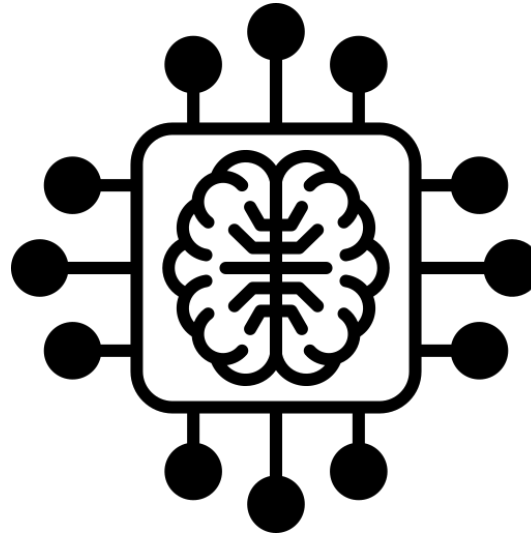# SBS4115 Fundamentals of AI & Data Analytics

# Machine Learning Fundamentals

**Lecturer:** Ir Dr Kelvin K. W. Siu

**email:** kelvinsiu@thei.edu.hk

Department of Construction, Environment and Engineering

# Intended Learning Outcomes

- By the end of this lecture, you will be able to…
  - Describe Artificial Intelligence (AI) and Machine Learning (ML)
  - Classify ML into three types
  - Understand some basic Python language
  - Appreciate the basic terminology in the ML algorithm
  - Discuss the four steps in the roadmap for building a ML system.
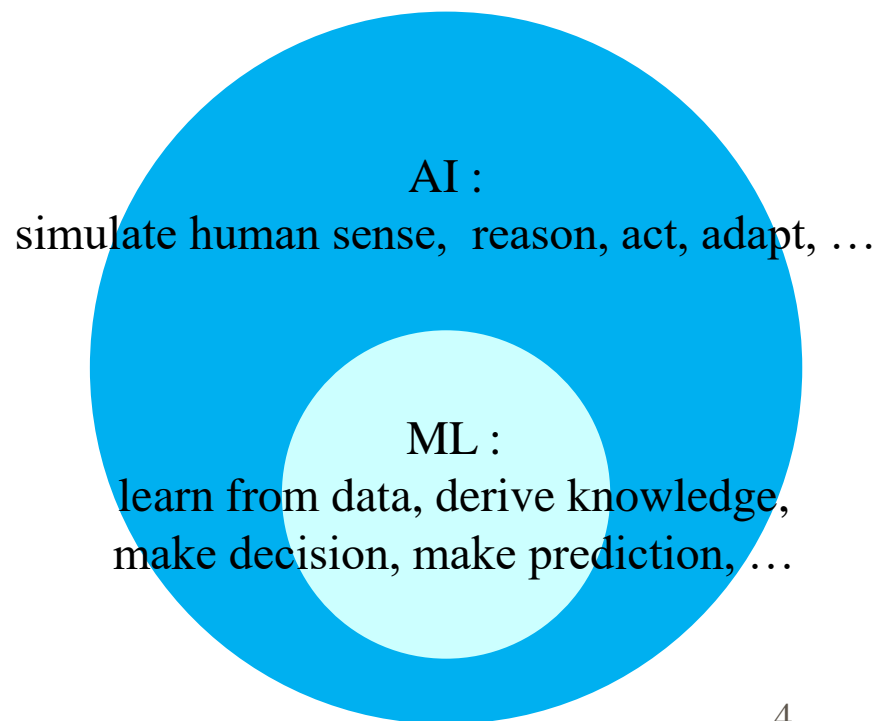
# Introduction to AI and ML

- **Artificial intelligence (AI)** refers to intelligence demonstrated by machine in opposite to human intelligence.

- The goal of AI is to **simulate human intelligence for reasoning and solving complex problems**.

- AI technology can be applied to various area such as healthcare, entertainment, finance, business, etc.



- One example of AI is AlphaGo which is a computer program that plays Go defeating the top human players in the world.

3

# Introduction to AI and ML

- **Machine learning (ML)** is a subfield of AI which enables a machine to learn from data and to derive knowledge from it.

- The goal of ML is to **make decision and prediction based on the given data**.

- In recent years, AI and ML have shown an increasing impact in applications and research areas including data science in particular.

- ML provides a more efficient way for capturing the knowledge in data to improve the performance of predictive models and make data-driven decisions.

AI :
simulate human sense, reason, act, adapt, …

ML :
learn from data, derive knowledge, make decision, make prediction, …

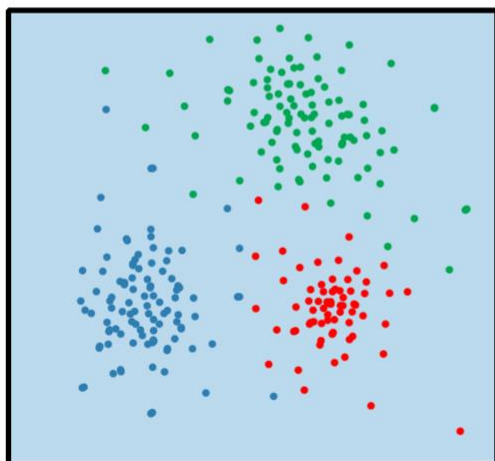# Introduction to AI and ML

- To summarize, ML can be divided into **three types**:
    1. **Supervised learning**: labelled data, direct feedback, predicted outcome, e.g. classification for predicting class labels or regression for predicting continuous outcomes
    2. **Unsupervised learning**: no labels, no feedback, find hidden structure in data, e.g. finding subgroups with clustering
    3. **Reinforcement learning**: decision process, reward system, learn series of actions, e.g. chess, computer games
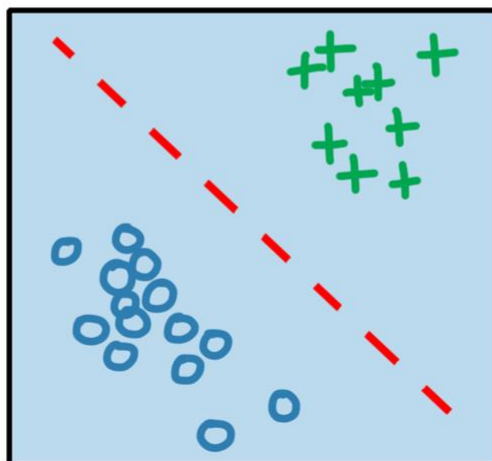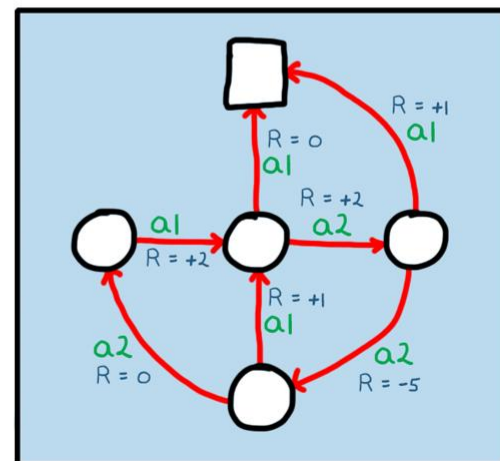
# Introduction to AI and ML



machine learning

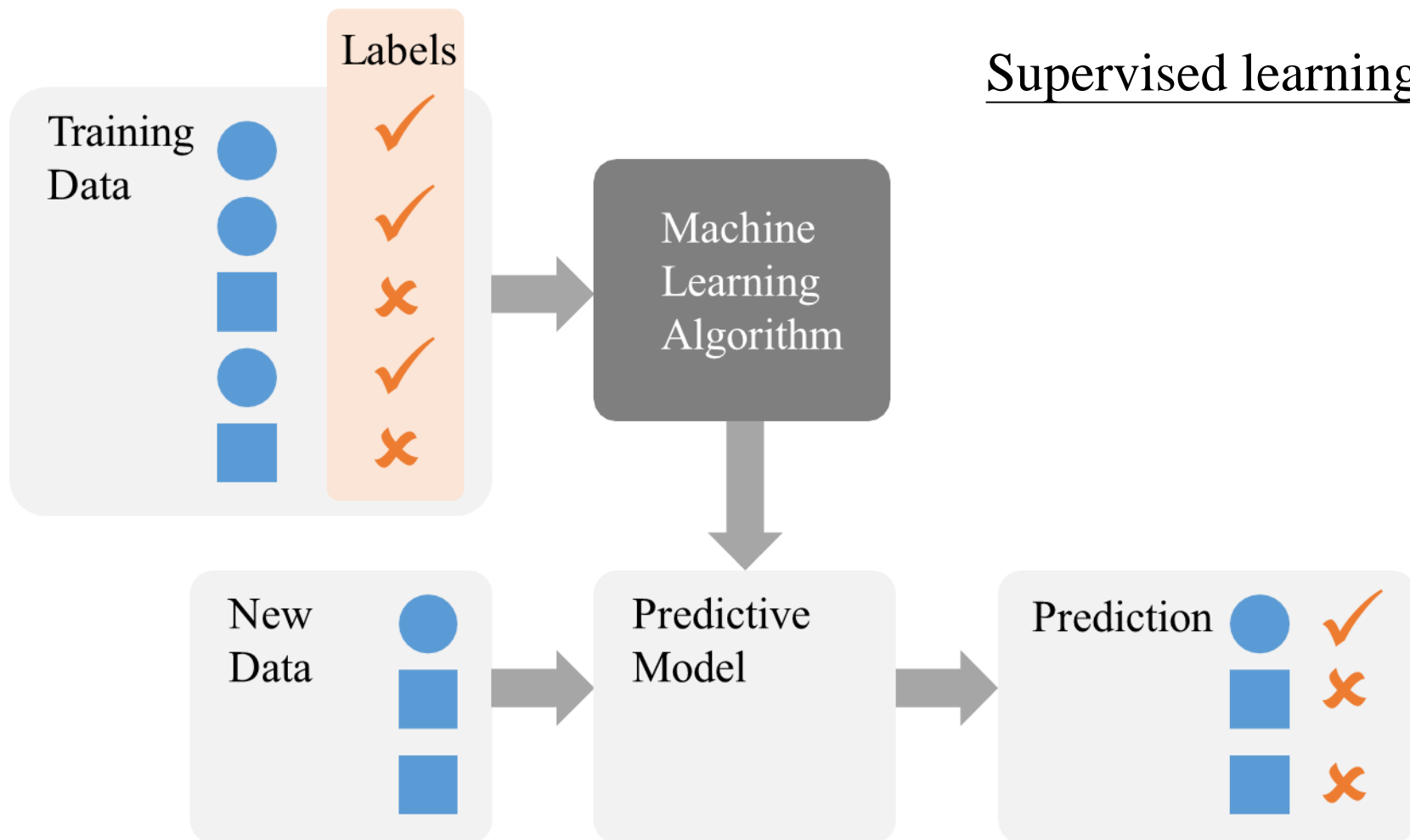unsupervised learning · supervised learning · reinforcement learning

# Introduction to AI and ML

1. **Supervised learning**

   - The main goal in supervised learning is to **learn a model from labelled training data** that allows us to **make predictions** about unseen or future data.

   - Here, the term "supervised" refers to **a set of training examples** (data inputs) where the desired output signals (labels) are already known.

   - The next figure summarizes a typical supervised learning workflow, where the labelled training data is passed to a machine learning algorithm for fitting a predictive model that can make predictions on new, unlabelled data inputs.
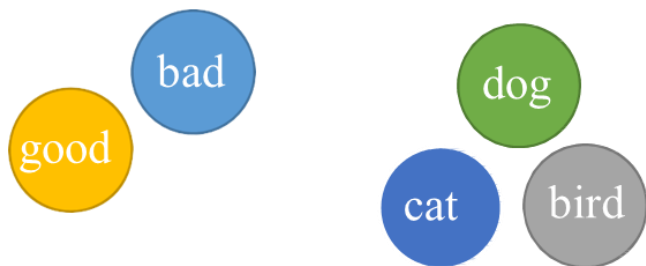
# Introduction to AI and ML

Labels

Supervised learning

Training Data

Machine Learning Algorithm

New Data

Predictive Model
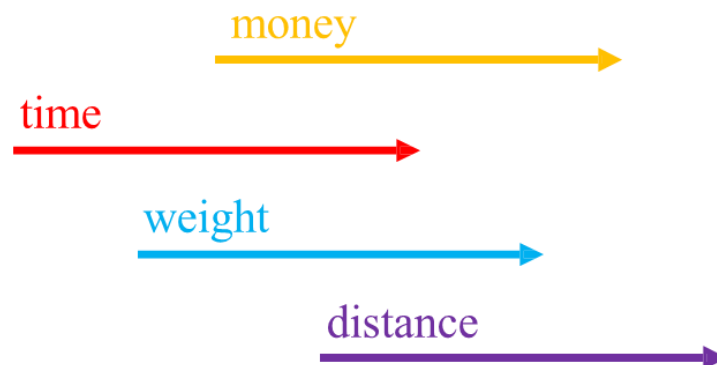
Prediction

# Introduction to AI and ML

1. **Supervised learning**
   - We will study into two major tasks of supervised learning:
     - A supervised learning task with discrete class labels is called a **classification** task.
     - Another subcategory of supervised learning is **regression**, where the outcome signal is a continuous value.
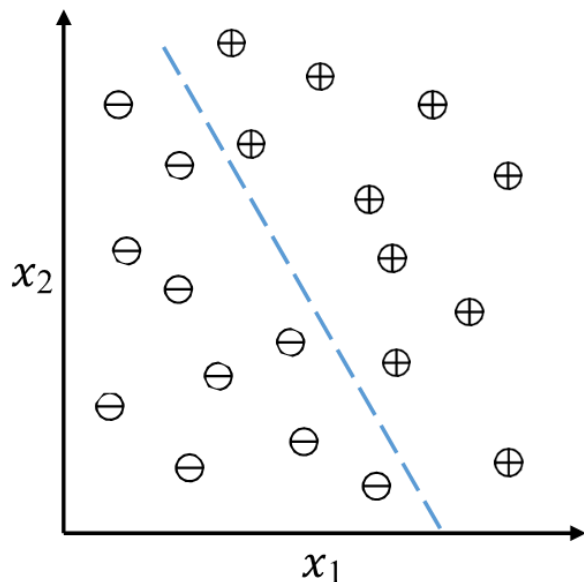
Classification - discrete labels

Regression - continuous values

bad

good

dog

cat

bird

money

time

weight

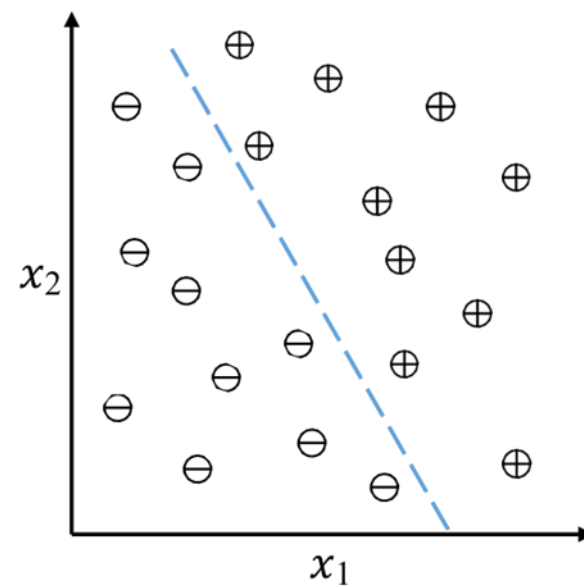distance

# Introduction to AI and ML

- **Classification** is a subcategory of supervised learning where the goal is to predict the categorical class labels of new instances, based on past observations.

- Those class labels are discrete, unordered values that can be understood as the group memberships of the instances.

- The figure on the left below illustrates the concept of a binary classification task given 20 training examples:

  - 10 training examples are labelled as the negative class (minus signs)

  - 10 training examples are labelled as the positive class (plus signs).

# Introduction to AI and ML

- In this scenario, our dataset is two-dimensional, which means that each example has two values associated with it: $x_1$ and $x_2$ along the horizontal and vertical axes respectively.

- Now, we can use a supervised machine learning algorithm to learn a rule – the decision boundary represented as a dashed line – that can separate those two classes and classify new data into each of those two categories given its $x_1$ and $x_2$ values:

# Introduction to AI and ML

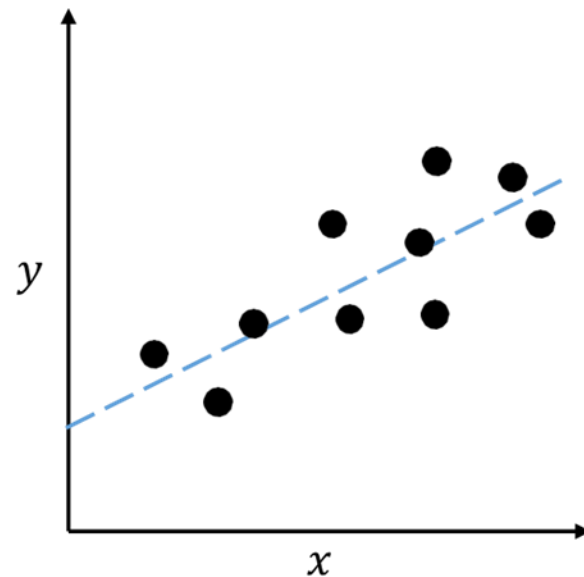- A second type of supervised learning is the prediction of continuous outcomes, which is also called **regression analysis**.

- In regression analysis, we are given a number of predictor variables (features) and a continuous response variable (outcome/target), and we try to find a relationship between those variables that allows us to predict an outcome.

# Introduction to AI and ML

- The figure on the right illustrates the concept of **linear regression**.

- Given a feature variable, $x$, and a target variable, $y$, we fit a straight line to this data that minimizes the distance between the data points and the fitted line.

- We can now use the intercept and slope learned from this data to predict the target variable of new data.

# Introduction to AI and ML

2. **Reinforcement learning**

- In reinforcement learning, the goal is to **develop a system (agent) that improves its performance based on interactions with the environment**.
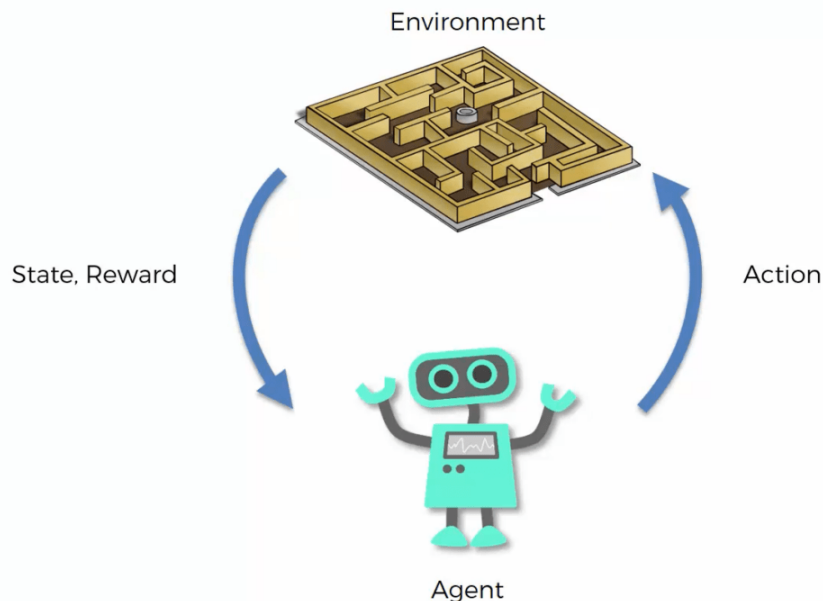

Reinforcement Learning

- Since the information about the current state of the environment typically also includes a so-called **reward signal**, we can think of reinforcement learning as a field related to supervised learning.

14

# Introduction to AI and ML

2. **Reinforcement learning**

- However, in reinforcement learning, this feedback is not the correct ground truth label or value, but a measure of how well the action was **measured by a reward function**.



- Through its interaction with the environment, an agent can then use reinforcement learning to learn a series of actions that maximizes this reward via an exploratory trial-and-error approach or deliberative planning.

# Introduction to AI and ML

2. **Reinforcement learning**

- A popular example of reinforcement learning is a chess engine.
- Here, the agent decides upon a series of moves depending on the state of the board (the environment), and the reward can be defined as win or lose at the end of the game.

# Introduction to AI and ML

2. **Reinforcement learning**

- There are many different subtypes of reinforcement learning.
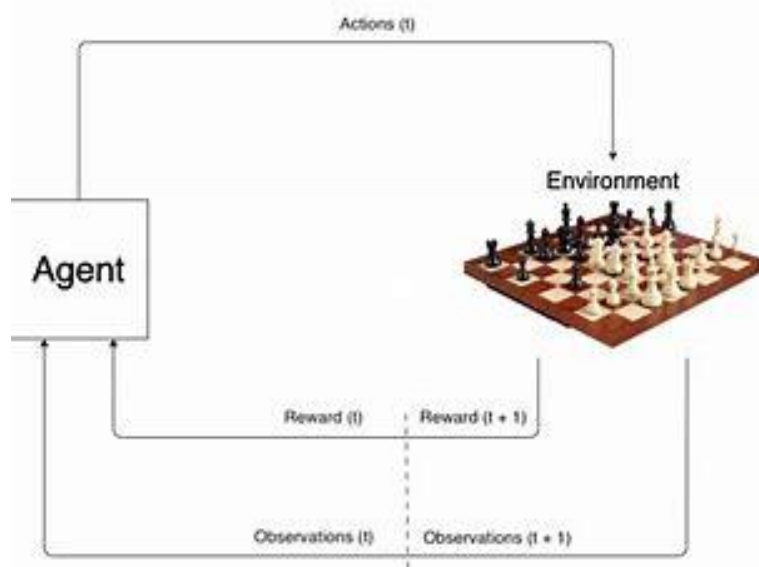- However, a general scheme is that the agent in reinforcement learning **tries to maximize the reward** through a series of interactions with the environment.



- Each state can be associated with a positive or negative reward, and a reward can be defined as accomplishing an overall goal, such as winning or losing a game of chess.
- For instance, in chess, the outcome of each move can be thought of as a different state of the environment.

17

# Introduction to AI and ML

3. **Unsupervised learning**

- In supervised learning, we know the right answer beforehand when we train a model, and in reinforcement learning, we define a measure of reward for particular actions carried out by the agent.

- In unsupervised learning, however, we are **dealing with unlabelled data or data of unknown structure**.

- Using unsupervised learning techniques, we are able to **explore the structure of our data to extract meaningful information without the guidance of a known outcome variable or reward function**.

# Introduction to AI and ML

# Introduction to AI and ML

3.  **Unsupervised learning**
    - One example is finding subgroups with **clustering**.
    - Clustering is an exploratory data analysis technique that allows us to organize a pile of information into meaningful subgroups (clusters) without having any prior knowledge of their group memberships.
    - Each cluster that arises during the analysis defines a group of objects that share a certain degree of similarity but are more dissimilar to objects in other clusters.
    - Clustering is also sometimes called **unsupervised classification**.

# Introduction to AI and ML

3.  **Unsupervised learning**

    - The following figure illustrates how **clustering** can be applied to organizing unlabelled data into three distinct groups based on the similarity of their features, $x_1$ and $x_2$:



    - Clustering is a great technique for **structuring information and deriving meaningful relationships from data**.

    - For example, it allows marketers to discover customer groups based on their interests, in order to develop distinct marketing programs.

# Introduction to AI and ML

- Consider the example below which illustrates sales of a fast-food shop.
- Each point represents a customer.
- $x_1$ represents the frequency of purchase (in a year).
- $x_2$ represents the average amount of purchase ($).
- Can you group these points into clusters and explain why?

# Introduction to Python

- High-level Interpreted Language
- No compilation needed, no gcc/g++ etc.
- Slow in general
- Not for (most) high-performance programming
- Trading time spent on programming for time spent on computation
- Popular for machine learning and data analysis

# Two ways to run Python

- Login to the below website with your Google Account
- https://colab.research.google.com/

- Use Anaconda
- https://www.anaconda.com/

# Google Colab

- Try these:
- print('Hello World')

- x = 3
- print(x, type(x))
- print(x + 1)   # Addition    Note: compare with print('x+1')
- print(x - 1)   # Subtraction
- print(x * 2)   # Multiplication
- print(x ** 2)  # Exponentiation

# A list is the Python equivalent of an array

- xs = [3, 1, 2]   # Create a list
- print(xs)
- print(xs[2])
- print(xs[-1])     # Negative indices count from the end of the list;

# Introduction to Python

- xs[2] = 'apple'    # Lists can contain elements of different types
- print(xs)


- xs.append('orange') # Add a new element to the end of the list
- print(xs)

# Introduction to Python

- num = float(input("Enter a number: "))
- print(num)

# Introduction to Python

- num1 = float(input("Enter the first number: "))
- num2 = float(input("Enter the second number: "))
- print(num1)
- print(num2)
- print('The second number you input is ', num2)

# Introduction to Python

- addition = num1 + num2
- subtraction = num1 - num2
- multiplication = num1 * num2
- division = num1 / num2

- print('Addition =', addition)

Try to do similar things for subtraction, multiplication and division…….

# Exercise 1:

- Write a Python program that asks the user to input a temperature in Celsius. The program should then convert the temperature to Fahrenheit using the formula and display the result. The conversion formula is:

$$F = \frac{9}{5} \times C + 32$$

Where F is the Fahrenheit temperature and C is the Celsius temperature.

# Array creation

import numpy as np

```python
my_array = np.array([1, 2, 3])
my_array
```

```python
my_2d_array = np.array([[1, 2, 3], [4, 5, 6]])
my_2d_array
```

```python
my_2d_array[1, 2]
```

What is the result? Why?

# Exercise 2:

- Write a Python program that:

Creates a 1-dimensional array of numbers (e.g. 1 x 5 entries)

e.g. [2, 5, 15, 35, 3]

Performs the following operations on the array:

Calculate the sum of all elements

Multiply all elements by 2

Find the maximum and minimum value in the array.

# Exercise 2:

- Some commands that you may use:

np.sum(numbers)
np.mean(numbers)

# Maximum and minimum value
np.max(numbers)
np.min(numbers)
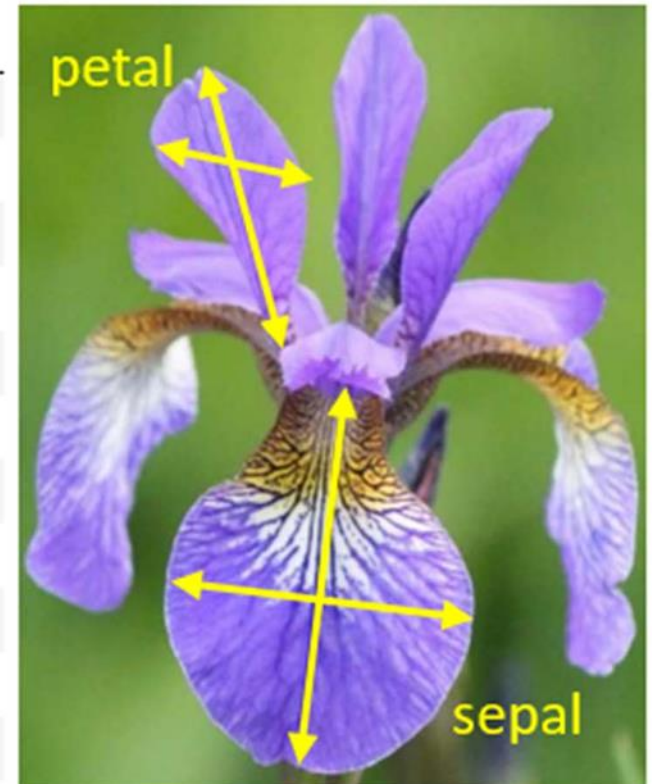
# Basic Terminology and Notation

- Before we study the algorithm of machine learning, we will look into the basic terminology, including the common terms referring to different aspects of a dataset and also the mathematical notations.

- To make it more precise, we will use a very classical and popular **dataset of Iris flower** in many data science or machine learning textbook.

- The dataset "iris.csv" contains the measurements of 150 Iris flowers from three different species – Setosa, Versicolor, and Virginica.

- Here, **each flower example represents one row in our dataset**, and the flower measurements in centimeters are stored as columns, which we also call the features of the dataset.

```
import pandas as pd
list1 = ["sepal length","sepal width","petal length","petal width","class label"]
df = pd.read_csv("iris.csv",header=None,names=list1)
```

| | sepal length | sepal width | petal length | petal width | class label |
|---|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |
| ... | ... | ... | ... | ... | ... |
| 145 | 6.7 | 3.0 | 5.2 | 2.3 | Iris-virginica |
| 146 | 6.3 | 2.5 | 5.0 | 1.9 | Iris-virginica |
| 147 | 6.5 | 3.0 | 5.2 | 2.0 | Iris-virginica |
| 148 | 6.2 | 3.4 | 5.4 | 2.3 | Iris-virginica |
| 149 | 5.9 | 3.0 | 5.1 | 1.8 | Iris-virginica |



150 rows × 5 columns

- These include sepal length, sepal width, petal length and petal width.
- The last column is the class labels which contain the species of each flower.

36

# Basic Terminology and Notation

- We can represent the features part of the dataset **using a matrix**.

- The number of rows refers to the number of samples.

- The number of columns refers to the number of attributes in the features.

- In the Iris example, we have 150 samples with 4 attributes (class label is not included).

- This gives us a $150 \times 4$ matrix:

  where an entry $x_j^{(i)}$ represents the $j$-th attribute of the $i$-th sample.

$$\begin{bmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} & x_4^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} & x_4^{(2)} \\ \vdots & \vdots & \vdots & \vdots \\ x_1^{(150)} & x_2^{(150)} & x_3^{(150)} & x_4^{(150)} \end{bmatrix}$$

# Basic Terminology and Notation

- In particular, all the attributes of the $i$-th sample is represented by the $i$-th row:

$$\boldsymbol{x}^{(i)} = \begin{bmatrix} x_1^{(i)} & x_2^{(i)} & x_3^{(i)} & x_4^{(i)} \end{bmatrix}$$

- Also, the $j$-th attribute of all the samples is represented by the $j$-th column:

$$\vec{x}_j = \begin{bmatrix} x_j^{(1)} \\ x_j^{(2)} \\ \vdots \\ x_j^{(150)} \end{bmatrix}$$

- Our target variable (the class labels) will be stored in a single column:

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(150)} \end{bmatrix}$$

# Basic Terminology and Notation

- Recall the python data analytic techniques in a previous chapter.
- You can extract any entry, row or column from the DataFrame as shown in the example below.

```
df.loc[3]
```

```
sepal length              4.6
sepal width               3.1
petal length              1.5
petal width               0.2
class label        Iris-setosa
Name: 3, dtype: object
```

```
df.loc[3,"petal length"]
```
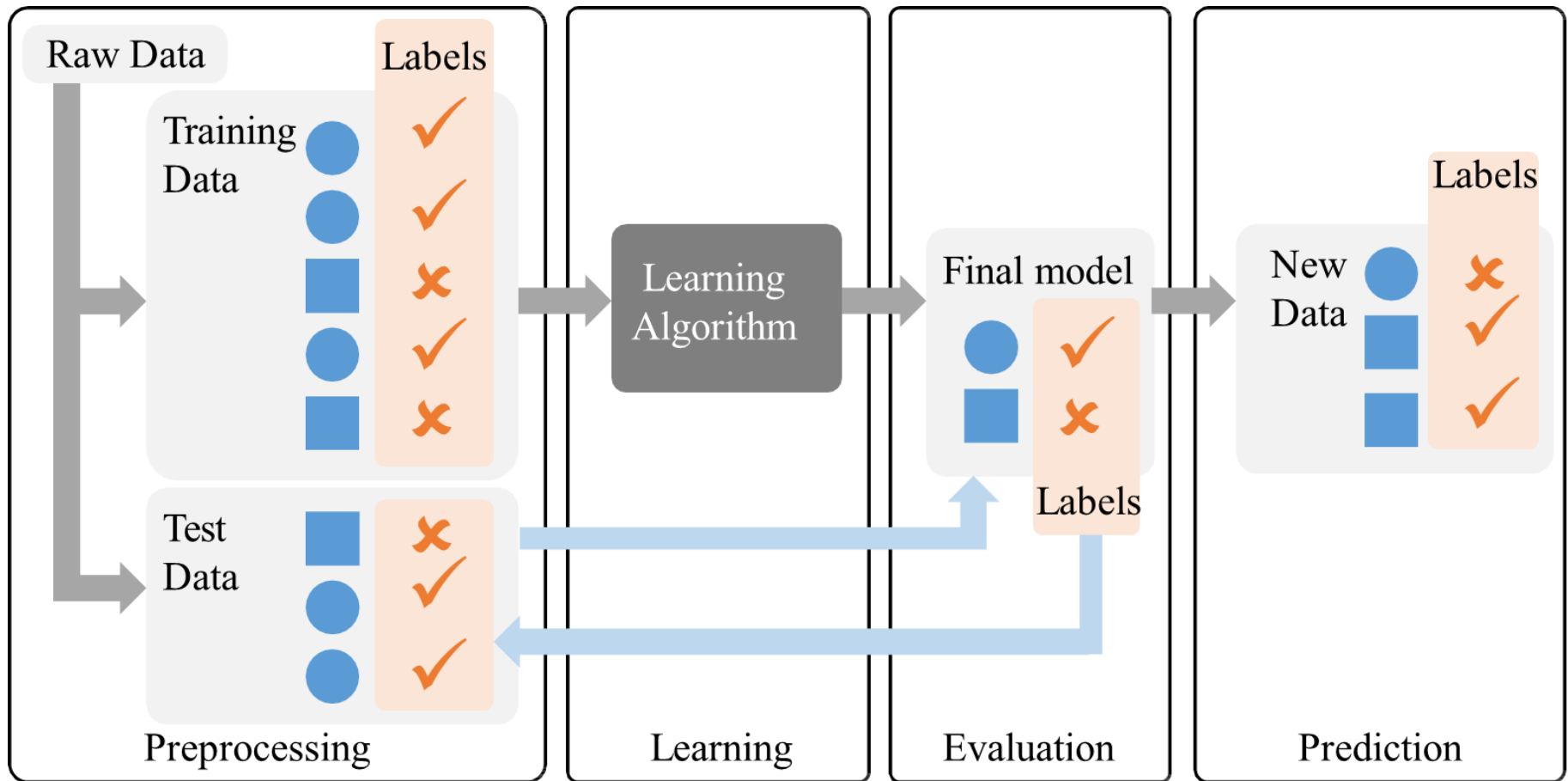
```
1.5
```

```
df["petal length"]
```

```
0        1.4
1        1.4
2        1.3
3        1.5
4        1.4
        ...
145      5.2
146      5.0
147      5.2
148      5.4
149      5.1
Name: petal length, Length: 150, dtype: float64
```

# Roadmap for Building ML System

- The Iris dataset is an example of labelled data.
- We will split it into:
    - a **training set** for training a machine learning model, and
    - a **testing set** for verifying the performance of our model.
- This will be introduced in the next chapter.
- Beforehand, we need to have **a roadmap of such building a machine learning system**, with a number of steps.
- The diagram on the next slide shows a typical workflow for using machine learning in **predictive modelling** which we are going to discuss in the remaining chapters.
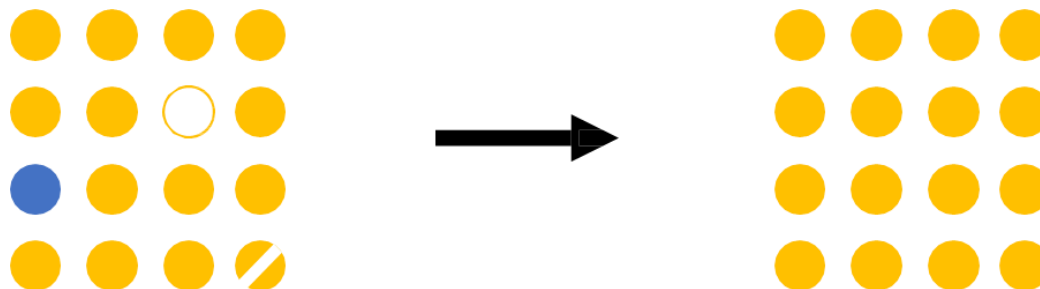
- This can be concluded into four steps:
  1. preprocessing
  2. learning
  3. evaluation
  4. prediction

# Roadmap for Building ML System

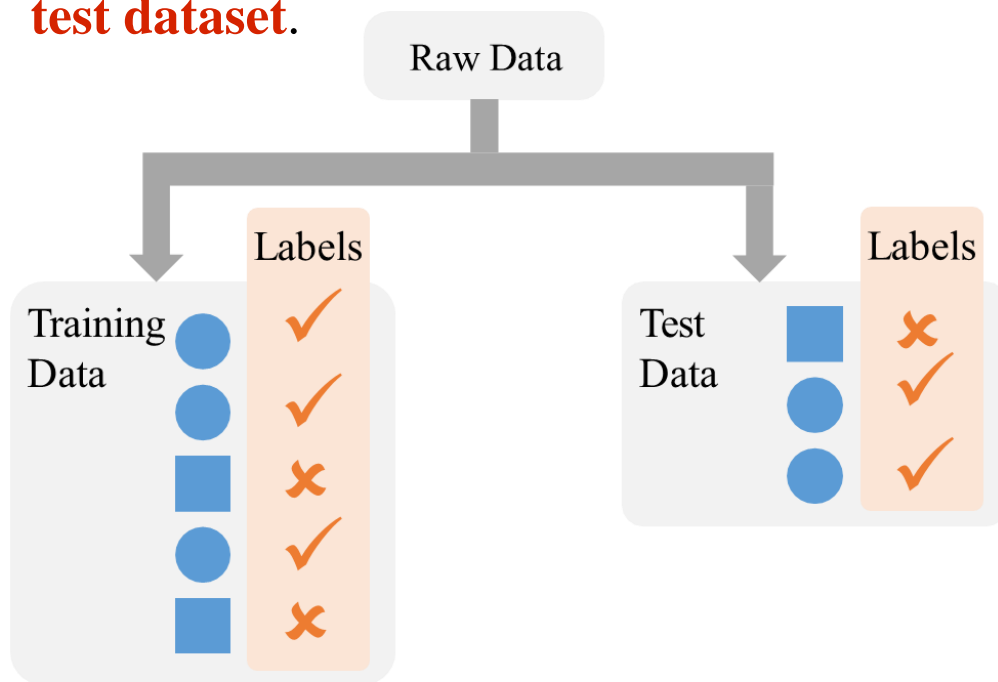1. **Preprocessing** – getting data into shape
   - Raw data might come in different forms and shapes that is not suitable for optimal performance of a learning algorithm.
   - Preprocessing refers to the step of **extracting useful features from the raw data and to convert it into desired form**.
   - This usually include data cleaning, standardization, dimension reduction, noise reduction, etc.
   - In many cases, this step is crucial to make the algorithm more efficient and also improves predictive performance of the model.

# Roadmap for Building ML System

1. **Preprocessing** – getting data into shape
   - To determine whether our machine learning algorithm not only performs well on the training dataset but also generalizes well to new data, we also want to **randomly divide the dataset into a separate training and test dataset**.
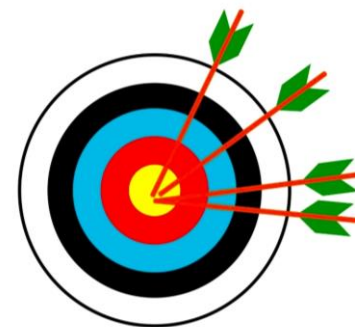


   - We use the training dataset to **train and optimize our machine learning model**, while we keep the test dataset until the very end to **evaluate the final model**.

# Roadmap for Building ML System

2. **Learning** – training and selecting a predictive model

- There are many different machine learning algorithms developed to solve different types of problems.

- Each of these algorithm is **based on some assumptions** and therefore has its **inherent biases**.

- In practice, it is essential to **compare the result of different algorithms** in order to train and select the best performing model.

- For such comparison, we need to **set up a metric to measure performance**.

- One commonly used metric is **classification accuracy**.

# Roadmap for Building ML System

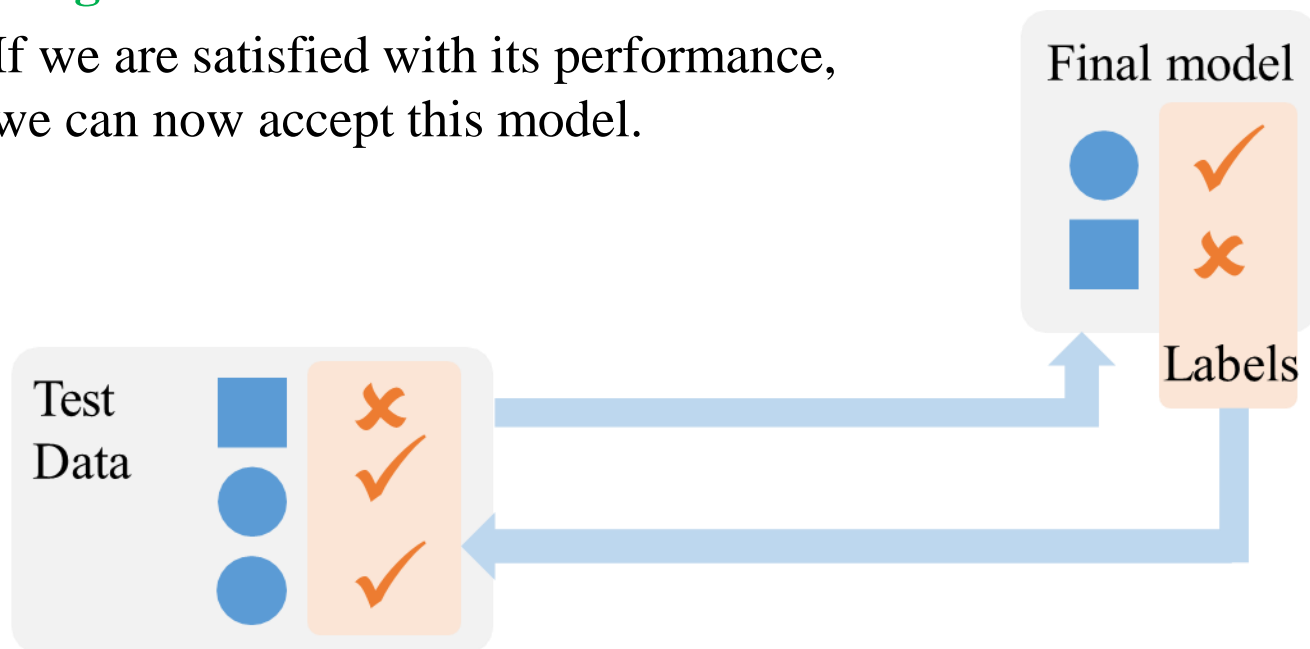2. **Learning** – training and selecting a predictive model
   - In order to make sure that the performance of the model is **not dependent to the selection of final test data**, various techniques of cross validation is used.
   - This refers to **dividing a dataset into training and validation subsets** in order to estimate the **generalization performance of the model**.

Training Data → Learning Algorithm → Final Model

# Roadmap for Building ML System

3.  **Evaluation** – evaluating models

    - After we have selected a model that has been fitted on the training dataset, we can **use the test dataset to estimate its performance and the generalization error**.

    - If we are satisfied with its performance, we can now accept this model.

# Roadmap for Building ML System

4. **Prediction** – predicting unseen data instances

- After all the previous steps, we come up with a model to predict new, future data.

- It is important to note that the parameters for the previously mentioned procedures, such as feature scaling and dimensionality reduction, are **solely obtained from the training dataset**, and the same parameters are later reapplied to transform the test dataset, as well as any new data instances.

- The performance measured on the test data may be overly optimistic otherwise.

# Checklist

- Can you:
    1. Describe Artificial Intelligence (AI) and Machine Learning (ML)?
    2. Classify ML into three types?
    3. Appreciate the basic terminology in the ML algorithm?
    4. Discuss the four steps in the roadmap for building a ML system?