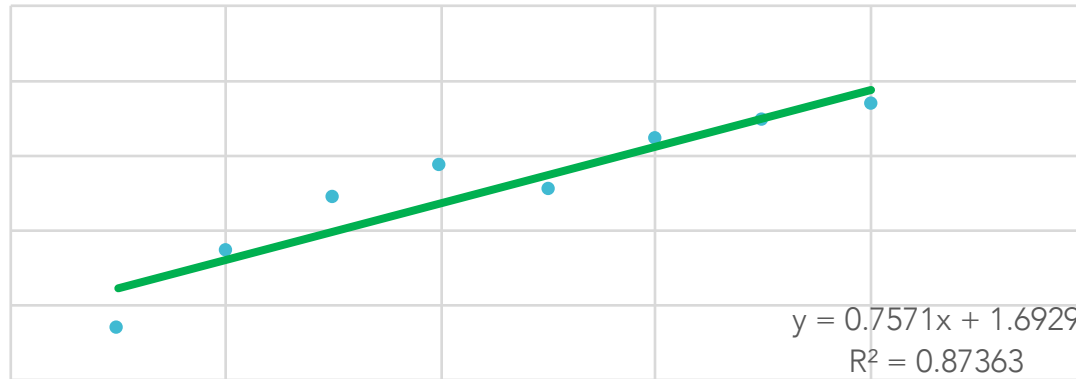
*How to prevent it?*

# Overfitting



Data

Normal Fitting

Overfitting

Serious Overfitting

# Overfitting



Underfitting
High Bias

$$y = 0.7571x + 1.6929$$
$$R^2 = 0.87363$$

Normal Fitting
Bias-Variance Tradeoff

$$y = -0.0988x^2 + 1.6464x + 0.2107$$
$$R^2 = 0.93314$$

Overfitting
High Variance

$$y = 0.0031x^6 - 0.0895x^5 + 1.019x^4 - 5.6404x^3 + 15.359x^2 - 17.289x + 7.95$$
$$R^2 = 0.98342$$
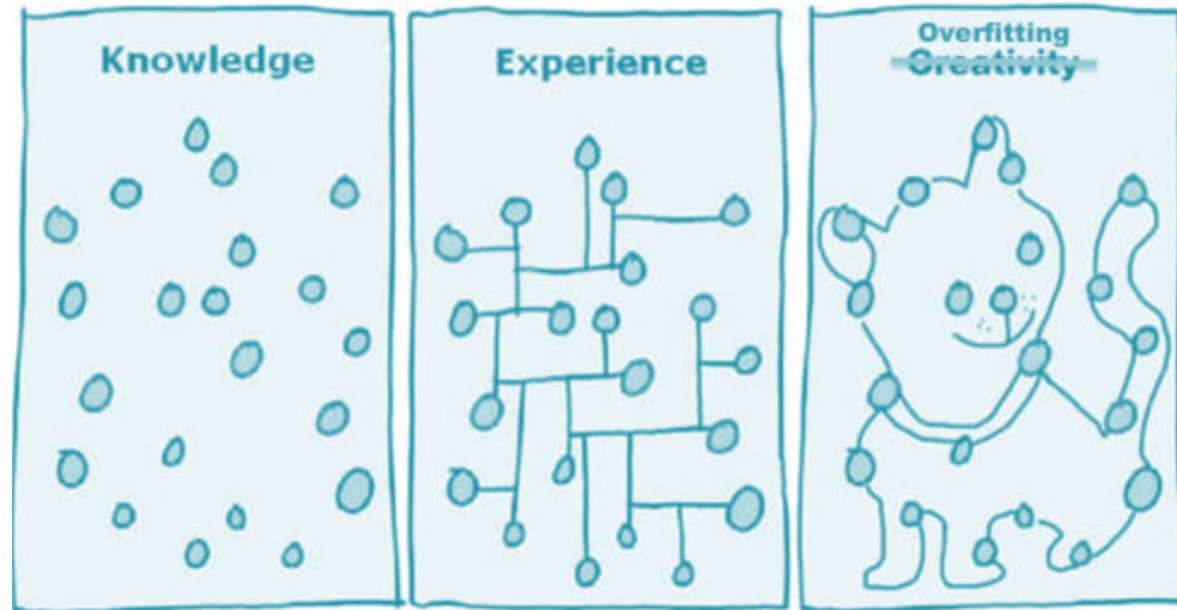
# Overfitting

- 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

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

## L1 & L2 Regularization

- L1 LASSO (Least absolute shrinkage and selection operator)

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

- L2 Ridge

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

# Elastic Net

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

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} \left[ \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^{n} |\theta_j| + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$
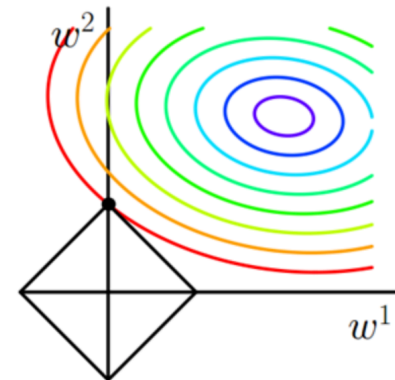
L1          L2

# Effect of L1 Regularization

- L1 LASSO (Least absolute shrinkage and selection operator)

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

L1 Regularization encourages sparsity

$$\theta \,\text{-->}\, \theta' = \theta - \eta\lambda \cdot sgn(\theta) - \eta \frac{\partial C0}{\partial \theta}$$
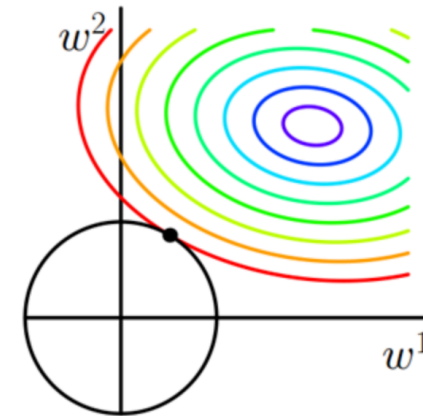
# Effect of L2 Regularization

- L2 Ridge

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

As λ increases, sum of squares decreases
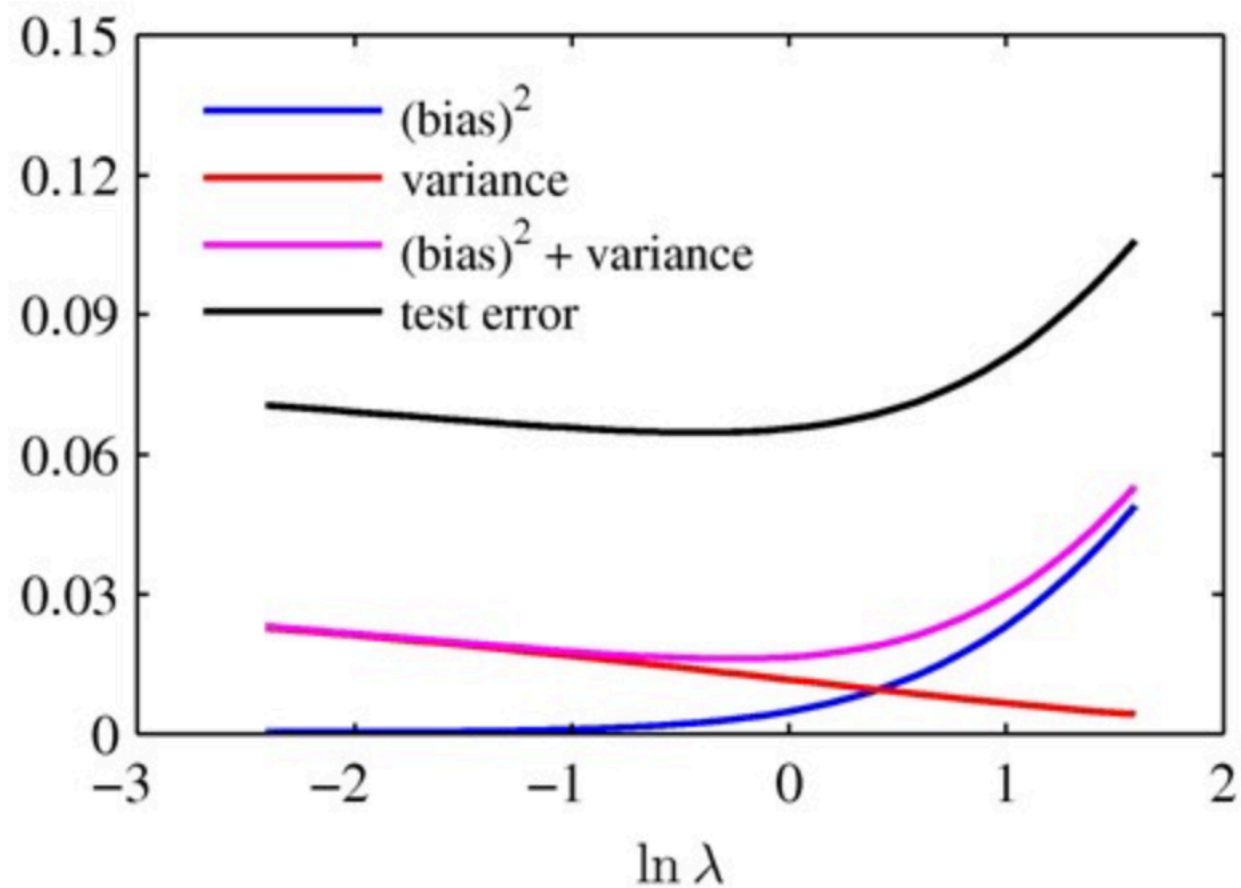
Weight decay

$$\theta \dashrightarrow \theta' = (1-\eta\lambda)\theta - \eta\frac{\partial C0}{\partial\theta}$$

# Regularization

- Reduces overfitting
- Reduces variance
- Minimizes the test-set error

- Minimizes the $R_{in}$

# Regularization

## Only Regularization?

- "Optimal" $\lambda$ ?
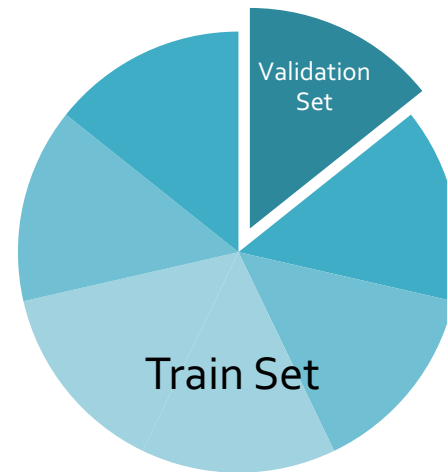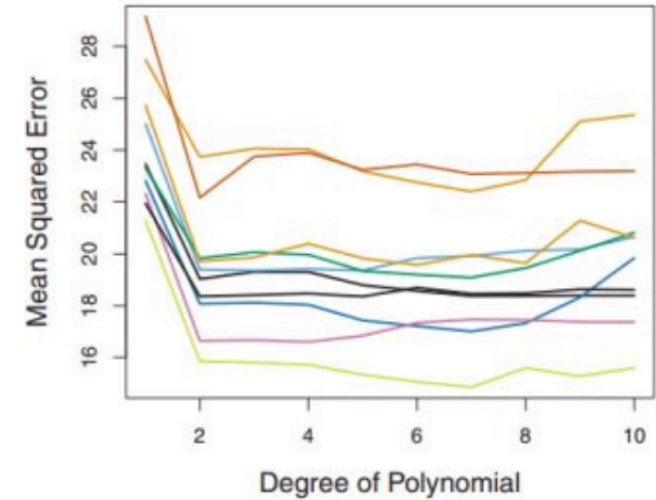- We need a **validation** set
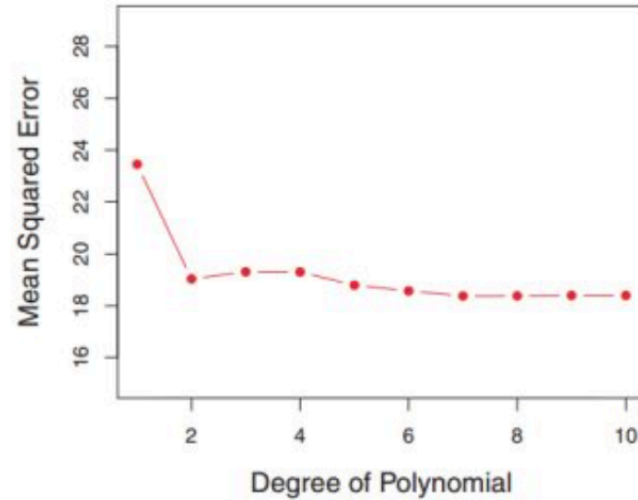
$\lambda_1, \lambda_2, \lambda_3\ldots$

$\downarrow$

learn, test

$\downarrow$

$R_{in1}, R_{in2}, R_{in3}\ldots$
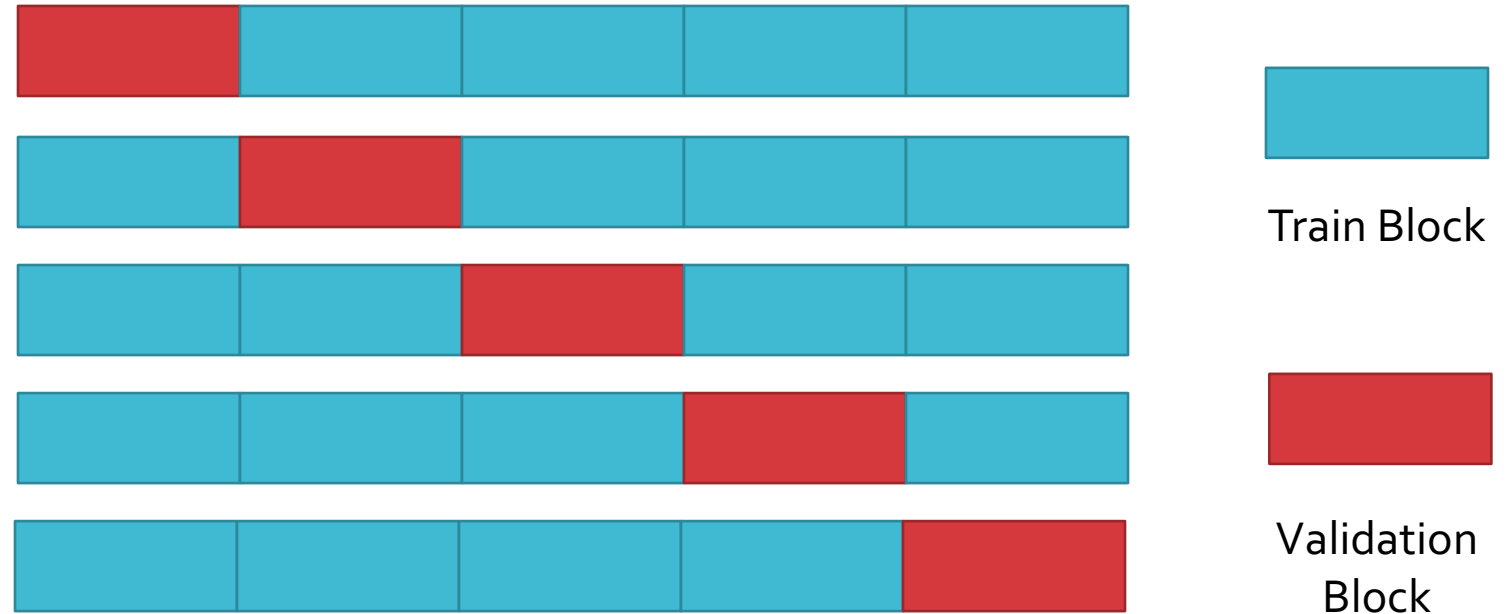
# Why do we need Cross Validation?

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

# Cross Validation
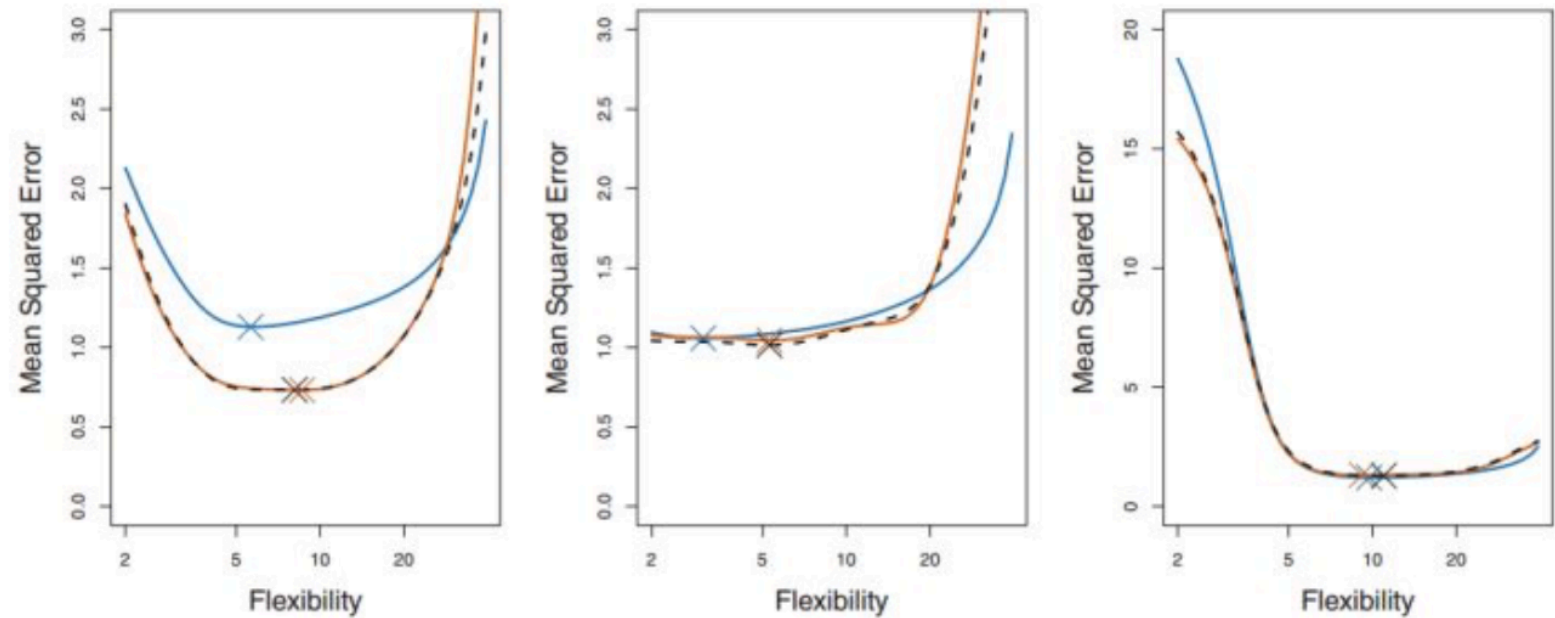
*Estimate the "optimal" λ by using it*

# K-fold cross validation

Divide the dataset into k blocks
for k = 1 to k
train on blocks except kth block, test on kth block
average the results, choose best λ.



Train Block

Validation Block

# K-fold cross validation

Common cases: K = 5, 10 or K = N (LOOCV)
High computation cost: K folds × many choices of model or λ



- - - - -  LOOCV          ——————  10-fold CV

# Summary

*Regularization*

*&*

*Cross Validation*

# Regularization & Cross Validation

- Trading off bias and variance is hard.
  - Degree of Polynomial ↗ Bias ↘ Variance ↗

- 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*