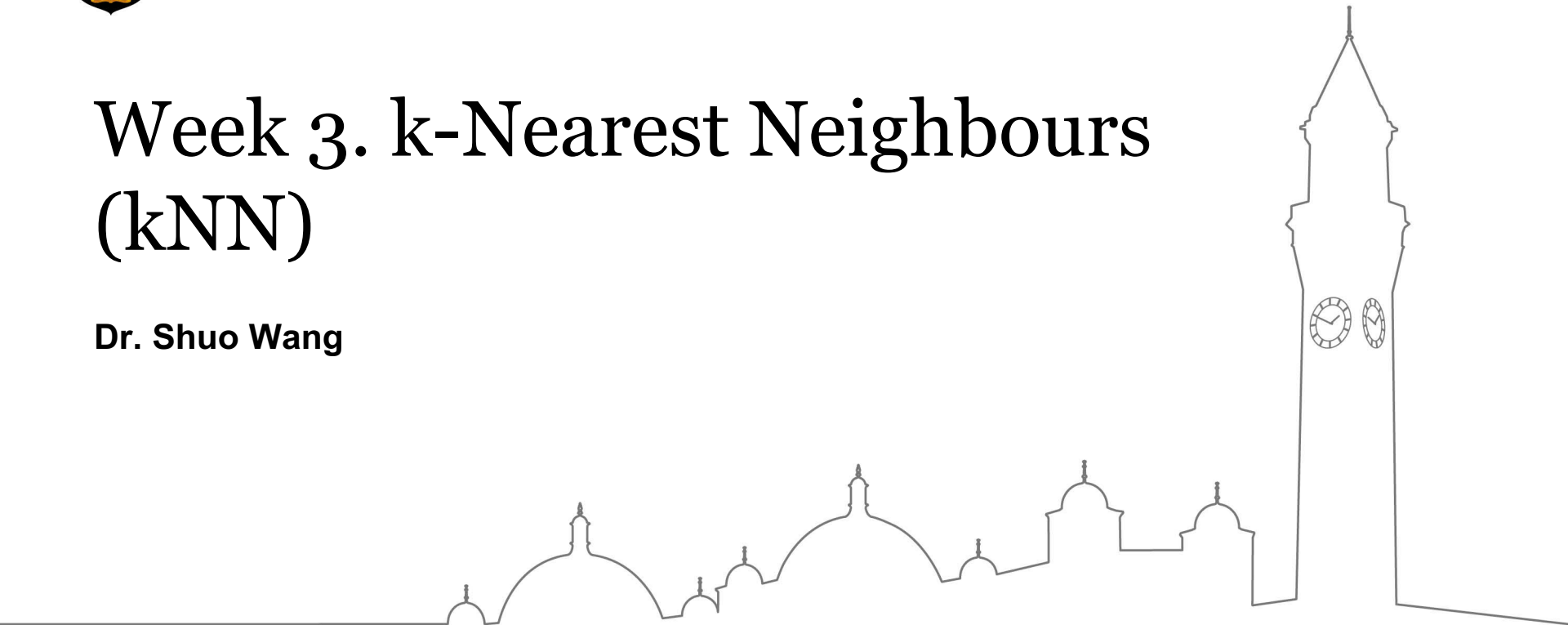




UNIVERSITY OF  
BIRMINGHAM

# Week 3. k-Nearest Neighbours (kNN)

**Dr. Shuo Wang**



# Overview

- Intuitive understanding
- The kNN algorithm
- Pros/cons

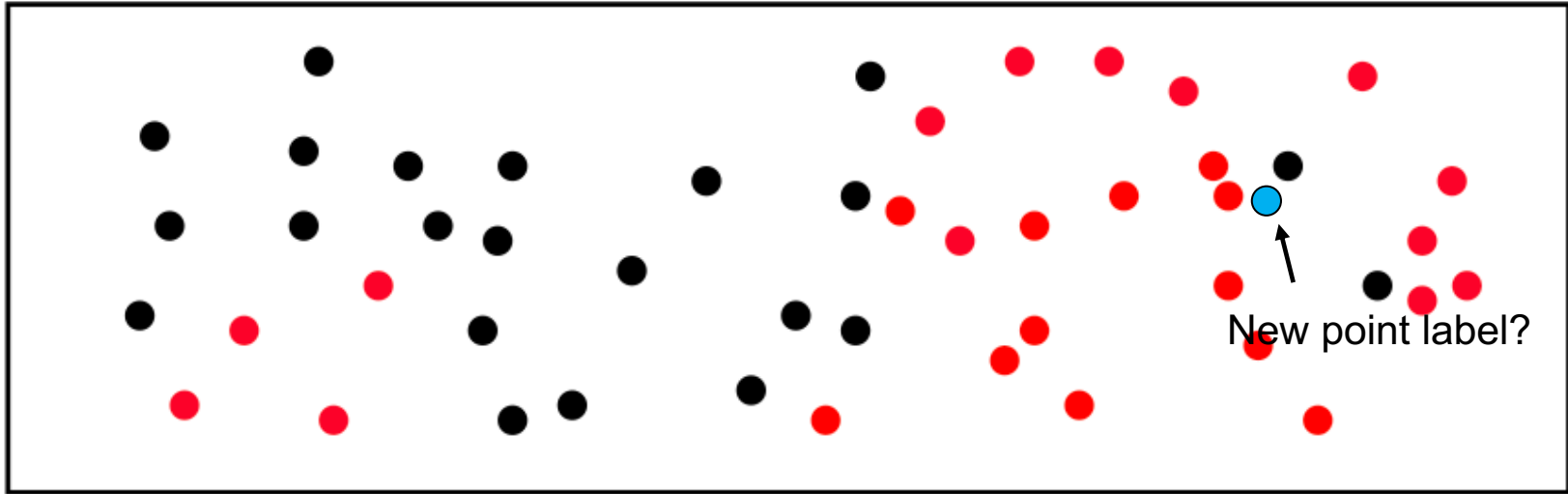


# kNN Basics

- Full name: k-Nearest Neighbours (kNN, or k-NN).
- It is **nonparametric**.
- It is **instance-based**.
- It is a **lazy** algorithm.

# Intuitive Understanding

Instead of approximating a model function  $f(x)$  globally, kNN approximates the label of a new point based on its **nearest** neighbours in training data.



Q1: How to choose  $k$ ? e.g. let  $k = 3$  to avoid issues.

Q2: how to we measure the distance between examples?

# Distance metrics (or similarity metrics)

Given two points  $\mathbf{x}^{(1)} = (x_1^{(1)}, x_2^{(1)}, \dots, x_d^{(1)})$ ,  $\mathbf{x}^{(2)} = (x_1^{(2)}, x_2^{(2)}, \dots, x_d^{(2)})$  in a d-dimensional space:

- Minkowski distance (or  $L^p$  norm)

$$D(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = \sqrt[p]{\sum_{i=1}^d |x_i^{(1)} - x_i^{(2)}|^p}$$

- When  $p=1$ , it becomes Manhattan distance

$$D(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = \sum_{i=1}^d |x_i^{(1)} - x_i^{(2)}|$$

- When  $p=2$ , it becomes Euclidean distance

$$D(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = \sqrt{\sum_{i=1}^d |x_i^{(1)} - x_i^{(2)}|^2}$$

# Distance metrics in kNN (common choice)

- Euclidean distance for real values (also called  $L^2$  distance).

$$D(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = \sqrt{\sum_{i=1}^d |x_i^{(1)} - x_i^{(2)}|^2}$$

- Hamming distance for discrete/categorical values, e.g.  $x \in \{rainy, sunny\}$ .

$$D(x^{(1)}, x^{(2)}) = \begin{cases} 0, & \text{if } x^{(1)} = x^{(2)} \\ 1, & \text{otherwise} \end{cases}$$

# kNN algorithm

Input: neighbour size  $k > 0$ , training set  $\{(\mathbf{x}^{(n)}, y^{(n)}): n = 1, 2 \dots N\}$ , a new unlabelled data  $\mathbf{x}^{(j)}$

for  $n = 1, 2 \dots N$  // each example in the training set

    Calculate  $D(\mathbf{x}^{(j)}, \mathbf{x}^{(n)})$  // distance between  $\mathbf{x}^{(j)}$  and  $\mathbf{x}^{(n)}$

    Select  $k$  training examples closest to  $\mathbf{x}^{(j)}$

Return  $y^{(j)}$  = the plurality vote of labels from the  $k$  examples.

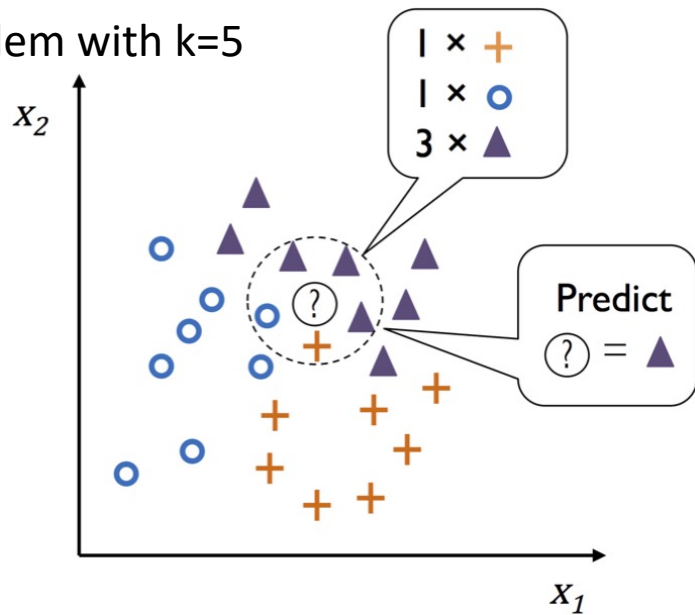
(classification) or

$y^{(j)}$  = average/median of the  $y$  values of the  $k$  examples.

(regression)

# Another visual example

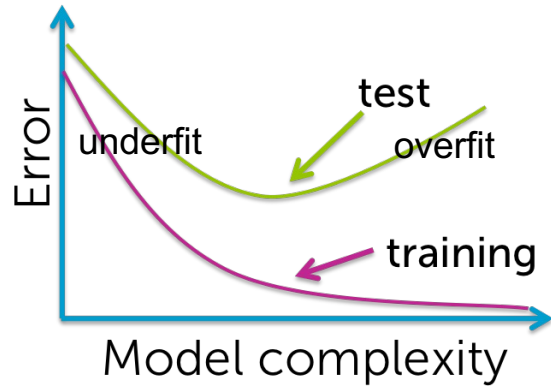
kNN for a 3-class problem with  $k=5$





# How to choose k?

- Recall: Overfitting and Underfitting



# The issue in numeric attribute ranges

- Attributes  $\mathbf{x} = (x_1, x_2, \dots, x_d)$  may have different ranges.
- What is the problem?
- Solutions?

# Normalisation and Standardization

- Method 1 Normalisation: Linearly scale the range of each attribute to be, e.g. in  $[0,1]$ .

$$x_{j\_new}^{(n)} = \frac{x_j^{(n)} - \min x_j}{\max x_j - \min x_j}$$

- Method 2 Standardization: Linearly scale each dimension to have 0 mean and variance 1 (by computing mean  $\mu$  and variance  $\sigma^2$ ).

$$x_{j\_new}^{(n)} = \frac{x_j^{(n)} - \mu_j}{\sigma_j}, \text{ where } \mu_j = \frac{1}{N} \sum_{n=1}^N x_j^{(n)}, \sigma_j = \sqrt{\frac{1}{N} \sum_{n=1}^N (x_j^{(n)} - \mu_j)^2}$$

# kNN algorithm with normalisation

Input: neighbour size  $k > 0$ , training set  $\{(\mathbf{x}^{(n)}, y^{(n)}): n = 1, 2 \dots N\}$ , a new unlabelled data  $\mathbf{x}^{(j)}$

Normalise/standardize  $\mathbf{x}^{(j)} \rightarrow \mathbf{x}_{new}^{(j)}$

for  $n = 1, 2 \dots N$  // each example in the training set

Normalise/standardize  $\mathbf{x}^{(n)} \rightarrow \mathbf{x}_{new}^{(n)}$

Calculate  $D(\mathbf{x}_{new}^{(j)}, \mathbf{x}_{new}^{(n)})$  // normalized/standardized distance

Select  $k$  training examples closest to  $\mathbf{x}^{(j)}$

Return  $y^{(j)}$  = the plurality vote of labels from the  $k$  examples.

(classification) or

$y^{(j)}$  = average/median of the  $y$  values of the  $k$  examples.

(regression)

# Fun project using kNN: where on earth is this photo from?

- Problem: where was this picture taken (country or GPS)?
- <http://graphics.cs.cmu.edu/projects/im2gps/>



- Get images from Flickr with gps info.
- Represent each image with meaningful features
- Apply kNN.



UNIVERSITY OF  
BIRMINGHAM

# Q/A

**Teams Channel:** [www.birmingham.ac.uk/](http://www.birmingham.ac.uk/)

**Office Hour:** [\[faculty or individual email\]@bham.ac.uk](mailto:[faculty or individual email]@bham.ac.uk)

