

k-nearest neighbours

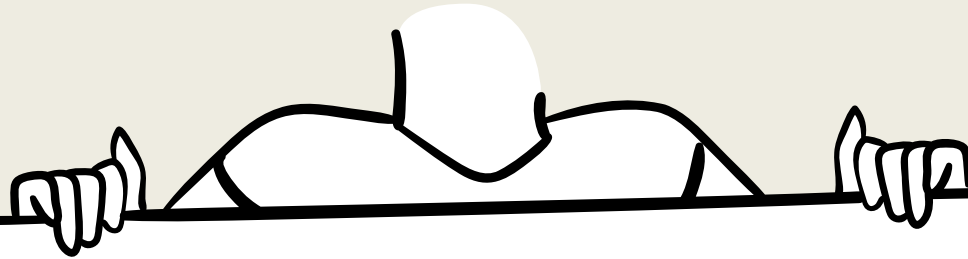
MACHINE LEARNING

Dr. Temitayo Olugbade

Recall from Week 1

1. The most basic element of **machine learning** is a **model** that learns from data.
2. Training a ML model involves optimizing the model parameters based on a **loss function**.

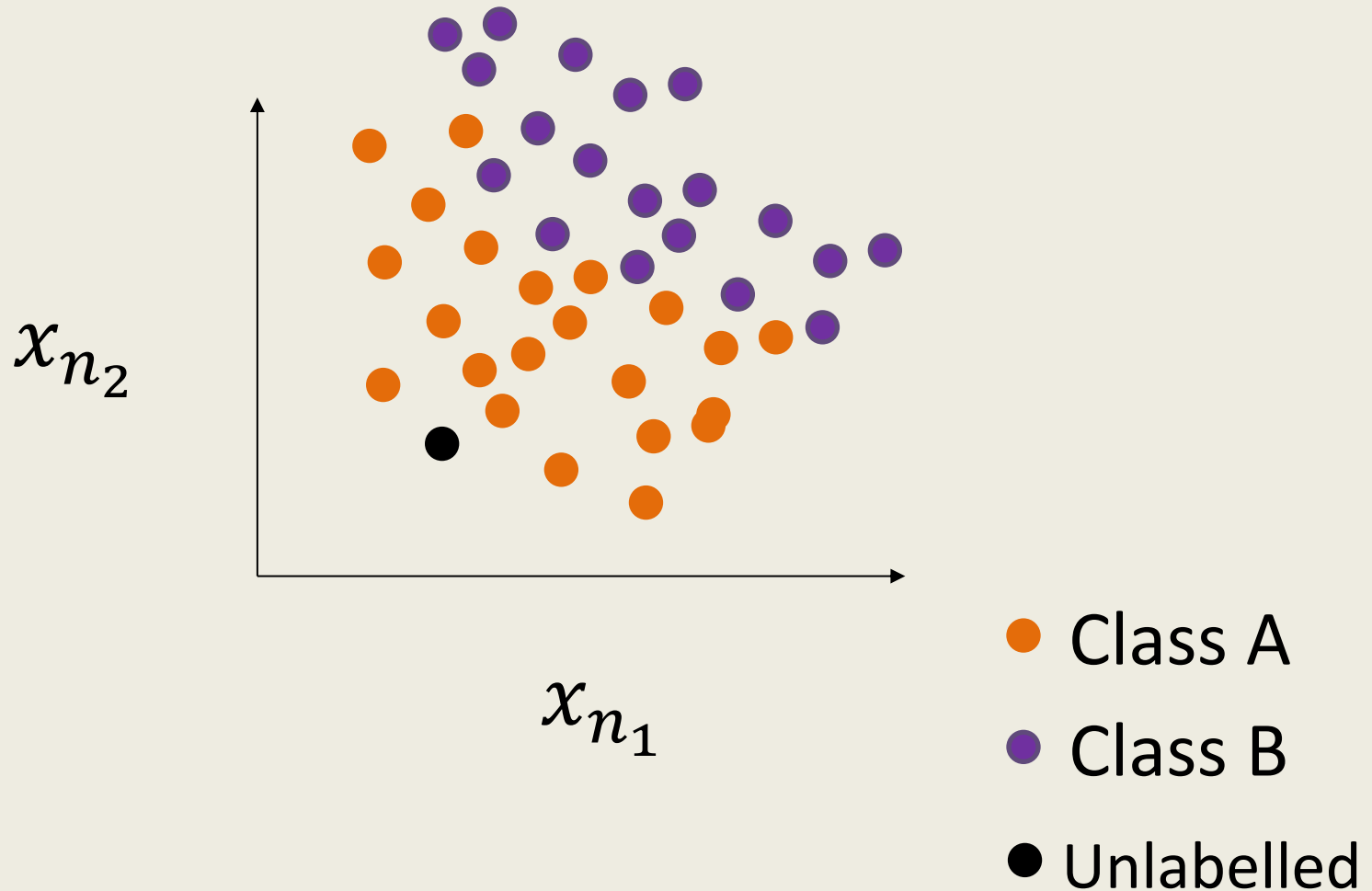
Learning outcome



After working through this mini-video,
you'll see how

- ❑ k-nearest neighbours (kNNs) work.

A toy data for illustration



k-nearest neighbours (kNN)

- Consider that you can and desire to automatically obtain some answer y_l
- You would need to use N training data instances to ~~find a model/function $f(\cdot)$ such that...~~
- With kNN, you use the class(es) of k nearest neighbours of x_i to determine its output \hat{y}_i

No explicit model is defined

How a kNN inference works

- **Step 1**

For each given test instance x_i , find its k nearest neighbours based on a **distance metric** $dist(x_i, x_n) \forall n = 1, 2, \dots, N$

- **Step 2**

Determine the class \hat{y}_i of x_i from the labels y_1, y_2, \dots, y_k of these k nearest neighbours using a **voting strategy**



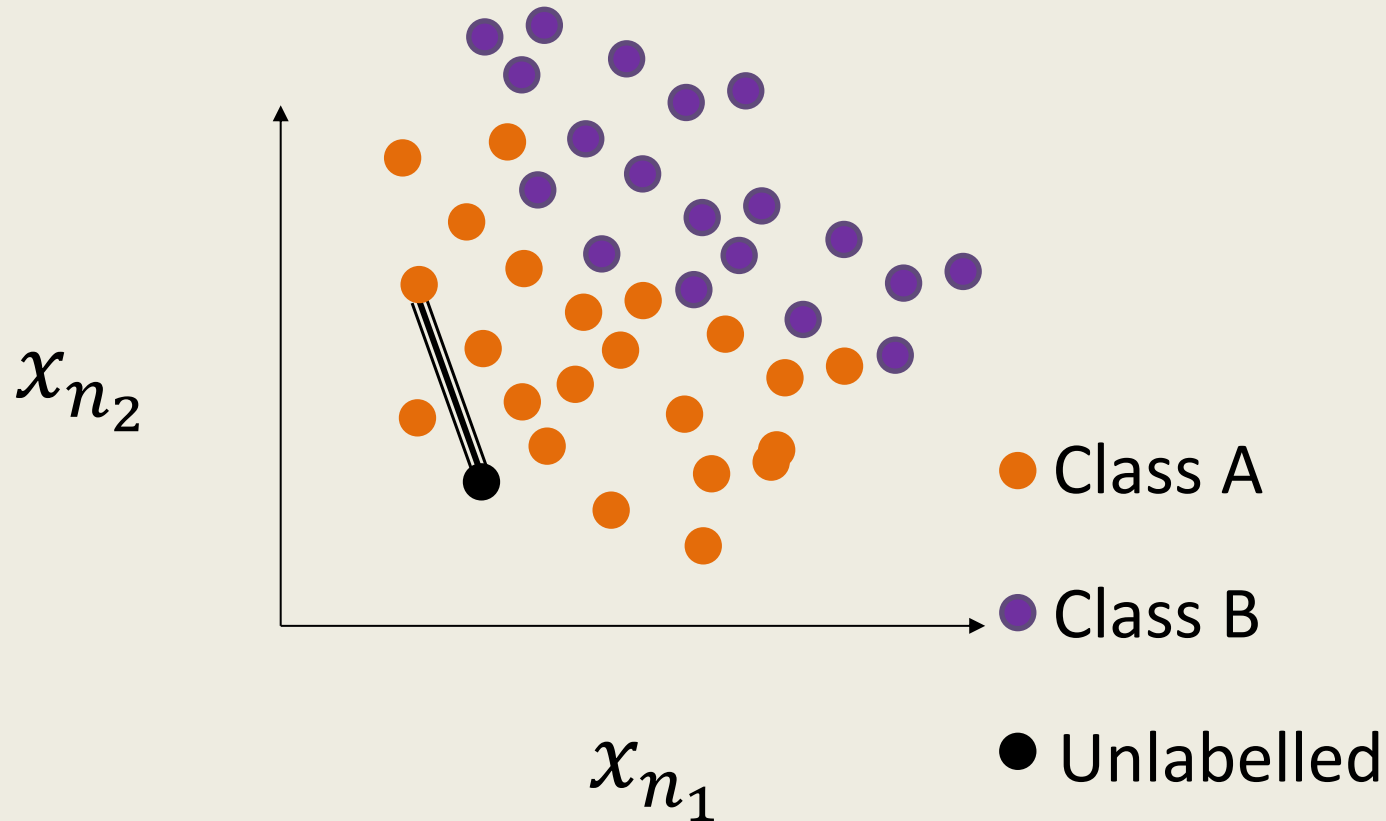
How a kNN inference works

- **Step 1**

For each given test instance x_i , find its k nearest neighbours based on a **distance metric** $dist(x_i, x_n) \forall n = 1, 2, \dots, N$



Euclidean distance



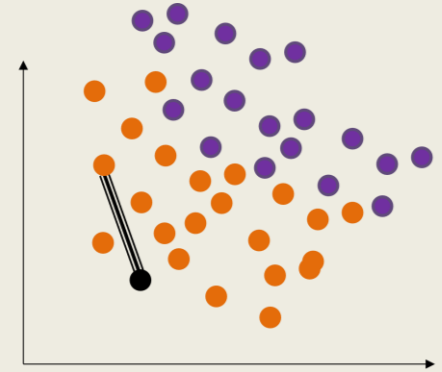
$$\text{dist}(x_i, x_n) = \sqrt{\sum_{d=1}^D (x_{i_d} - x_{n_d})^2} = \|x_i - x_n\|$$

see math proof in last slide pages

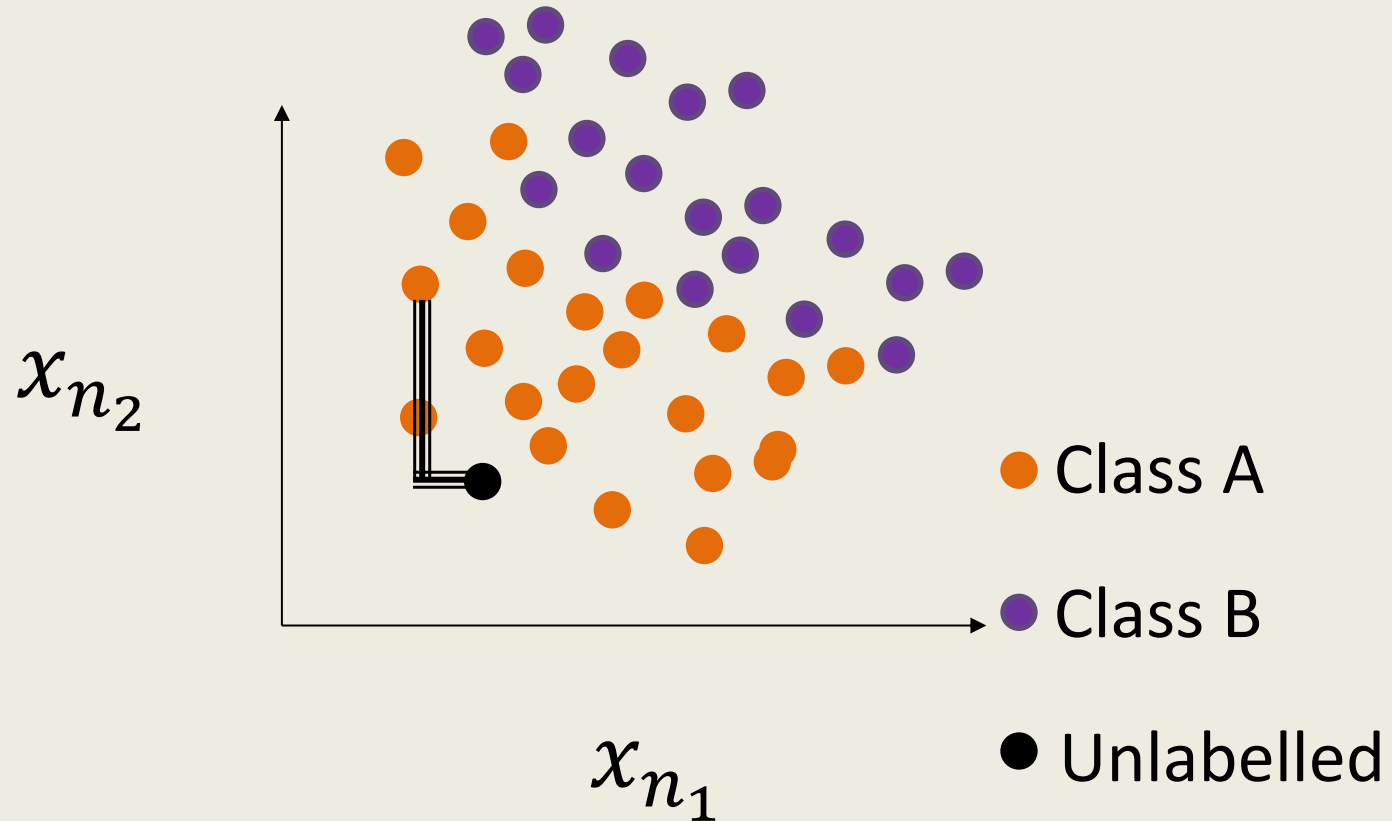
Euclidean distance use cases

Differentiating:

- objects in the sky
- football fans in a pub
- sediment particles vs floating particles in water

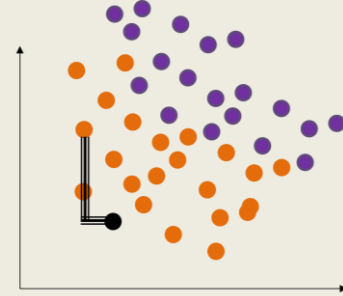


Manhattan (city block) distance



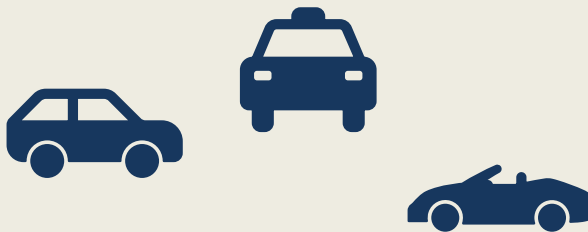
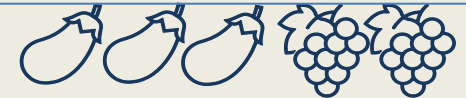
$$\text{dist}(x_i, x_n) = \sum_{d=1}^D |x_{i_d} - x_{n_d}| = |x_i - x_n|$$

Manhattan distance use cases

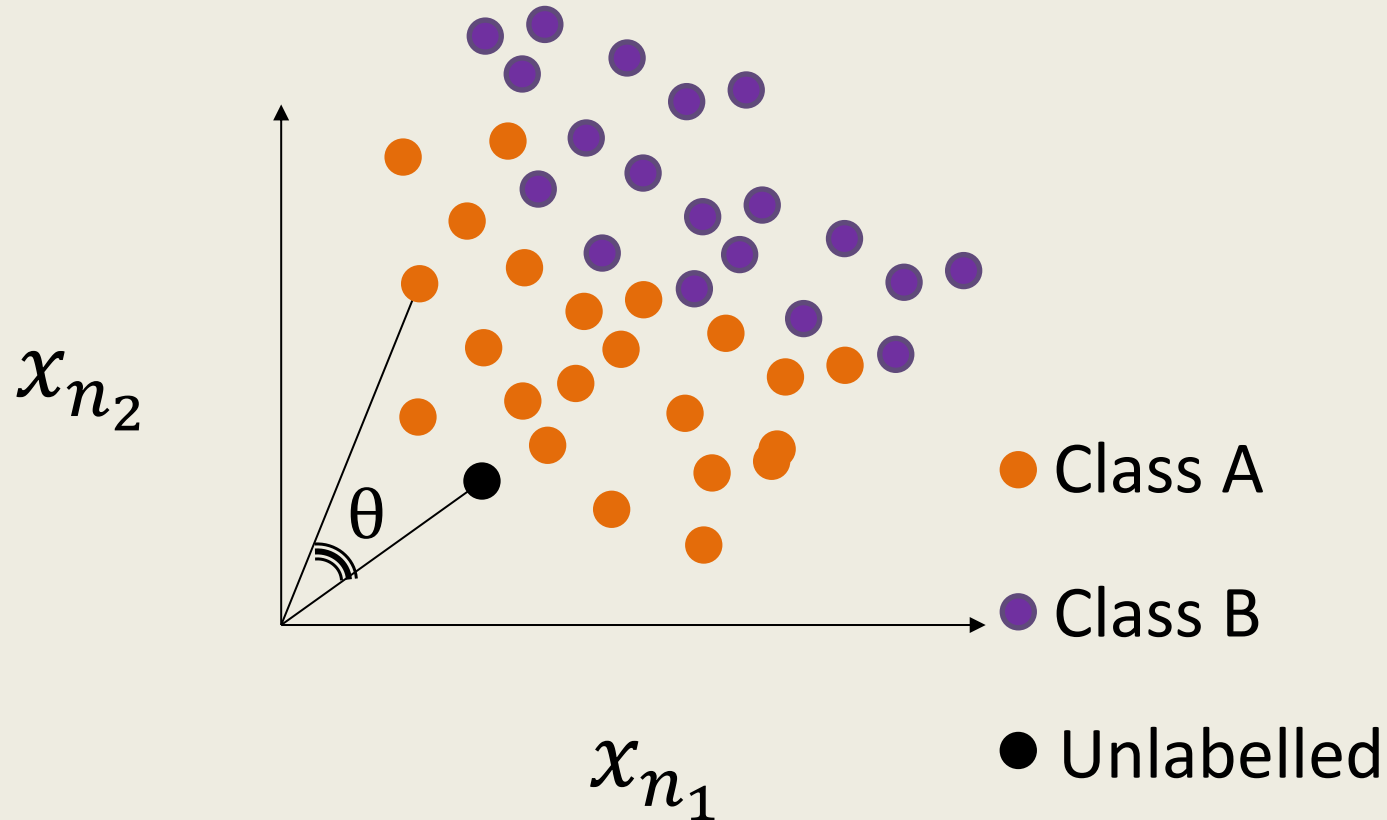


Predicting:

- prices of houses in a region
- categories of items in a supermarket
- trajectories of cars in traffic



Cosine similarity



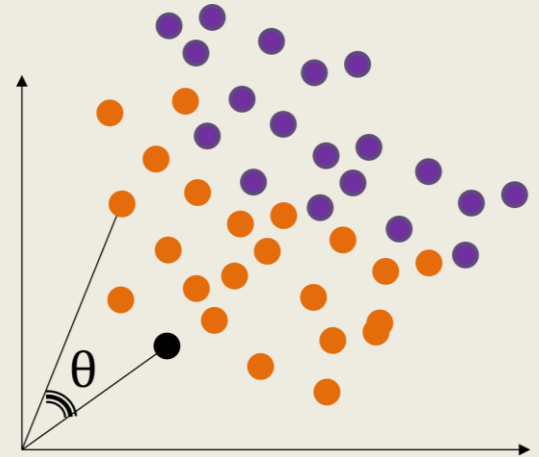
$$\text{dist}(x_i, x_n) = 1 - \frac{\sum_{d=1}^D x_{i_d} \cdot x_{n_d}}{\sqrt{\sum_{d=1}^D x_{i_d}^2} \cdot \sqrt{\sum_{d=1}^D x_{n_d}^2}}$$

see math proof in last slide pages

Cosine similarity use cases

Similarity between

- text documents (represented by the number of occurrences of given words)



Other 'distance' metrics

- Correlation (trend, e.g. Pearson's for normally distributed data, Spearman's otherwise)
e.g. between house prices over two time periods
- Tchebychev (maximum difference)
e.g. between the longest subtask for each of two processes
- Metric learning (optimization, e.g. based on Linear Discriminant Analysis, or Generalized Mahalanobis Distance)

How a kNN inference works

- **Step 1**

For each given test instance x_i , find its k nearest neighbours based on a **distance metric** $dist(x_i, x_n) \forall n = 1, 2, \dots, N$

- **Step 2**

Determine the class \hat{y}_i of x_i from the labels y_1, y_2, \dots, y_k of these k nearest neighbours using a **voting strategy**



Voting strategy

- Majority voting (for classification)

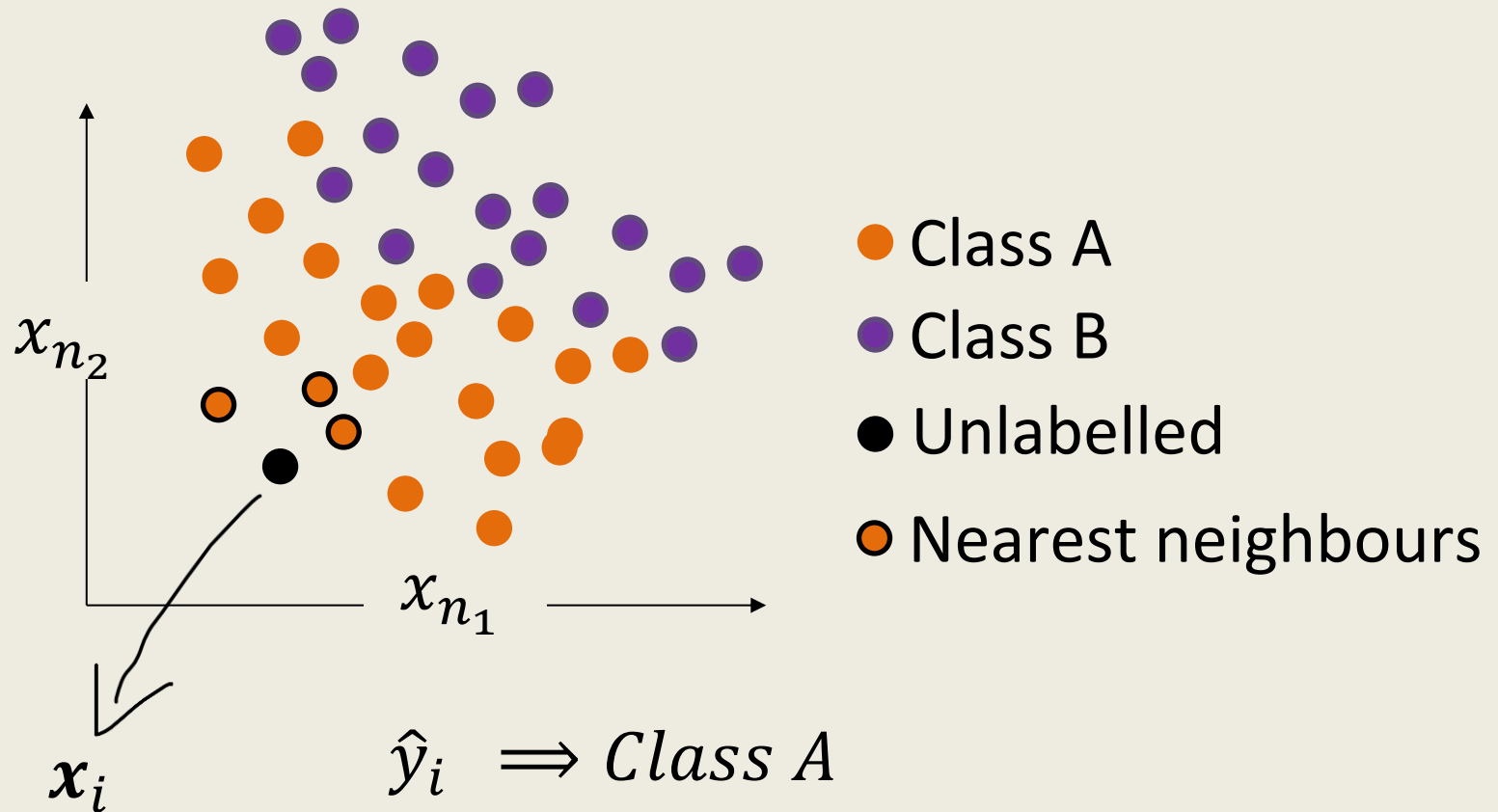
$$\hat{y}_i = \max_c \left\{ \sum_{j=1}^k 1_{y_j=c}, \quad \forall c \in \mathcal{C} \right\}$$

- Averaging (for regression)

$$\hat{y}_i = \frac{1}{k} \sum_{j=1}^k y_j$$

Majority voting example

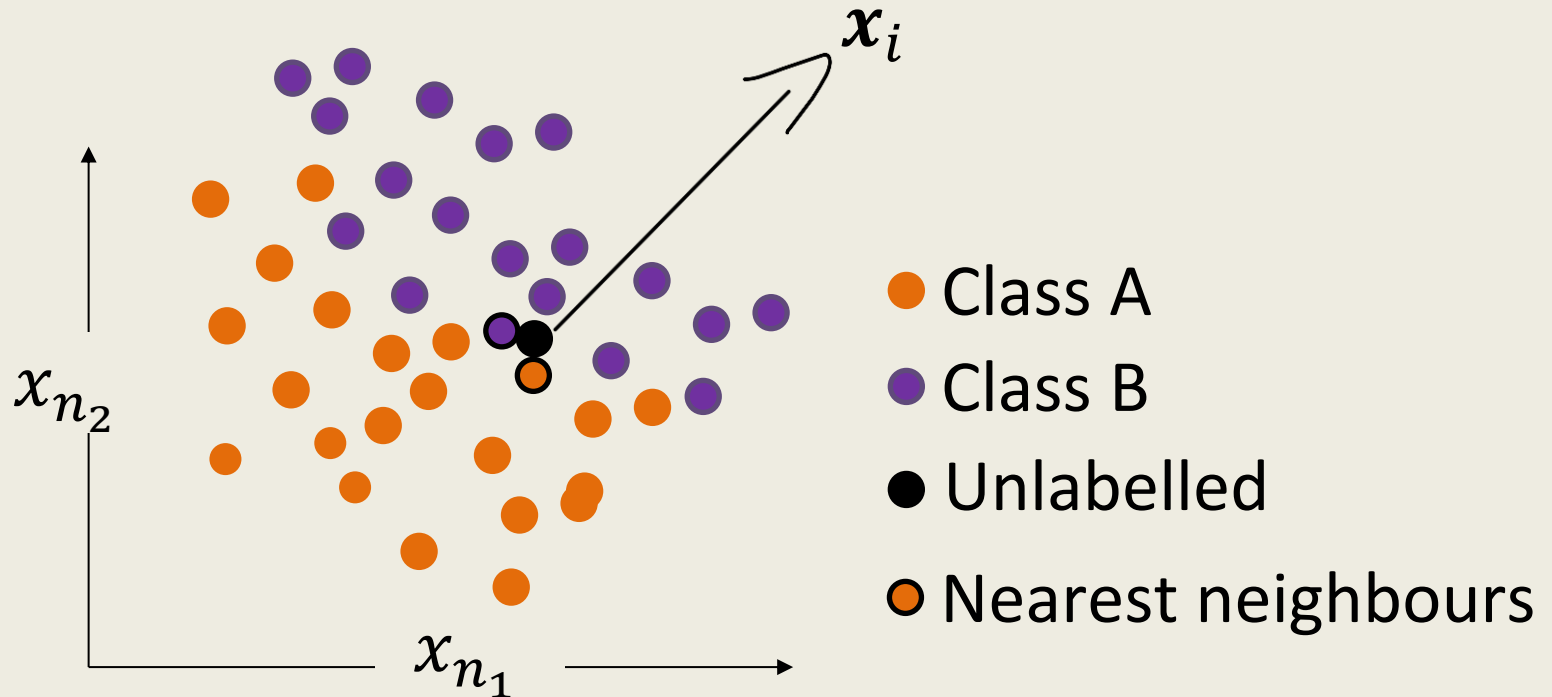
Consider number of neighbours $k=3$,



see math details in last slide pages

Another example

Consider $k=2$



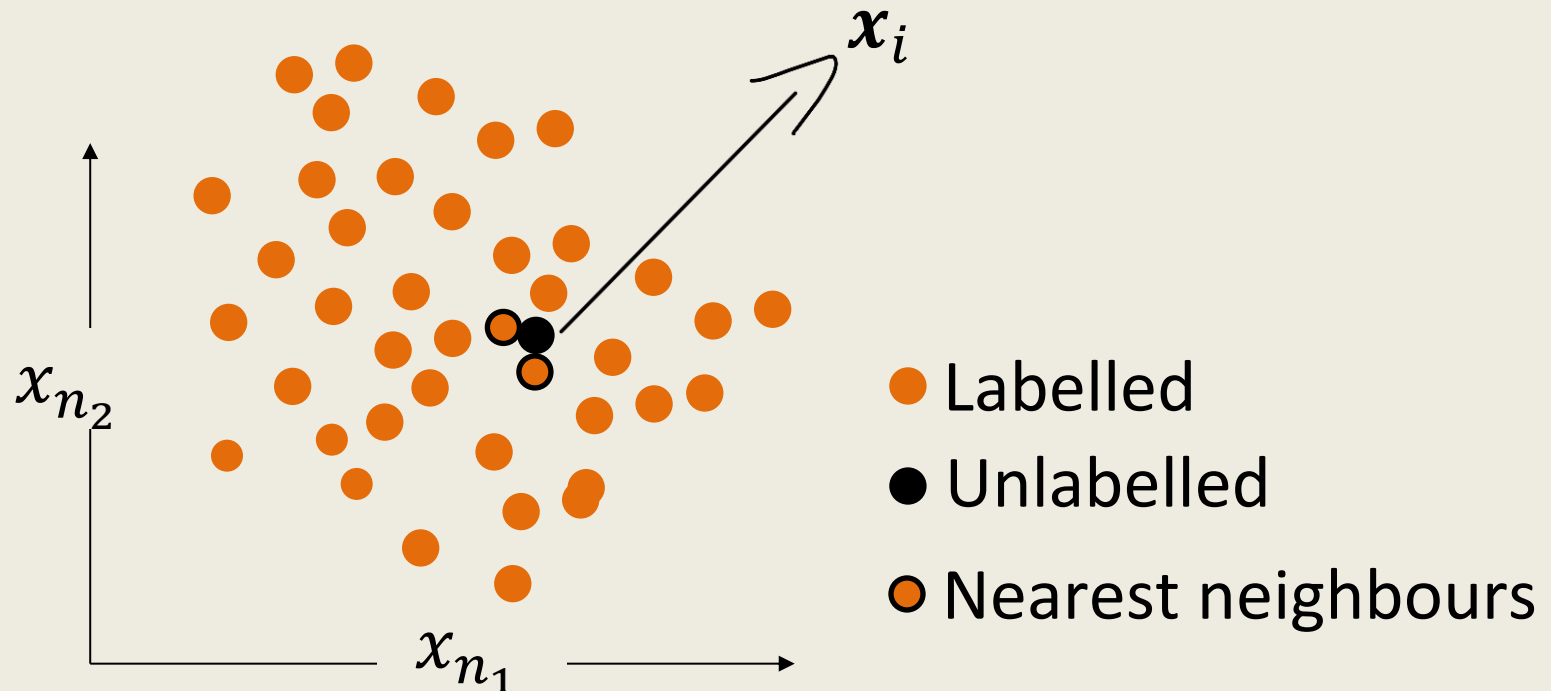
$\hat{y}_i \Rightarrow$ a tie between Class A & Class B

see math details on last slide pages

Averaging example

Consider a regression task with $k=2$

$$\begin{aligned}\Rightarrow \hat{y}_i &= \frac{1}{k} \sum_{j=1}^k y_j \\ &= \frac{1}{2} (y_{n.neighbour1} + y_{n.neighbour2})\end{aligned}$$



Other voting strategies

- Distance-weighted voting (for classification)

$$\hat{y}_i = \max_c \left\{ \sum_{j=1}^k \frac{1}{d(x_i, x_j)^r} \cdot 1_{y_j=c}, \quad \forall c \in \mathcal{C} \right\}, \quad r \geq 1$$

- Distance-weighted averaging (for regression)

$$\hat{y}_i = \frac{1}{k} \sum_{j=1}^k \frac{1}{d(x_i, x_j)^r} \cdot y_j, \quad r \geq 1$$

Inference time complexity for kNN

$$N^2 D$$

N for the total number of data instances

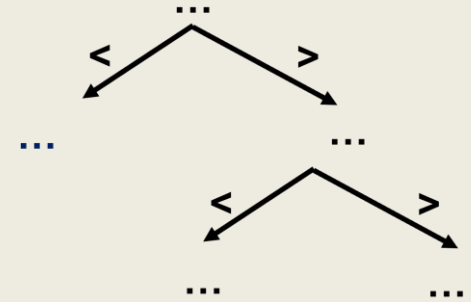
D for the total number of features

see math details in last slide pages



Strategy for minimizing time complexity

- Indexing for faster search
e.g. k-dimension tree (kd tree)



- kd tree – Create a tree-based indexing of the data
 - Repeatedly partition the data region
 - Similar to creating a decision tree
 - The partition median is the preferred split point
 - Leaf nodes are data instances
 - Sorting of values per feature – time complexity $N \log N$
 - Overall time complexity given D features – $ND \log N$
- kNN search – Goes through the tree from root node
 - Values to the left of a node are smaller than the node
 - Aim is to find the leaf node closest to x_i

Other strategies: Dimensionality reduction

- Time complexity of the ordinary kNN (i.e. kNN with brute force search) is

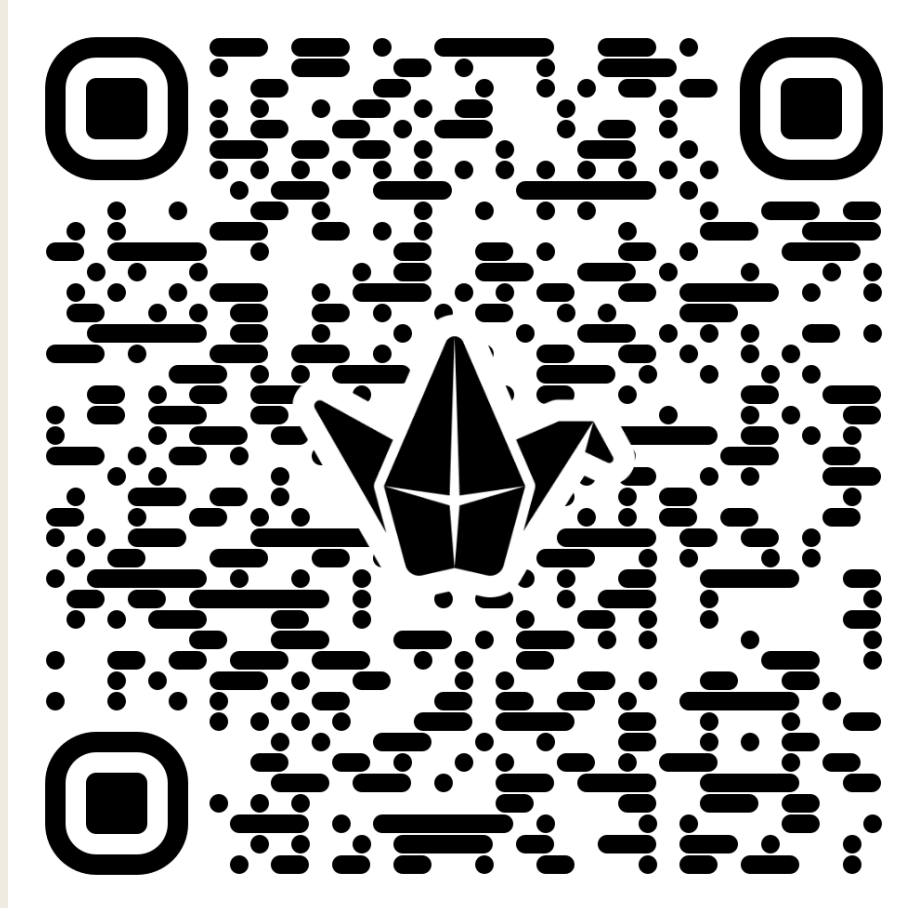
$$N^2 D$$

- **Dimensionality reduction** – Transforming the feature set into a smaller number of features, i.e. from a larger number of dimensions D to a smaller number of dimensions D^*
- Dimensionality reduction methods to be discussed in more detail in Week 6

Summary

1. At inference time, the kNN finds **k-nearest neighbours** for obtaining predictions for a given data instance.
2. It has three main hyperparameters: the number of neighbours **k** , the **distance metric**, the **voting strategy**.
3. Brute force search for nearest neighbours has high time complexity for inference time.

Any questions???



scan the QR code to ask questions

Math details and proofs

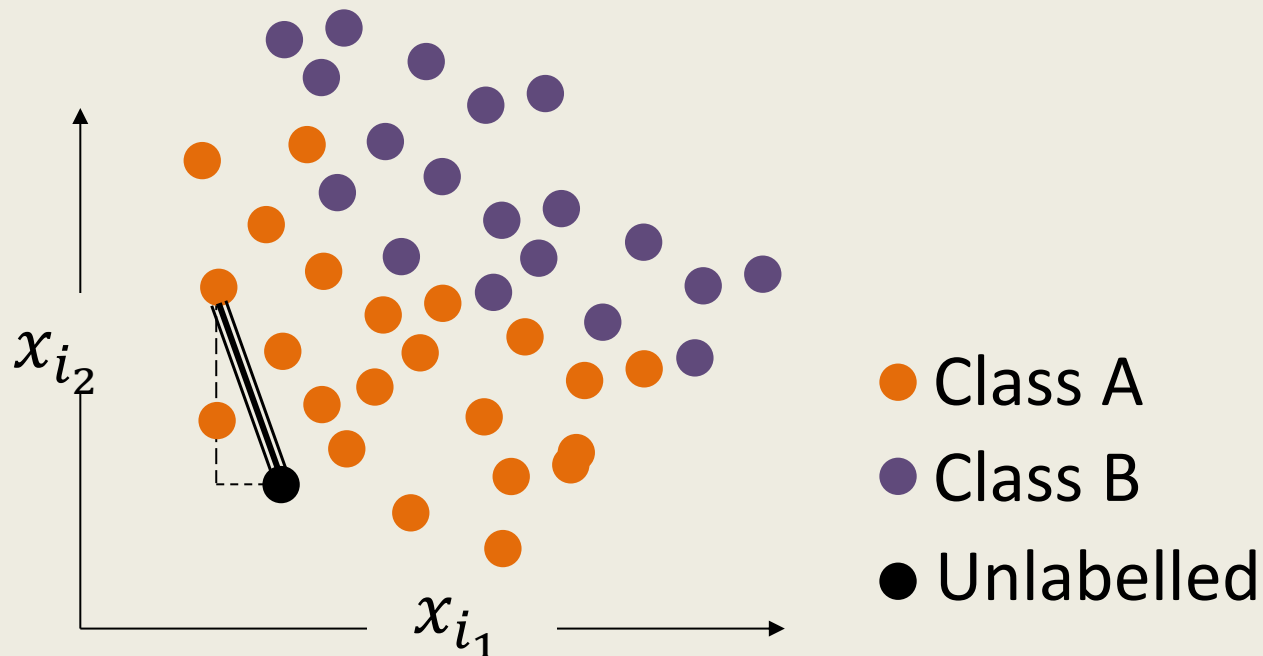
Derivation of Euclidean distance – MATH

From Pythagoras' theorem,

$$(dist(x_i, x_n))^2 = (x_{i_1} - x_{n_1})^2 + (x_{i_2} - x_{n_2})^2$$

taking square root of both sides

$$\Rightarrow dist(x_i, x_n) = \sqrt{(x_{i_1} - x_{n_1})^2 + (x_{i_2} - x_{n_2})^2}$$



Derivation of cosine similarity – MATH

$$\vec{x_i} \cdot \vec{x_j} = \|\vec{x_i}\| \|\vec{x_j}\| \cos \theta$$

making $\cos \theta$ the subject of the formula

$$\cos \theta = \frac{\vec{x_i} \cdot \vec{x_j}}{\|\vec{x_i}\| \|\vec{x_j}\|}$$

expanding each term on the right hand side

$$\cos \theta = \frac{\sum_{d=1}^D x_{i_d} \cdot x_{j_d}}{\sqrt{\sum_{d=1}^D x_{i_d}^2} \cdot \sqrt{\sum_{d=1}^D x_{j_d}^2}}$$

Distance is actually dissimilarity

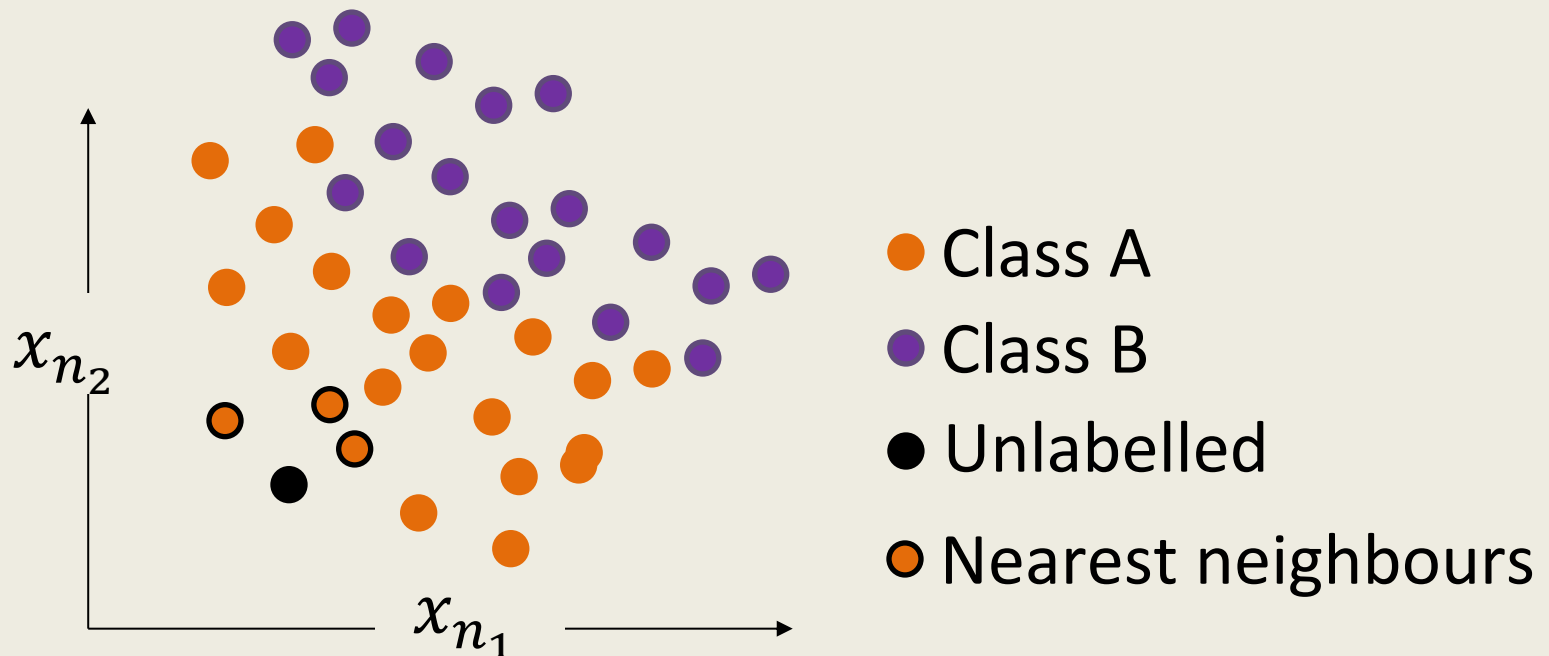
$$\Rightarrow 1 - \cos \theta = 1 - \frac{\sum_{d=1}^D x_{i_d} \cdot x_{j_d}}{\sqrt{\sum_{d=1}^D x_{i_d}^2} \cdot \sqrt{\sum_{d=1}^D x_{j_d}^2}}$$

Majority voting example – MATH DETAILS

Consider number of neighbours $k=3$,

$$\Rightarrow \hat{y}_i = \max_c \left\{ \sum_{j=1}^3 1_{y_j=c}, \quad \forall c \in C \right\}$$

$$\Rightarrow \hat{y}_i = \max_c \{3, 0\} \Rightarrow \text{Class A}$$

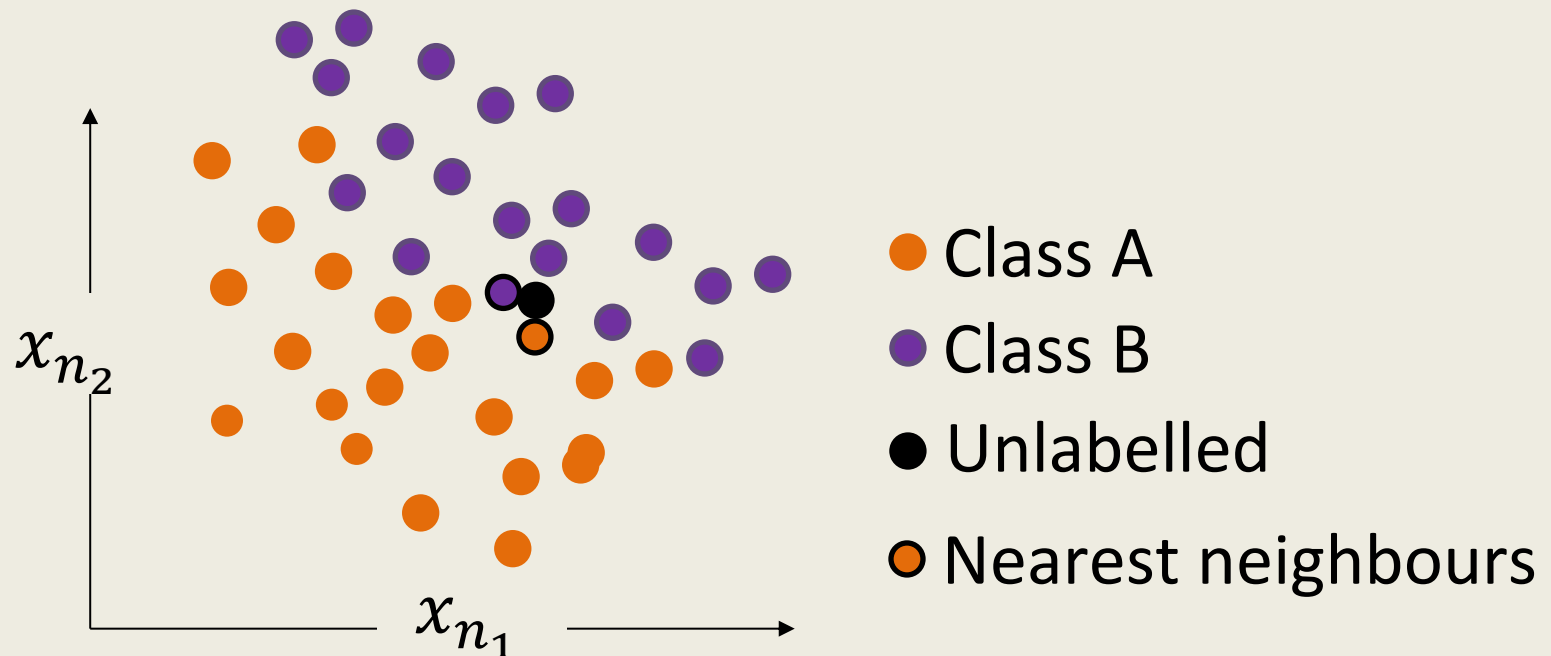


Another example – MATH DETAILS

Consider $k=2$

$$\Rightarrow \hat{y}_i = \max_c \left\{ \sum_{j=1}^2 1_{y_j=c}, \quad \forall c \in \mathcal{C} \right\}$$

$\Rightarrow \hat{y}_i = \max_c \{1, 1\} \Rightarrow$ a tie between Class A & Class B



Inference time complexity for kNN – MATH

- Time complexity

$$O(N^2 D)$$

N for the total number of data instances

D for the total number of features

- Recall that

$$\text{dist}(x_i, x_j) =$$

$$\sqrt{\sum_{d=1}^D (x_{i_d} - x_{j_d})^2}, \forall i, i = 1, 2, \dots, N, j \neq i, \forall j, j = 1, 2, \dots, M$$

i.e. each of M data instances in the test dataset is compared to each of N data instances in the training dataset, and each data instance in both sets has D dimensions that need to be included in the comparison