

# CSC311 Introduction to Machine Learning

## Supervised Learning and Nearest Neighbours

---

Alice Gao

# Learning Outcomes (Supervised Learning)

By the end of this lecture, you will be able to

- Distinguish artificial intelligence, machine learning, and deep learning.
- Distinguish supervised, unsupervised, and reinforcement learning.
- Formulate a problem as a supervised learning problem in formal notation.
- Explain the purposes and uses of the training, validation, and test sets.

# Learning Outcomes (KNN)

By the end of this lecture, you will be able to

- Classify a new data point by applying KNN to a small data set.
- Draw the decision boundaries of KNN for a small data set.
- Analyze the impact of varying  $k$  on KNN's behaviour.
- Explain the common pitfalls of KNN.

# Outline

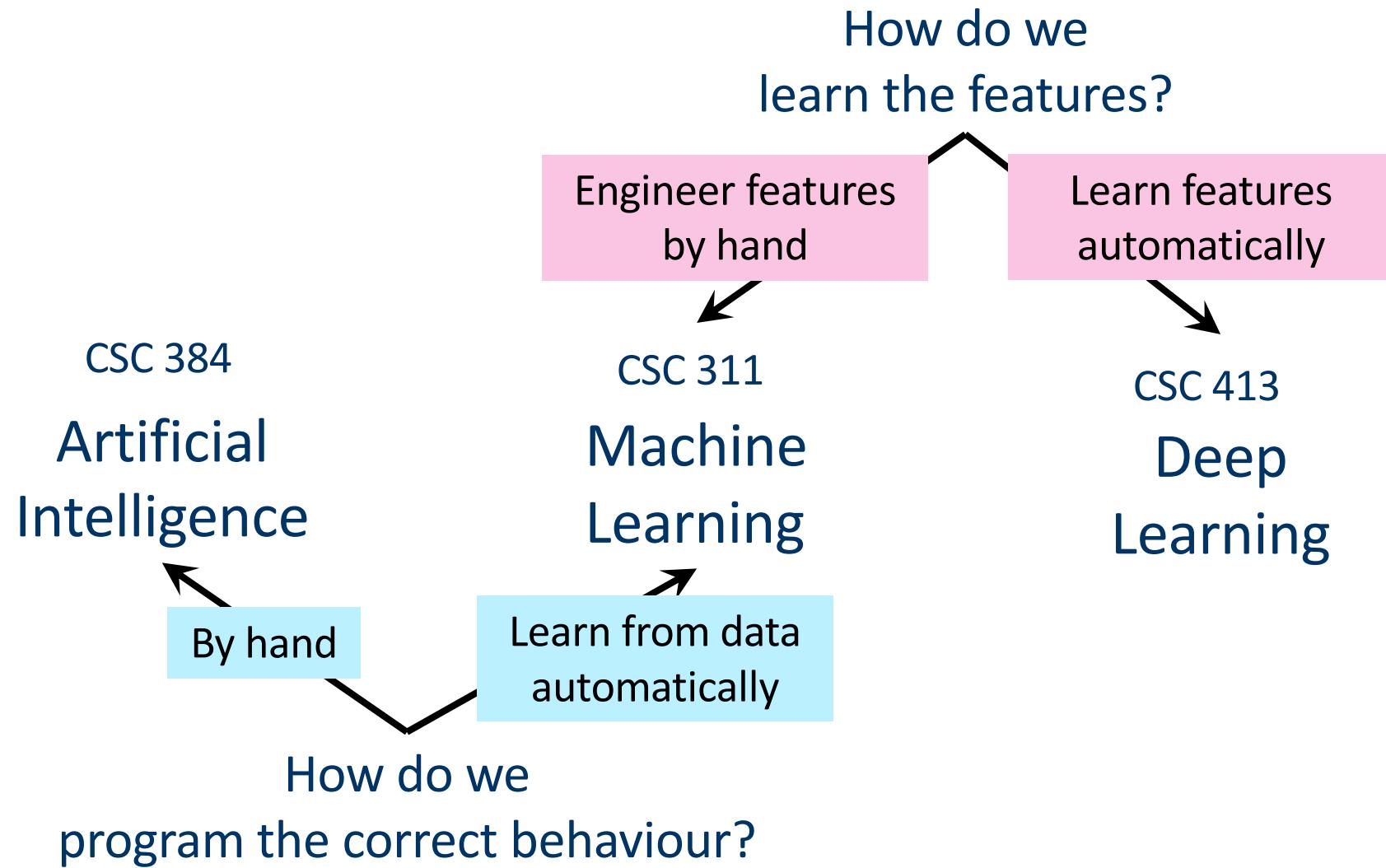
- Introduction to Machine Learning
- Nearest Neighbours
- Decision Boundaries
- K Nearest Neighbours
- Training, Validation and Test Sets
- Limitations of KNN



# Introduction to Machine Learning

---

# Artificial Intelligence vs. Machine Learning vs. Deep Learning



Supervised learning → predict ground truth labels of new data  
Unsupervised learning → trying to find patterns in unlabeled data

## Supervised vs. Reinforcement vs. Unsupervised Learning

### Supervised Learning

- data with ground truth labels
- predict labels given inputs

### Reinforcement Learning

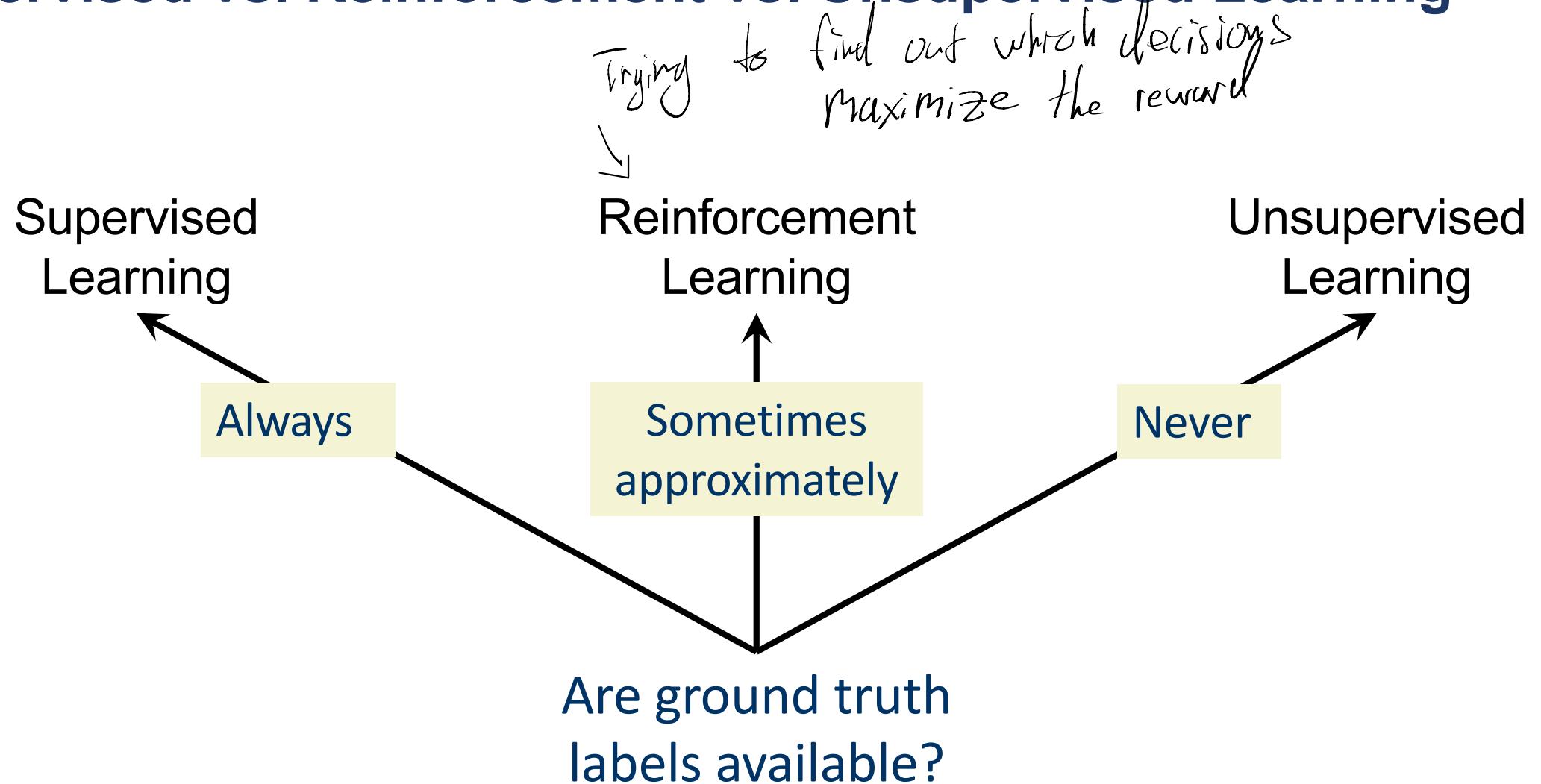
- receives reward signal and
- tries to maximize the reward signal

An agent is trying to maximize a reward signal

### Unsupervised Learning

- data without ground truth labels
- find pattern in data

# Supervised vs. Reinforcement vs. Unsupervised Learning





# Supervised Learning

---

# Supervised Learning

## Objective:

hypothesis / model |  $\rightarrow$   
function mapping input data to  
a label

# Learn a function that maps an input to an output

based on given input-output pairs

training set, test set based on given

Given the **training set** (a set of labeled examples)

Produce a hypothesis or model (a function that maps inputs to outputs)

To perform predictions in the **test set** (unlabeled examples)



# Supervised Learning Setup: Input

Input is a vector  $\mathbf{x} \in \mathbb{R}^D$ :

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = [x_1 \quad x_2 \quad \cdots \quad x_D]^T$$

**Feature**: each scalar  $x_i$  in  $\mathbf{x}$  describes a characteristic of the input.

Many types of data (e.g., text, image, audio) can be represented as vectors.

Vectors are useful because we can **use linear algebra**.

## Supervised Learning Setup: Output

Prediction output  
Target output

Input: a vector  $\mathbf{x} \in \mathbb{R}^D$

Goal of Supervised learning: Create a hypothesis / model

Target output:  $t \in \mathbb{R}$ , the ground-truth value that we want to predict.

Prediction output:  $y \in \mathbb{R}$ , the value generated by our model.

Types of output:

- Classification:  $t$  is discrete (categorical, a value in a finite set  $\{1, \dots, C\}$ ).
- Regression:  $t$  is continuous (a real number).
- Structured prediction:  $t$  is structured (e.g., text, image).

types of labels: classification, regression,  
structured

# Supervised Learning Setup: Training Set

The **training data set** has  $N$  labeled examples:

$$\{(\mathbf{x}^{(1)}, t^{(1)}), (\mathbf{x}^{(2)}, t^{(2)}), \dots, (\mathbf{x}^{(N)}, t^{(N)})\}$$

The  $i^{\text{th}}$  data point has feature vector  $\mathbf{x}^{(i)}$  and target  $t^{(i)}$ .

Notation:

- Superscript = index of the data point in the dataset.
- Subscript = index of the feature in the input vector.

Question: What does  $x_{12}^{(35)}$  represent?

the 12<sup>th</sup> feature of input vector  
at index 35

## Example: Handwritten Digit Recognition

Famous Dataset: MNIST

Problem: Given an image of a handwritten digit, determine which digit the image represents.

- Dataset: MNIST, a large collection of handwritten digits commonly used to benchmark machine learning methods
- Input: a greyscale 28×28 pixel image
- Output: a digit 0-9



# An Image is a Matrix



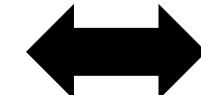
UNIVERSITY OF  
TORONTO

Alice Gao, CSC311 Intro to ML, University of Toronto.

# From a Matrix to a Vector

Image

0	0	0
0	2	53
24	123	252



Vector

0
0
0
0
2
53
24
123
252

# Input Dimension for MNIST Dataset

Handwritten Digit Recognition Problem:

- Input: a greyscale  $28 \times 28$  pixel image
- Output: a digit 0-9

Question: What is the dimension of an input vector  $x$ ?

$$X \in \mathbb{R}^{784}$$

$X \in \mathbb{R}^D$ , vector of Scalar features  
 target:  $t \in \mathbb{R}$  ground truth value we want to predict

Input	$\mathbf{x} \in \mathbb{R}^D$ , a vector of input features
Feature	$x_j \in \mathbb{R}$ , a scalar value describing a characteristic of the input
Target	$t \in \mathbb{R}$ , a scalar ground-truth value that we aim to predict
Prediction	$y \in \mathbb{R}$ , a scalar predicted value generated by a model
Data Point	$(\mathbf{x}, t)$ , a pair of input-output values
Data Set	$(\mathbf{x}^{(1)}, t^{(1)}), \dots, (\mathbf{x}^{(i)}, t^{(i)}), \dots, (\mathbf{x}^{(N)}, t^{(N)})$ , a set of labeled examples

Learning: takes  $\{x^{(i)}, t^{(i)}\}_{i=1}^N$   
→ ground truth labels

## Supervised Learning: Learning versus Inference

A hypothesis or model is a function  $f$  mapping input  $x$  to output prediction  $y$ .

Inference → using a model to predict value  $y$

### Learning

Find a model  $f$  that maps from inputs to outputs given a labeled training dataset.

Training set  $\{(x^{(i)}, t^{(i)})\}_{i=1}^N \rightarrow$  Model  $f$

testing the  
model  
make a  
to prediction

### Inference

Use a model  $f$  to predict an output  $y$  given a new example  $x$

Input  $x$ , Model  $f \rightarrow$  Prediction  $y$

# Supervised Learning: Model Family

A **model** is a function  $f$  mapping input  $x$  to output prediction  $y$ .

*Model family → Set of models that share same structure*

A **model family** is a set of models that share the same structure

(e.g., linear models, decision trees, neural networks)

but differ in their **parameters**.

*Parameters, features of the data themselves*

## Supervised Learning: Evaluation

$$\text{Accuracy} = \frac{\sum_{i=1}^N \mathbb{I}[t^{(i)} = y^{(i)}]}{N}$$

- Evaluation is assessing how well a model performs.
- Evaluation is generally independent of the model family.

Accuracy is measured as fraction of correct predictions in a data set

For classification, we can measure the model accuracy (or error) as the fraction of correct (or incorrect) predictions in a data set.

$$\text{Accuracy} = \frac{\sum_{i=1}^N \mathbb{I}[t^{(i)} = y^{(i)}]}{N}, \quad \text{Error} = 1 - \text{Accuracy}$$

where  $\mathbb{I}[e] = 1$  if condition  $e$  is true, and 0 otherwise.



# Nearest Neighbours

---

# Training Set and New Point

Training Data

$x^{(i)}$	$t^{(i)}$
8	8
5	5
2	2
0	0
7	7
1	1

New Data

$x$	$y$	$+$
1	?	1

# Distance Between First Training Point and New Point

Training Data

$\mathbf{x}^{(i)}$	$t^{(i)}$
8	8
5	5
2	2
0	0
7	7
1	1

distance( $\mathbf{x}^{(1)}, \mathbf{x}$ )

New Data

$\mathbf{x}$	$y$
1	?

# Compute Distances to All Training Points

Training Data

$\mathbf{x}^{(i)}$	$t^{(i)}$
8	8
5	5
2	2
0	0
7	7
1	1

$\text{distance}(\mathbf{x}^{(1)}, \mathbf{x})$

New Data

$\mathbf{x}$	$y$
1	?

# Identify Closest Neighbour

Training Data

$\mathbf{x}^{(i)}$	$t^{(i)}$
8	8
5	5
2	2
0	0
7	7
1	1

New Data

$\mathbf{x}$	$y$
1	?

distance( $\mathbf{x}^{(1)}, \mathbf{x}$ )

distance( $\mathbf{x}^{(2)}, \mathbf{x}$ )

distance( $\mathbf{x}^{(3)}, \mathbf{x}$ )

distance( $\mathbf{x}^{(4)}, \mathbf{x}$ )

distance( $\mathbf{x}^{(5)}, \mathbf{x}$ )

distance( $\mathbf{x}^{(6)}, \mathbf{x}$ )

# Predict by Copying Closest Neighbour's Label

Training Data

$x^{(i)}$	$t^{(i)}$
8	8
5	5
2	2
0	0
7	7
1	1

New Data

$x$	$y$
1	7

Copied label

Closest training  
data point

Copy its label

Pick the closest training  
data point, and  
Copy its label (target)

## Nearest Neighbours (Formal Description) 1/2

1. Calculate the distance between every training data point  $\mathbf{x}^{(i)}$  and the new data point  $\mathbf{x}$ .

$$\text{distance}(\mathbf{x}^{(i)}, \mathbf{x})$$

2. Find training data point  $(\mathbf{x}^{(c)}, t^{(c)})$  that is closest to the new data point  $\mathbf{x}$ .

$$\text{distance}(\mathbf{x}^{(c)}, \mathbf{x}) = \min_{i \in 1, \dots, N} \text{distance}(\mathbf{x}^{(i)}, \mathbf{x})$$

3. Predict  $y = t^{(c)}$  as the label for the new data point  $\mathbf{x}$ .

distance ( $\mathbf{x}^{(i)}, \mathbf{x}$ )

## Nearest Neighbours (Formal Description) 2/2

1. Calculate the distance between every training data point  $\mathbf{x}^{(i)}$  and the new data point  $\mathbf{x}$ .

$$\text{distance}(\mathbf{x}^{(i)}, \mathbf{x})$$

2. Find training data point  $(\mathbf{x}^{(c)}, t^{(c)})$  that is closest to the new data point  $\mathbf{x}$ .

$$\mathbf{x}^{(c)} = \underset{\mathbf{x}^{(i)} \in \text{-training set}}{\arg \min} \text{distance}(\mathbf{x}^{(i)}, \mathbf{x})$$

(compare with previous notation:  $\text{distance}(\mathbf{x}^{(c)}, \mathbf{x}) = \min_{i \in 1, \dots, N} \text{distance}(\mathbf{x}^{(i)}, \mathbf{x})$ )

3. Predict  $y = t^{(c)}$  as the label for the new data point  $\mathbf{x}$ .

## How Do We Measure Distance?

Must choose a distance measure, depending on the problem we are solving.

### Euclidean distance ( $L^2$ distance)

a reasonable default choice.

other  $L^p$  distances also reasonable.

$$x^{(b)} = \arg \min_{x^{(b)}} \text{distance}(x^{(b)}, \{x^{(i)}\}_{\text{training set}})$$

$$\|x^{(a)} - x^{(b)}\|_2 = \sqrt{\sum_{j=1}^D (x_j^{(a)} - x_j^{(b)})^2}$$

Manhattan distance  
 $\|x^{(a)} - x^{(b)}\|_1$

### Cosine similarity

measures angle between two vectors.

$$\cosine(x^{(a)} - x^{(b)}) = \frac{x^{(a)} \cdot x^{(b)}}{\|x^{(a)}\| \|x^{(b)}\|}$$

$$\cosine(x^{(a)} - x^{(b)}) = \frac{\mathbf{x}^{(a)} \cdot \mathbf{x}^{(b)}}{\|\mathbf{x}^{(a)}\|_2 \|\mathbf{x}^{(b)}\|_2}$$

*data point*  $X \in \mathbb{R}^2$

## Example: Citrus Fruit Prediction

Classify whether a fruit is  
an **orange** or a **lemon**  
based on its **weight** and **height**  
using the training data.

Weight	Height	Label
6.8	7.4	Orange
7.1	7.5	Orange
7.6	8.2	Orange
7.2	7.2	Orange
7.2	9.3	Lemon
7.3	9.5	Lemon
7.2	9.0	Lemon
6.9	7.0	Lemon

# Classifying a New Fruit

Weight	Height	Label
6.8	7.4	Orange
7.1	7.5	Orange
7.6	8.2	Orange
7.2	7.2	Orange
7.2	9.3	Lemon
7.3	9.5	Lemon
7.2	9.0	Lemon
6.9	7.0	Lemon

A new fruit has weight 7.4 and height 8.5.  
How should we classify it using Nearest Neighbours?

The new data point is  $\mathbf{x} = \begin{bmatrix} 7.4 \\ 8.5 \end{bmatrix}$

Calculate distance to each training point:

$$\text{distance}(\mathbf{x}^{(1)}, \mathbf{x}) = \sqrt{(7.4 - 6.8)^2 + (7.4 - 8.5)^2}$$

Similarly, calculate  $(\mathbf{x}^{(i)}, \mathbf{x})$ , for  $i = 2, \dots, 8$ .

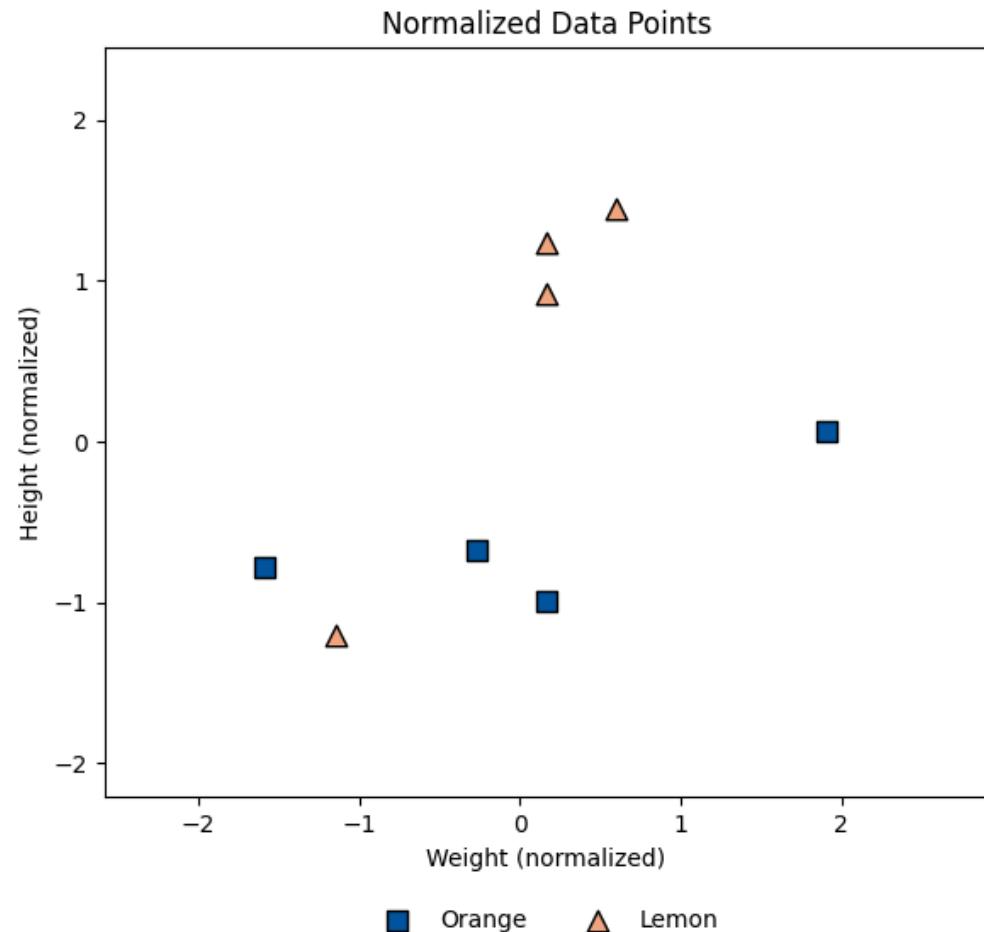
Copy the label of the closest training example.



# Nearest Neighbours Decision Boundaries

---

# Visualize Decision Boundaries



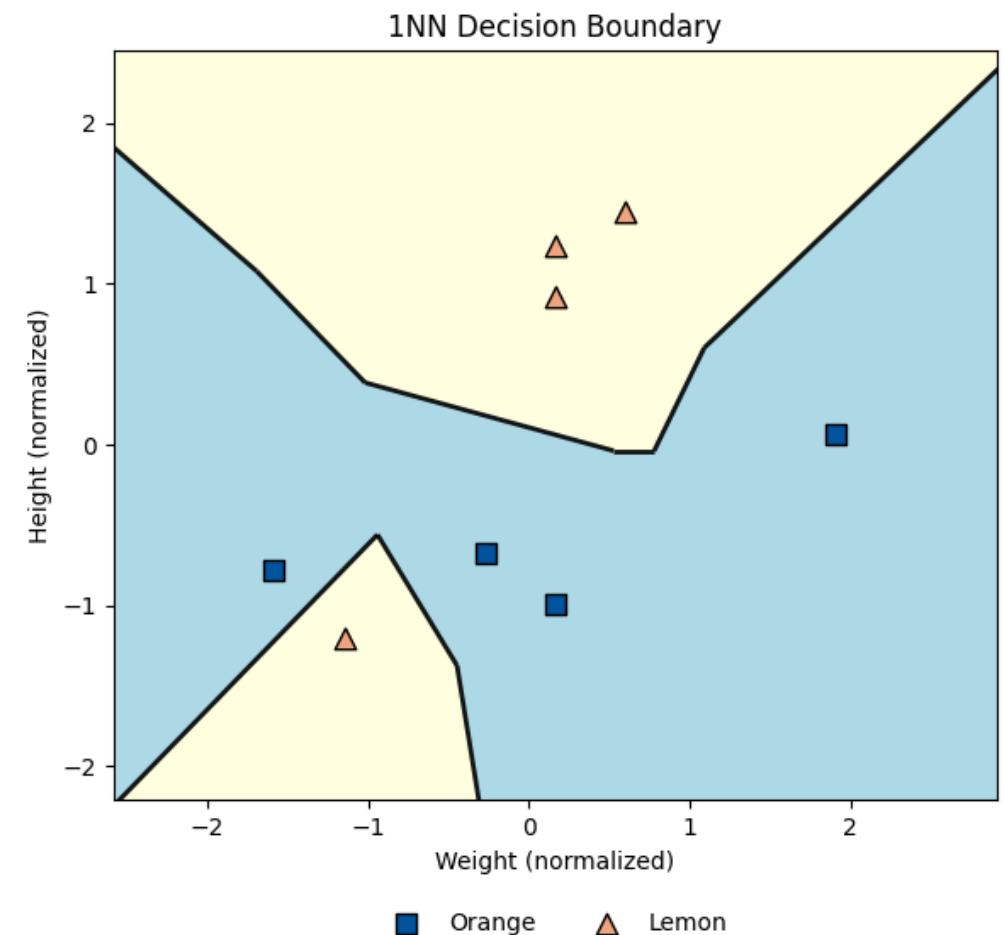
height  
is typically  
lemon

Weight	Height	Fruit
6.8	7.4	Orange
7.1	7.5	Orange
7.6	8.2	Orange
7.2	7.2	Orange
7.2	9.3	Lemon
7.3	9.5	Lemon
7.2	9.0	Lemon
6.9	7.0	Lemon

# Nearest Neighbours Decision Boundary

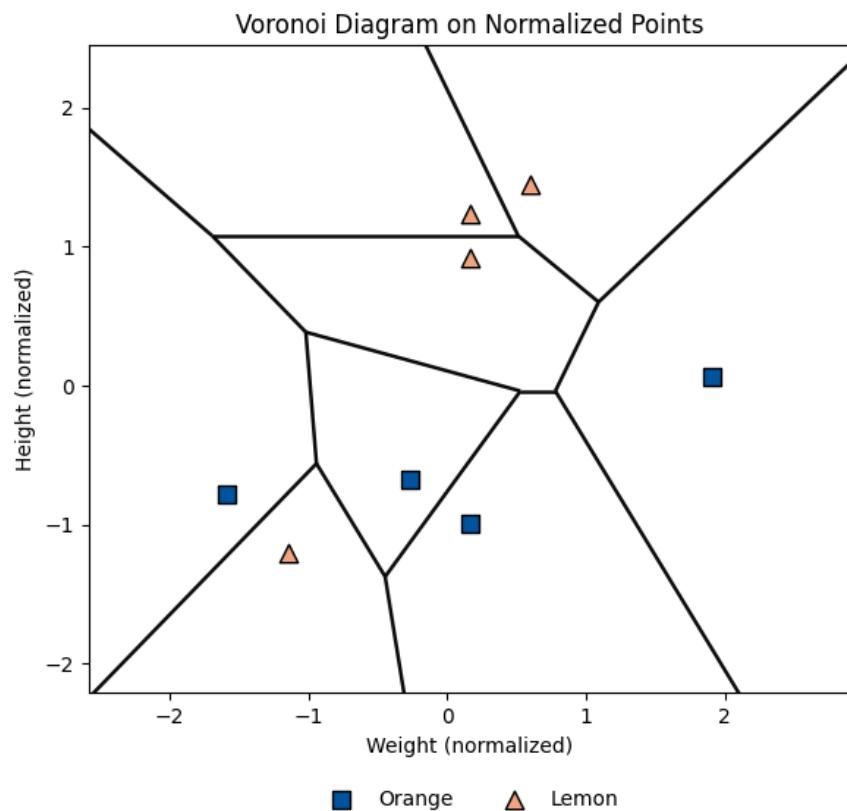
A decision boundary is

- a region where the classifier output is **ambiguous**
- a dividing line/surface where the classifier **switches** from predicting one class to another

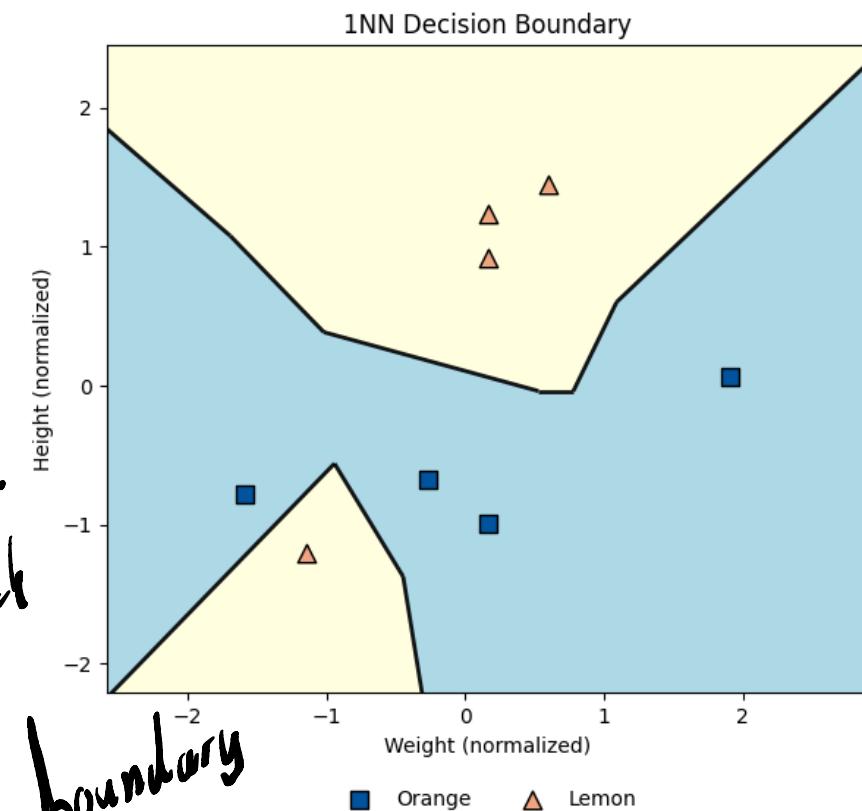


# 1NN Decision Boundary via Voronoi Diagram

Voronoi Diagram  
*Perpendicular bisectors*

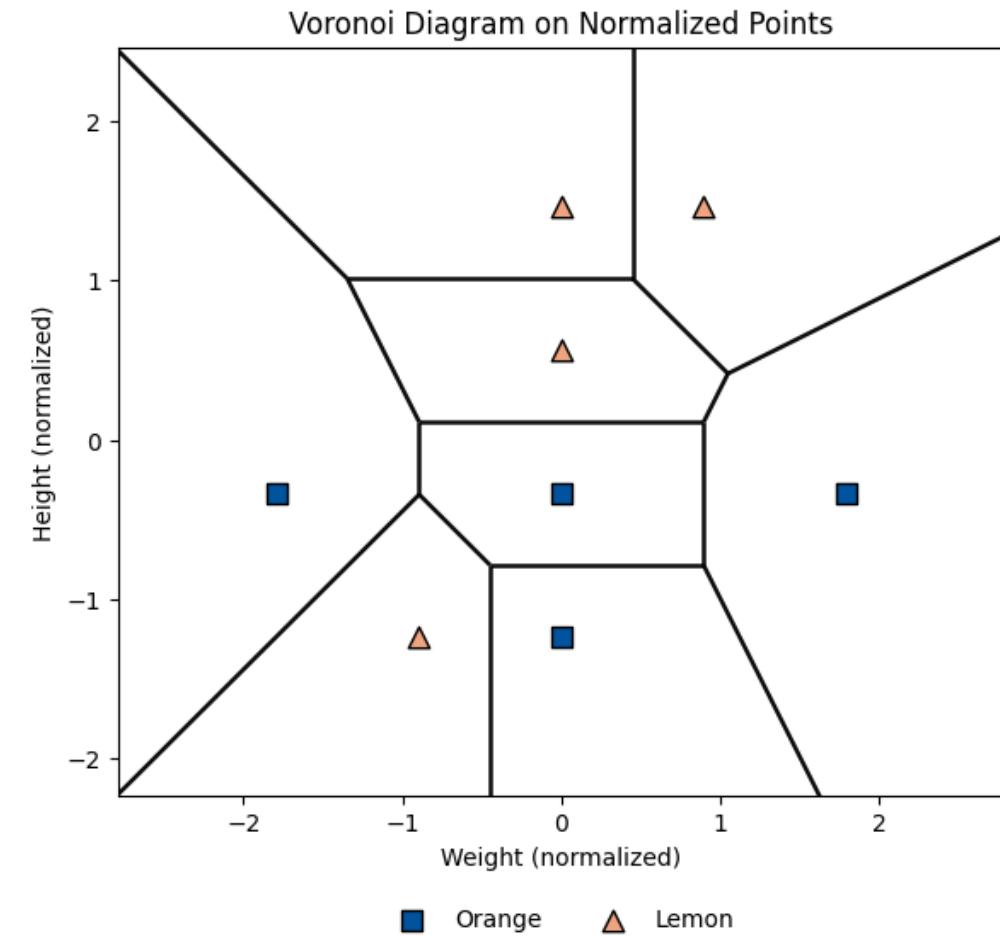
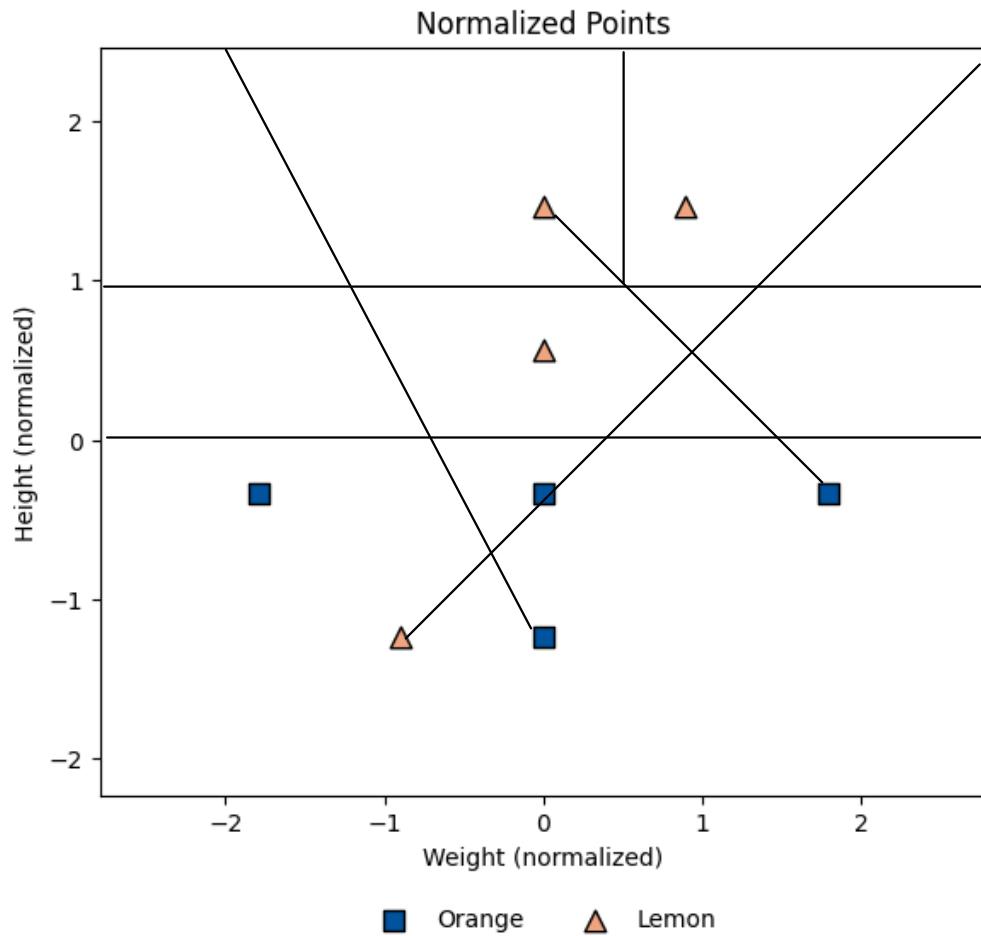


## 1NN Decision Boundary



Difference between Voronoi and 1NN:  
Voronoi creates regions for each point, 2 NN creates decision boundary where classifiers change

# Sketch Voronoi Diagram

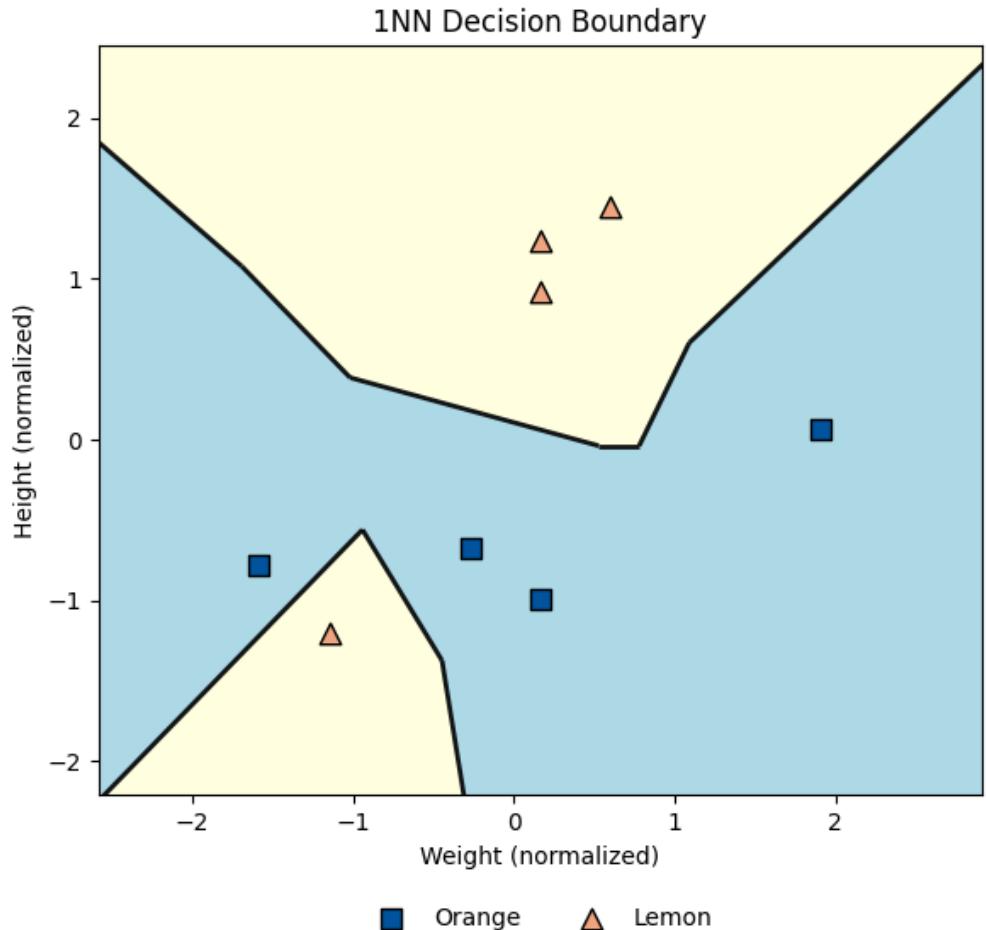




# *K* Nearest Neighbours

---

# Problem with Nearest Neighbours



The NN model is **sensitive** to noise or mis-labeled data.

Which area may have unreasonable predictions?

$K$  nearest neighbours

## $K$ Nearest Neighbours (KNN)

Instead of using the **closest** point, use  $K$  **closest** points for prediction.

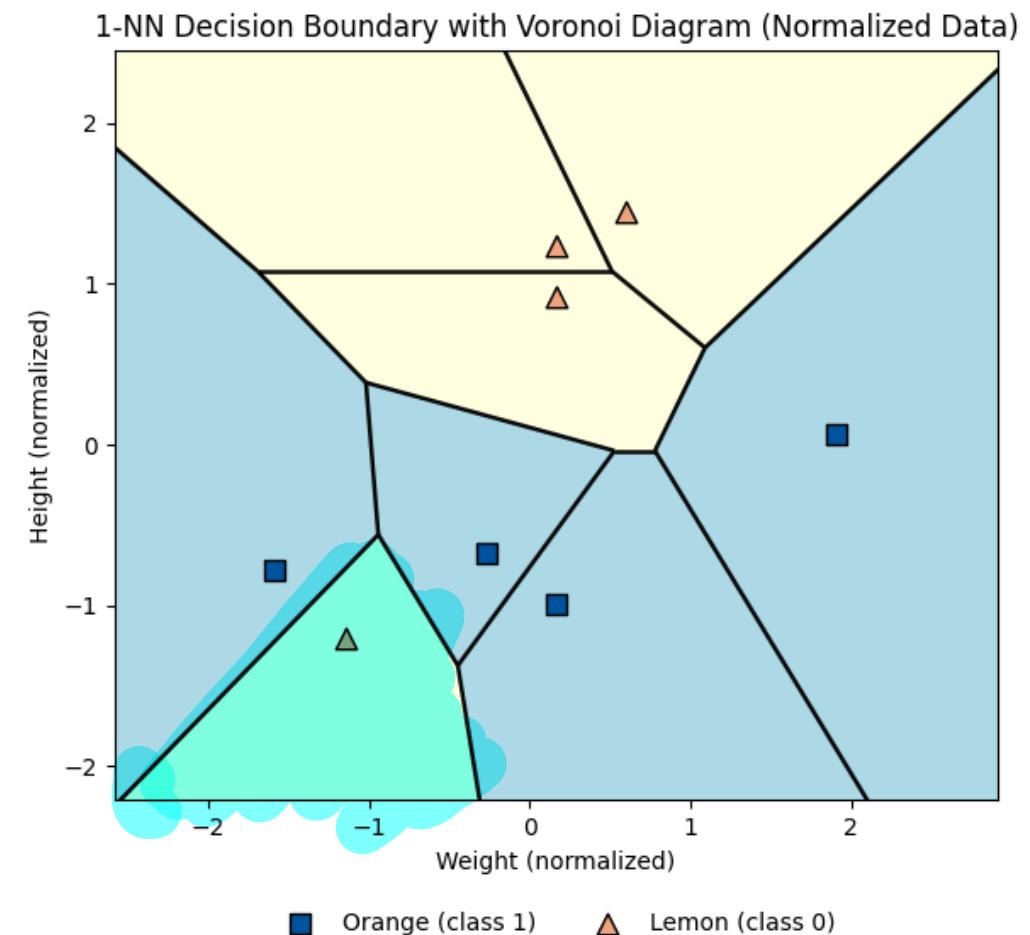
$K$  Nearest Neighbours:

- Find the  $K$  closest training points to new data point  $x$ .
- Predict the label by **majority vote** among  $K$  nearest neighbours.

# 3NN Decision Boundary (Optional Exercise)

Find a region in the data space where **1NN** and **3NN** models predict **different labels**.

A.



## **K is a Hyperparameter of KNN**

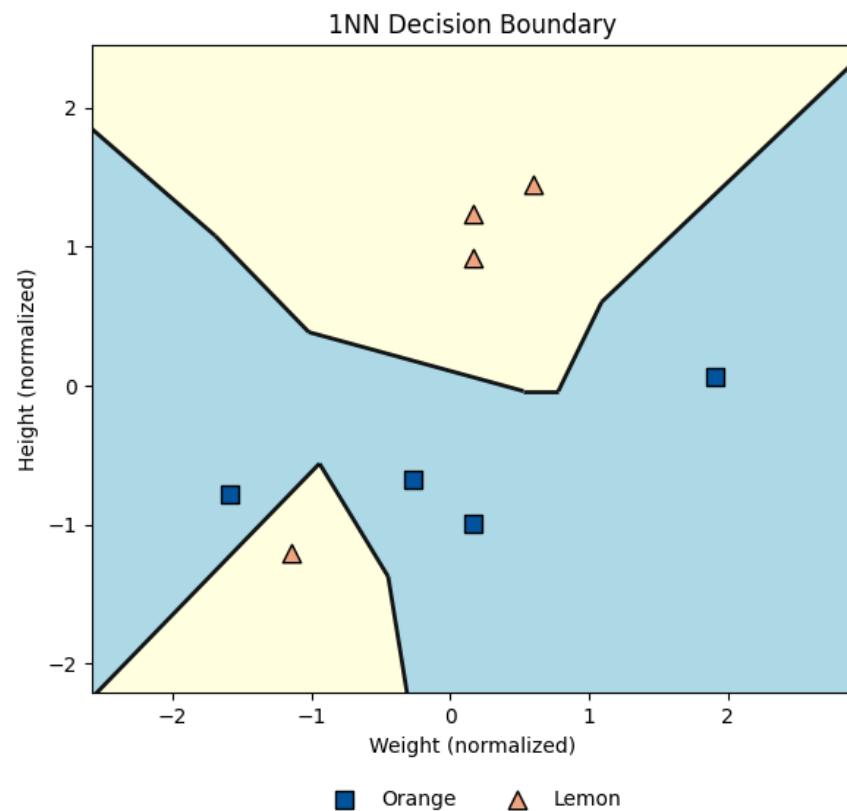
A hyperparameter is

- A setting or configuration **external** to the model
- Chosen **before training**
- Governs how the model is learned
- Strongly influences model **behaviour** and **performance**

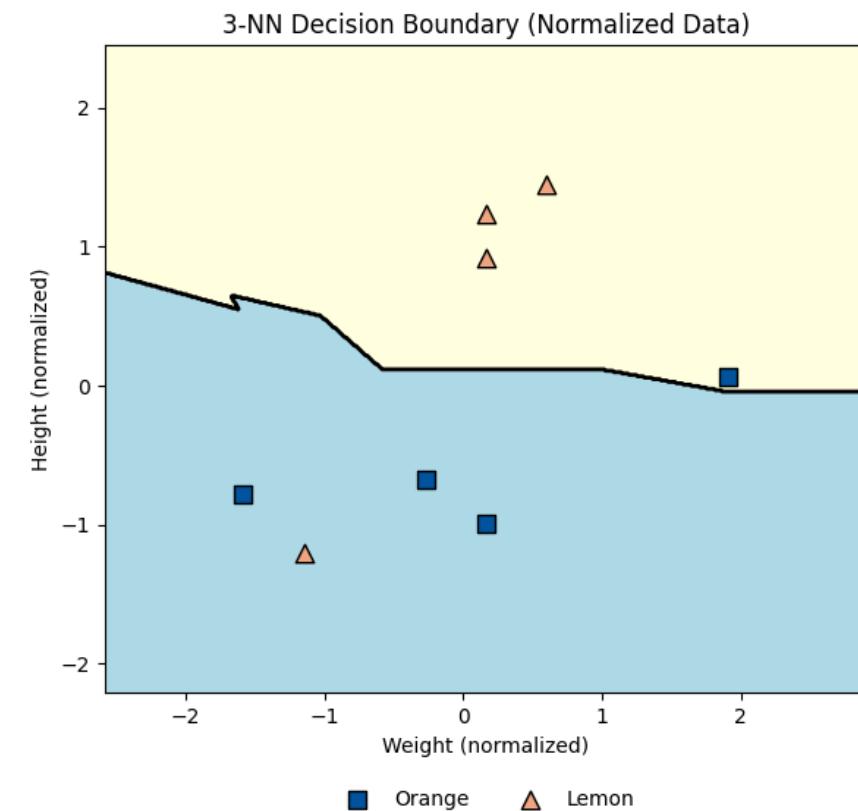
model example: for KNN, choosing K determines how the function works

# Effects of Varying $K$

$K = 1$

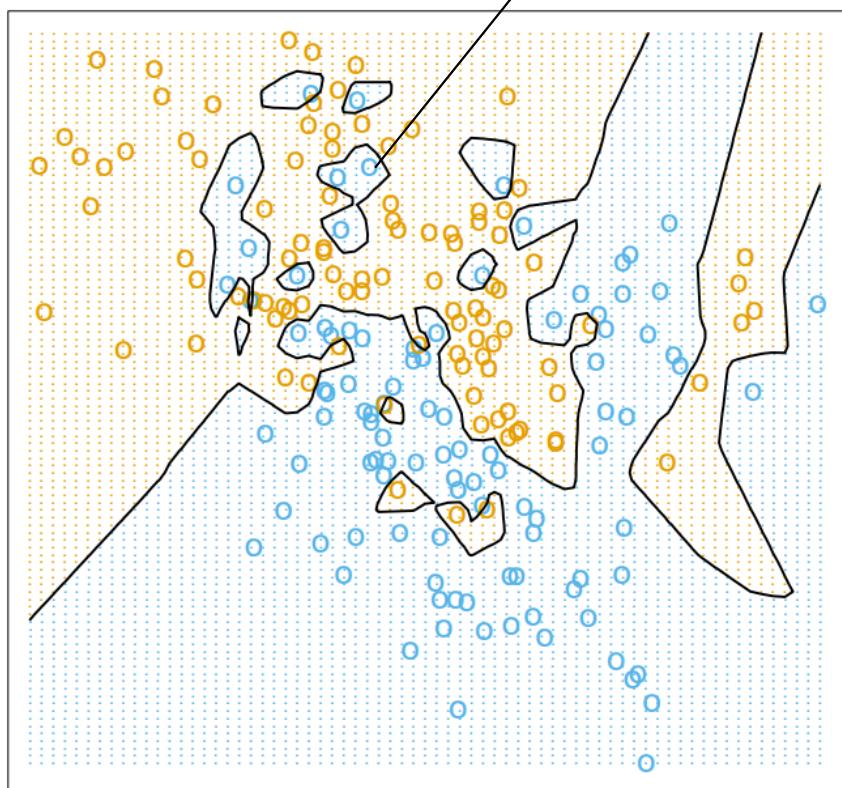


$K = 3$

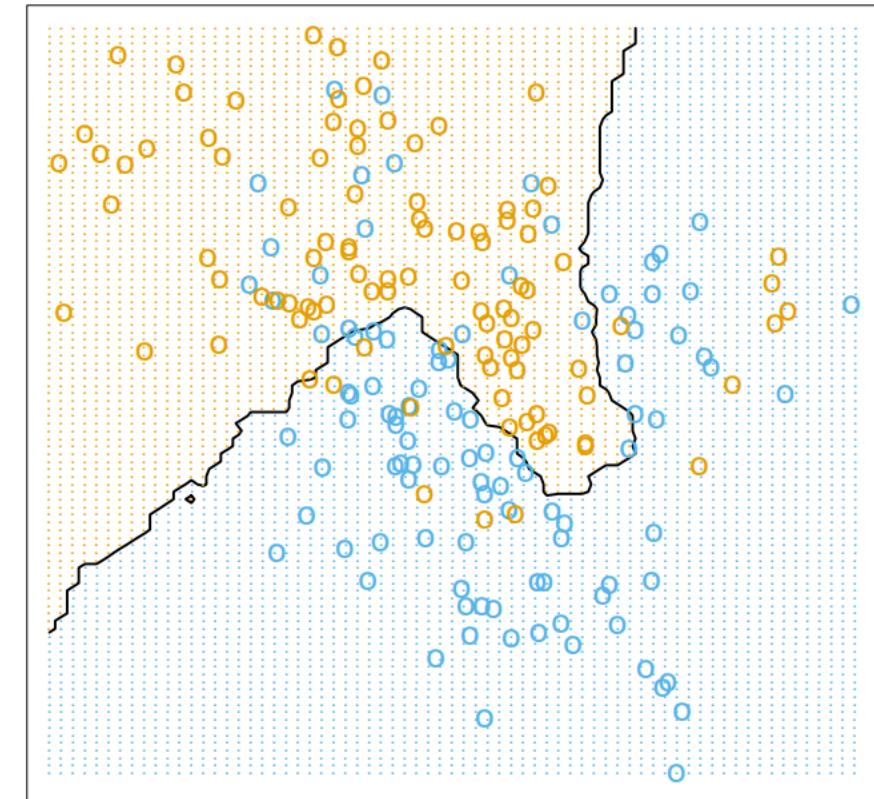


# Effects of Varying $K$

$K = 1$  Overfits



$K = 15$



## Tradeoffs in Choosing $K$

$K$  is small  $\Rightarrow$  overfitting (high variance)  
 $K$  is large  $\Rightarrow$  underfitting (high bias)

$K$  is too small

- Decision boundary very complex
- Captures fine-grained patterns
- Highly sensitive to noise

$K$  is too large

- Decision boundary very smooth
- Misses important/local patterns
- Produces stable predictions by averaging many neighbours

**May overfit (high variance)**

**May underfit (high bias)**

# Tuning Hyperparameters

Choosing  $K$  in KNN

$$k < \sqrt{N}$$

General approach

No single “best” choice of  $K$

Depends on dataset size  $N$

Common heuristic:  $K < \sqrt{N}$

Use a validation set  
(or cross-validation)  
to select hyperparameters



# Training, Validation, and Test Sets

---

# Why Split the Data?

What do we need to do with our data?

- Train a model.
- Choose hyperparameters before training a model.
- Evaluate our model on unseen data.

Solution:

- Split available data into training, validation, and test sets.

$$\sum_i \ell_i C + \frac{1}{N} = y_i$$

## Roles of Training, Validation, and Test Sets

- Training set is used to **learn model parameters**.
- Validation set is used to **tune hyperparameters** and choose between models.
- Test set is used **only once** at the end to evaluate how well our model generalizes to unseen data.

Citrus fruit example:  
width and height of  
various fruit

Training Set	Validation Set	Test Set
Learn model parameters	Tune hyperparameters	Measure generalization error once, at the end

## Workflow with KNN

$$k \in \{1, 2, 3\}$$

Divide data into training, validation, and test sets.  
Train KNN model on the training set  
Find out which K hyperparameter is most accurate  
Find the accuracy of the most accurate on new training set

1. Divide the dataset into training, validation and test sets.
2. For each  $K \in \{1, 2, 3\}$ , **train** a KNN model on the **training** set.

What happens at step 2?

3. Compute accuracy of each KNN model on the **validation** set.

What happens at step 3?

4. Select value of  $K$  with best validation accuracy. Assume  $K = 3$  is best.
5. Compute accuracy of the KNN model with  $K = 3$  on the **test** set.

## Exercise: Which should you not do with the test set?

- A. Perform data exploration.
- B. Perform feature engineering.
- C. Learn model parameters.
- D. Tune hyperparameters.
- E. Report performance on the set as the final unbiased result.



# Limitations of KNN

---

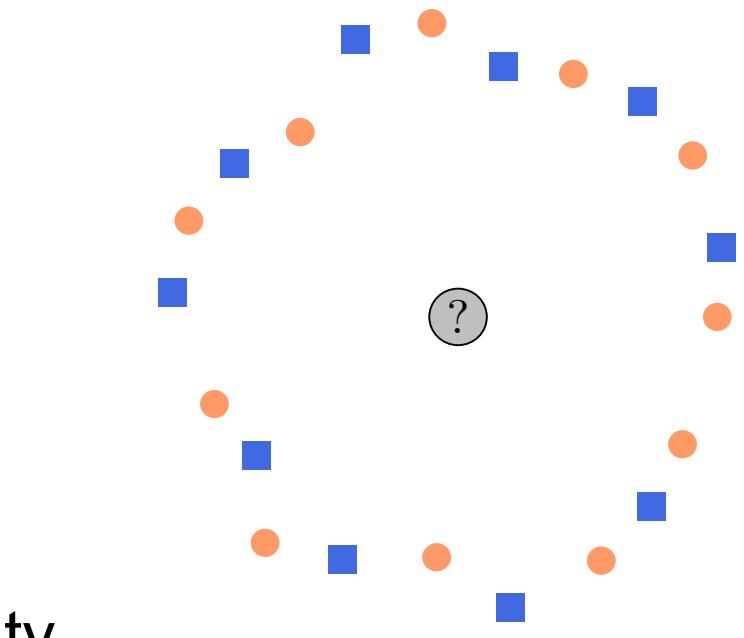
# Curse of Dimensionality

As dimensions grow, the data space expands **exponentially**,  
and points becomes **very far apart**.

If **the nearest neighbour** is **as far as any other point**,  
it is not very useful for prediction.

Problem: KNN is not suitable for high-dimensional data.

Possible solutions: add more data or reduce dimensionality.



# The Problem - Sensitivity to Scale

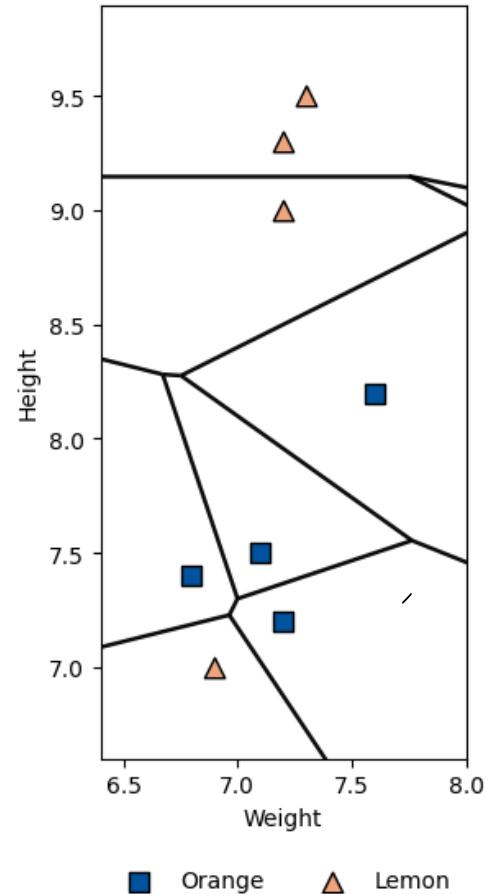
KNN is sensitive to ranges of features.

Why:

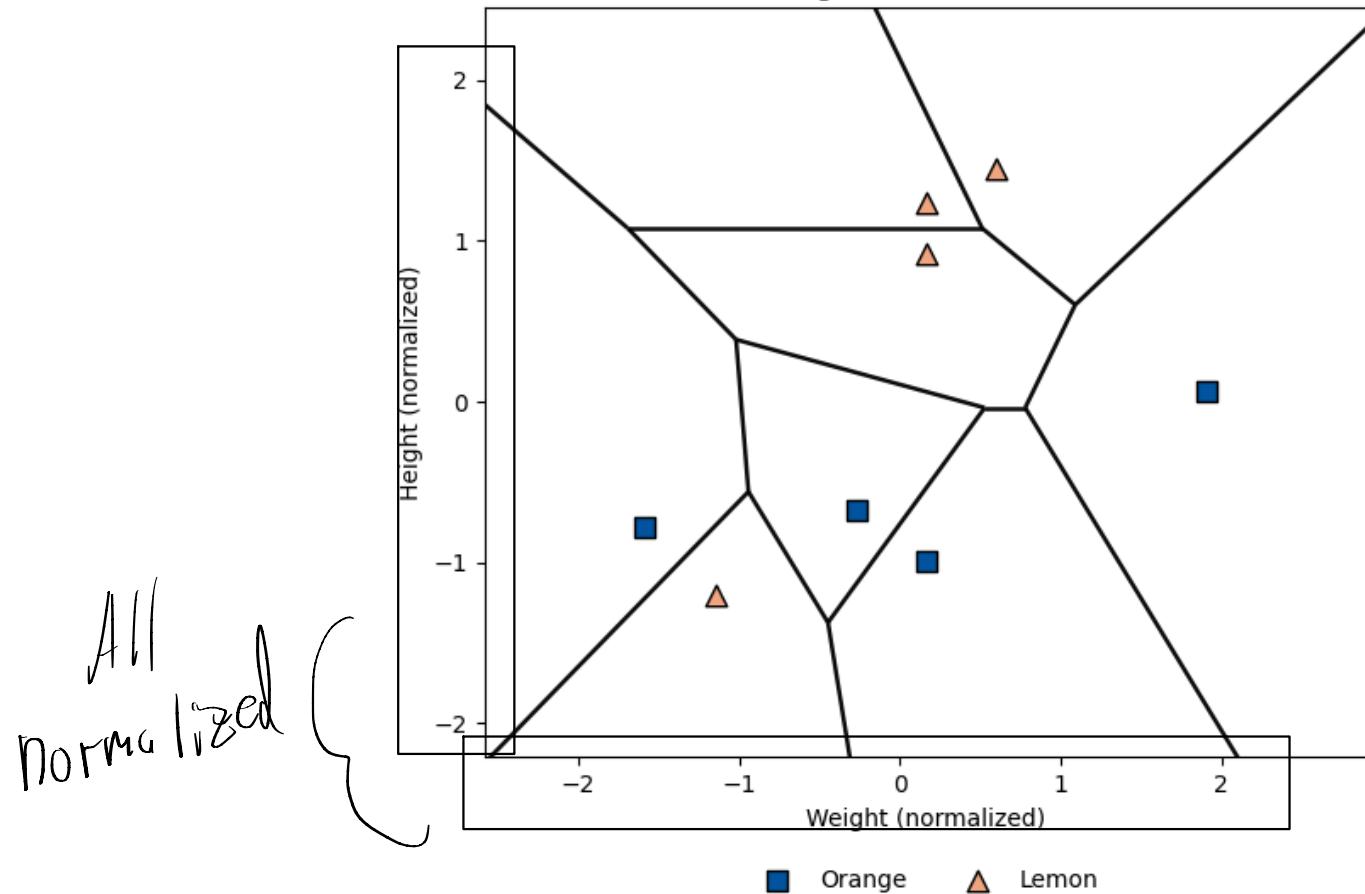
- KNN relies on distance calculations.
- A feature with a larger range dominates distance calculations.

# Predictions for Unnormalized vs. Normalized Data

Voronoi Diagram on Unnormalized Points



Voronoi Diagram on Normalized Points



$$N(\mu=0, \sigma^2=1)$$

$$x_{\text{norm}} = \frac{x - \mu}{\sigma}$$

## The Solution - Normalization

Problem: A feature with a larger range dominates distance calculations.

Solution: Normalize features to mean = 0 and variance = 1.

$$x_{\text{norm}} = \frac{x - \mu}{\sigma}$$

where  $\mu$  = mean and  $\sigma$  = the standard deviation.

# KNN: Advantages and Limitations

## Advantages

- Simple and intuitive
- No training required
- Works for both classification and regression
- Adapts easily to new data

## Limitations

- Struggles in high dimensions
- Sensitive to feature scaling
- Slow for large datasets
- Must store all training data