# DS 4400

# Machine Learning and Data Mining I
# Spring 2024

David Liu

Khoury College of Computer Science

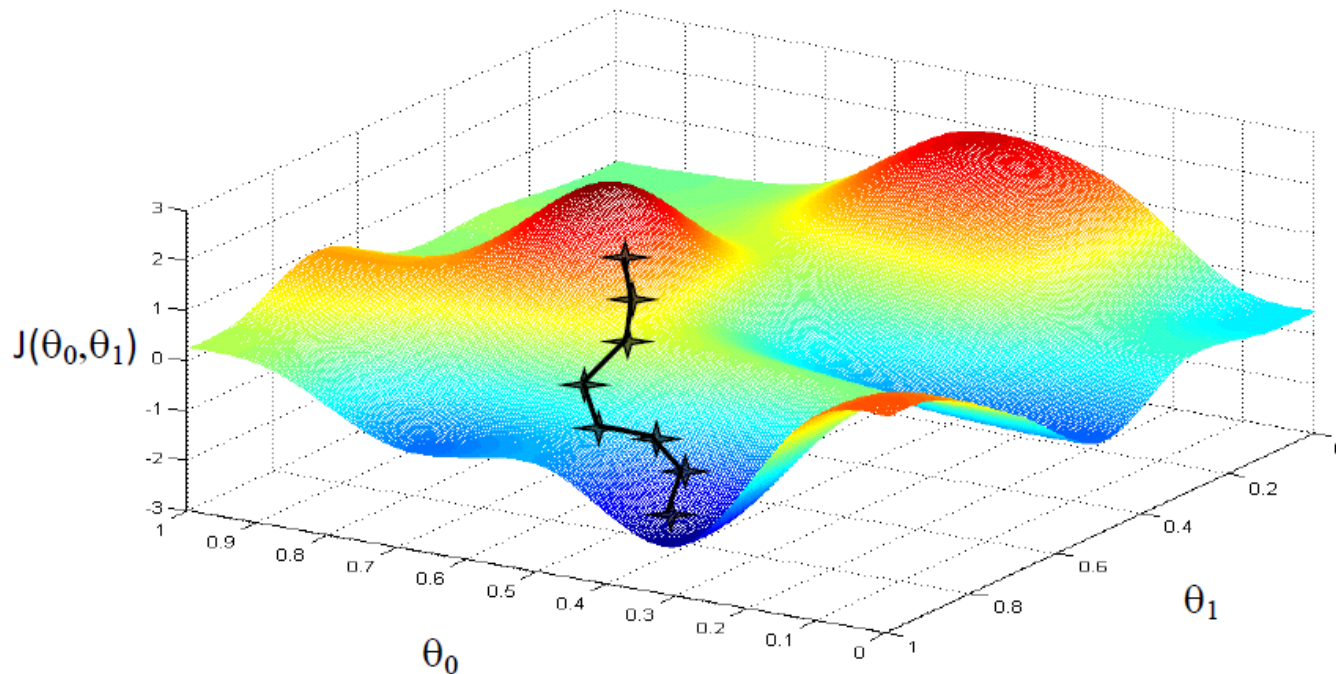Northeastern University

January 30 2024

# Announcements

- Will start grading HW 1 after tonight's deadline

- Will release HW 2 this week

- Midterm exam on Friday Feb 23

- Will release further final project guidance this week.

# Outline

- Review of Gradient Descent
- Non-linear regression
  - Polynomial regression
  - Cubic, spline regression
- Regularization
  - Ridge regression
  - Lasso regression
- Classification
  - K Nearest Neighbors (kNN)
  - Bias-Variance tradeoff

# How to optimize $J(\theta)$?

- Choose initial value for $\theta$
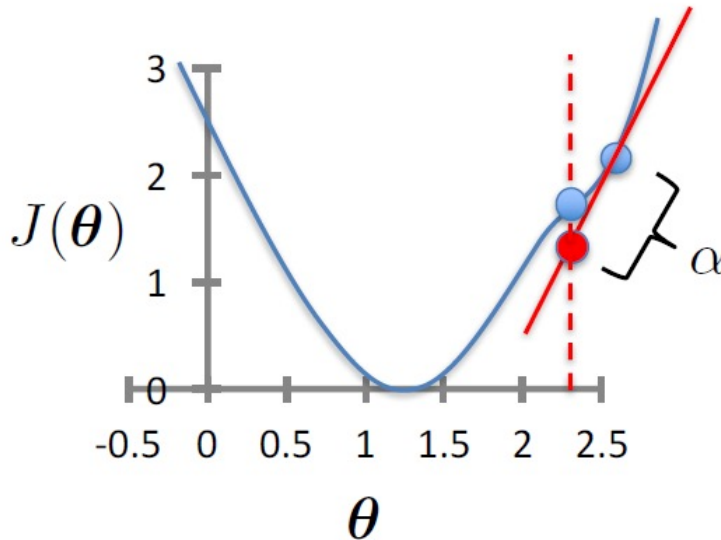- Until we reach a minimum:
  - Choose a new value for $\theta$ to reduce $J(\theta)$

# Gradient Descent

- Initialize $\boldsymbol{\theta}$
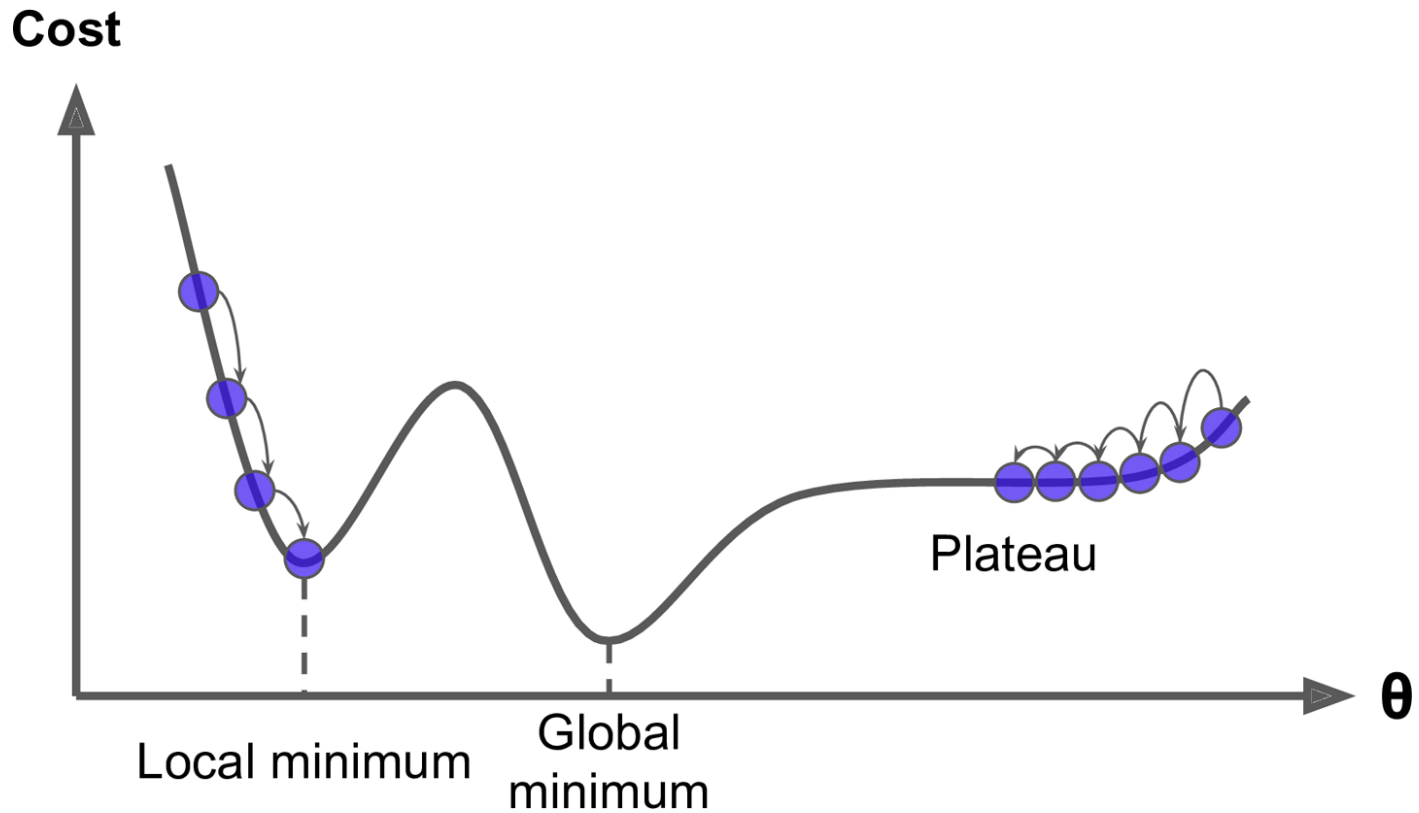- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$
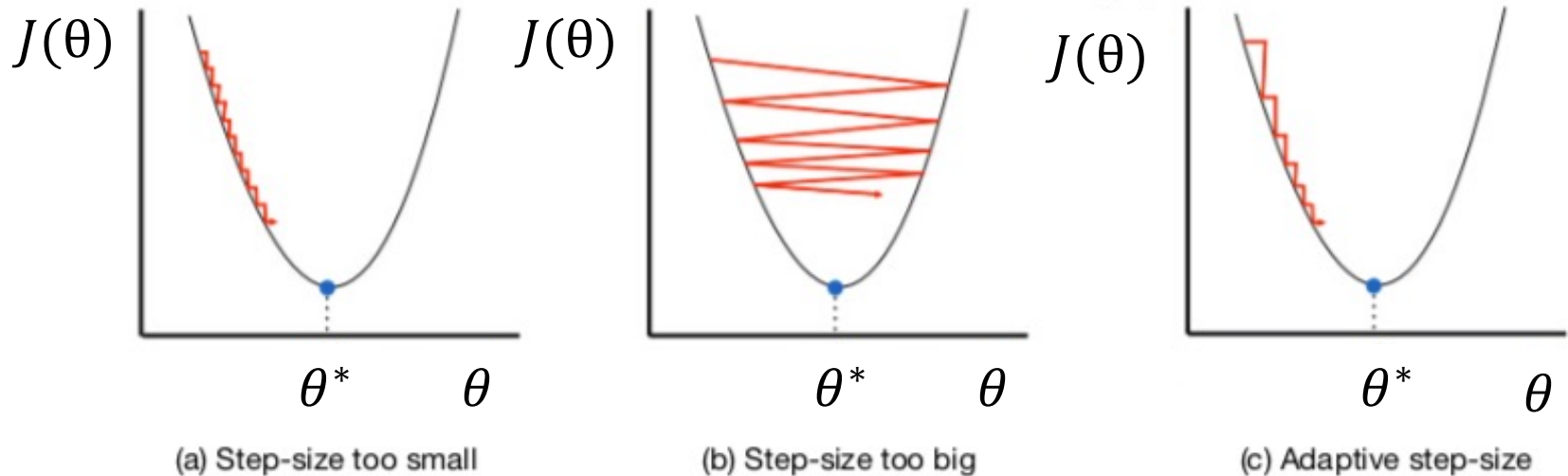
simultaneous update
for j = 0 ... d

learning rate (small)
e.g., α = 0.05

# GD Convergence Issues

# Adaptive step size

$J(\theta)$

$J(\theta)$

$J(\theta)$

(a) Step-size too small

(b) Step-size too big

(c) Adaptive step-size

$\theta^*$  $\theta$

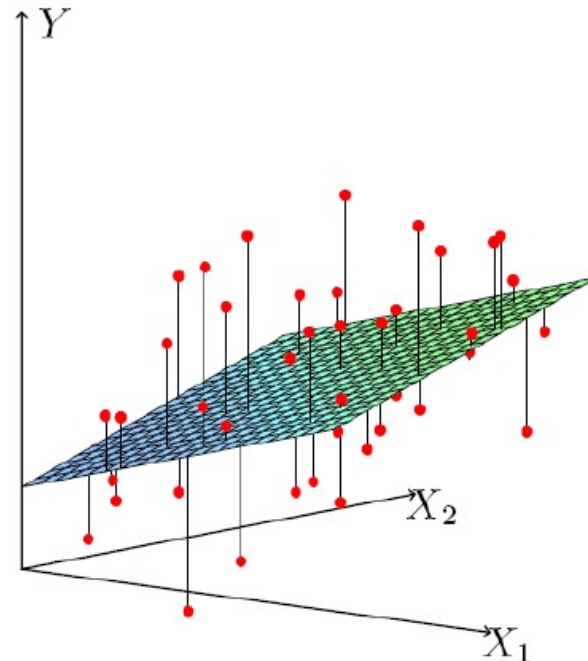$\theta^*$  $\theta$

$\theta^*$  $\theta$

- Start with large step size and reduce over time, adaptively
- Line search method
- Measure how objective decreases

# NON-LINEAR REGRESSION

# Multiple Linear Regression

- Dataset: $x_i \in R^d, y_i \in R$
- Hypothesis $h_\theta(x) = \theta^T x$
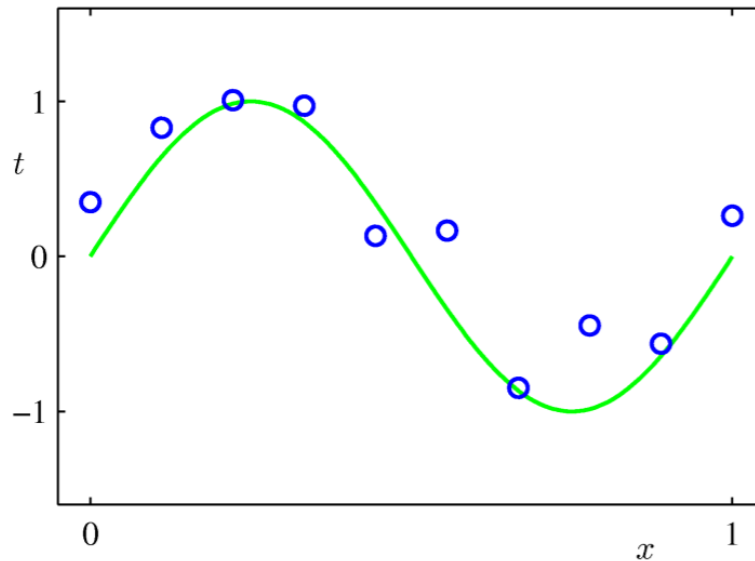- MSE $= \frac{1}{N} \sum (\theta^T x_i - y_i)^2$  <span style="color:red">Loss / cost</span>

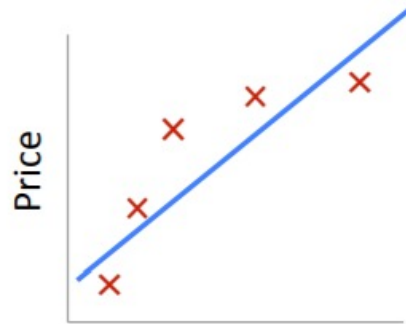$$\theta = (X^\intercal X)^{-1} X^\intercal y$$

# Polynomial Regression

- Polynomial function on single feature
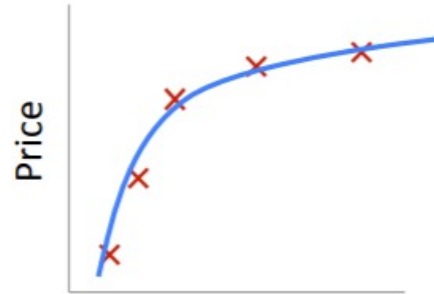
$$- h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_p x^p$$

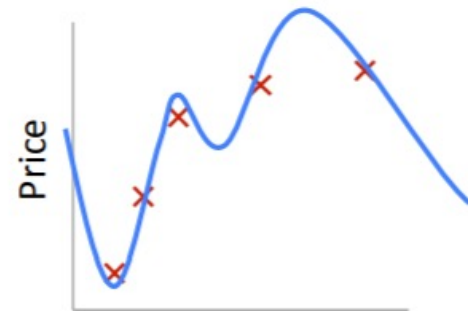# Polynomial Regression



$$\theta_0 + \theta_1 x$$
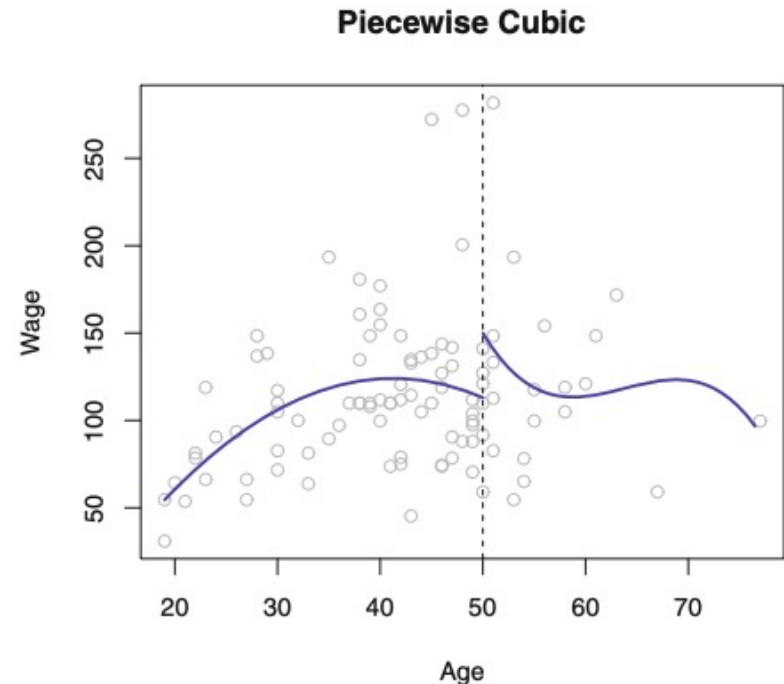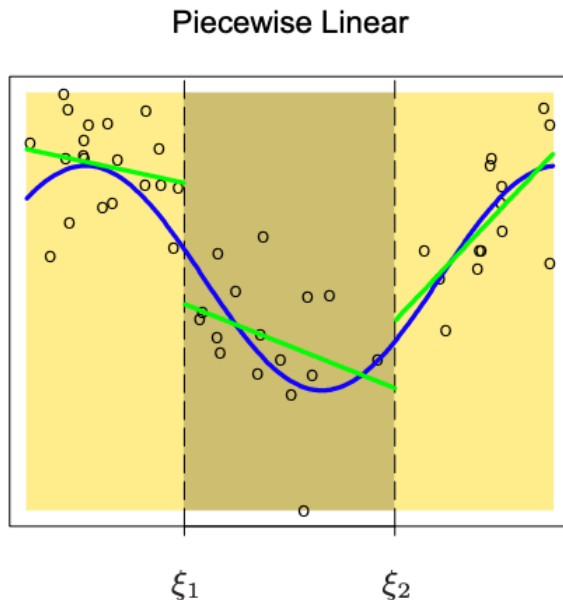
$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$
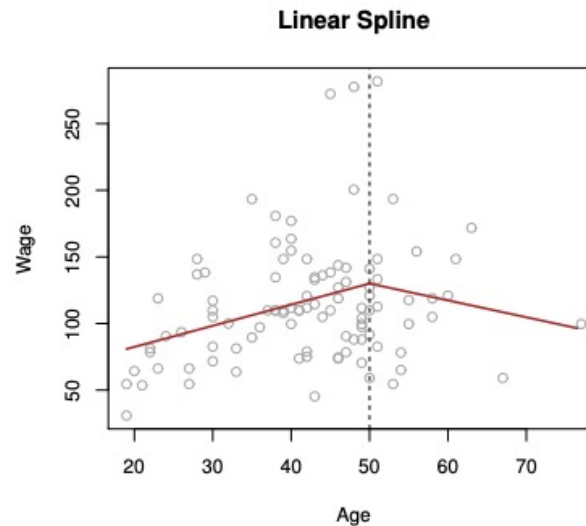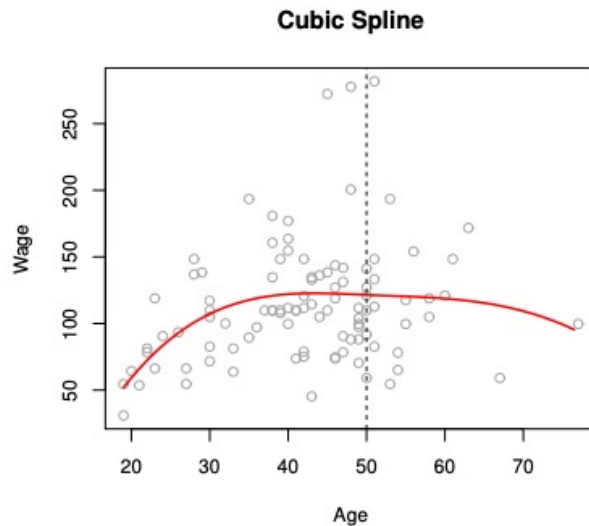
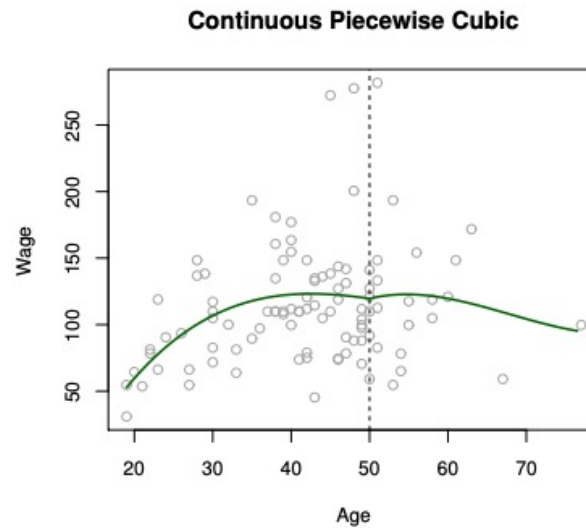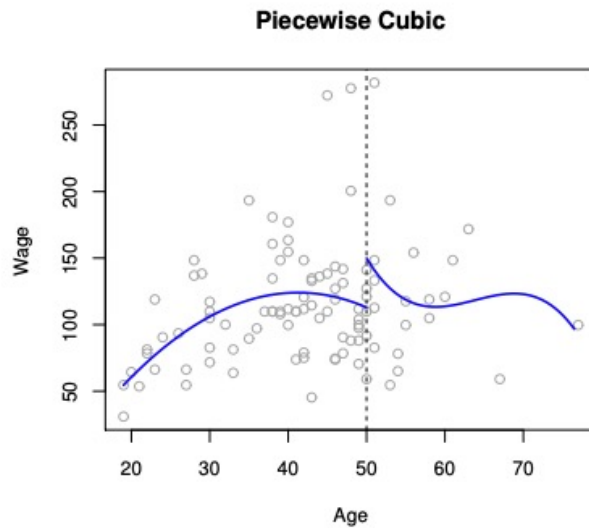# Polynomial Regression Training

- Simple Linear Regression
- $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_p x^p$
- How to train model?

# Piecewise Polynomial

- Divide the space into regions
- Polynomial regression on each region
  - Linear piecewise (degree 1), quadratic piecewise (degree 2), cubic piecewise (degree 3)



Piecewise Linear



Piecewise Cubic

# Piecewise and spline regression

# Piecewise polynomial vs Regression spline



**Piecewise Cubic**

**Cubic Spline**

1 break at Age $= 50$

1 knot at Age $= 50$

**Definition: Cubic spline**

A cubic spline with knots at $x$-values $\xi_1, \ldots, \xi_K$ is a continuous piecewise cubic polynomial with continuous derivates and continuous second derivatives at each knot.

# Cubic splines



**Cubic Spline**

- Turns out, cubic splines are sufficiently flexible to *consistently* estimate smooth regression functions $f$
- You can use higher-degree splines, but *there's no need to*
- To fit a cubic spline, we just need to pick the knots

# Additive Models

- Multiple Linear Regression Model

  $- y_i = \theta_0 + \theta_1 x_1 + \cdots + \theta_d x_d$

- Additive Models

  $- y_i = \theta_0 + \textcolor{red}{f_1}(x_1) + \cdots + \textcolor{red}{f_d}(x_d)$

- Can instantiate functions f with:

  - Linear functions:

  - Quadratic:

  - Cubic:

# Generalization in ML



Underfitting ⟷ Overfitting

Simple model                    Complex model
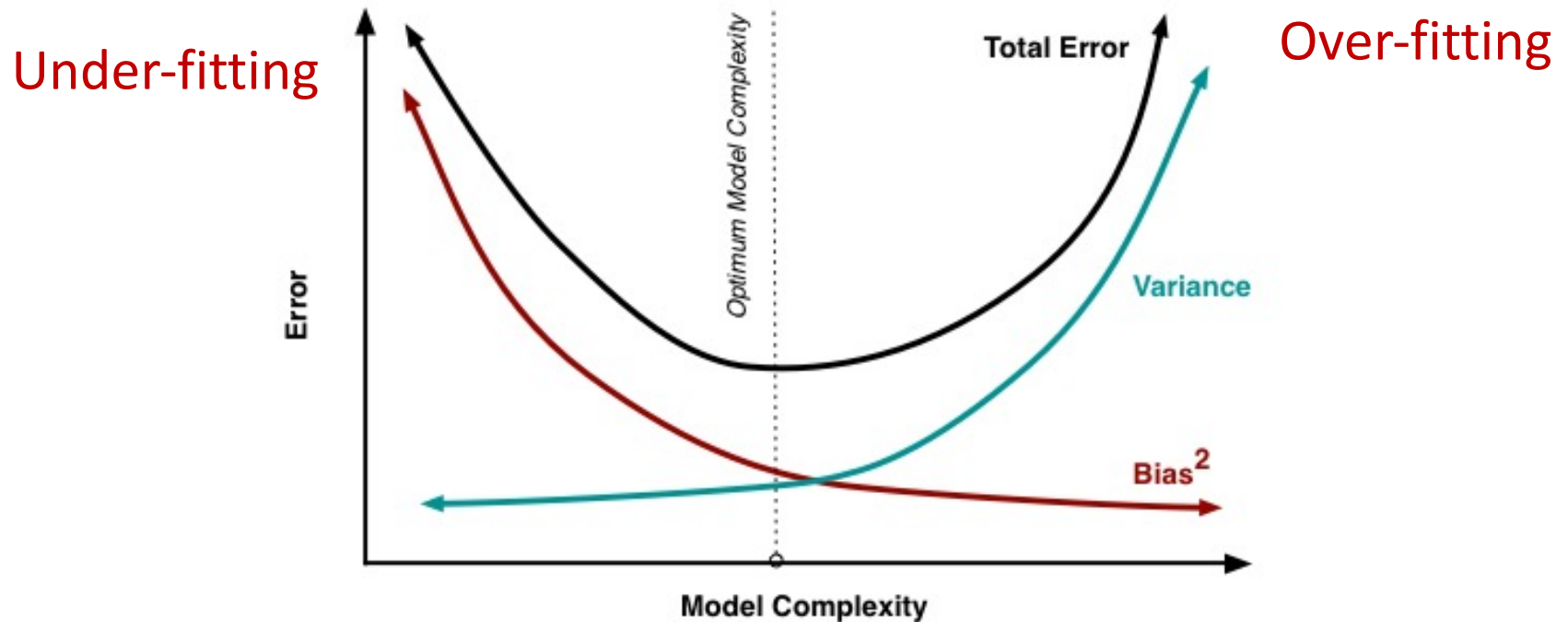
- Goal is to generalize well on new testing data
- Risk of overfitting to training data

# Bias-Variance Tradeoff



Under-fitting

Over-fitting

- Bias = Difference between estimated and true models
- Variance = Model difference on different training sets
  MSE is proportional to Bias + Variance

# REGULARIZATION
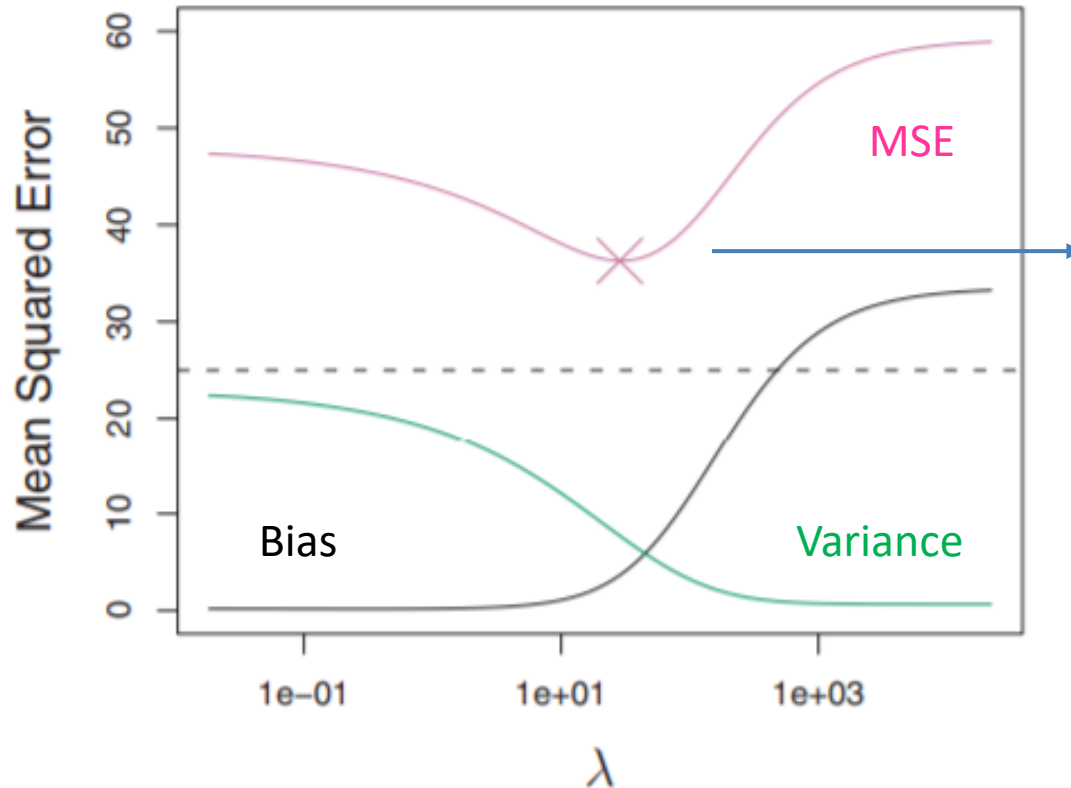
# Regularization

- A method for automatically controlling the complexity of the trained model
- Goals
  - Reduce model complexity
  - Reduce variance
  - Mitigate the bias-variance tradeoff
- Main techniques
  - Modify loss function to account for regularization term (Ridge, Lasso)
  - Perform feature selection and fit model on subset of features

# Ridge regression
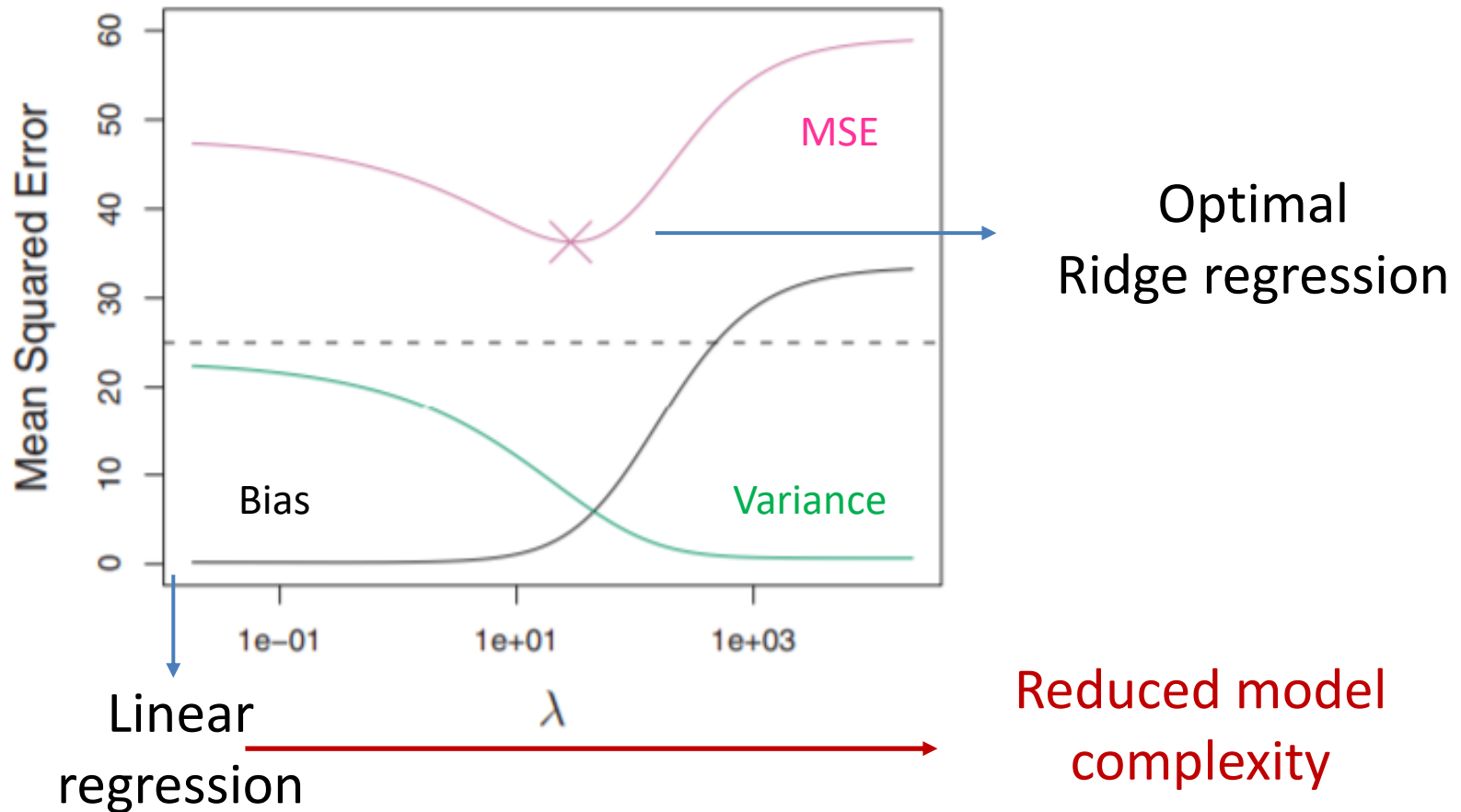
- Linear regression objective function

$$J(\theta) = \sum_{i=1}^{N} [h_\theta(x_i) - y_i]^2 + \lambda \sum_{j=1}^{d} \theta_j^2$$
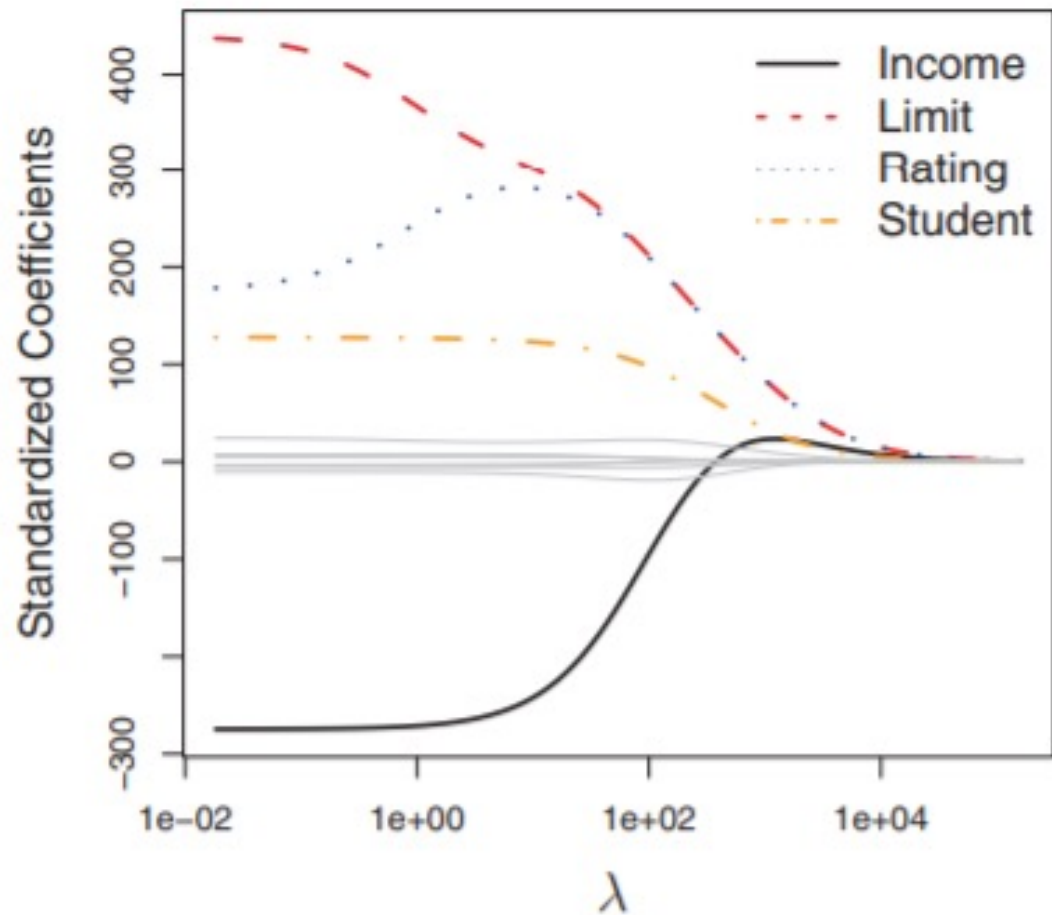
# Bias-Variance Tradeoff



Optimal
Ridge regression

# Bias-Variance Tradeoff

# Coefficient shrinkage



Predict credit card balance

# GD for Ridge Regression

Min Loss

$$J(\theta) = \sum_{i=1}^{N} [h_\theta(x_i) - y_i]^2 + \lambda \sum_{j=1}^{d} \theta_j^2$$

# GD for Ridge Regression

Min MSE

$$J(\theta) = \sum_{i=1}^{N}[h_\theta(x_i) - y_i]^2 + \lambda \sum_{j=1}^{d}\theta_i^2$$

Gradient update: $\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^{N}(h_\theta(x_i) - y_i)$

$$\theta_j \leftarrow \theta_j - \alpha \sum_{i=1}^{N}(h_\theta(x_i) - y_i)x_{ij} - \underbrace{\alpha\lambda\theta_j}_{}$$

Regularization

$$\theta_j \leftarrow \theta_j(1 - \alpha\lambda) - \alpha \sum_{i=1}^{N}(h_\theta(x_i) - y_i)x_{ij}$$

# Lasso Regression

$$J(\theta) = \sum_{i=1}^{N}[h_\theta(x_i) - y_i]^2 + \lambda \sum_{j=1}^{d} |\theta_j|$$

- L1 norm for regularization
- Results in sparse coefficients
- Small issue: gradients cannot be computed around 0
  - Can use sub-gradient at 0

# Lasso Regression

$$J(\theta) = \sum_{i=1}^{N} (h_\theta(x_i) - y_i)^2 + \lambda \sum_{j=1}^{d} |\theta_j|$$

Squared Residuals

Regularization

- L1 norm for regularization

- Results in sparse coefficients

- Issue: gradients cannot be computed around 0

- Method of sub-gradient optimization

# Alternative Formulations

- Ridge
  - L2 Regularization
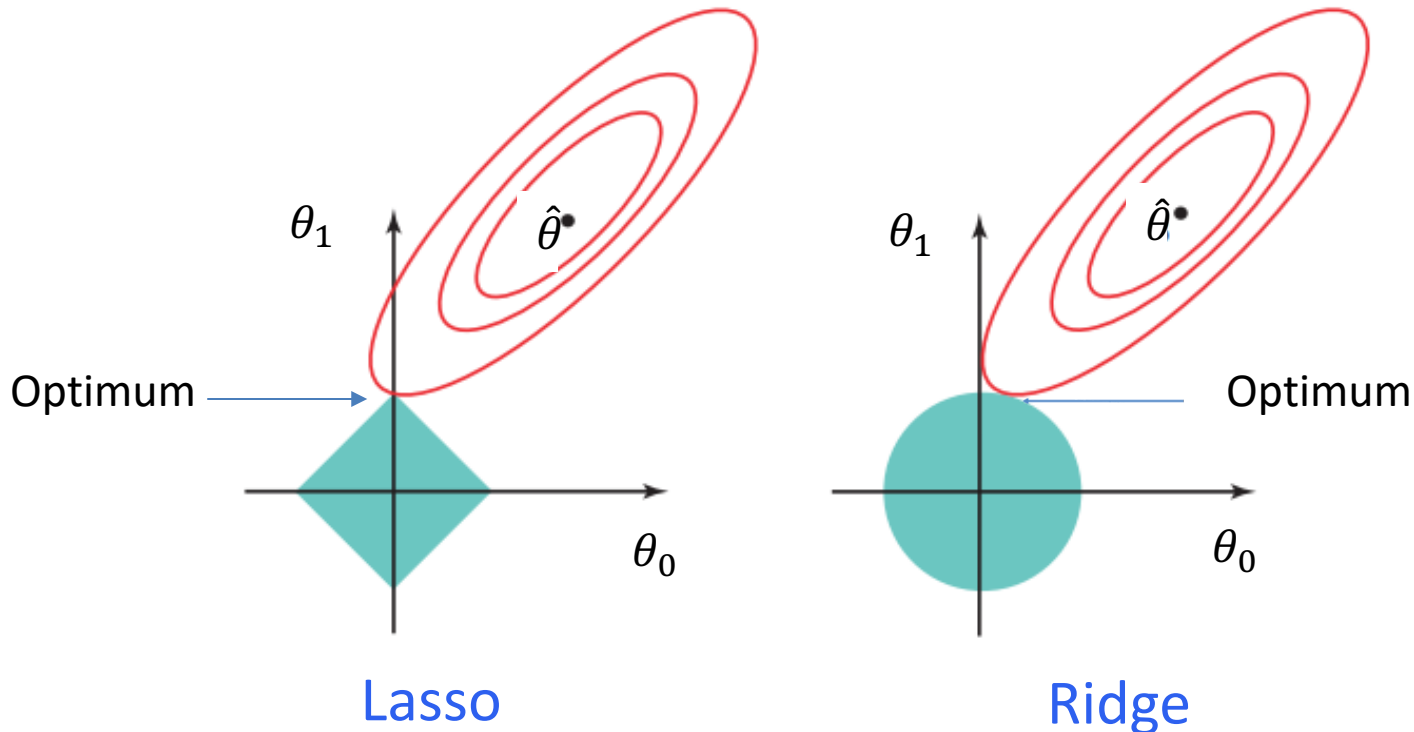  - $\min_\theta \sum_{i=1}^N [h_\theta(x_i) - y_i]^2$ subject to $\sum_{j=1}^d |\theta_j|^2 \le \epsilon$

- Lasso
  - L1 regularization
  - $\min_\theta \sum_{i=1}^N [h_\theta(x_i) - y_i]^2$ subject to $\sum_{j=1}^d |\theta_j| \le \epsilon$

# Lasso vs Ridge

- Ridge shrinks all coefficients
- Lasso sets some coefficients at 0 (sparse solution)
  - Perform feature selection



Lasso        Ridge

# Ridge vs Lasso

- Both methods can be applied to any loss function (regression or classification)

| Ridge | Lasso |
|-------|-------|
| • Ridge | • Lasso |

# Ridge vs Lasso

- Both methods can be applied to any loss function (regression or classification)
- In both methods, value of regularization parameter $\lambda$ needs to be adjusted
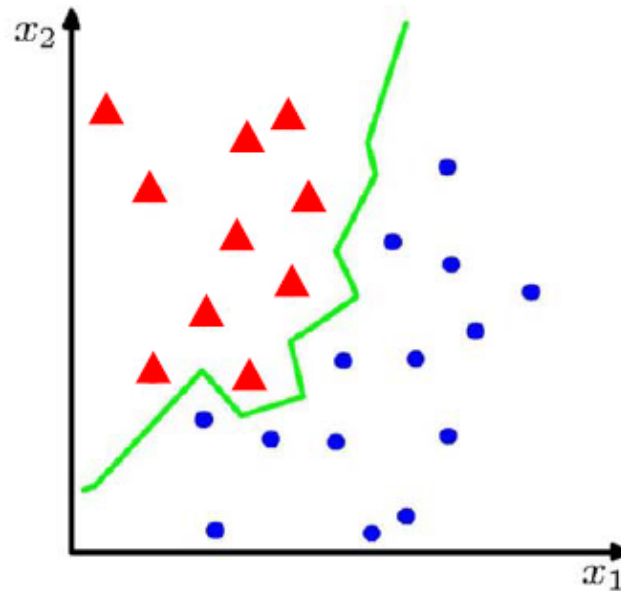- Both reduce model complexity

| - Ridge | - Lasso |
|---|---|
| + Differentiable objective | - Gradient descent needs to be adapted |
| + Gradient descent converges to global optimum | + Results in sparse model |
| - Shrinks all coefficients | + Can be used for feature selection in large dimensions |

# Classification
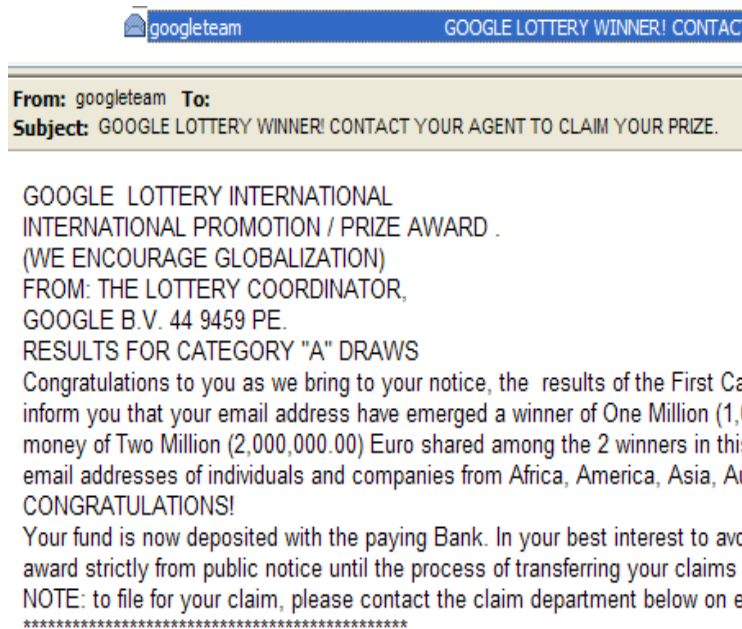


- Suppose we are given a training set of N observations

$$\{x_1, \ldots, x_N\} \text{ and } \{y_1, \ldots, y_N\}, x_i \in R^d, y_i \in \{0, 1\}$$

Binary or discrete

- Classification problem is to estimate f(x) from this data such that
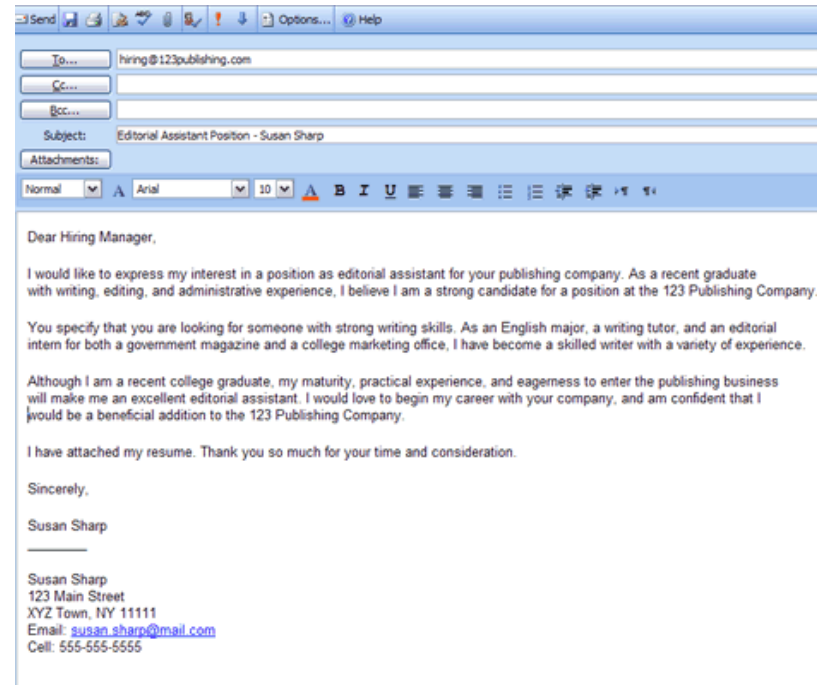
$$f(x_i) = y_i$$

# Example 1: Binary classification

## Classifying spam email





**Content-related features**
- Use of certain words
- Word frequencies
- Language
- Sentence

**Structural features**
- Sender IP address
- IP blacklist
- DNS information
- Email server
- URL links (non-matching)

## Binary classification: SPAM or HAM

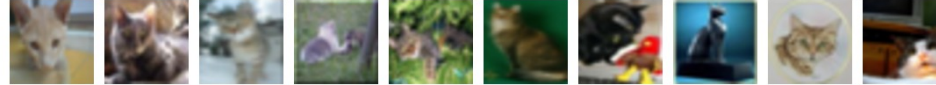# Example 2: Multi-class classification

Image classification
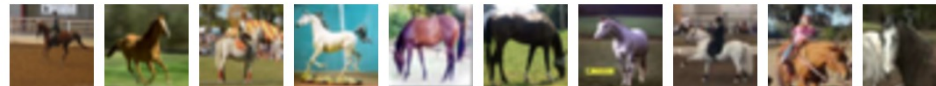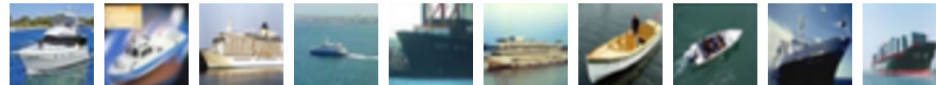


airplane
automobile
bird
cat
deer
dog
frog
horse
ship
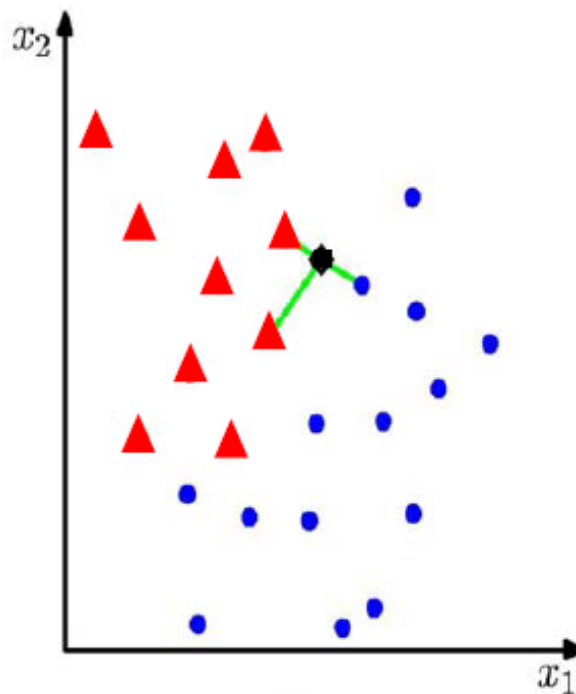truck

Multi-class classification

# K Nearest Neighbour (K-NN) Classifier

## Algorithm

- For each test point, x, to be classified, find the K nearest samples in the training data

- Classify the point, x, according to the majority vote of their class labels

e.g. K = 3

• applicable to multi-class case

# Distance Metrics

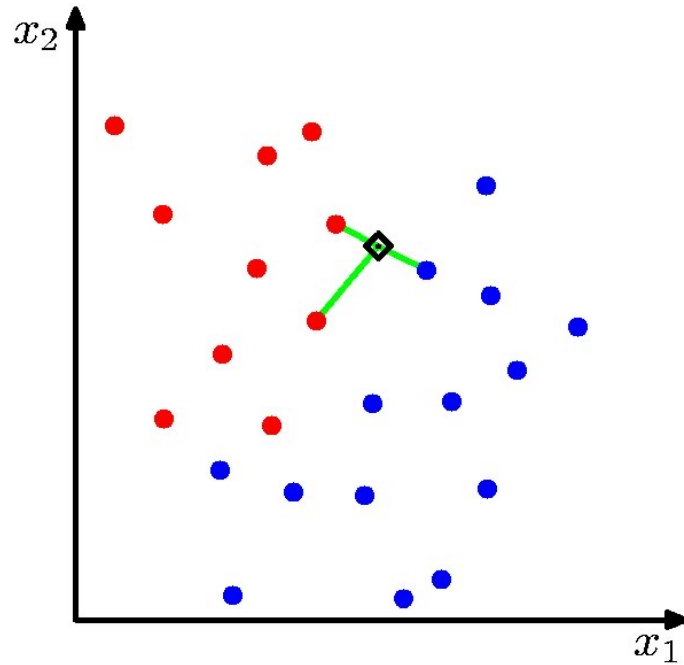- Euclidean Distance $$\sqrt{\left(\sum_{i=1}^{k}(x_i - y_i)^2\right)}$$

- Manhattan Distance $$\sum_{i=1}^{k}|x_i - y_i|$$
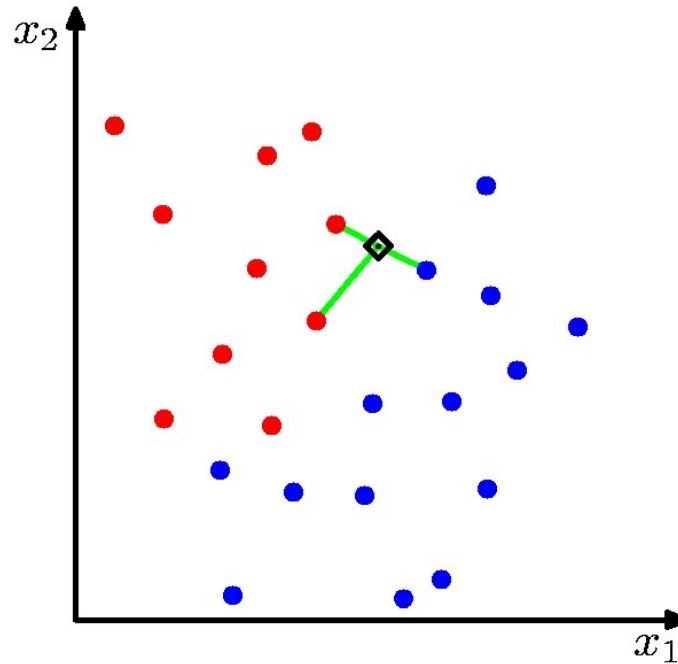
- Minkowski Distance $$\left(\sum_{i=1}^{k}(|x_i - y_i|)^q\right)^{\frac{1}{q}}$$
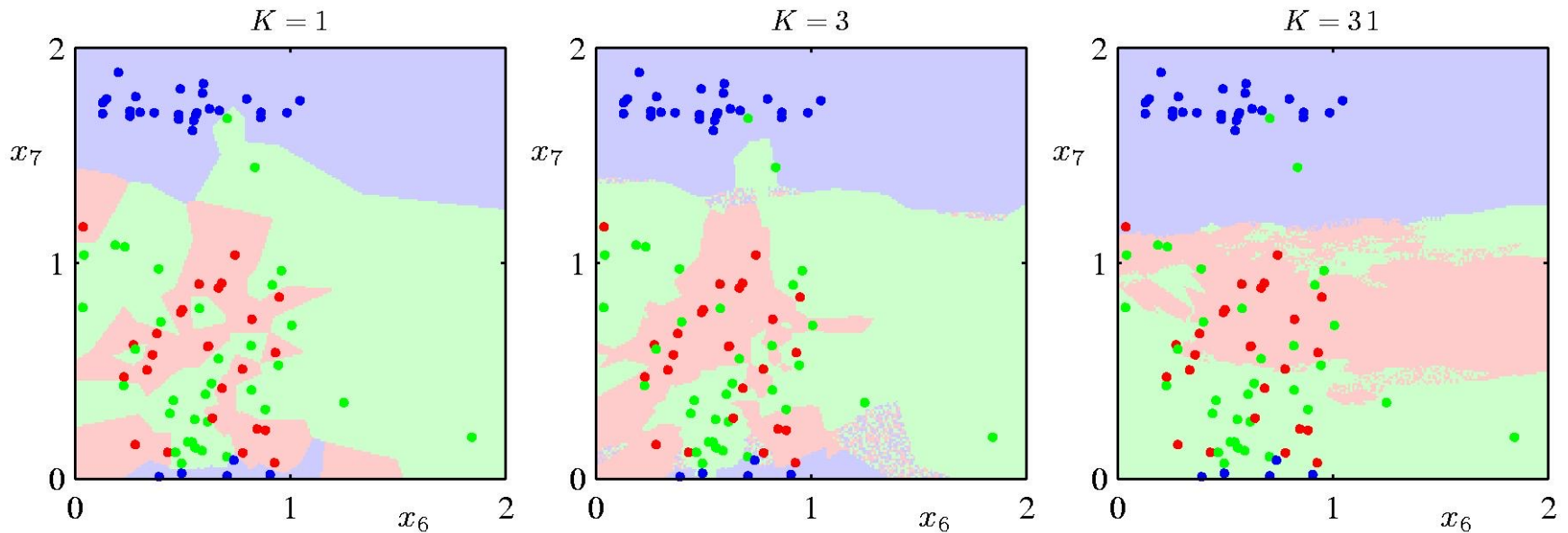
# kNN



- Algorithm (to classify point $x$)
  - Find $k$ nearest points to $x$ (according to distance metric)
  - Perform majority voting to predict class of $x$

# kNN



- Algorithm (to classify point $x$)
  - Find $k$ nearest points to $x$ (according to distance metric)
  - Perform majority voting to predict class of $x$
- Properties
  - Does not learn any model in training!
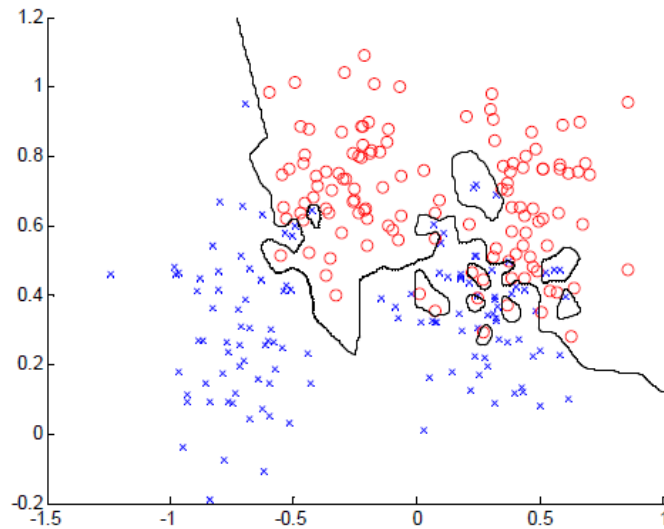  - Instance learner (needs all data at testing time)

# K-Nearest-Neighbours for Multi-class Classification
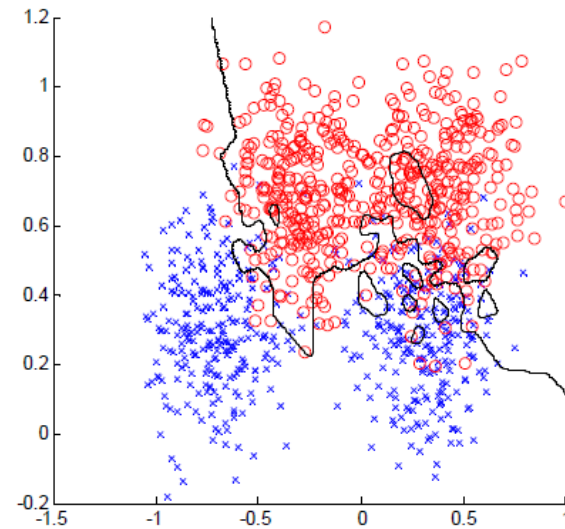


Vote among multiple classes

# K = 1



How to choose k (hyper-parameter)?

# K = 3



Training data

Testing data
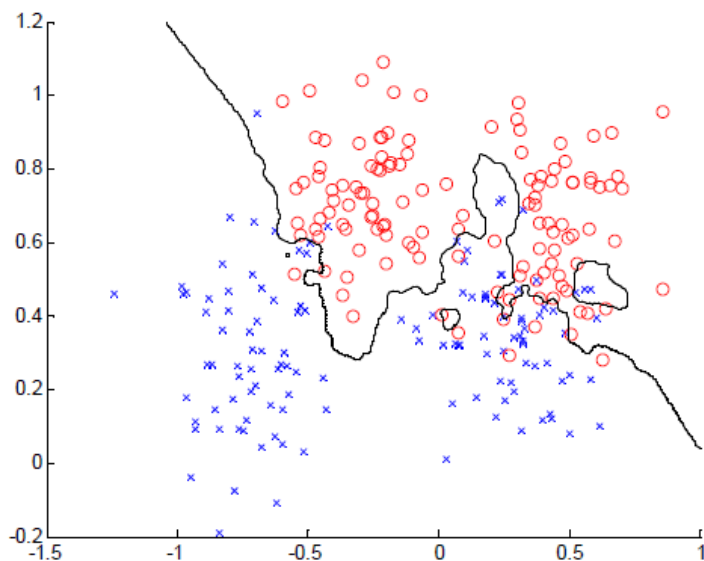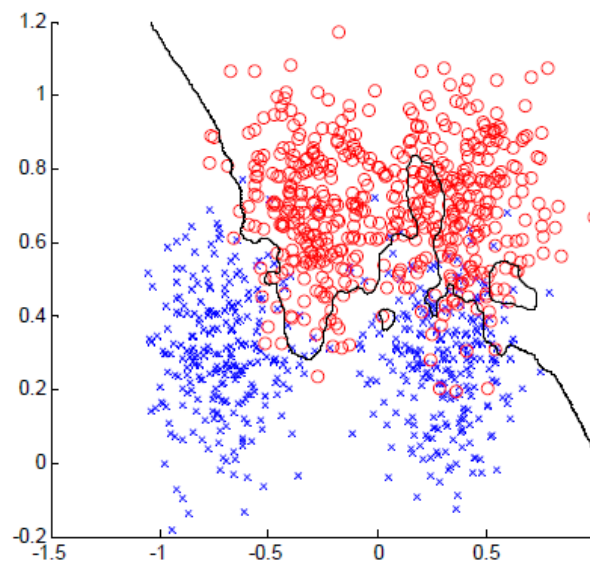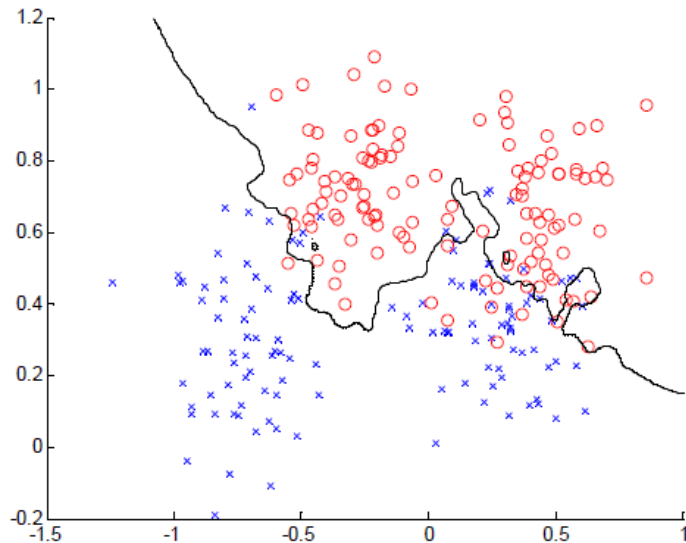
error = 0.0760

error = 0.1340

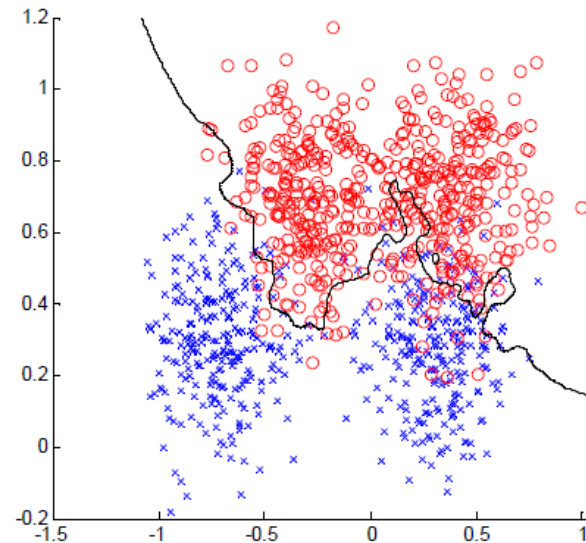How to choose k (hyper-parameter)?
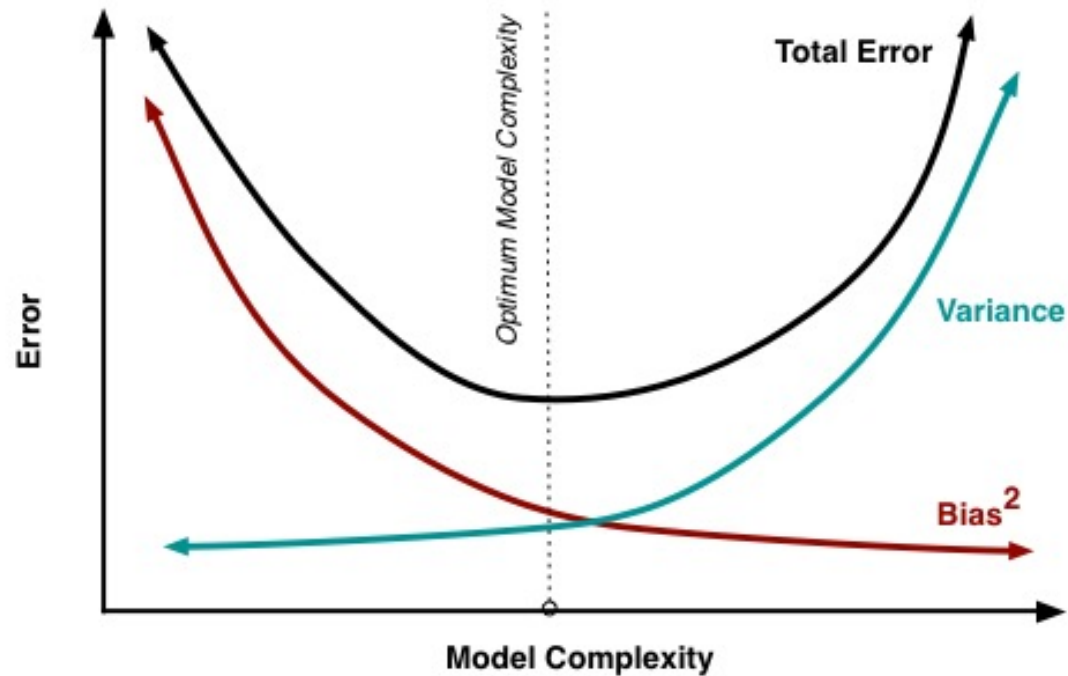
# K = 7



Training data
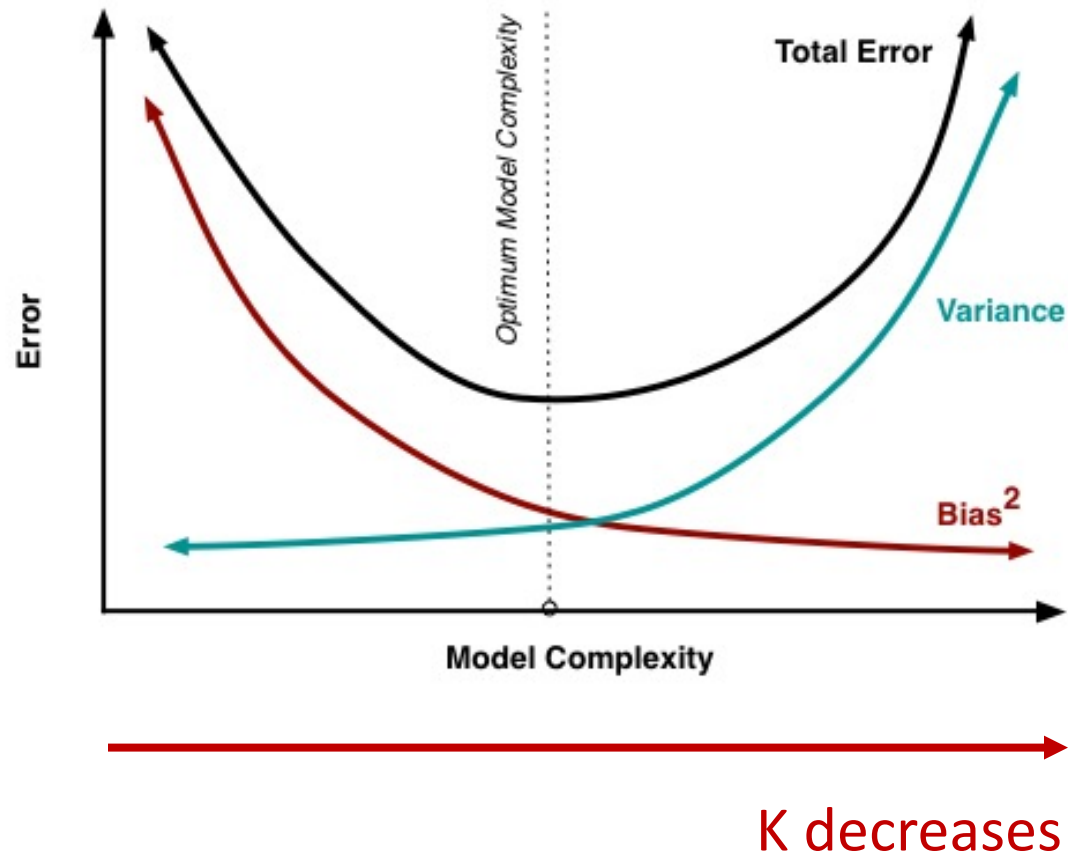
error = 0.1320

Testing data

error = 0.1110

How to choose k (hyper-parameter)?

# Bias-Variance Tradeoff for kNN

# Bias-Variance Tradeoff for kNN



K decreases

# How Overfitting Affects Prediction



How can we avoid over-fitting without having access to testing data?
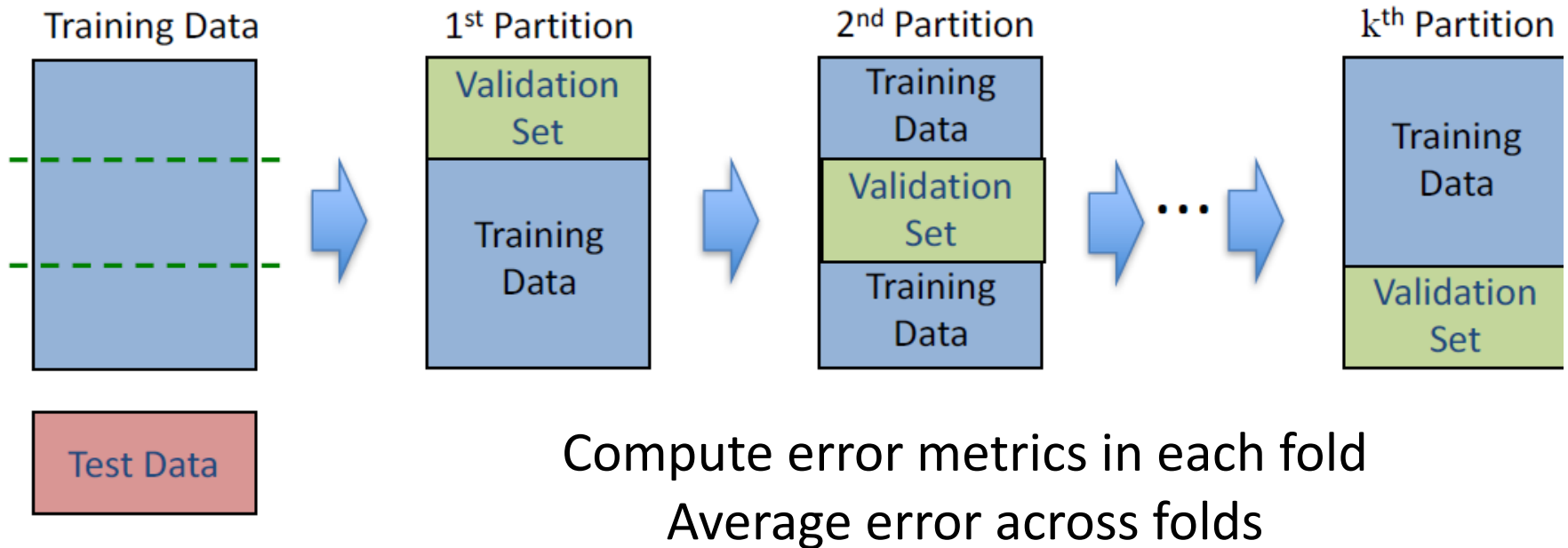
# Cross Validation

As K increases:

- Classification boundary becomes smoother
- Training error can increase

Choose (learn) K by cross-validation

- Split training data into training and validation
- Hold out validation data and measure error on this

# Cross Validation



**Training Data**

**1st Partition**
- Validation Set
- Training Data

**2nd Partition**
- Training Data
- Validation Set
- Training Data

$\mathbf{k^{th}}$ **Partition**
- Training Data
- Validation Set

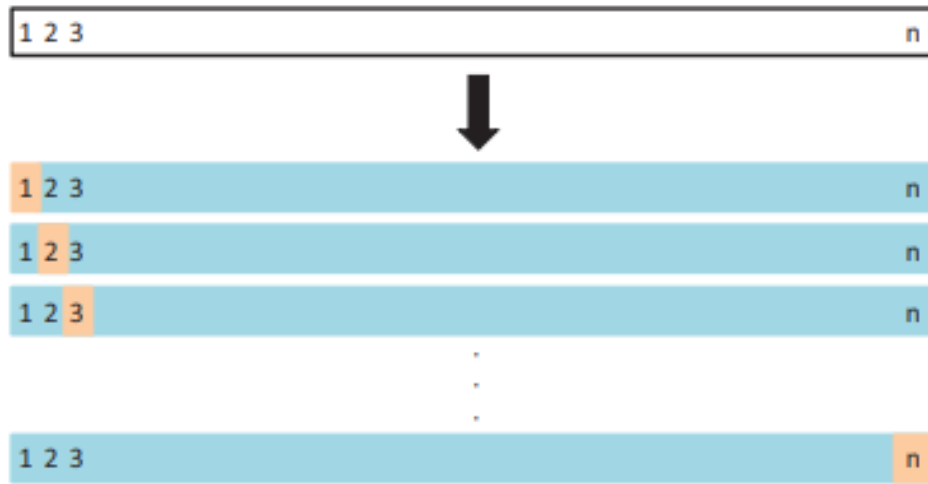**Test Data**

Compute error metrics in each fold
Average error across folds

## 1. k-fold CV

- Split training data into k partitions (folds) of equal size
- Pick the optimal value of hyper-parameter according to error metric averaged over all folds

# Cross Validation



## 2. Leave-one-out CV (LOOCV)

- k=n (validation set only one point)
- Pros: Less bias
- Cons: More expensive to implement, higher variance
- Recommendation: perform k-fold CV with k=5 or k=10

# Cross-Validation Takeaways

- General method to estimate performance of ML model at testing and select hyper-parameters
  - Improves model generalization
  - Avoids overfitting to training data
- Techniques for CV: k-fold CV and LOOCV
- Compare to regularization
  - Regularization works when training with GD
  - Cross-validation can be used for hyper-parameter selection
  - The two methods can be combined (Ridge, Lasso)