# DS 4400

# Machine Learning and Data Mining I
# Spring 2022

Alina Oprea

Associate Professor

Khoury College of Computer Science
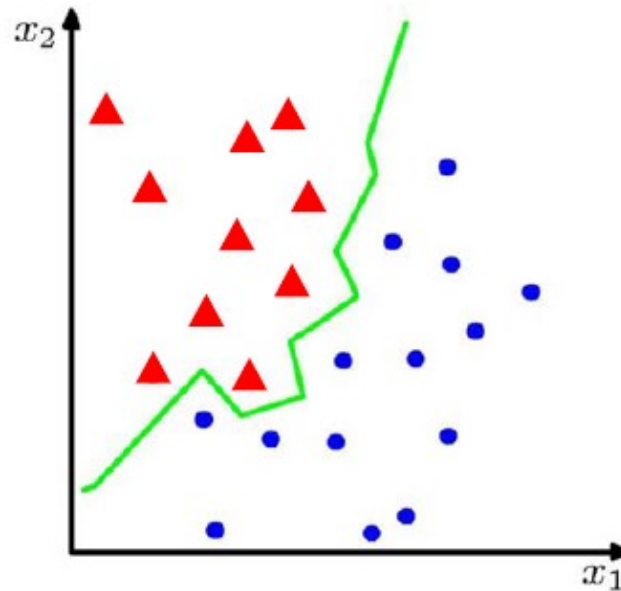
Northeastern University

February 14 2022

# Outline

- Classification
  - K Nearest Neighbors (kNN)
- Cross validation
  - K-fold cross validation
  - Leave one out cross validation
- Linear classifiers
- Logistic regression
  - Classification based on probability
  - Objective for logistic regression

# 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
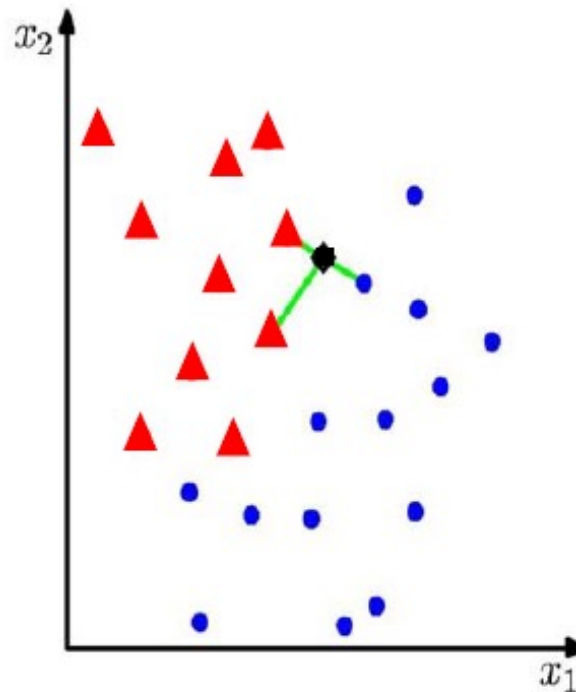
$$f(x_i) = y_i$$
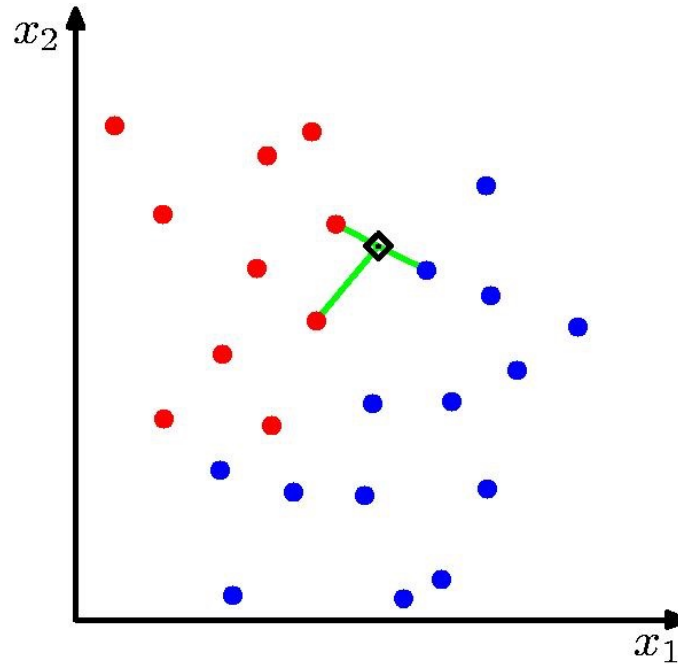
# K Nearest Neighbour (K-NN) Classifier

## 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
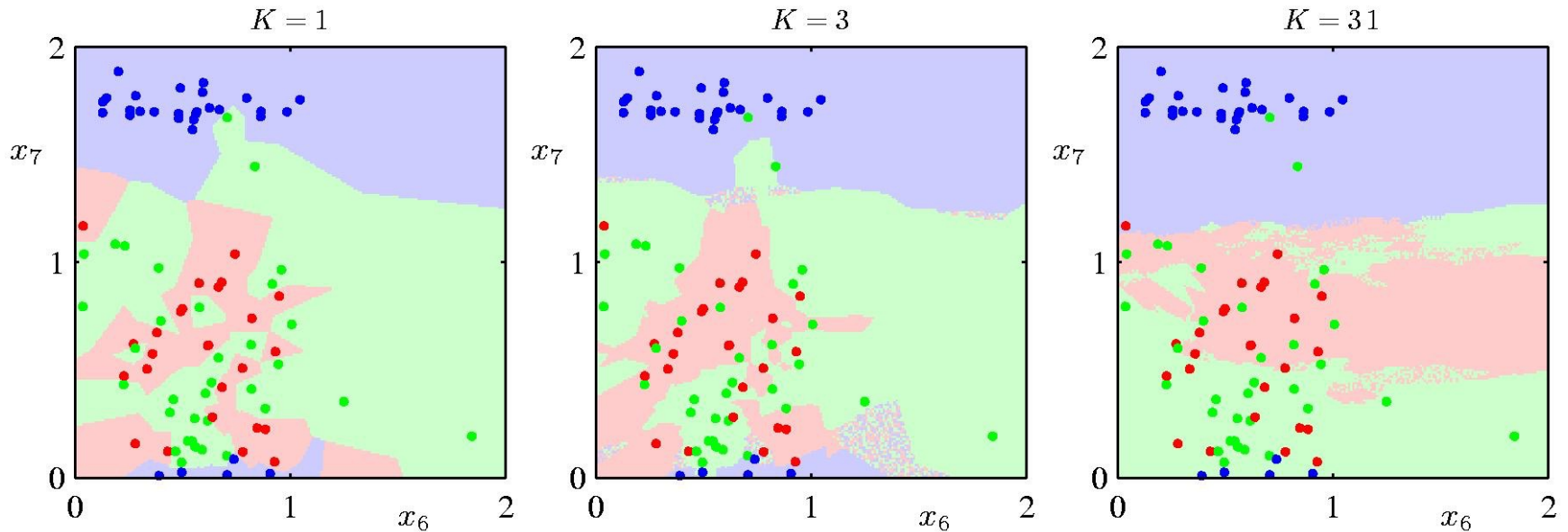
e.g. K = 3

• applicable to multi-class case

# 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
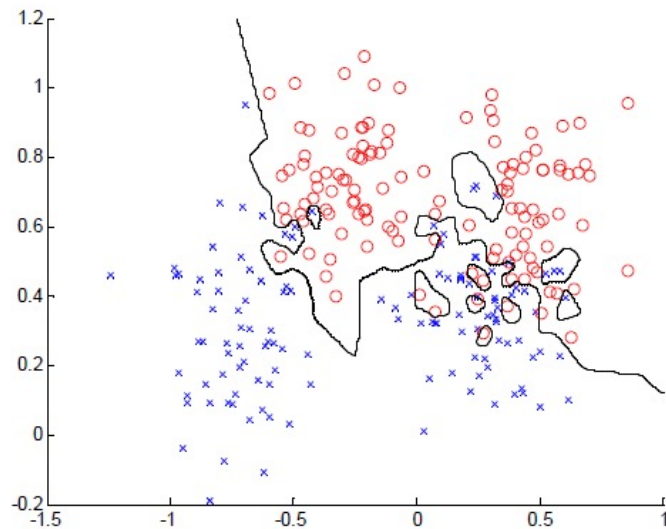
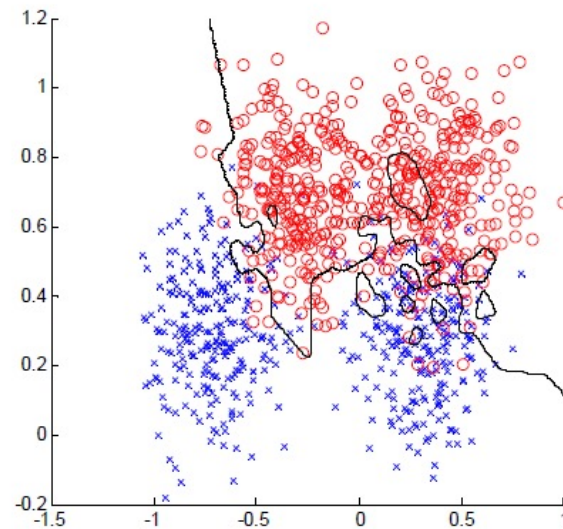

Vote among multiple classes
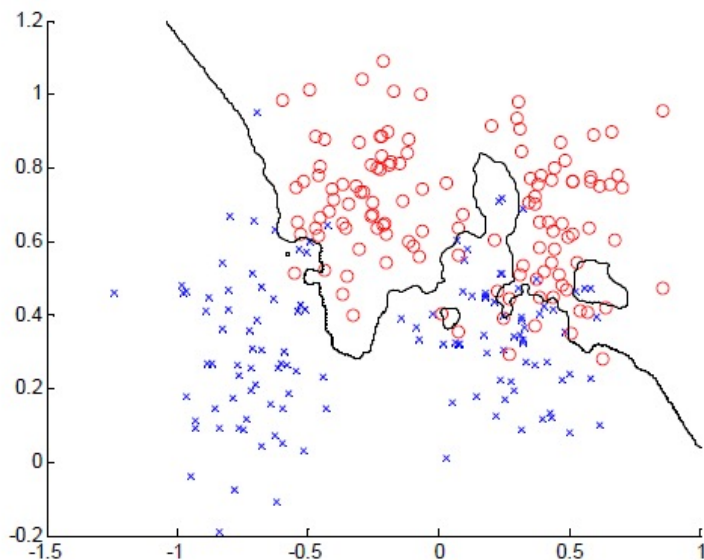
# K = 1



Training data

error = 0.0

Testing data

error = 0.15

How to choose k (hyper-parameter)?

# K = 3



Training data

error = 0.0760

Testing data

error = 0.1340

How to choose k (hyper-parameter)?

# K = 7



Training data

Testing data

error = 0.1320

error = 0.1110

How to choose k (hyper-parameter)?

# Bias-Variance Tradeoff for kNN



K decreases

# How Overfitting Affects Prediction



- How to pick hyper-parameters without access to testing data?
- Goal: Reduce overfitting and variance

Important: Do not use testing data for hyper-parameter selection even if it is available

# 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

K=1

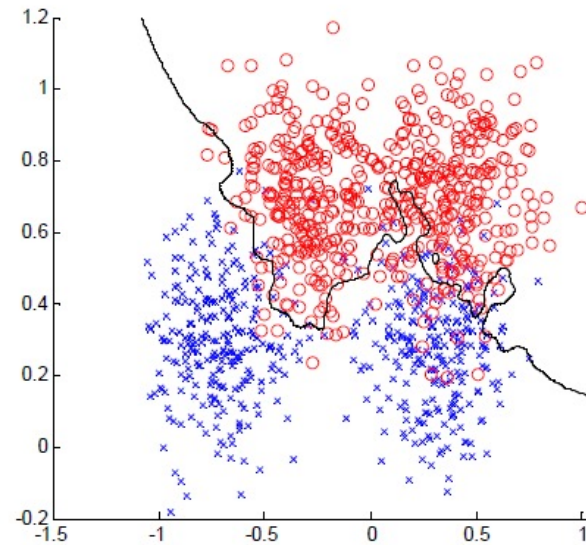| Train set 1 | Validation set 1 | Error |

...

| Train set 10 | Validation set 10 | Error |

Avg Validation Error

K=7

| Train set 1 | Validation set 1 | Error |

...

| Train set 10 | Validation set 10 | Error |

Avg Validation Error

Minimum error!

# Cross Validation



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
- Used for small training sets

   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

# Linear classifiers

- A **hyperplane** partitions space into 2 half-spaces
  - Defined by the normal vector $\theta \in \mathbb{R}^{d+1}$
    - $\theta$ is orthogonal to any vector lying on the hyperplane



- Consider classification with +1, -1 labels …

# Linear Classifiers

- **Linear classifiers**: represent decision boundary by hyperplane

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad x^{\mathsf{T}} = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$$



$h_\theta(x) = f(\theta^T x)$ linear classifier
- If $\theta^T x > 0$ classify "Class 1"
- If $\theta^T x < 0$ classify "Class 0"

All the points x on the hyperplane satisfy: $\theta^T x = 0$

# Linear vs Non-Linear Classifiers

# Classification Based on Probability

- Instead of just predicting the class, give the *probability of the instance being in that class*
  - Learn $P(Y|X)$
- Consider binary classifier with classes 0 and 1
  - $P(Y = 1|X) + P(Y = 0|X) = 1$
  - Sufficient to learn $P(Y = 1|X)$
- Advantages: interpretability and confidence of output

# Logistic Regression

- Setup
  - Training data: $\{x_i, y_i\}$, for $i = 1, \ldots, N$
  - Labels: $y_i \in \{0,1\}$
- Goals
  - Learn $P(Y = 1|X = x)$
- Highlights
  - Probabilistic output
  - At the basis of more complex models (e.g., neural networks)
  - Supports regularization (Ridge, Lasso)
  - Can be trained with Gradient Descent

# Interpretation of Model Output

$h_{\boldsymbol{\theta}}(\boldsymbol{x})$ = estimated $\quad P(Y = 1 | X; \theta)$

Example: Cancer diagnosis from tumor size

$$\boldsymbol{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$$

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = 0.7$$

→ Tell patient that 70% chance of tumor being malignant

Note that: $\quad P(Y = 0 | X; \theta) + P(Y = 1 | X; \theta) = 1$

Therefore, $\quad P(Y = 0 | X; \theta) = 1 - P(Y = 1 | X; \theta)$

# Logistic Regression

- Takes a probabilistic approach to learning discriminative functions (i.e., a classifier)

- $h_{\boldsymbol{\theta}}(\boldsymbol{x})$ should give $P(Y = 1 | X; \theta)$
  - Want $0 \leq h_{\boldsymbol{\theta}}(\boldsymbol{x}) \leq 1$

- Logistic regression model:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}\right)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}}}$$

Logistic / Sigmoid Function



$g(z)$

# LR is a Linear Classifier!

- Predict $Y = 1$ if:
  - $\text{P}[Y = 1 | X = x; \theta] > \text{P}[Y = 0 | X = x; \theta]$
  - $\text{P}[Y = 1 | X = x; \theta] > \frac{1}{2}$

$$\frac{1}{1 + e^{-\theta^T x}} > \frac{1}{2}$$

- Equivalent to:

  - $e^{\theta_0 + \sum_{j=1}^{d} \theta_j x_j} > 1$

  - $\theta_0 + \sum_{j=1}^{d} \theta_j x_j > 0$

Logistic Regression is a linear classifier!

# Logistic Regression

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}\right)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$g(z)$

$\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}$ should be large <u>negative</u> values for negative instances

$\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}$ should be large <u>positive</u> values for positive instances

- Assume a threshold and...
  - Predict $Y = 1$ if $h_{\boldsymbol{\theta}}(\boldsymbol{x}) \geq 0.5$
  - Predict $Y = 0$ if $h_{\boldsymbol{\theta}}(\boldsymbol{x}) < 0.5$

y = 1

$\theta$

y = 0

Logistic Regression is a linear classifier!

# Maximum Likelihood Estimation (MLE)

Given training data $X = \{x_1, \dots, x_N\}$ with labels
$Y = \{y_1, \dots, y_N\}$

What is the likelihood of training data for parameter $\theta$?

Define likelihood function

$$Max_\theta \, L(\theta) = P[Y|X; \theta]$$

Assumption: training points are independent

$$L(\theta) = \prod_{i=1}^{N} P[Y = y_i | X = x_i; \theta]$$

General probabilistic method for classifier training

# Log Likelihood

- Max likelihood is equivalent to maximizing log of likelihood

$$L(\theta) = \prod_{i=1}^{N} P[Y = y_i | X = x_i; \theta]$$

$$\log L(\theta) = \sum_{i=1}^{N} \log P[Y = y_i | X = x_i; \theta]$$

- They both have the same maximum $\theta_{MLE}$

# Maximum Likelihood Estimation (MLE)

Given training data $X = \{x_1, \dots, x_N\}$ with labels
$Y = \{y_1, \dots, y_N\}$

What is the likelihood of training data for parameter $\theta$?

Define likelihood function

$$Max_\theta \, L(\theta) = P[Y|X; \theta]$$

Assumption: training labels are conditionally independent

$$L(\theta) = \prod_{i=1}^{N} P[Y = y_i | X = x_i; \theta]$$

General probabilistic method for parameter estimation

# MLE for Logistic Regression

$$P(Y = y_i | X = x_i; \theta) = h_\theta(x_i)^{y_i} \left(1 - h_\theta(x_i)\right)^{1-y_i}$$

$$\theta_{MLE} = \text{argmax}_\theta \sum_{i=1}^{N} \log P[Y = y_i | X = x_i; \theta]$$

$$= \text{argmax}_\theta \sum_{i=1}^{N} y_i \log h_\theta(x_i) + (1 - y_i) \log \left(1 - h_\theta(x_i)\right)$$

Logistic regression objective

$$\min_\theta J(\theta)$$

$$J(\theta) = -\sum_{i=1}^{N} [y_i \log h_\theta(x_i) + (1 - y_i) \log \left(1 - h_\theta(x_i)\right)]$$

# Cross-Entropy Objective

$$J(\theta) = -\sum_{i=1}^{N} [y_i \log h_\theta(x_i) + (1 - y_i) \log(1 - h_\theta(x_i))]$$

- Cost of a single instance:

$$\text{cost}(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases}$$

- Can re-write objective function as

$$J(\boldsymbol{\theta}) = \sum_{i=1}^{n} \text{cost}\left( h_\theta(x_i), y_i \right)$$

Cross-entropy loss

# Review

- K nearest neighbors is the first example of classifier
  - Instance learner
- Cross-validation should be performed to
  - Improve generalization and avoid over-fitting
  - Choose hyper parameters (k in kNN)
- Logistic regression is a linear classifier that predicts class probability
  - Cross-entropy objective derived with MLE
  - MLE: probabilistic method of maximum likelihood

# Acknowledgements

- Slides made using resources from:
  - Andrew Ng
  - Eric Eaton
  - David Sontag
- Thanks!