# Predictive Modelling

## Classification - K Nearest Neighbours

### Jonathan Mwaura

Khoury College of Computer Sciences

July 23, 2024

# Introduction

## Textbook

Reading: Chapter 4 of: Gareth James et al (2021) . An Introduction to Statistical Learning (2nd Edition) .

`https://www.statlearning.com/`

## Acknowledgements

These slides have been adapted from the following Professors:

1) Andrew Ng - Stanford
2) Eric Eaton - UPenn
3) David Sontag - MIT
4) Alina Oprea - Northeastern
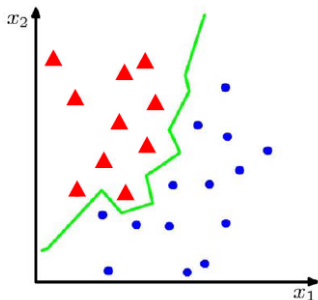
# Supervised Learning

**Problem Setting**

- Set of possible instances $\mathcal{X}$
- Set of possible labels $\mathcal{Y}$
- Unknown target function $f : \mathcal{X} \rightarrow \mathcal{Y}$
- Set of function hypotheses $H = \{h \mid h : \mathcal{X} \rightarrow \mathcal{Y}\}$

**Input**: Training examples of unknown target function f

$$\{x_i, y_i\}, \text{for } i = 1, \ldots, N$$

**Output**: Hypothesis $\hat{f} \in H$ that best approximates f

$$\hat{f}(x_i) \approx y_i$$

# Classification



Binary or discrete
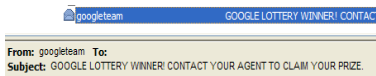
- Suppose we are given a training set of N observations

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

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

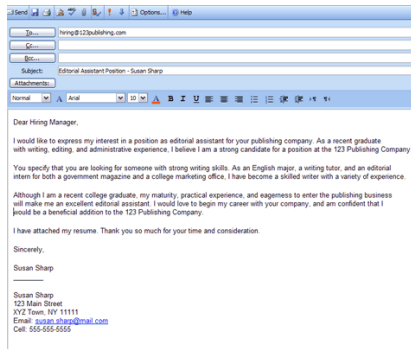$$f(x_i) = y_i$$

# 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)

The Roux Institute
at Northeastern University
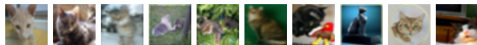
# Example 2: Multi-class classification

Image classification



airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

Multi-class classification
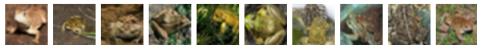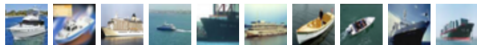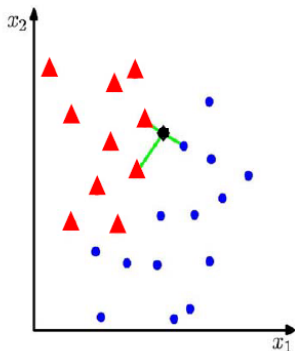
# K Nearest Neighbour (K-NN) Classifier

**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

The Roux Institute
at Northeastern University

# Distance Metrics

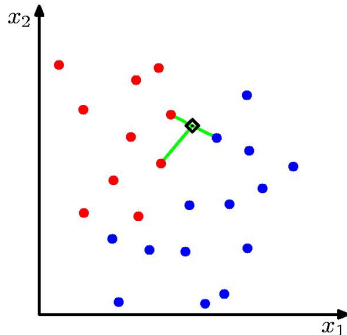- Euclidean Distance $\sqrt{\left(\sum_{i=1}^{k}(x_i - y_i)^2\right)}$

- 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 $x$
- Properties
  - Does not learn any model in training!
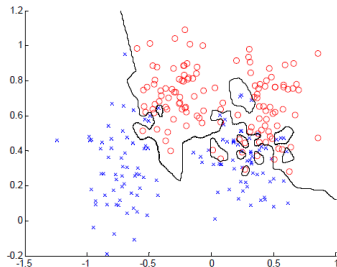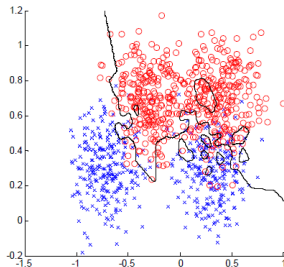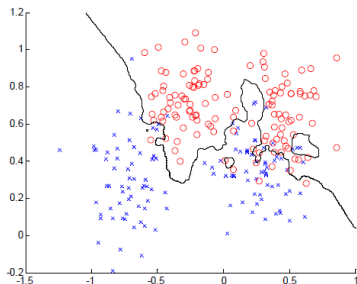  - Instance learner (needs all data at testing time)

# K = 1



Overfitting!

How to choose k (hyper-parameter)?

# K = 3



Training data

Testing data

error = 0.0760

error = 0.1340

How to choose k (hyper-parameter)?

The Roux Institute
at Northeastern University

22

# K = 7

Training data

Testing data



error = 0.1320

error = 0.1110

How to choose k (hyper-parameter)?

# Bias-Variance Tradeoff for kNN

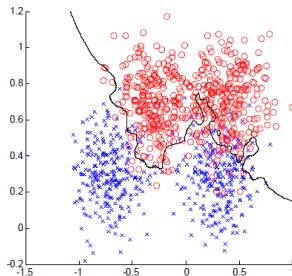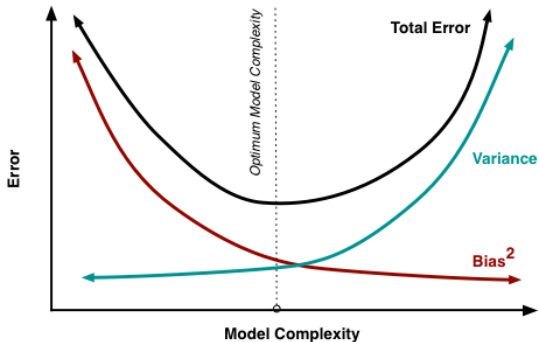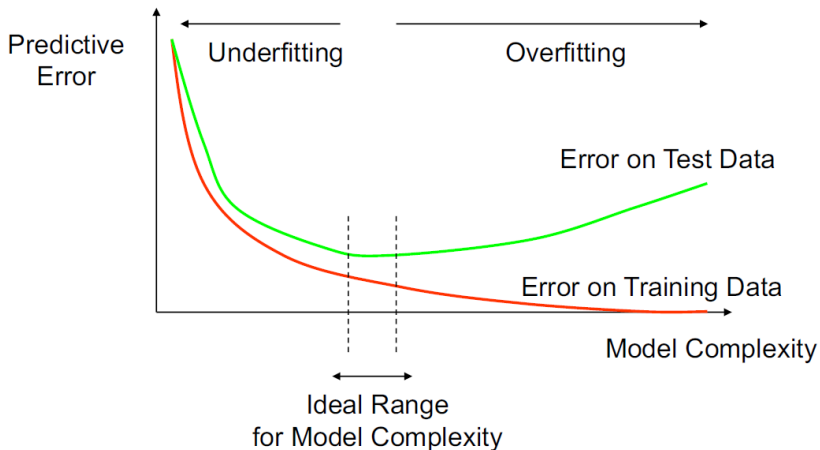# How Overfitting Affects Prediction



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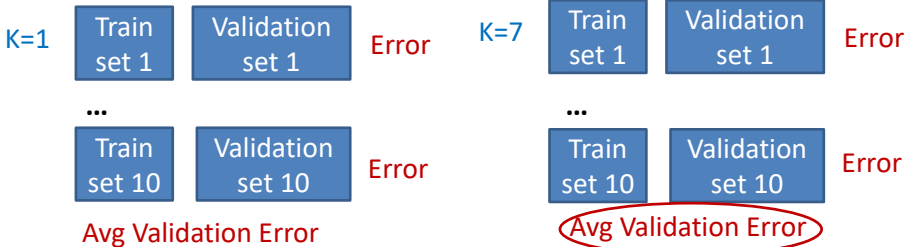
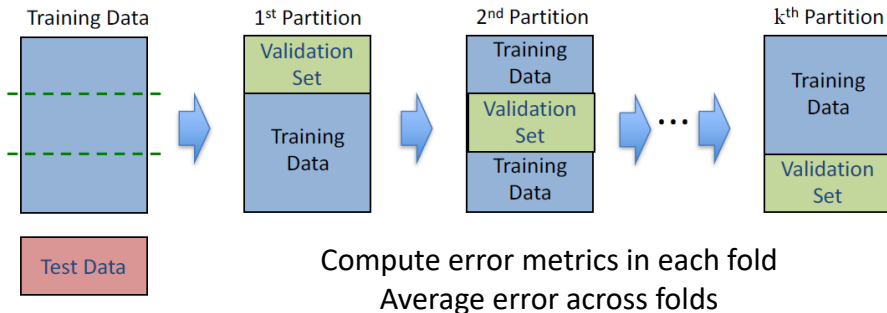# Cross Validation

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
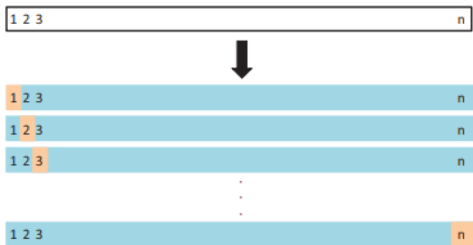


Compute error metrics in each fold
Average error across folds

## 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
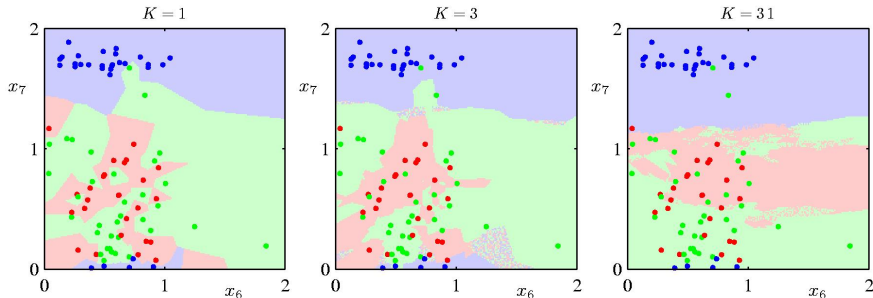
# Cross Validation



## 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

The Roux Institute
at Northeastern University

# K-Nearest-Neighbours for Multi-class Classification



Vote among multiple classes