

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. 9, 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 θ
- Repeat until convergence

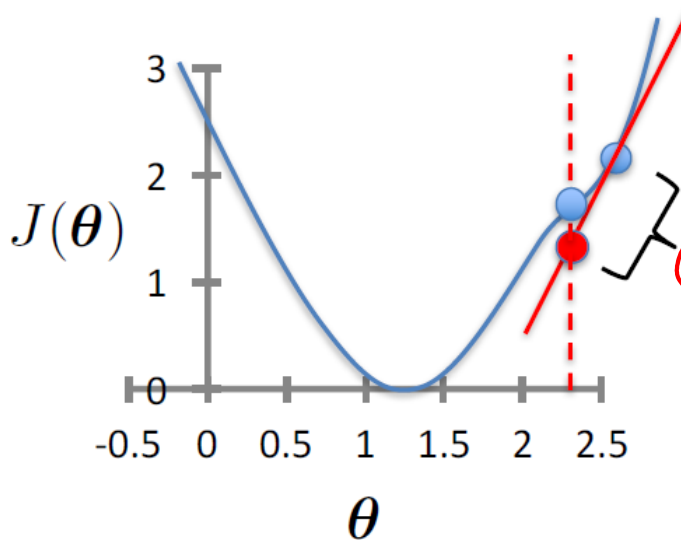
$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update
for $j = 0 \dots d$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

STOPPING:

- NUMBER STEPS
- HOW MUCH θ CHANGE



STEP SIZE
LEARNING RATE

GD for Linear Regression

- Initialize θ
- 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
 - 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

$$O(NdT)$$

Gradient Descent vs Closed Form

Gradient Descent

- Initialize θ
- 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

- + PERFORMANCE
- + MEMORY
- + GENERALIZABLE
- + GLOBAL OPT. FOR CONVEX OBJ.
- NO GUARANTEES ON OPT.
- PARAM TUNING

• Closed Form

- + GLOBAL OPTIMUM
- + NO PARAM TUNING
- NOT ALWAYS INVERTIBLE

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

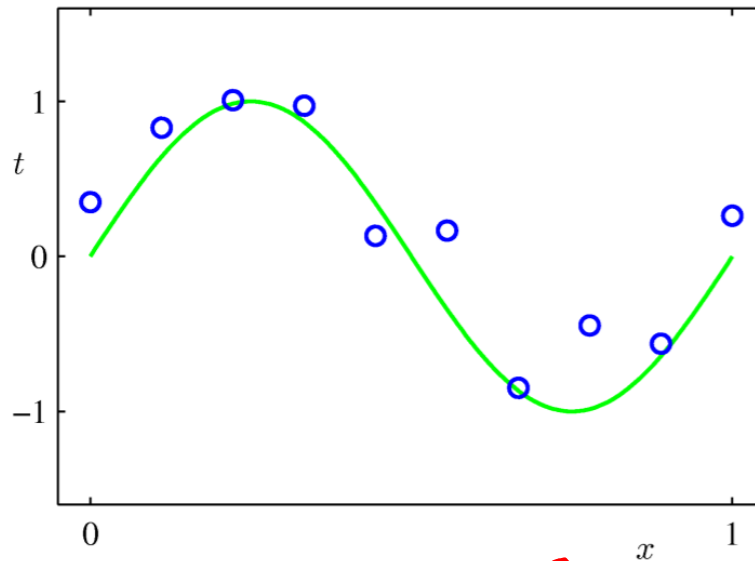
- Gradient descent is an efficient algorithm for optimization and training ML models
 - The most widely used algorithm in ML!
 - Much faster than using closed-form solution for linear regression
 - Main issues with Gradient Descent is convergence and getting stuck in local optima (for neural networks)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

Polynomial Regression

- Polynomial function on single feature

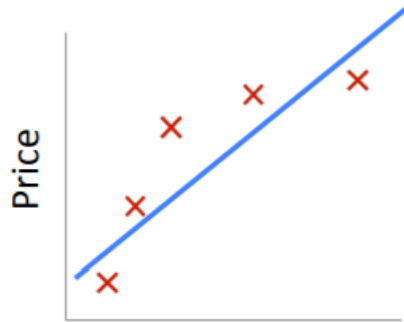
$$- h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_d x^p$$

degree p



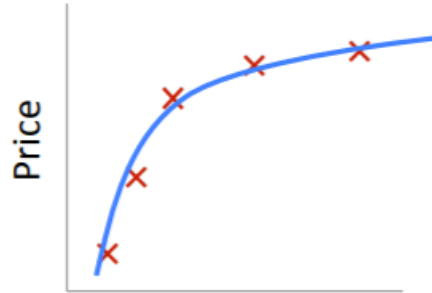
$x \rightarrow x, x^2, \dots, x^p$

Polynomial Regression



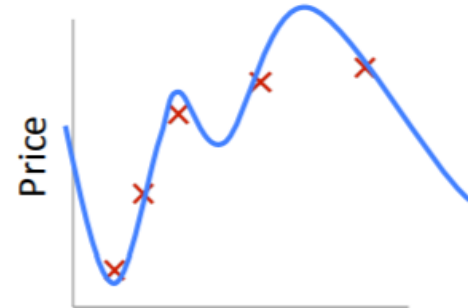
Size
 $\theta_0 + \theta_1 x$

UNDERFIT



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

CORRECT FIT

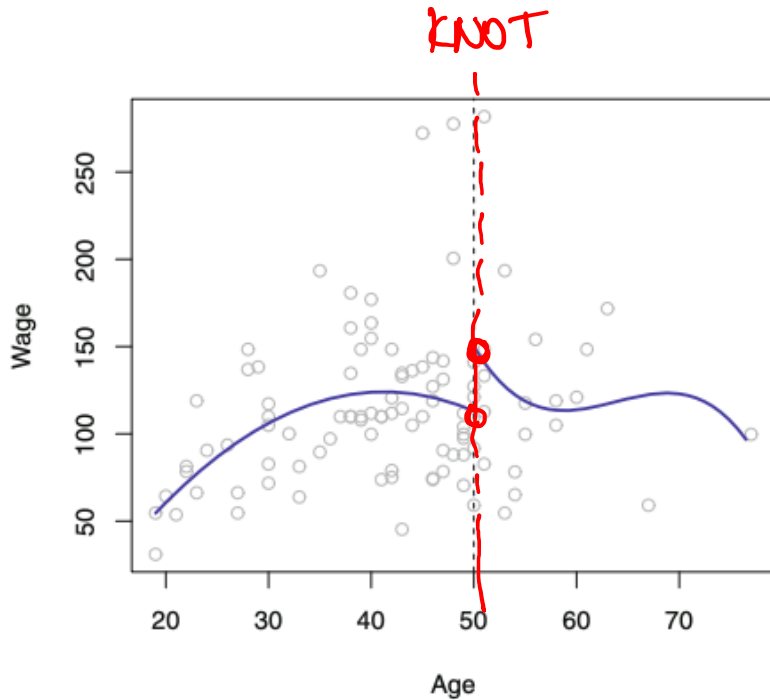


Size
 $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

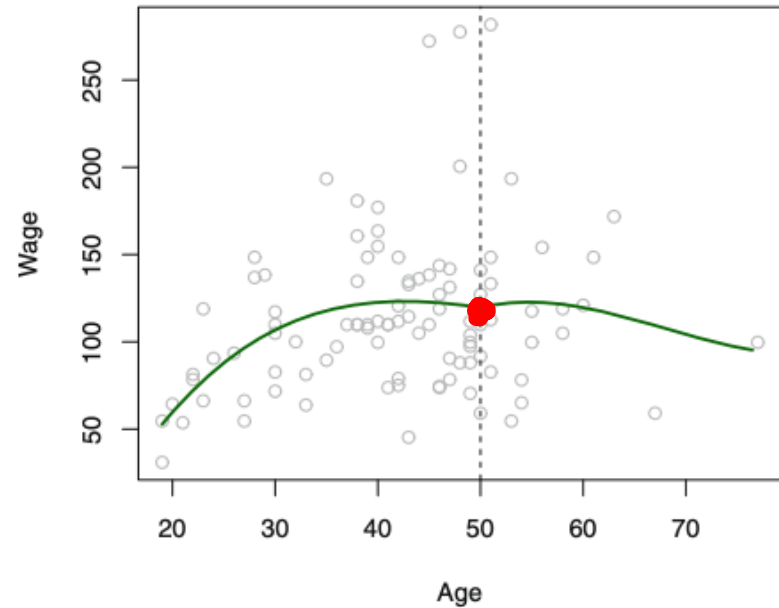
OVERFIT
TOO COMPLEX

TO AVOID OVERFITTING: $D \leq 4$.

Non-Linear Regression



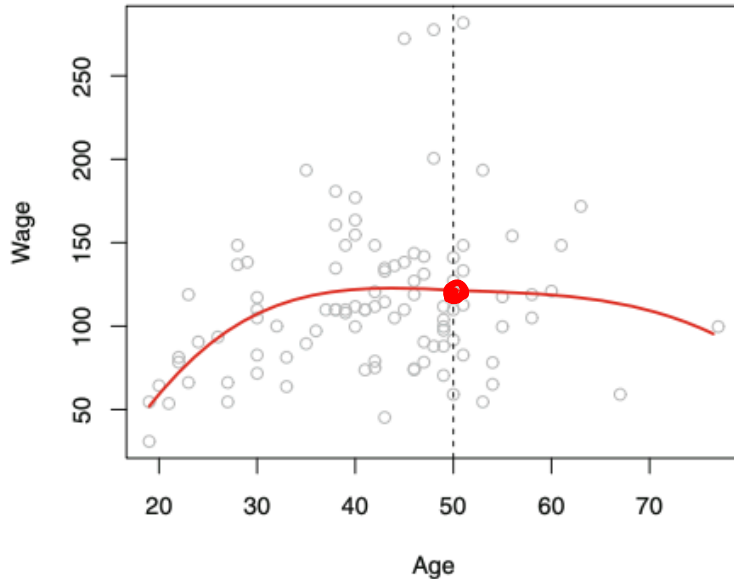
PIECEWISE CUBIC
DEG 3



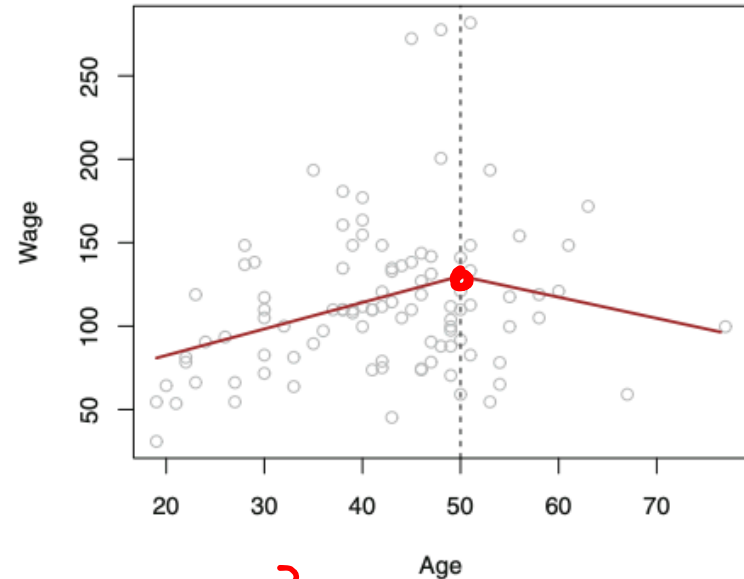
CONTINUOUS ...

Spline Regression

CUBIC SPLINE



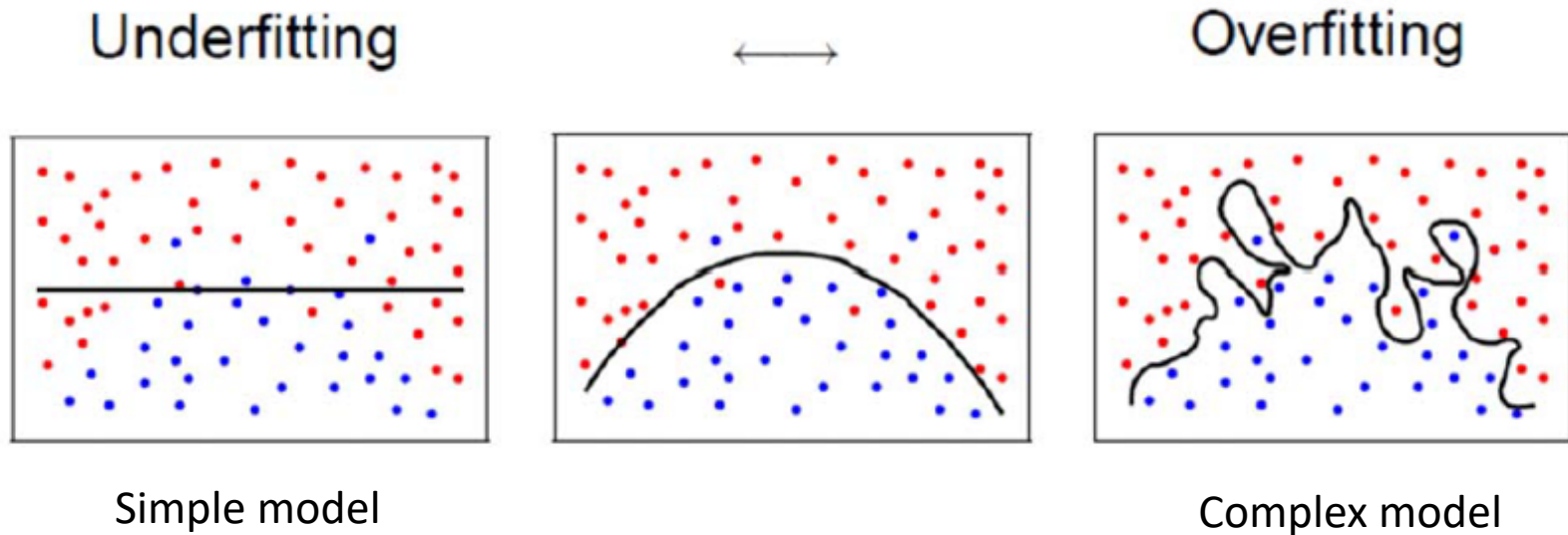
LINEAR SPLINE



How PICK # OF KNOTS?

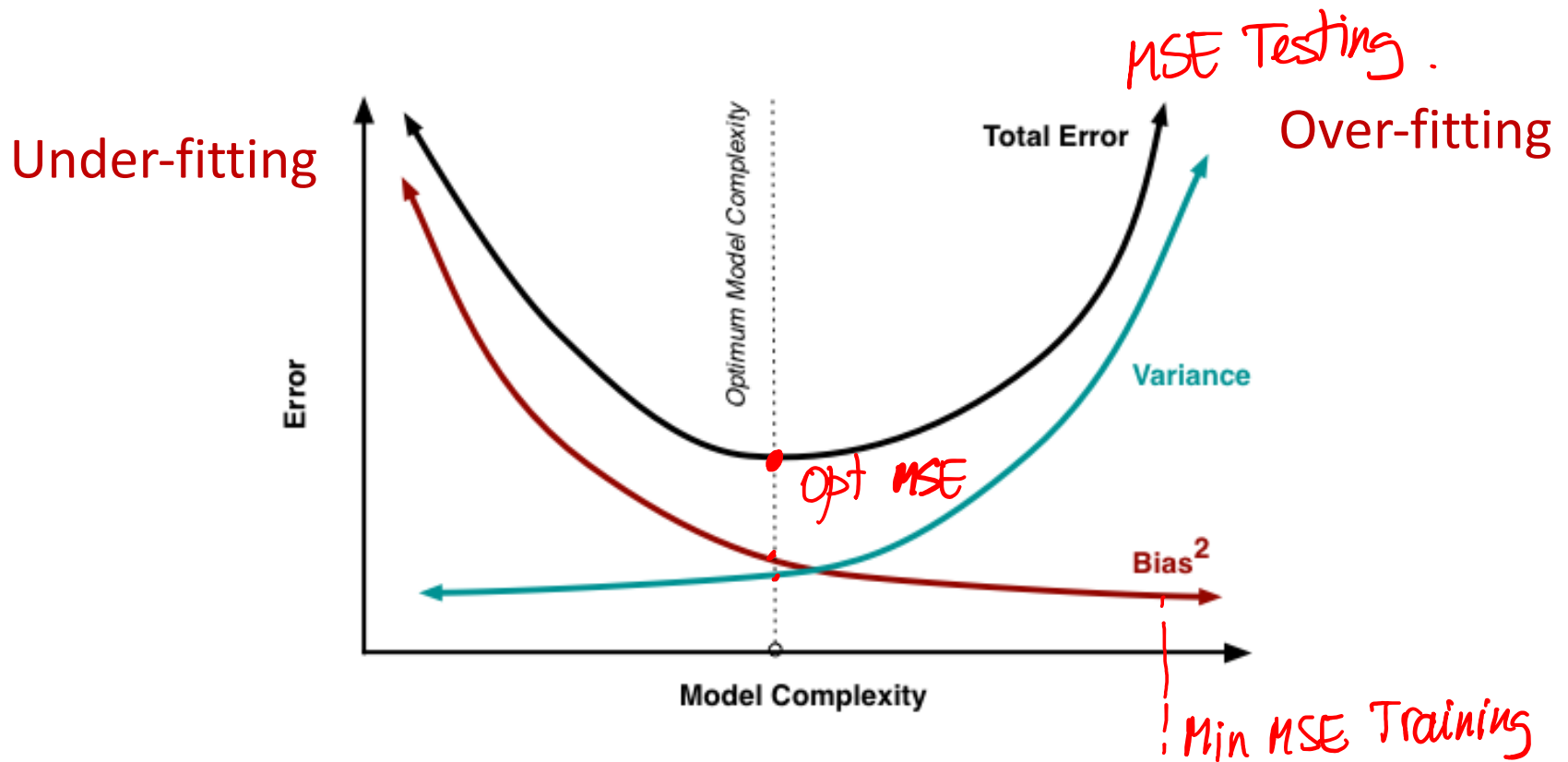
- Fit polynomial regression on each region (knot)
- Spline - Continuous and differentiable function at boundary

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 $\text{Bias}^2 + \text{Variance}$

Regularization

- A method for automatically controlling the complexity of the learned hypothesis
- **Idea:** penalize for large values of θ_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

L2 Regularization

- Linear regression objective function

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

$\|\theta\|_2^2$

λ - REG. PARAM

$\lambda = 0 \Rightarrow$ Linear reg.

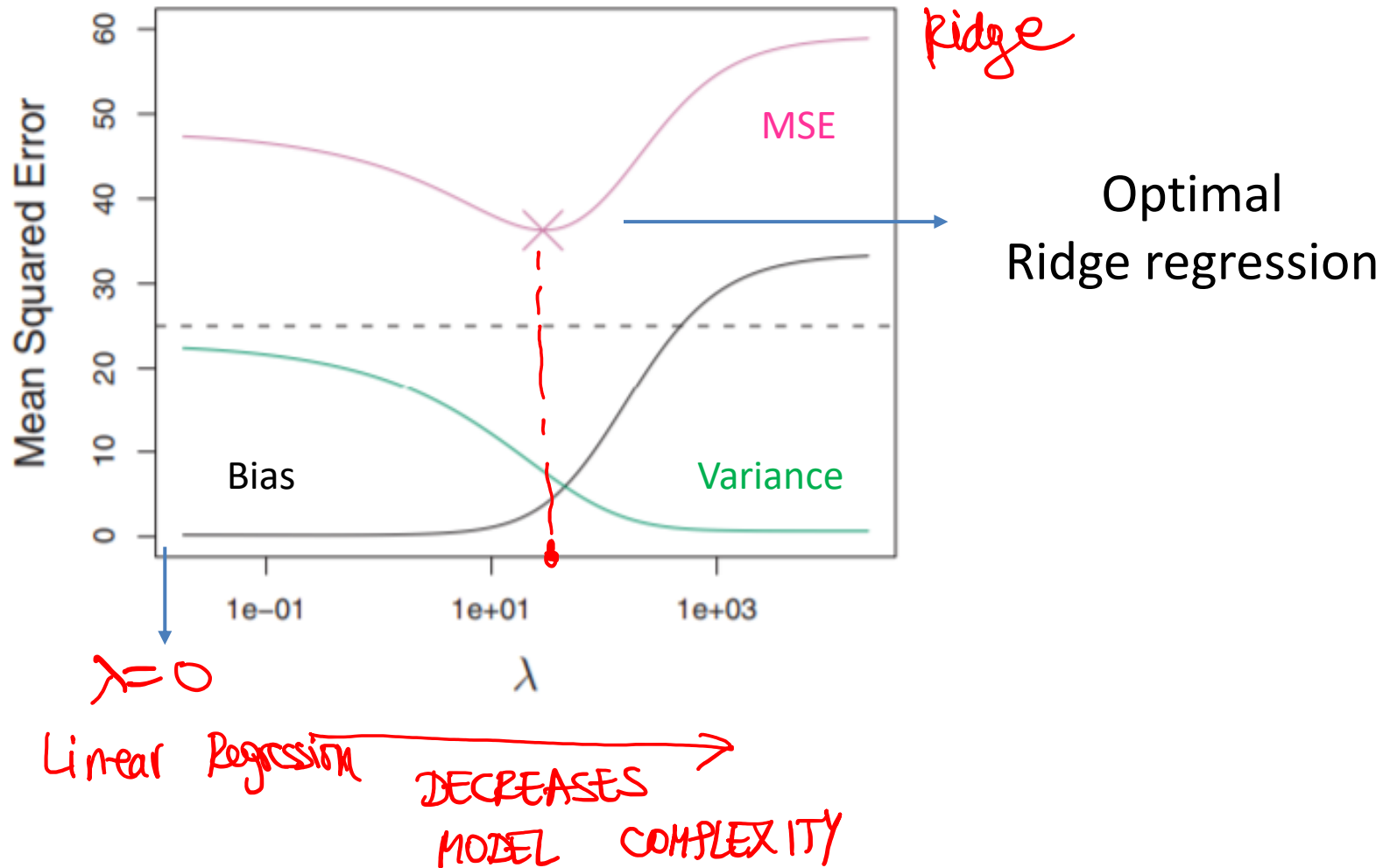
$\lambda \uparrow \Rightarrow \theta_j \downarrow$

DECREASE MODEL COMPLEXITY

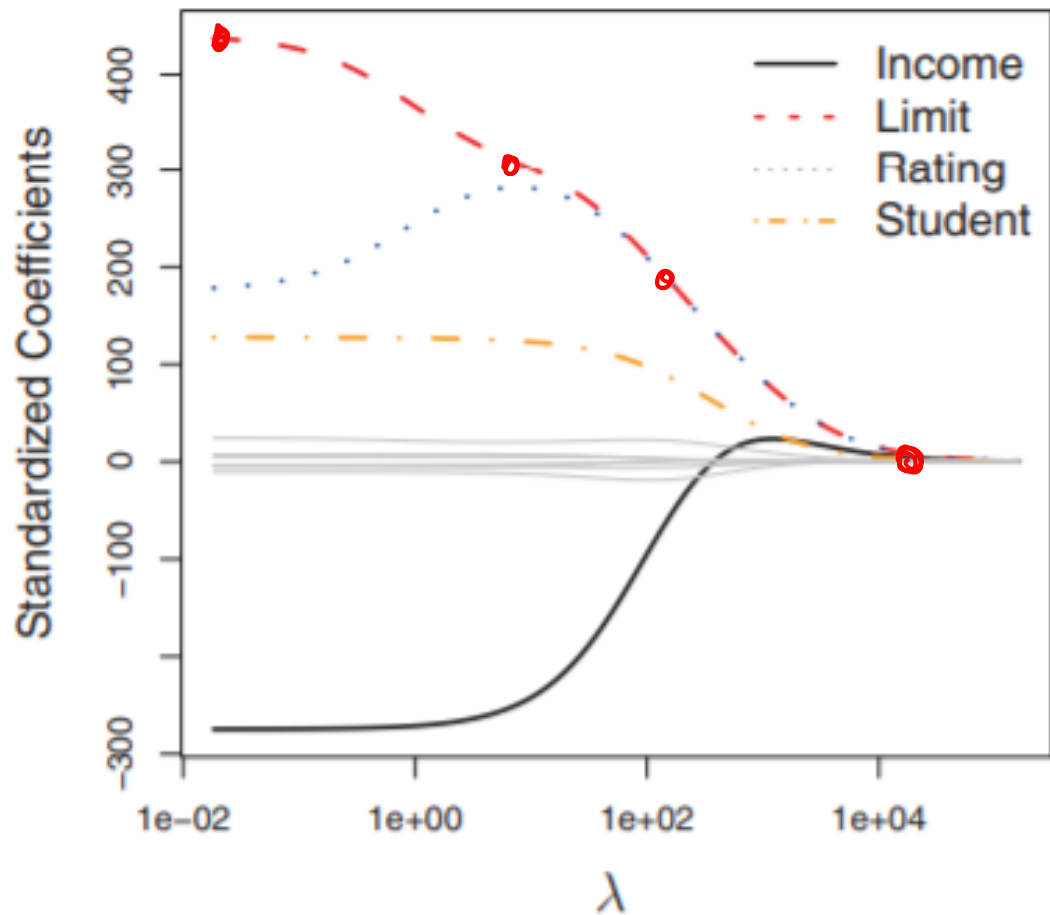
$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

$J(\theta)$ convex \Rightarrow HELPS GRADIENT DESCENT

Bias-Variance Tradeoff



Coefficient shrinkage



Predict credit card balance

GD for Ridge Regression

Min MSE

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

Handwritten notes: $\frac{\partial \theta_j^2}{\partial \theta_j} = 2\theta_j$ and $\theta_1^2 + \theta_2^2 + \dots + \theta_j^2 + \dots + \theta_d^2$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^N (h_{\theta}(x_i) - y_i) \cdot x_{ij} + \lambda \theta_j$$

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

$$= \theta_j (1 - \alpha \lambda) - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i) x_{ij}, \quad j=1, \dots, d$$

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

TRAINING EX

$$x_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{id} \end{bmatrix}$$

$$y_i$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

Lasso Regression

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

\sim MSE

REG
L1 norm $\|\theta\|_1$
FEATURE SELECTION
METHOD

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

Alternative Formulations

- Ridge

- L2 Regularization

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

MSE

ball of radius ϵ

- Lasso

- L1 regularization

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

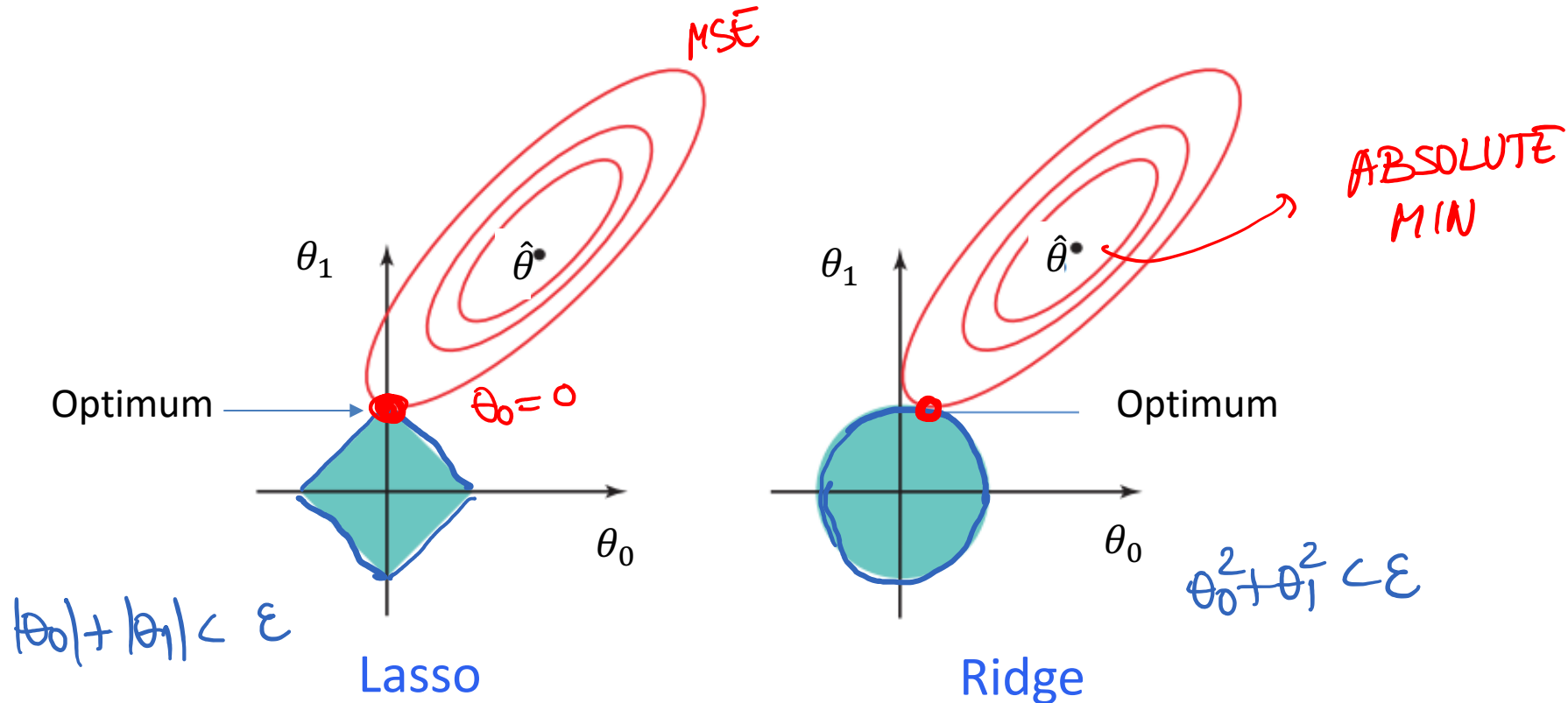
MSE

polytope

$|\theta_1| + |\theta_2| < \epsilon$

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 *MSE* (regression or classification)

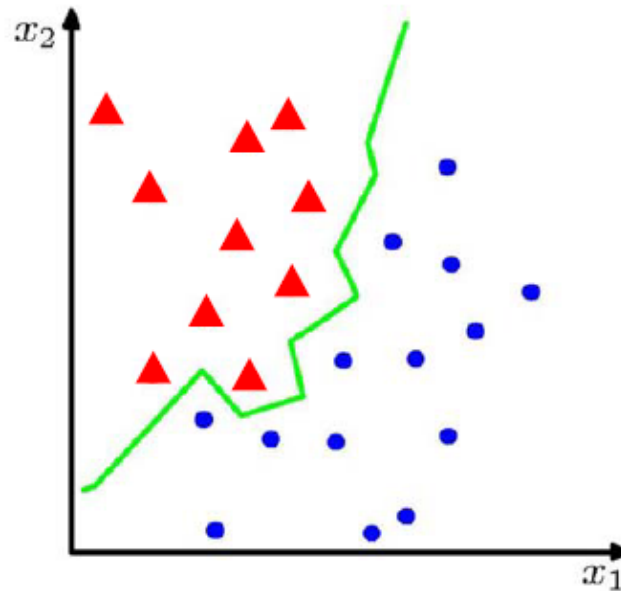
- Ridge

- + CONVEX OBJECTIVE
- + GLOBAL OPT
- + CAN USE GRADIENT DESCENT
- + HAS CLOSED FORM

- Lasso

- + LARGE FEATURE SPACE
(FEATURE SELECTION)
- NEED ADAPT GRADIENT
DESCENT

Classification



Binary or
discrete

- Suppose we are given a training set of N observations

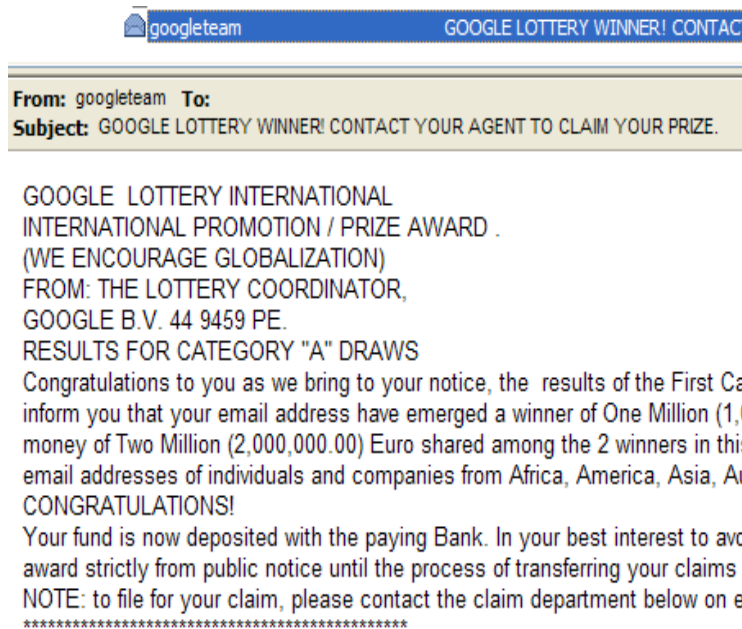
$$\{x_1, \dots, x_N\} \text{ and } \{y_1, \dots, y_N\}, x_i \in R^d, y_i \in \{0, 1\}$$

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

$$f(x_i) = y_i$$

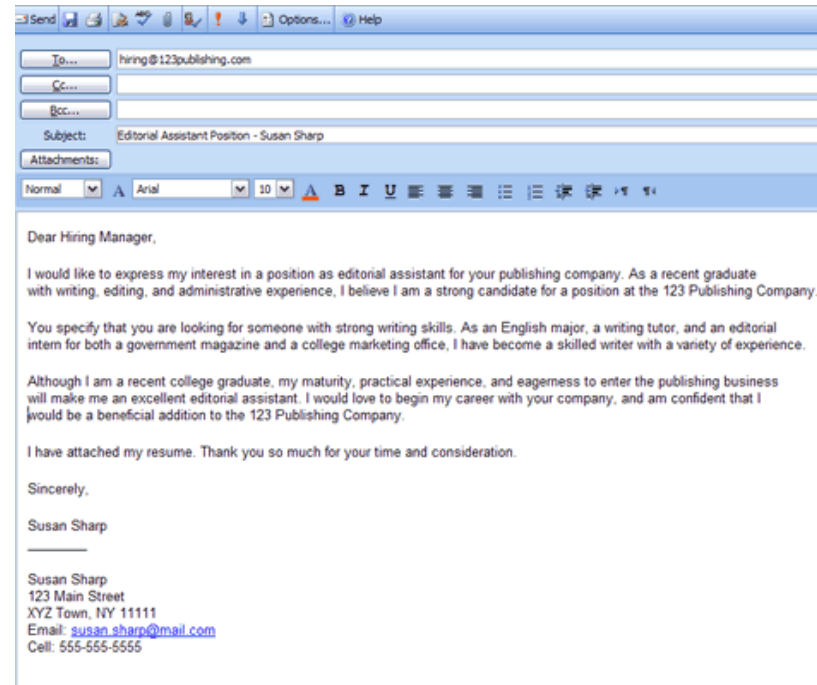
Example 1: Binary classification

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

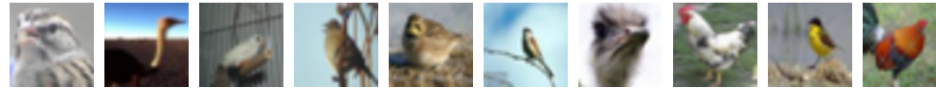
airplane



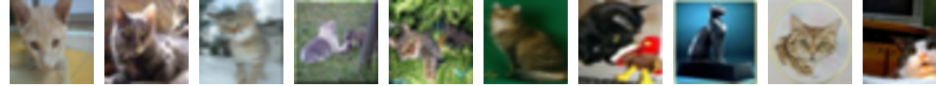
automobile



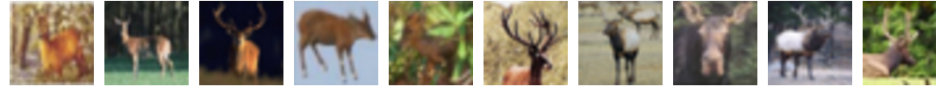
bird



cat



deer



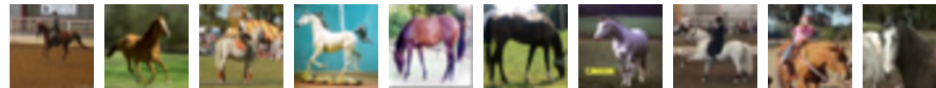
dog



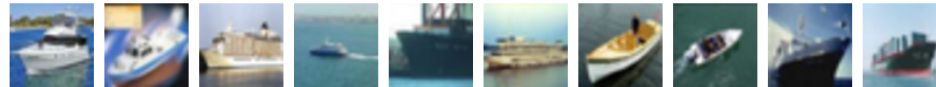
frog



horse



ship



truck



Multi-class classification

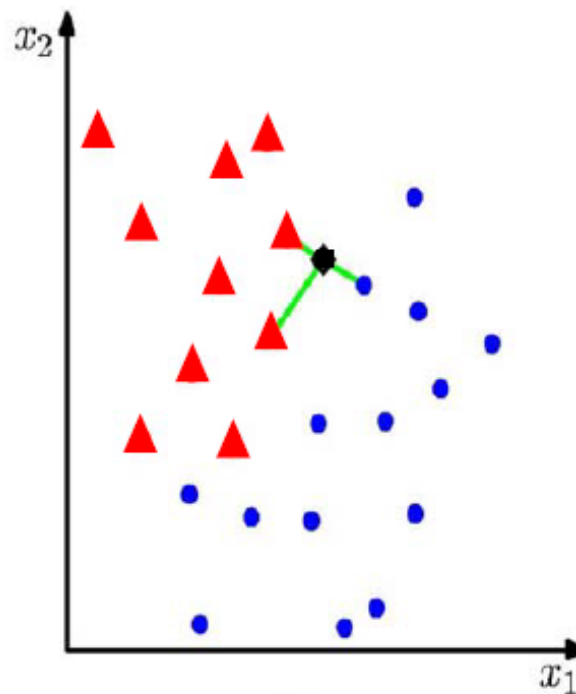
K Nearest Neighbour (K-NN) Classifier

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

$$\sqrt{\sum_{i=1}^k (x_i - y_i)^2} \quad \mathcal{L}_2$$

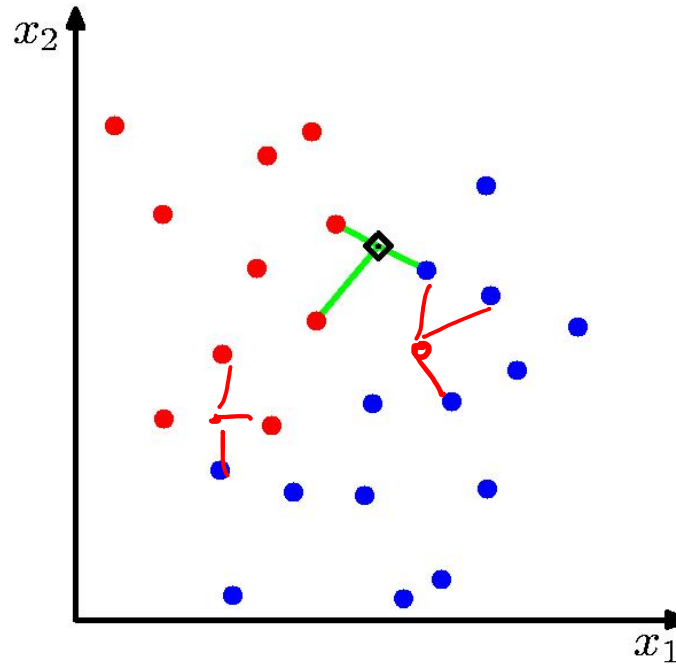
- Manhattan Distance

$$\sum_{i=1}^k |x_i - y_i| \quad \mathcal{L}_1$$

- Minkowski Distance

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

kNN

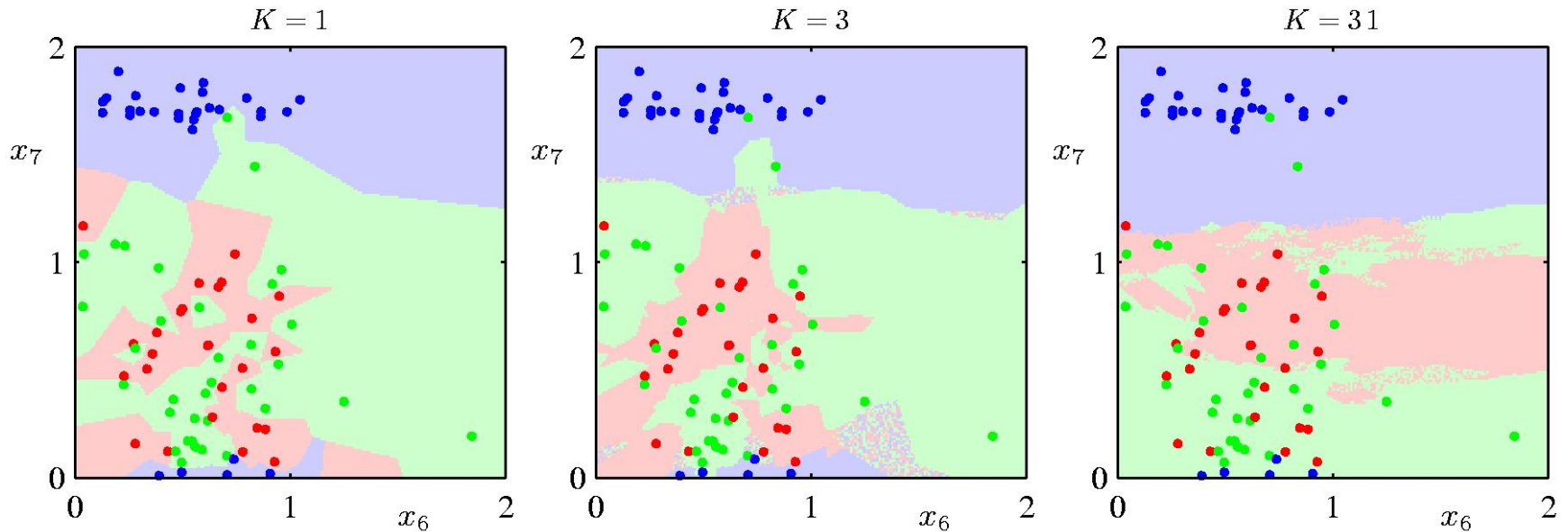


- Algorithm (to classify point x)
 - Find k nearest points to x (according to distance metric)
 - Perform majority voting to predict class of x

ISSUES:

- DOES NOT LEARN
- INSTANCE LEARNER

K-Nearest-Neighbours for Multi-class Classification



Vote among multiple classes