

# DS 4400

## Machine Learning and Data Mining I

Alina Oprea  
Associate Professor  
Khoury College of Computer Science  
Northeastern University

October 1 2020

# Announcements

- Thanks for submitting HW 1
- HW 2 is posted on Piazza and Gradescope
  - It is due on Monday, Oct. 12, at midnight
- Start thinking about class projects
  - Will post guidelines soon
  - Find a project partner and dataset you are interested in

# Outline

- Gradient Descent comparison with closed-form solution
- Regularization
  - Ridge and Lasso regression
- Classification
  - K Nearest Neighbors (kNN)
  - Cross-validation

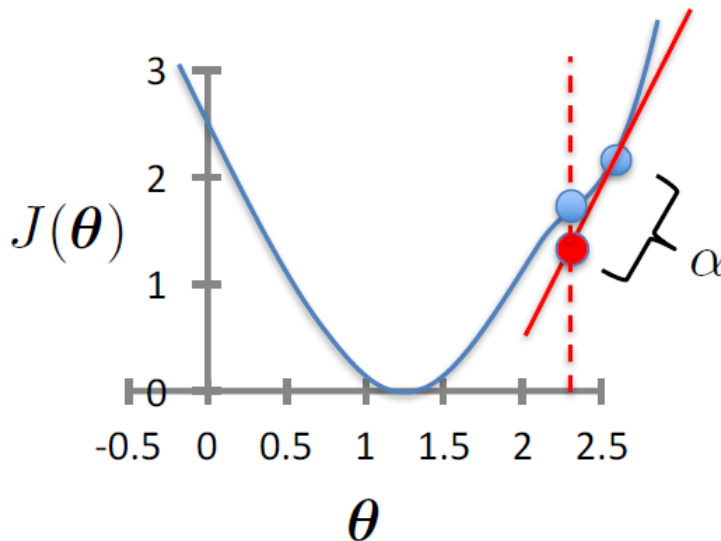
# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update  
for  $j = 0 \dots d$

learning rate (small)  
e.g.,  $\alpha = 0.05$



- Gradient = slope of line tangent to curve
- Function decreases faster in negative direction of gradient
- Larger learning rate => larger step

# GD for Linear Regression

- Initialize  $\theta$
- Repeat until convergence  $\|\theta_{new} - \theta_{old}\| < \epsilon$  or  $iterations == MAX\_ITER$

$$\theta_j \leftarrow \theta_j - \alpha \frac{2}{N} \sum_{i=1}^N (h_{\theta}(x_i) - y_i) x_{ij}$$

simultaneous  
update  
for  $j = 0 \dots d$

- To achieve simultaneous update
  - At the start of each GD iteration, compute  $h_{\theta}(x_i)$
  - Use this stored value in the update step loop
- Assume convergence when  $\|\theta_{new} - \theta_{old}\|_2 < \epsilon$

$$\text{L}_2 \text{ norm: } \|v\|_2 = \sqrt{\sum_i v_i^2} = \sqrt{v_1^2 + v_2^2 + \dots + v_{|v|}^2}$$

Can also bound number of iterations

# Gradient Descent in Practice

- Asymptotic complexity
  - $O(NTd)$ ,  $N$  is size of training data,  $d$  is feature dimension, and  $T$  is number of iterations
- Most popular optimization algorithm in use today
- At the basis of training
  - Linear Regression
  - Logistic regression
  - SVM
  - Neural networks and Deep learning
  - Stochastic Gradient Descent variants

# Gradient Descent vs Closed Form

## Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update  
for  $j = 0 \dots d$

## Closed form

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

### • Gradient Descent

- + Linear increase in  $d$  and  $N$
- + Generally applicable
- Need to choose  $\alpha$  and stopping conditions
- Might get stuck in local optima

### • Closed Form

- + No parameter tuning
- + Gives the global optimum
- Not generally applicable
- Slow:  $O(Nd^2) + O(d^3)$

# Issues with Gradient Descent

- Might get stuck in local optimum and not converge to global optimum
  - Restart from multiple initial points
- Only works with differentiable loss functions
- Small or large gradients
  - Feature scaling helps
- Tune learning rate
  - Can use line search for determining optimal learning rate



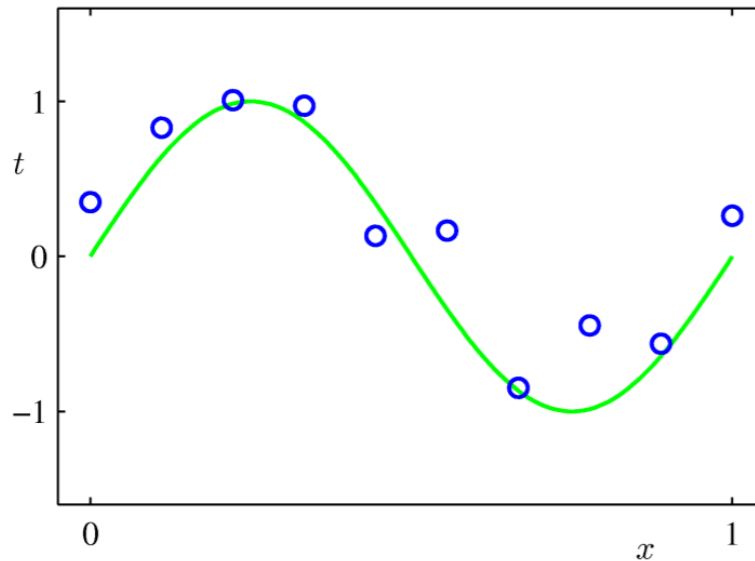
# Review Gradient Descent

- Gradient descent is an efficient algorithm for optimization and training ML models
  - The most widely used algorithm in ML!
  - Much faster than using closed-form solution for linear regression
  - Main issues with Gradient Descent is convergence and getting stuck in local optima (for neural networks)
- Gradient descent is guaranteed to converge to optimum for strictly convex functions if run long enough

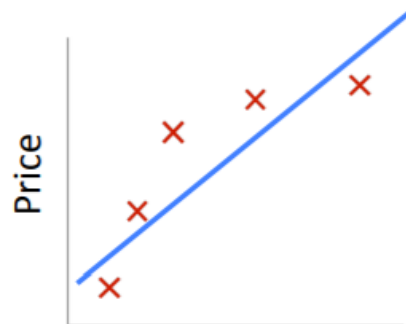
# Polynomial Regression

- Polynomial function on single feature

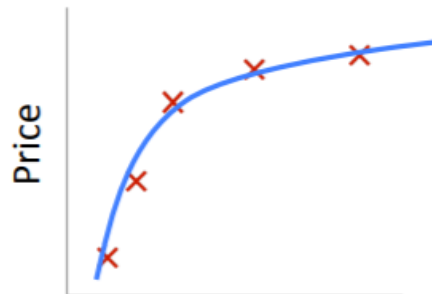
$$- h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_d x^p$$



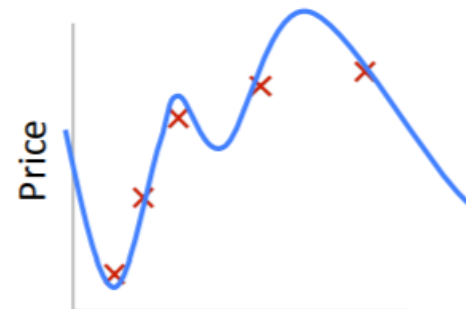
# Polynomial Regression



Size  
 $\theta_0 + \theta_1 x$   
Underfitting  
(high bias)



Size  
 $\theta_0 + \theta_1 x + \theta_2 x^2$   
Correct fit

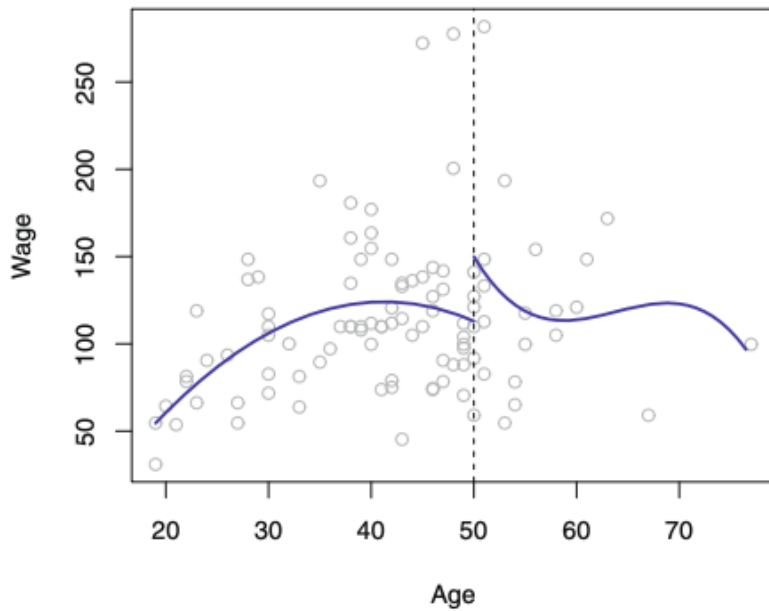


Size  
 $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$   
Overfitting  
(high variance)

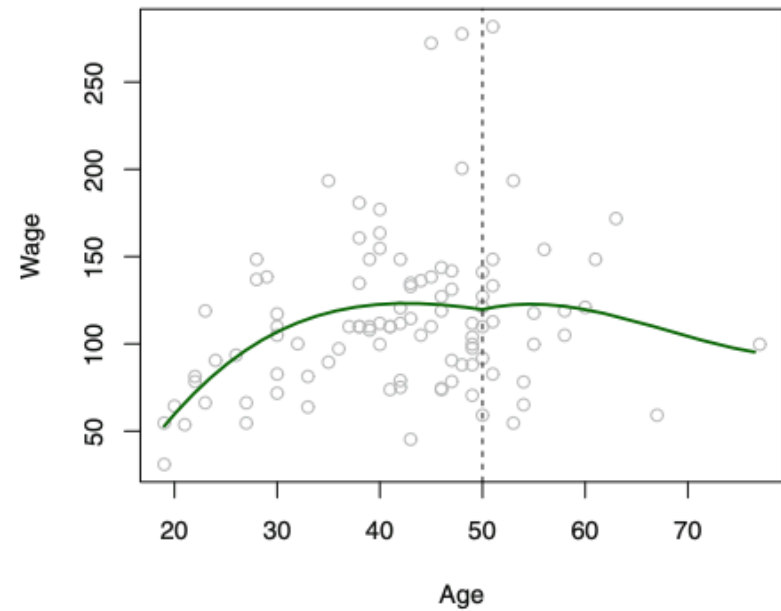
- Typically to avoid overfitting  $p \leq 4$

# Non-Linear Regression

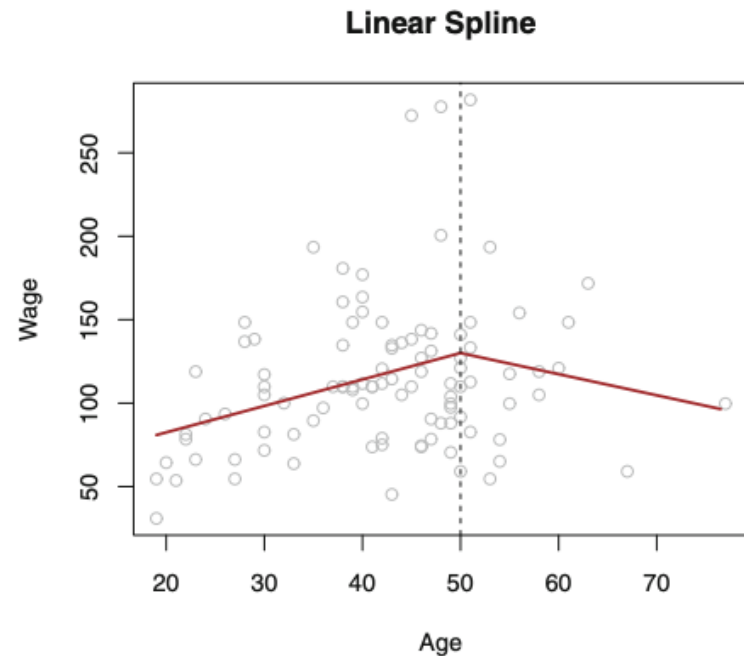
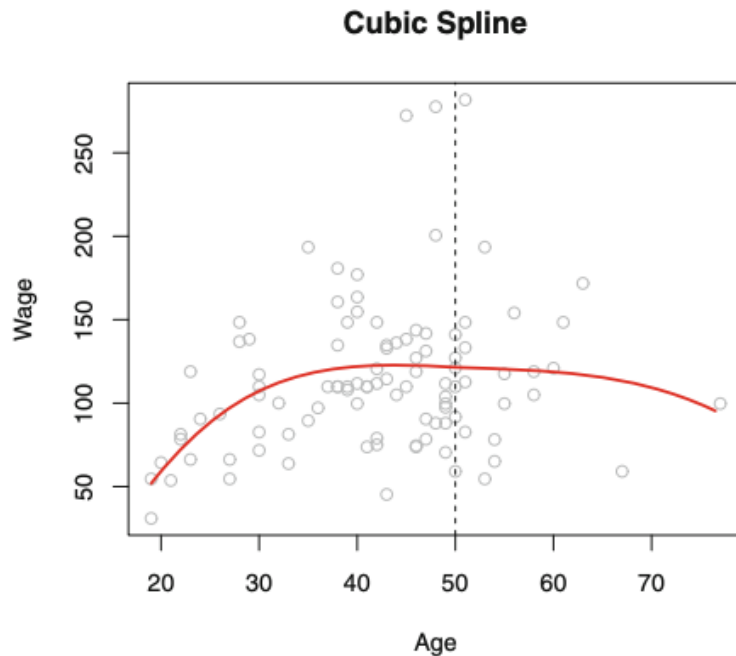
**Piecewise Cubic**



**Continuous Piecewise Cubic**

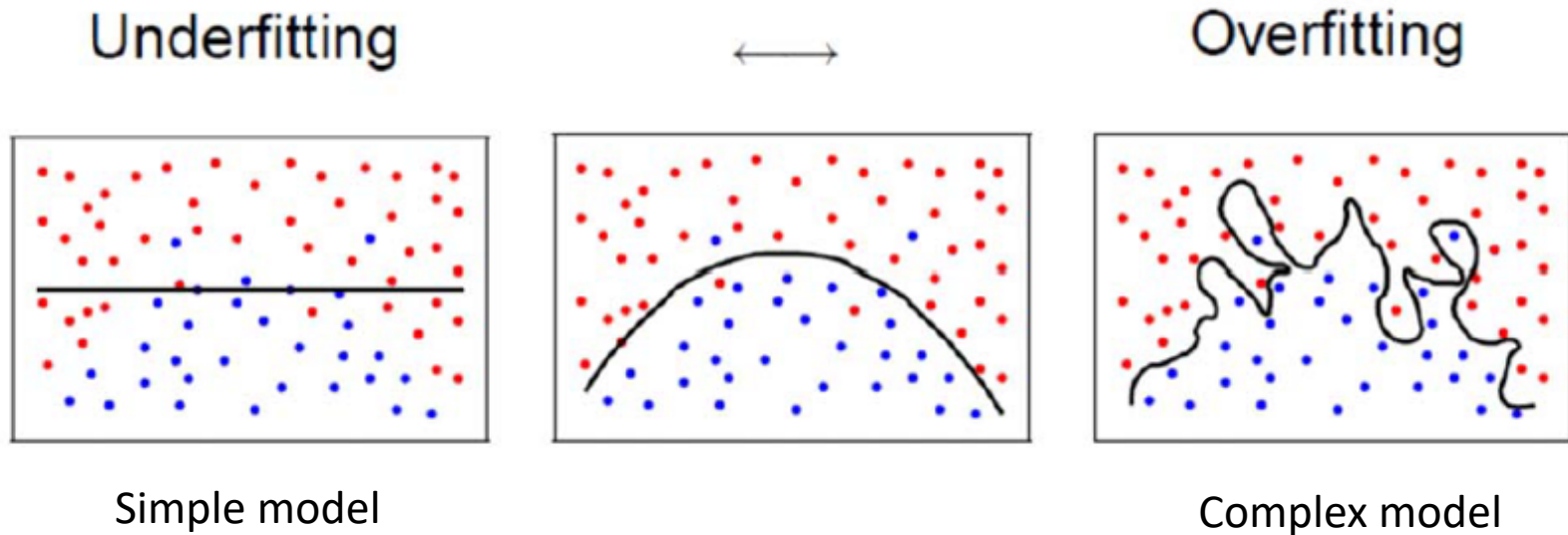


# Spline Regression



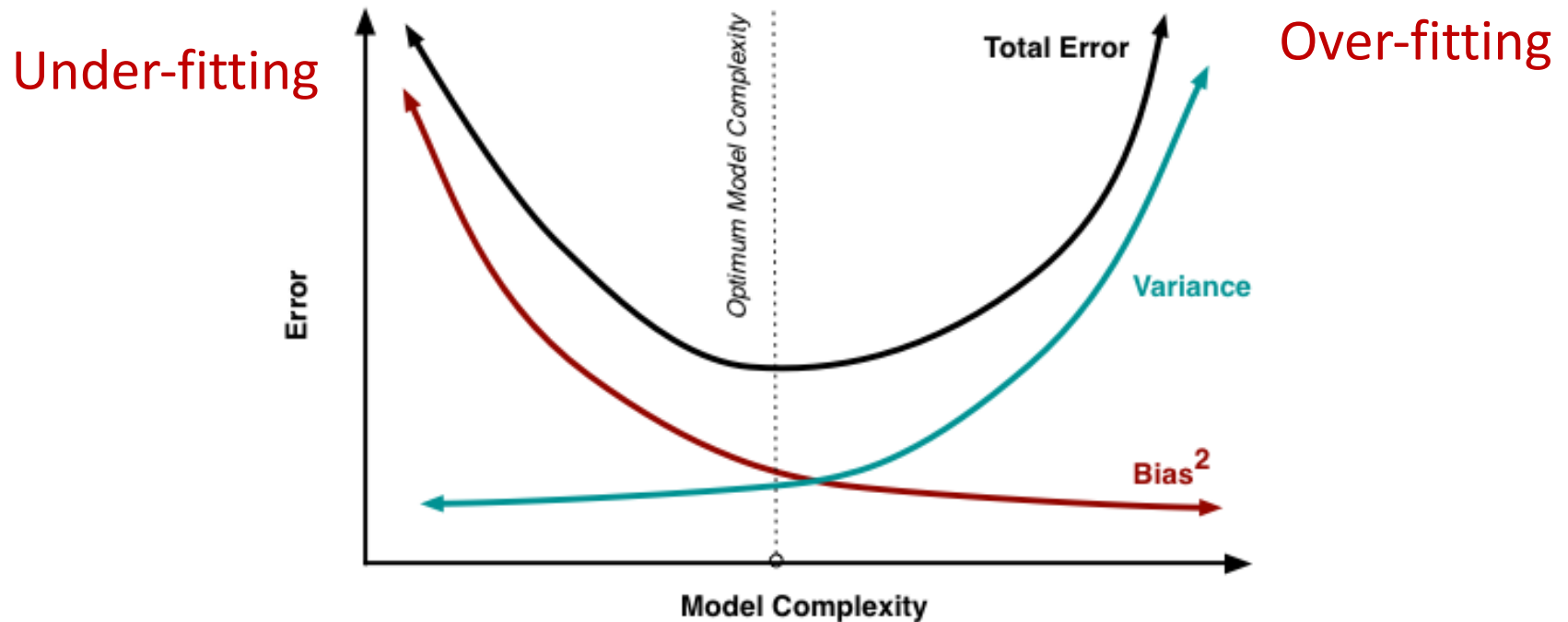
- Fit polynomial regression on each region (knot)
- Spline - Continuous and differentiable function at boundary

# Generalization in ML



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

# Bias-Variance Tradeoff



- Bias = Difference between estimated and true models
- Variance = Model difference on different training sets

**MSE is proportional to Bias + Variance**

# Regularization

- A method for automatically controlling the complexity of the learned hypothesis
- **Idea:** penalize for large values of  $\theta_j$ 
  - Can incorporate into the cost function
  - Works well when we have a lot of features, each that contributes a bit to predicting the label
- Can also address overfitting by eliminating features (either manually or via model selection)

Reduce model complexity

Reduce model variance



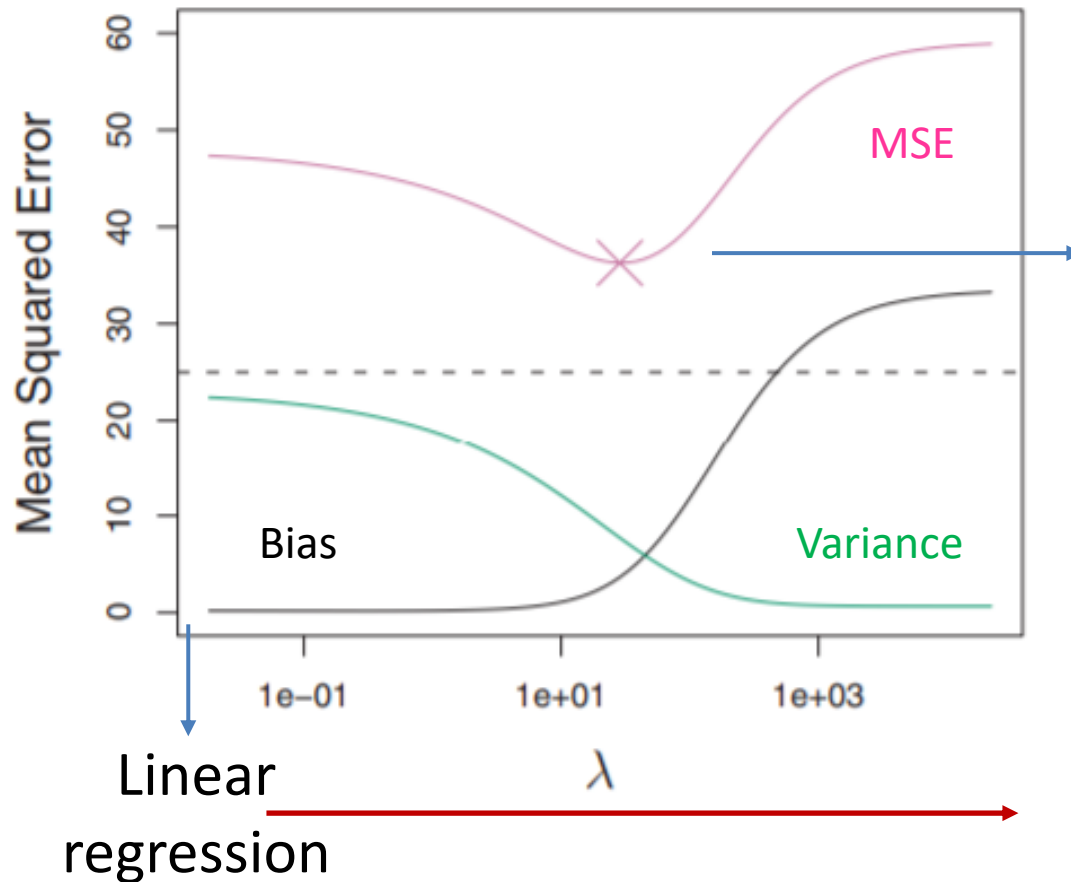
# Ridge regression

- Linear regression objective function

$$J(\theta) = \underbrace{\frac{1}{2} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2}_{\text{model fit to data}} + \underbrace{\frac{\lambda}{2} \sum_{j=1}^d \theta_j^2}_{\text{regularization}}$$

- $\lambda$  is the regularization parameter (  $\lambda \geq 0$  )
  - No regularization on  $\theta_0$ !
- If  $\lambda = 0$ , we train linear regression
  - If  $\lambda$  is large, the coefficients will shrink close to 0

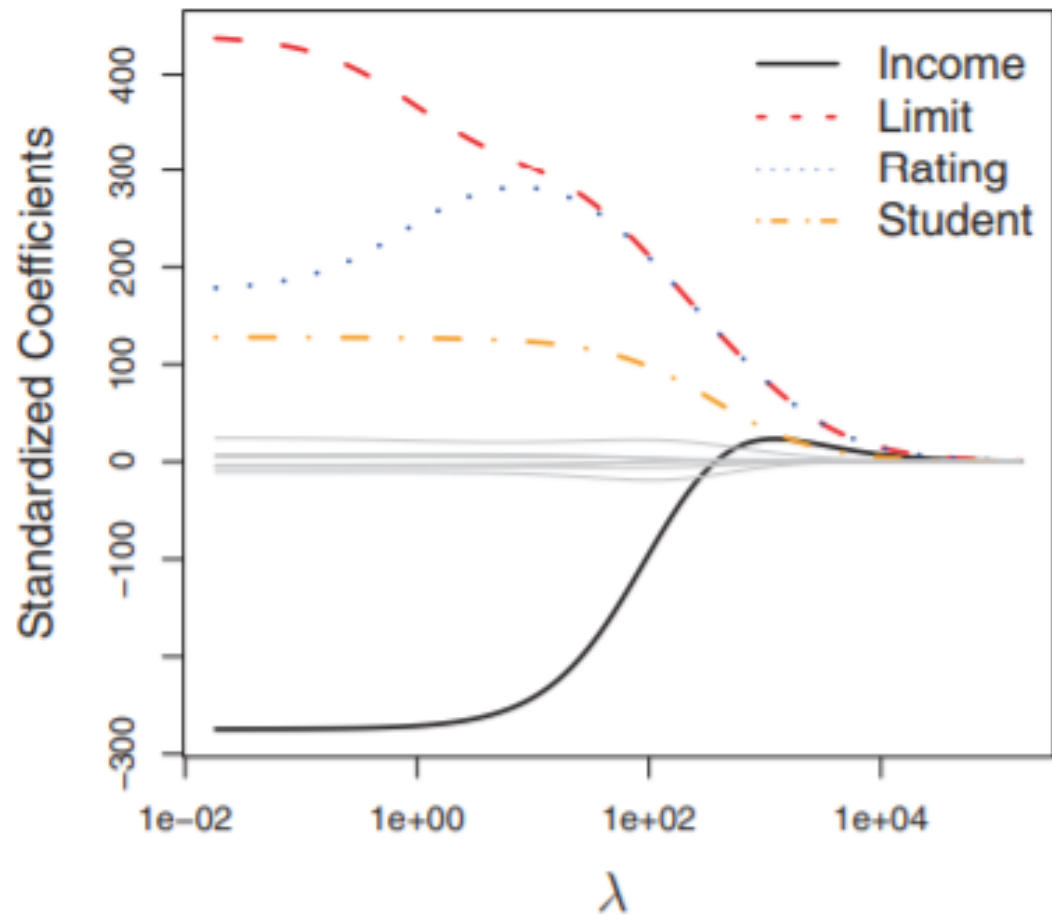
# Bias-Variance Tradeoff



Optimal  
Ridge regression

Reduced model  
complexity

# Coefficient shrinkage



Predict credit card balance

# GD for Ridge Regression

Min MSE

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

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

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

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

# Lasso Regression

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

- 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  $\sum_{j=1}^d |\theta_j|^2 \leq \epsilon$

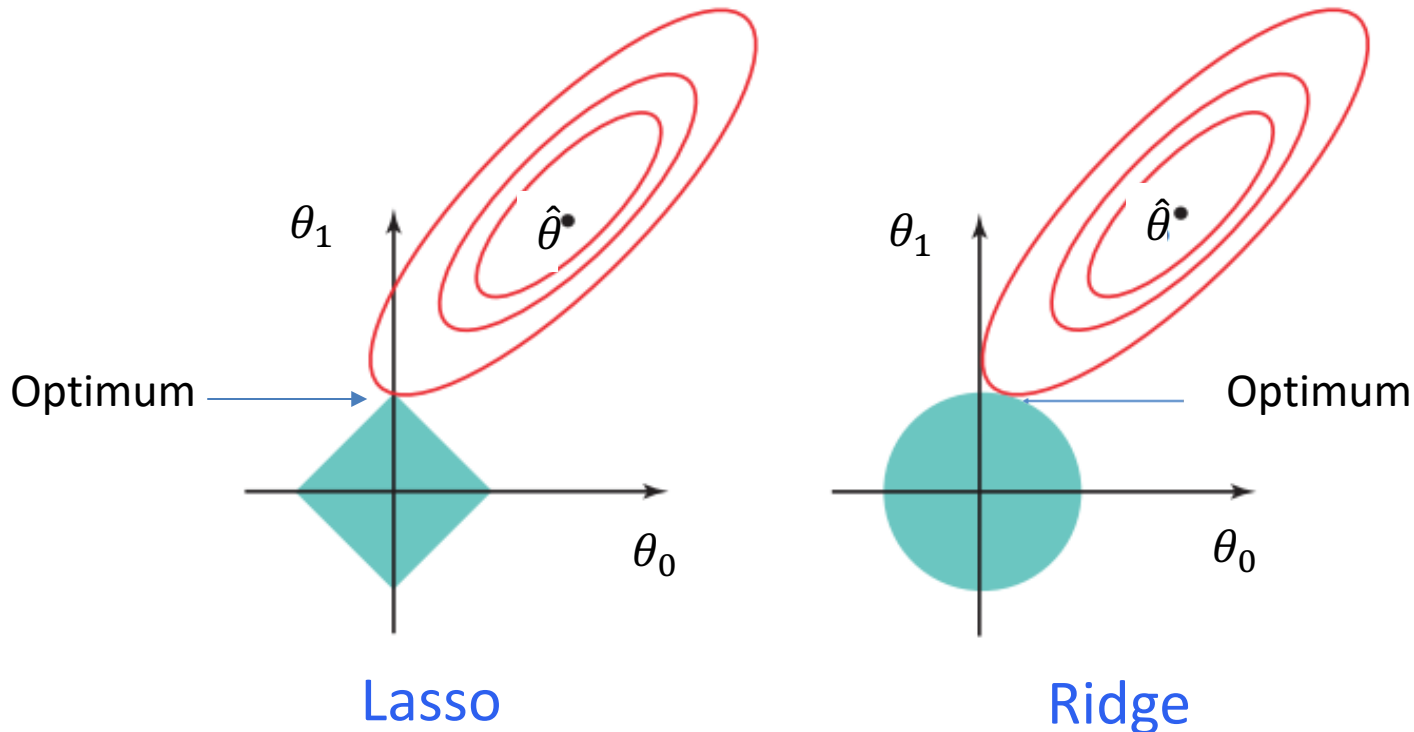
- Lasso

- L1 regularization

- $\min_{\theta} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2$  subject to  $\sum_{j=1}^d |\theta_j| \leq \epsilon$

# Lasso vs Ridge

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



# 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**

- + Differentiable objective
- + Gradient descent converges to global optimum
- Shrinks all coefficients

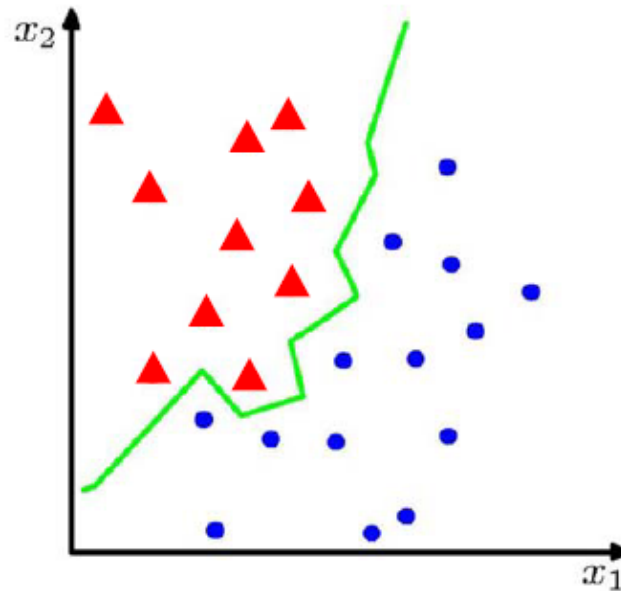
- **Lasso**

- Gradient descent needs to be adapted
- + Results in sparse model
- + Can be used for feature selection in large dimensions



# Classification

---



Binary or  
discrete

- Suppose we are given a training set of  $N$  observations

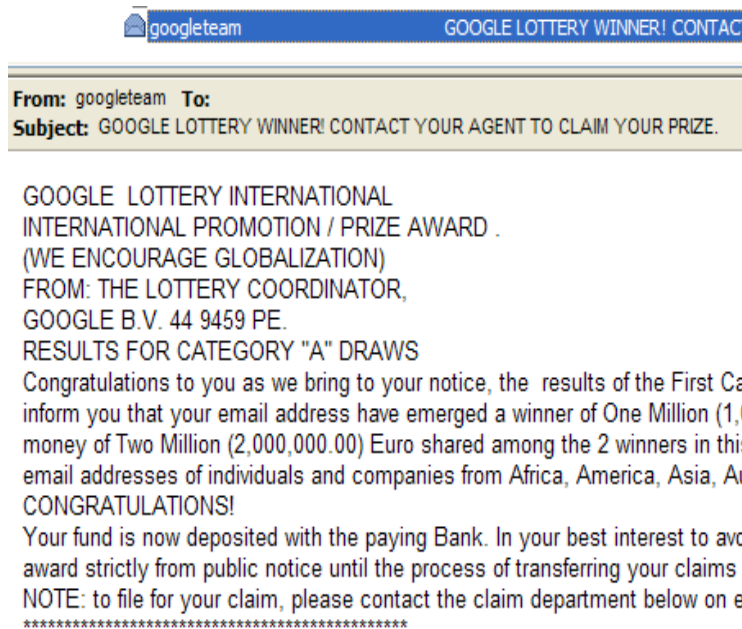
$$\{x_1, \dots, x_N\} \text{ and } \{y_1, \dots, y_N\}, x_i \in R^d, y_i \in \{0, 1\}$$

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

$$f(x_i) = y_i$$

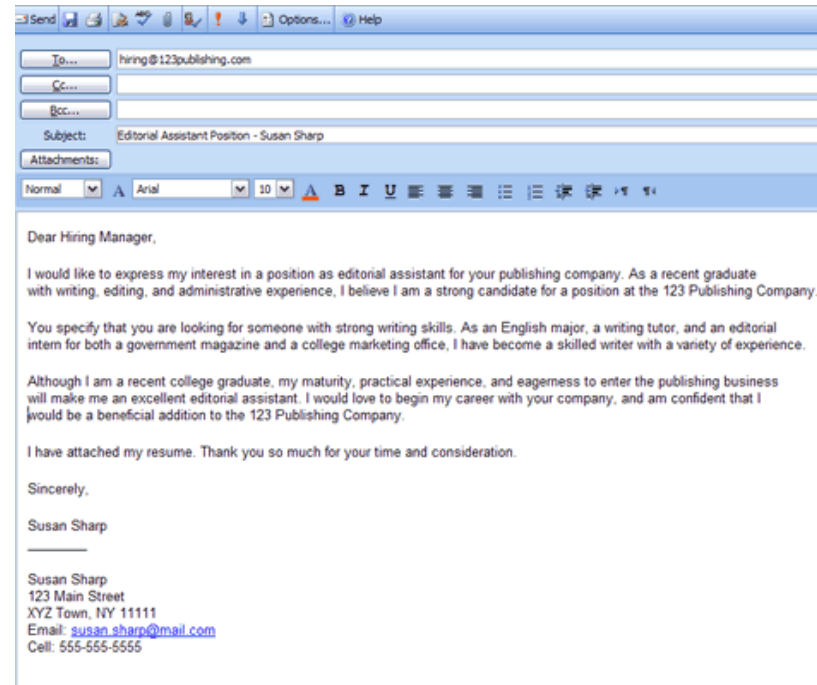
# Example 1: Binary classification

## 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

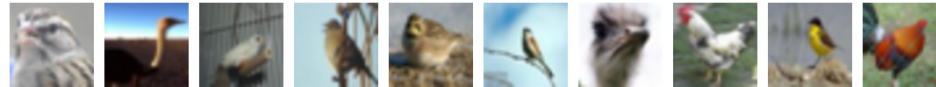
**airplane**



**automobile**



**bird**



**cat**



**deer**



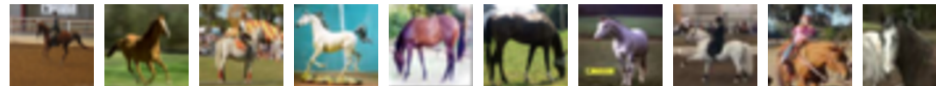
**dog**



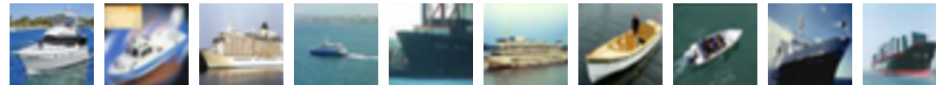
**frog**



**horse**



**ship**



**truck**



**Multi-class classification**

# K Nearest Neighbour (K-NN) Classifier

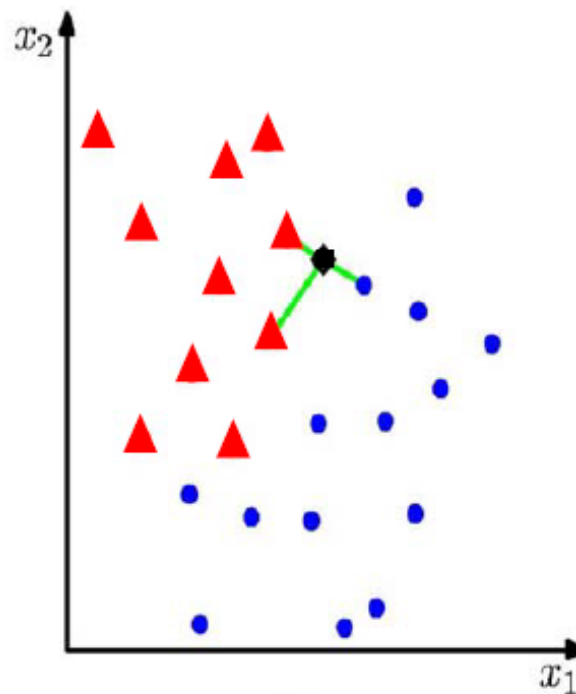
---

## 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

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

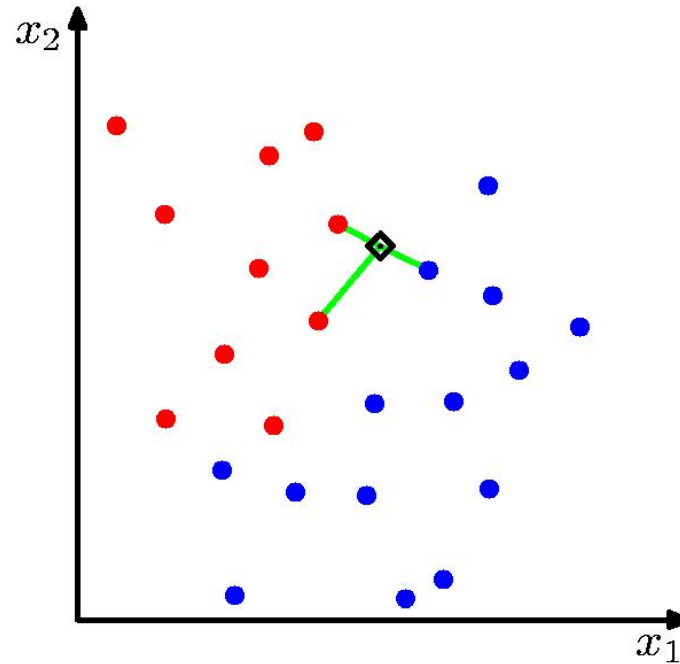
- Manhattan Distance

$$\sum_{i=1}^k |x_i - y_i|$$

- Minkowski Distance

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

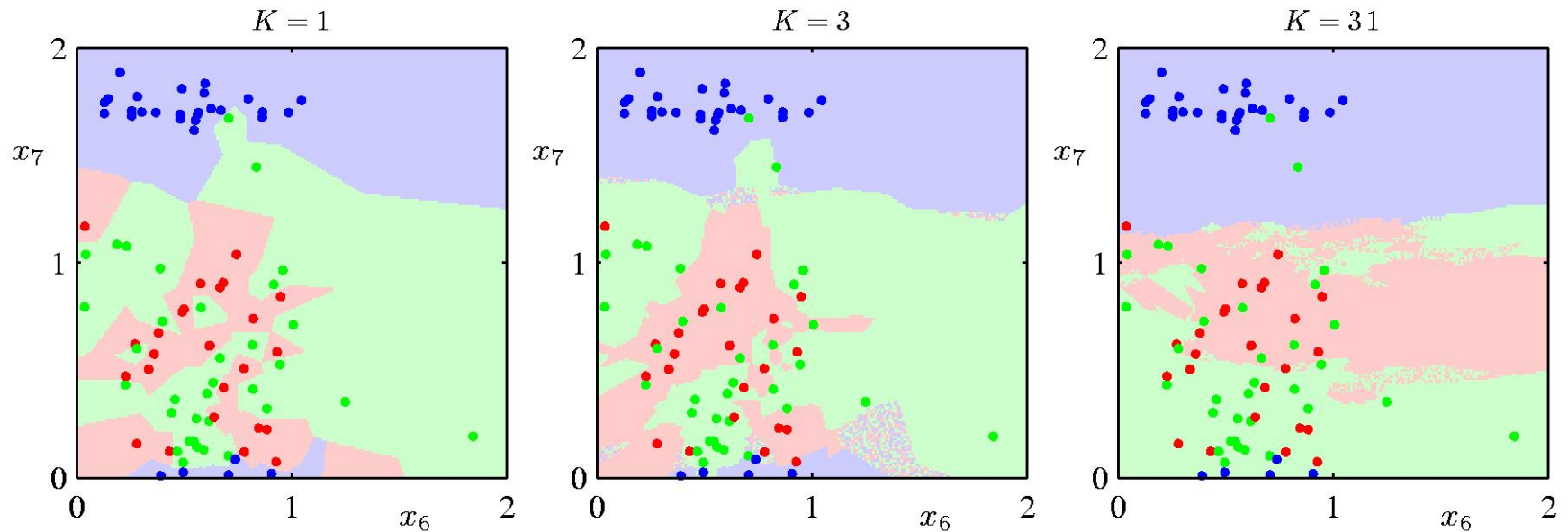
# kNN



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