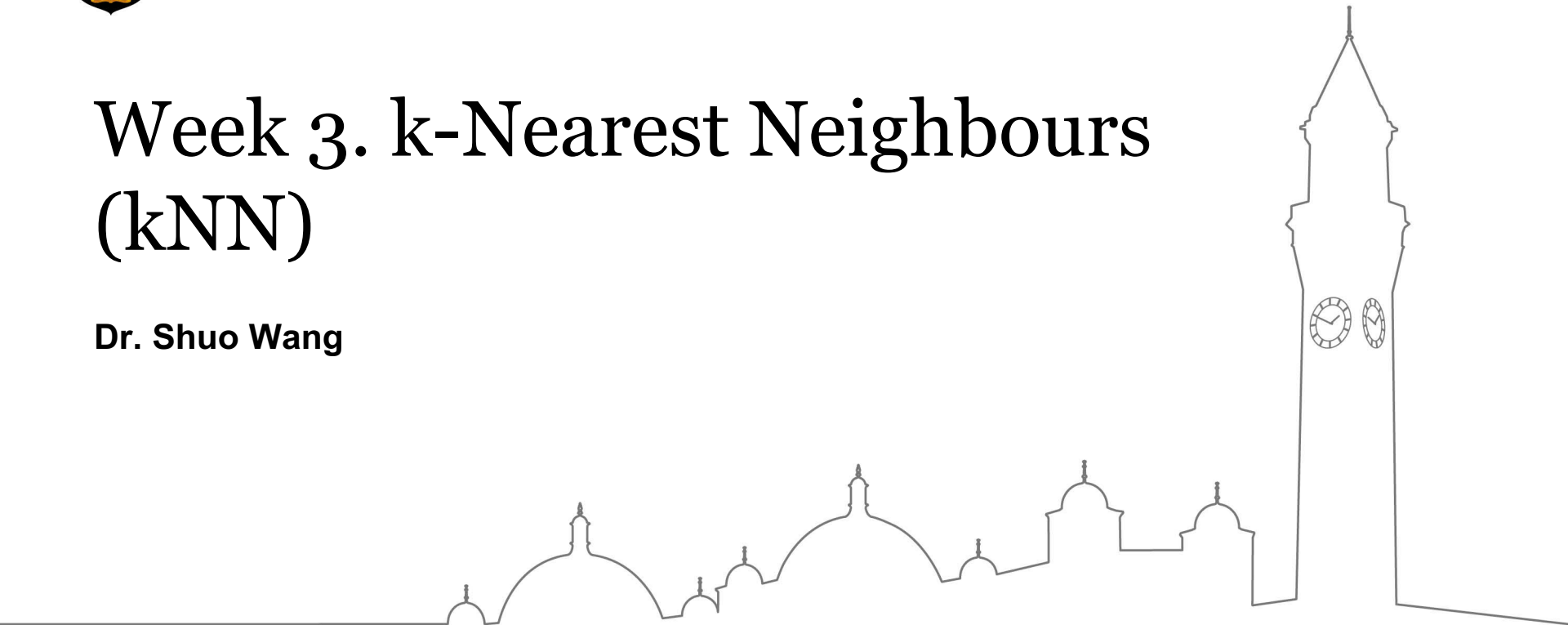




UNIVERSITY OF  
BIRMINGHAM

# Week 3. k-Nearest Neighbours (kNN)

**Dr. Shuo Wang**



# Overview

- Intuitive understanding
- The kNN algorithm
- Pros/cons

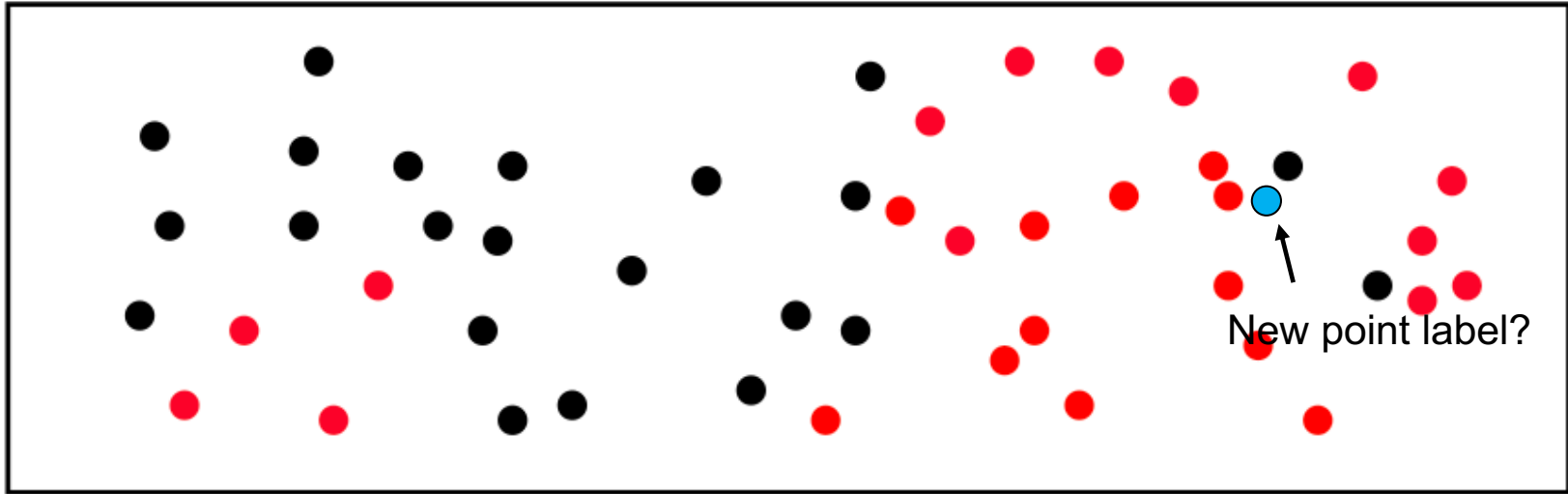


# kNN Basics

- Full name: k-Nearest Neighbours (kNN, or k-NN).
- It is **nonparametric**.  
No assumption about the functional form of the model.
- It is **instance-based**.  
The prediction is based on a comparison of a new point with data points in the training set, rather than a model.
- It is a **lazy** algorithm.  
No explicit training step. Defers all the computation until prediction.
- Can be used for both classification and regression problems.

# 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)

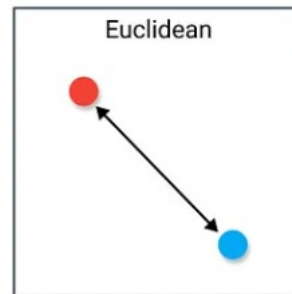
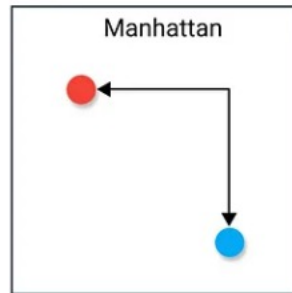
$$D(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = \sqrt[p]{\sum_{i=1}^d |x_i^{(1)} - x_i^{(2)}|^p}$$

- 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

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



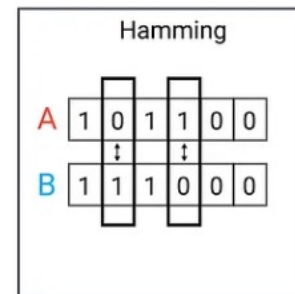
# Distance metrics in kNN (common choice)

- 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)

# Check your understanding

Consider a binary problem (lemon or orange) with 2 dimensions (height and width) with following training examples:

- $\mathbf{x}^{(1)} = (6, 6)$ ,  $y^{(1)} = \text{orange}$
- $\mathbf{x}^{(2)} = (8, 10)$ ,  $y^{(2)} = \text{lemon}$
- $\mathbf{x}^{(3)} = (7, 6)$ ,  $y^{(3)} = \text{orange}$

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

New example

- $\mathbf{x}^{(4)} = (8, 7)$ ,  $y^{(4)} = ?$  Using  $k=1$  nearest neighbour, and Euclidean distance
- $D(\mathbf{x}^{(4)}, \mathbf{x}^{(1)}) = \sqrt{\sum_{i=1}^d |x_i^{(4)} - x_i^{(1)}|^2} = \sqrt{(x_1^{(4)} - x_1^{(1)})^2 + (x_2^{(4)} - x_2^{(1)})^2} = \sqrt{(8 - 6)^2 + (7 - 6)^2} = \sqrt{5}$
- Can you calculate  $D(\mathbf{x}^{(4)}, \mathbf{x}^{(2)})$  and  $D(\mathbf{x}^{(4)}, \mathbf{x}^{(3)})$ , and see which point  $\mathbf{x}^{(4)}$  is closest to? 3 and  $\sqrt{2}$



# Check your understanding

Consider a regression problem (lemon' weight) with 2 dimensions (height and width) with following training examples:

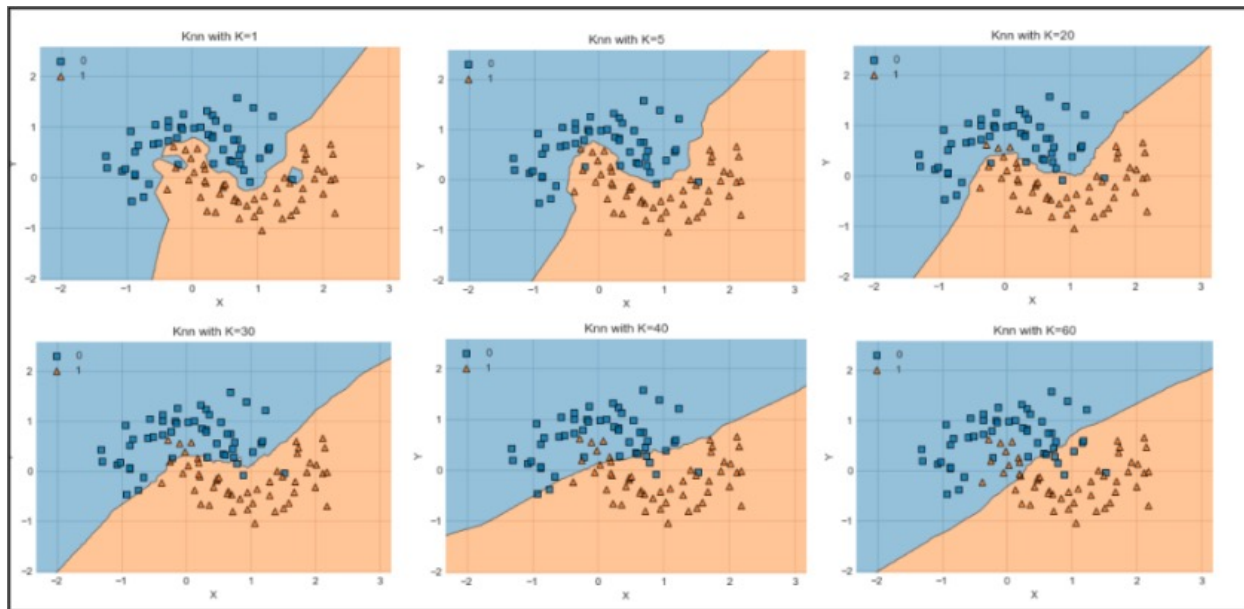
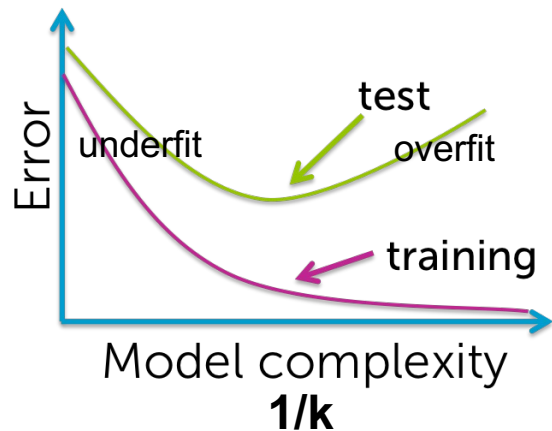
- $\mathbf{x}^{(1)} = (6,6), y^{(1)} = 10$
- $\mathbf{x}^{(2)} = (8,10), y^{(2)} = 20$
- $\mathbf{x}^{(3)} = (7,6), y^{(3)} = 15$

New example

- $\mathbf{x}^{(4)} = (8,7), y^{(4)} = ?$
- $D(\mathbf{x}^{(4)}, \mathbf{x}^{(1)}) = \sqrt{5}, D(\mathbf{x}^{(4)}, \mathbf{x}^{(2)}) = 3, D(\mathbf{x}^{(4)}, \mathbf{x}^{(3)}) = \sqrt{2}$
- If  $k = 2$ , what is the label of  $\mathbf{x}^{(4)}$ ?
- $y^{(4)} = (10+15)/2 = 12.5$

# How to choose k?

- Recall: Overfitting and Underfitting
- k changes model complexity: smaller k  $\rightarrow$  higher complexity



# How to choose k?

- Small k -> small neighborhood -> high complexity -> may overfit
- Large k -> large neighborhood -> low complexity -> may underfit
- Practitioners often choose k between 3 – 15, or  $k < \sqrt{N}$  (N is the number of training examples).
- Refer to “model selection/evaluation” to be learnt next week.

# The issue in numeric attribute ranges

- Attributes  $x = (x_1, x_2, \dots, x_d)$  may have different ranges.
- The attribute with a larger range is treated as more important by the kNN algorithm (*some learning bias is embedded!*)
- It can affect the performance if you don't want to treat attributes differently.
- For example, if  $x_1$  is in  $[0, 2]$  (e.g. height), and  $x_2$  is in  $[0, 100]$  (e.g. age),  $x_2$  will affect the distance more.
- 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  $\sigma^2$ ).

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

# Example

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

- Consider a dataset with 2 dimensions (ie. Attributes), where  $x_1$  represents the age of a patient and  $x_2$  represents the body weight. The output  $y \in \{normal, abnormal\}$ .

| Patient            | $x_1$ | $x_2$ | $y$ |
|--------------------|-------|-------|-----|
| $\mathbf{x}^{(1)}$ | 14    | 70    | n   |
| $\mathbf{x}^{(2)}$ | 12    | 90    | a   |
| $\mathbf{x}^{(3)}$ | 15    | 66    | n   |

- Normalize each attribute of  $\mathbf{x}^{(1)}$  to  $[0,1]$ .

- $$x_{1\_new}^{(1)} = \frac{x_1^{(1)} - 12}{15 - 12} = \frac{14 - 12}{15 - 12} = 0.667$$

- $$x_{2\_new}^{(1)} = \frac{x_2^{(1)} - 66}{90 - 66} = \frac{70 - 66}{90 - 66} = 0.167$$

- $\mathbf{x}^{(1)} : (14, 70) \rightarrow (0.667, 0.167)$ ,  $\mathbf{x}^{(2)}$  and  $\mathbf{x}^{(3)}$ ?

- $\mathbf{x}^{(2)} : (12, 90) \rightarrow (0, 1)$

- $\mathbf{x}^{(3)} : (15, 66) \rightarrow (1, 0)$

$$\mathbf{x}^{(4)} : (16, 64)?$$

# kNN algorithm with normalization/standardization

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)

# Pros/cons

- kNN is a nonparametric, instance-based, lazy algorithm.
- Need to specify the distance function and pre-define k value.
- Easy to implement and interpret.
- It can approximate complex functions, so it has very good accuracy.
- It has to store all training data (large memory space), and calculate distance of each training example to the new example.

There are smarter ways to store and use training data, e.g. KD-trees, remove redundant data.

- It can be sensitive to noise, especially when k is small.
- Its performance is degraded greatly as data dimension increases. (curse of dimensionality)

As the volume grows larger, the “neighbors” become further apart and not so close anymore. The prediction thus becomes less accurate.



# 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)

