

# DDA3020 Machine Learning: Lecture 01 Introduction

Baoyuan Wu  
School of Data Science, CUHK-SZ

Jan 08/10, 2024

# Outline

- 1 About this course
- 2 What is machine learning
- 3 Supervised learning
- 4 Unsupervised learning
- 5 Some basic concepts in machine learning

- 1 About this course
- 2 What is machine learning
- 3 Supervised learning
- 4 Unsupervised learning
- 5 Some basic concepts in machine learning

# Instructors and teaching assistants (session 1)

## Time and venue:

- Lecture: Tuesday/Thursday 01:30–02:50 pm, TCx C201
- Tutorial: **TBD**

## Instructor:

- **Haizhou Li:** Professor, IEEE Fellow, ISCA Fellow, SAEng Fellow
- Email: [haizhouli@cuhk.edu.cn](mailto:haizhouli@cuhk.edu.cn) Homepage: <https://colips.org/~eleliha/>
- Office hour (OH): Wednesday biweekly (4:00-5:30pm, 06 Mar, 20 Mar, 03 Apr, 17 Apr, 4:30-6:00pm 10 May), DY 512

## TAs:

- **Haoying Li:** [222043004@link.cuhk.edu.cn](mailto:222043004@link.cuhk.edu.cn)  
OH: Thu. 2:30-3:30pm, 4F-38, Zhixin Building
- **Jiawei Xu:** [223040240@link.cuhk.edu.cn](mailto:223040240@link.cuhk.edu.cn)  
OH: Tue. 2:00-3:00pm, 4F-86, Zhixin Building
- **Zihan Zhou:** [zzh1282260738@gmail.com](mailto:zzh1282260738@gmail.com)  
OH: Wed. 4:00-5:00pm, 4F-88, Zhixin Building
- **Zhijun Liu:** [zhijunliu1@link.cuhk.edu.cn](mailto:zhijunliu1@link.cuhk.edu.cn)  
OH: Thu. 1:30-2:30pm, 4F-115, Zhixin Building

## To request for letters of recommendation:

- Expectation: Regular participation in lectures, internship or research experience
- Preparation: Self-evaluation, CV, transcript...
- Limits: 20 requests per student-year in 2 batches
- Processing time: 4-5 weeks after all materials have been submitted \*university links should be sent at one time, please provide a summary of all requests

# Instructors and teaching assistants (session 2)

## Time and venue:

- Lecture: Monday/Wednesday 08:30–09:50 am, AB-E205
- Tutorial: **TBD**

## Instructor:

- **Baoyuan Wu:** Associate Professor, SDS (wubaoyuan@cuhk.edu.cn)
- Personal homepage: <https://sites.google.com/site/baoyuanwu2015>
- Office hour (OH): 10:20-11:40 am, Monday biweekly (15 Jan, 29 Jan, 26 Feb, 11 Mar, 25 Mar, 08 Apr, 22 Apr, 8 May), DY 411

## TAs:

- **Haoying Li:** 222043004@link.cuhk.edu.cn  
OH: Thu. 2:30-3:30pm, 4F-38, Zhixin Building
- **Jiawei Xu:** 223040240@link.cuhk.edu.cn  
OH: Tue. 2:00-3:00pm, 4F-86, Zhixin Building
- **Zihan Zhou:** zzh1282260738@gmail.com  
OH: Wed. 4:00-5:00pm, 4F-88, Zhixin Building
- **Zhijun Liu:** zhijunliu1@link.cuhk.edu.cn  
OH: Thu. 1:30-2:30pm, 4F-115, Zhixin Building

# Agenda

Week	Content	
W1	Introduction	
W2	Review of Probability & Linear Algebra	
W3	Optimization	
W4	Linear Regression I	
W5	Linear Regression II & Logistic Regression	Homework 1
W6	Logistic Regression & Support Vector Machines I	
W7	Support Vector Machines II	Homework 2
W8	Decision Tree and Random Forest	
W9	Neural Networks	
W10	CNN & Over-fitting, Bias-Variance Tradeoff	Homework 3
W11	Performance Evaluation	
W12	K-Means, Mixture Models	Homework 4
W13	Expectation Maximization (EM)	
W14	Dimensionality Reduction & PCA	

# Learning materials

## Textbooks:

- Required: Andriy Burkov, “The Hundred-Page Machine Learning Book”, 2019. (read first, buy later: [Link](#))
- Required: K. Murphy. Machine Learning: A Probabilistic Perspective. MIT Press, 2012. ([Link](#))
- Recommended: Andreas C. Muller and Sarah Guido, “Introduction to Machine Learning with Python: A Guide for Data Scientists”, O’Reilly Media, Inc., 2017. (The PDF will be uploaded to the BB system, only used for learning)
- Recommended: Jeff Leek, “The Elements of Data Analytic Style: A guide for people who want to analyze data”, Lean Publishing, 2015. ([Link](#))
- Recommended: C. Bishop. Pattern Recognition and Machine Learning. Springer, 2011. ([Link](#))

**Lecture slides:** The slides of the following week will be uploaded to the BB system before the lecture.

# Grading policy

## Grading:

- 30% Written homework assignments (4 times)
- 30% Programming homework assignments (Python/Scikit-learn) (4 times)
- 40% Final exam

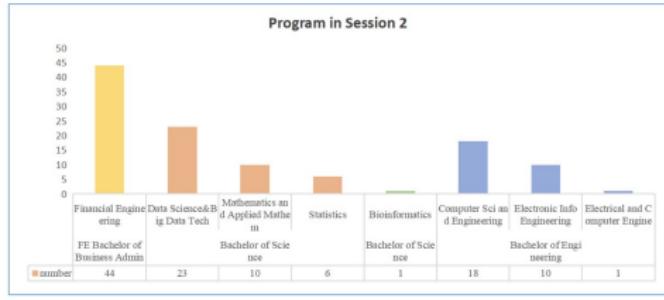
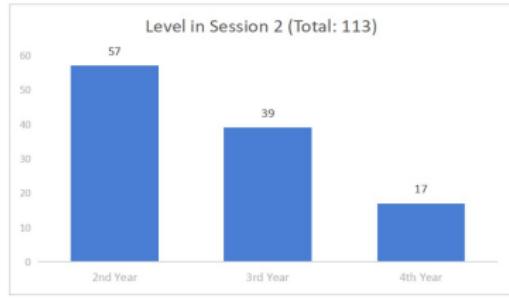
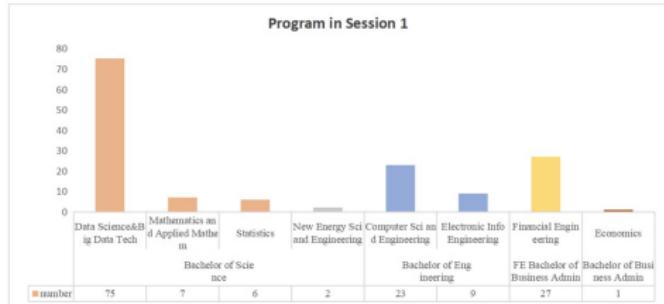
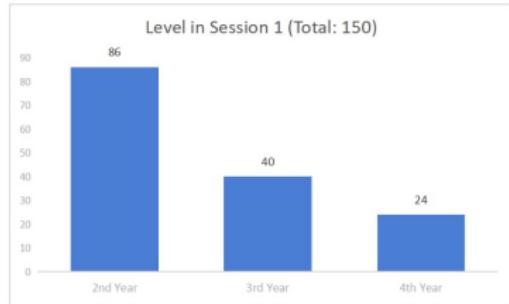
**Submission:** You have 2 weeks to **independently** complete each assignment (written + programming). Late submission will get discounted score:  $(0, 24]$  hours  $\rightarrow$  80%;  $(24, 120]$  hours  $\rightarrow$  50%;  $(120, \infty)$  hours  $\rightarrow$  0%.

**Requirements to finish the homework smoothly:** basic knowledge of probability and linear algebra, familiar with Python programming, correct understanding of the teaching content.

**Plagiarism:** Zero mark is given for the whole assignment (including written and programming) at the first plagiarism case. Students will **FAIL** the whole course for repeated plagiarism. **Note:** if there are heavy overlaps between two answers, then both will be identified as plagiarism (we don't time to distinguish). Thus, discussions are encouraged, but you must finish the assignment by yourself, and don't share your answer to others.

# Students' background

- The students come from **diverse years and majors**, thus, may have different expectations and experiences.



# Course difficulty

- Machine learning is a **complex** and **evolving** discipline. Our course will cover the **fundamental concepts and algorithms** in machine learning.
- **Fundamental**  $\neq$  **SIMPLE**. Machine learning is built on several other foundation subjects, including **probability**, **linear algebra**, **optimization**, **computer science**, *etc.*
- If you find it **difficult** to understand the teaching content, please don't hesitate to tell me or TAs, we will help you overcome the challenges. If you feel the content is **too easy**, no worry, we can provide further reading materials on advanced and emerging topics (*e.g.*, deep learning, adversarial machine learning). Wish each of you an enjoyable learning experience!

# Communications

**Smooth communications** between you and us are very important to build a successful course. If any question/difficulty/suggestion, you are welcome to

- Talk to me directly, after class or during office hour
- Email or talk to your TAs, questions will be collected and sent to me weekly for responses
- Questions about programming are handled by TAs
- You may also send me emails if you need my help

- 1 About this course
- 2 What is machine learning
- 3 Supervised learning
- 4 Unsupervised learning
- 5 Some basic concepts in machine learning

# Definition of machine learning

- Arthur Samuel: “the field of study that gives computers the ability to learn without being explicitly programmed.”
- Tom Mitchell: “A computer program is said to learn from **experience E** with respect to some class of **tasks T** and **performance measure P**, if its performance at tasks in T, as measured by P, improves with experience E.”

# Definition of machine learning

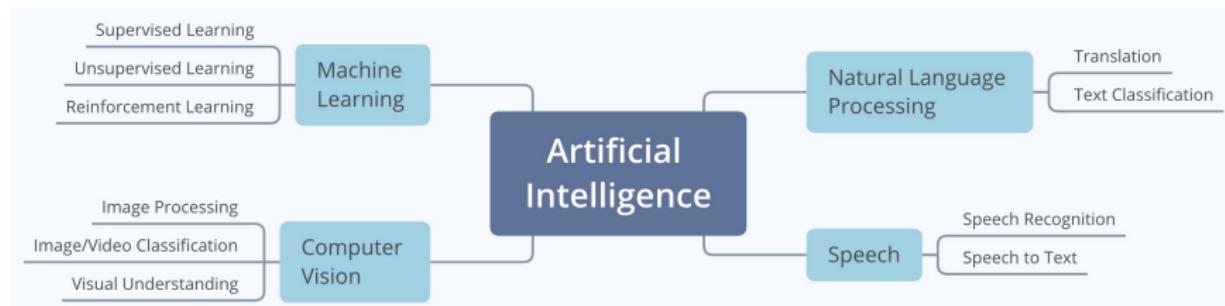
Tom Mitchell: “A computer program is said to learn from **experience E** with respect to some class of **tasks T** and **performance measure P**, if its performance at tasks in T, as measured by P, improves with experience E.”

Suppose your email program watches which emails you do or do not mark as spam, and based on that learns how to better filter spam. What is the task T in this setting?

- Classifying emails as spam or not spam.
- Watching you label emails as spam or not spam.
- The number (or fraction) of emails correctly classified as spam/not spam.
- None of the above – this is not a machine learning problem.

# ML is a branch of artificial intelligence

- **Artificial intelligence (AI)** is intelligence demonstrated by **machines** (e.g., computer, robots), unlike the natural intelligence displayed by humans and animals, which involves consciousness and emotionality.
- AI covers many topics, such as machine learning (ML), computer vision (CV), natural language processing (NLP) and speech processing.



**Figure:** Machine learning is one of the most important branches of artificial intelligence.

# Connection with other disciplines

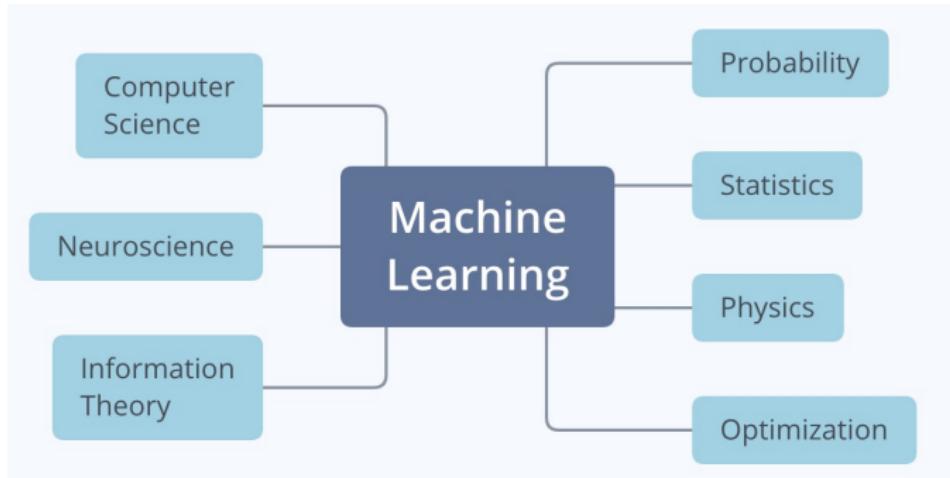
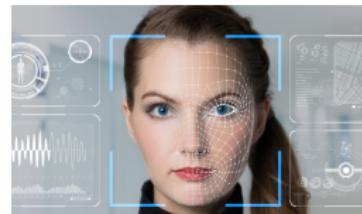
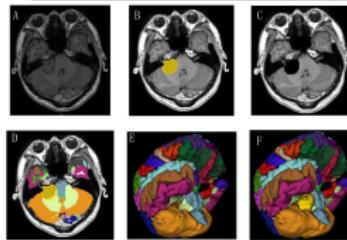
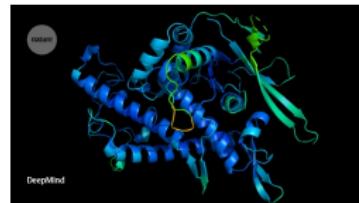
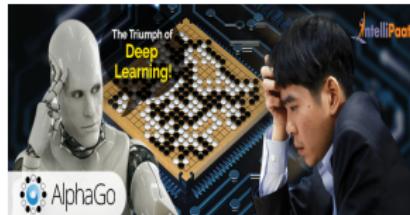


Figure: Machine learning can be seen as an interdisciplinary.

# Applications of machine learning



**Figure:** Machine learning has been widely used in many mission-critical tasks, such as Game (AlphaGo), protein structure prediction (AlphaFold2), EEG signal processing, medical image diagnosis, face recognition, etc..

Let's see a video!

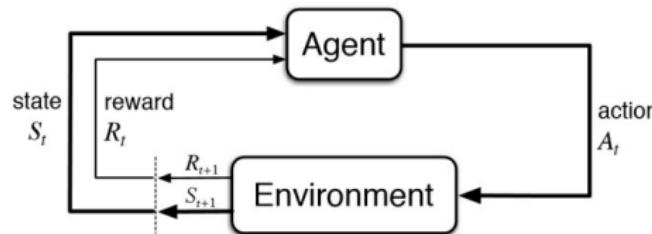
# Two basic paradigms of ML

Given some data  $x_1, x_2, x_3, \dots$ , you can do

- **Supervised learning:** you are also provided some human-labeled outputs  $y_1, y_2, y_3, \dots$ , and your task is to learn a mapping function from one input data  $x_i$  to one output  $y_i$ . **Learning from teacher**
- **Unsupervised learning:** your task is to build/learn a good model of  $x$ , such that some characteristics of the data could be revealed, such as clustering, dimensionality reduction, *etc.* **Learning by oneself**

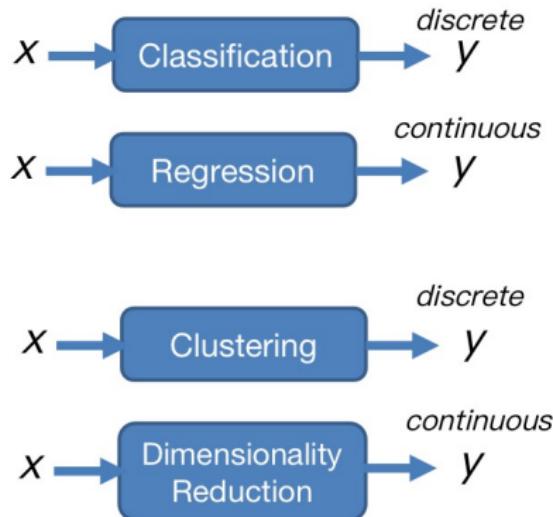
# Other learning paradigms

- **Reinforcement learning:** you can make some actions  $a_i$  to change the data  $x_i$  (the state of environment), and you will receive some rewards/punishments  $r_i$ . Gradually, you will automatically learn to make the suitable action for different data to get more rewards. It is one of three basic machine learning paradigms, alongside supervised learning and unsupervised learning.  
**Learning from rewards/punishments**



In this course, we will focus on **supervised learning** and **unsupervised learning**.

# Supervised vs. unsupervised learning



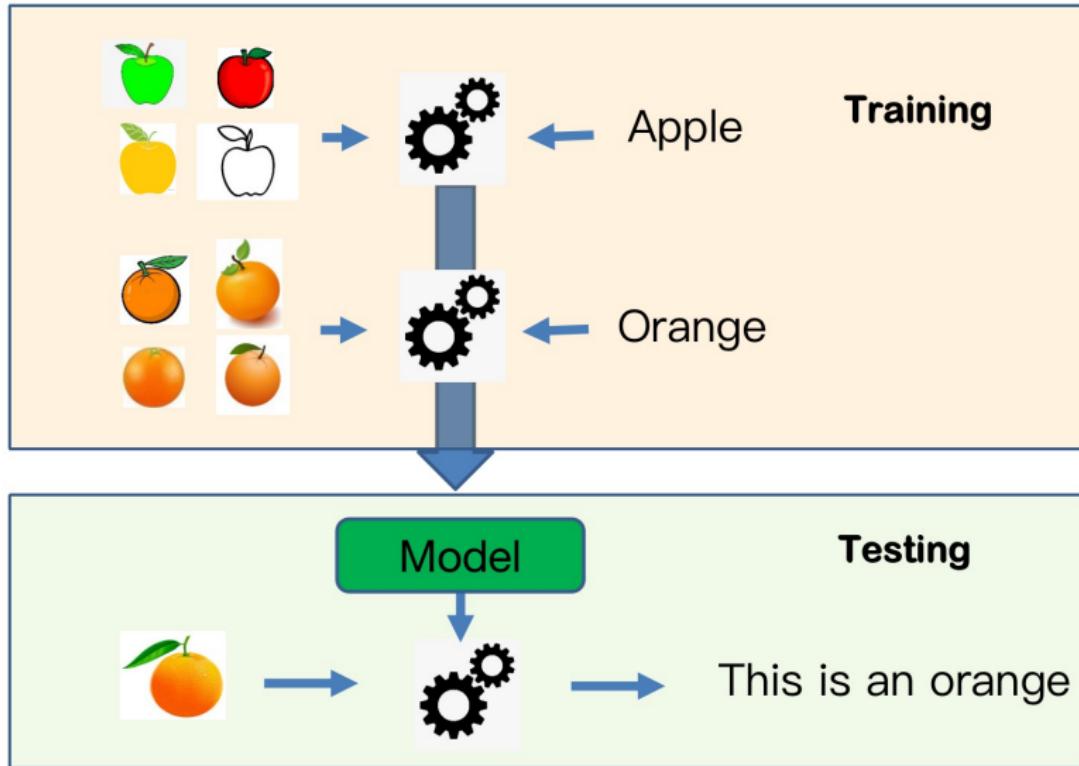
	Supervised Learning (with labels)	Unsupervised Learning (without labels)
Discrete	Classification	Clustering
Continuous	Regression	Dimensionality Reduction

Refer to:

<https://towardsdatascience.com/supervised-vs-unsupervised-learning-14f68e32ea8d>

- 1 About this course
- 2 What is machine learning
- 3 Supervised learning
- 4 Unsupervised learning
- 5 Some basic concepts in machine learning

# Supervised learning



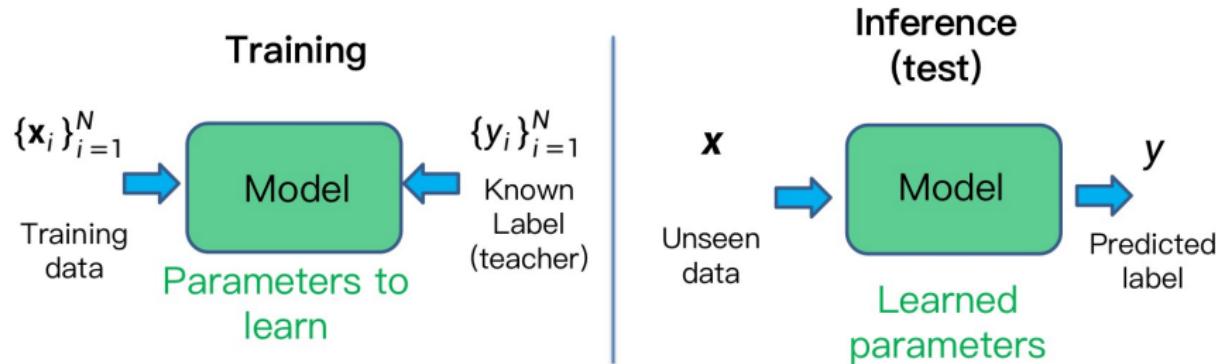
# Supervised learning

In **supervised learning**, the **dataset** is the collection of **labeled examples**, denoted as  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ ,

$$\mathbf{x}_i = [x_i^{(1)}, \dots, x_i^{(j)}, \dots, x_i^{(D)}]^\top, \quad i = 1, \dots, N$$

- Each element  $\mathbf{x}_i$  is called a **feature vector**: it is a vector in which each dimension  $j = 1, \dots, D$  contains a value that describes the example somehow.
- The **label**  $y_i$  can be either an element belonging to a **finite set of classes**  $\{1, 2, \dots, C\}$ , or a **real number**.

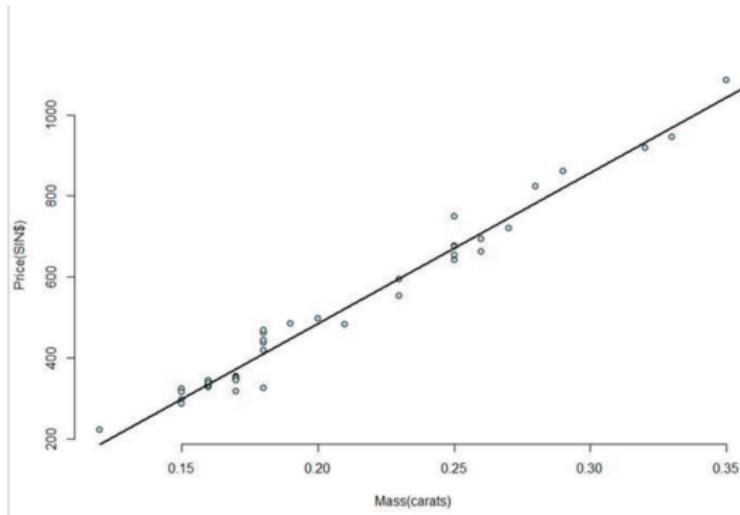
# How supervised learning works



General procedure of supervised learning:

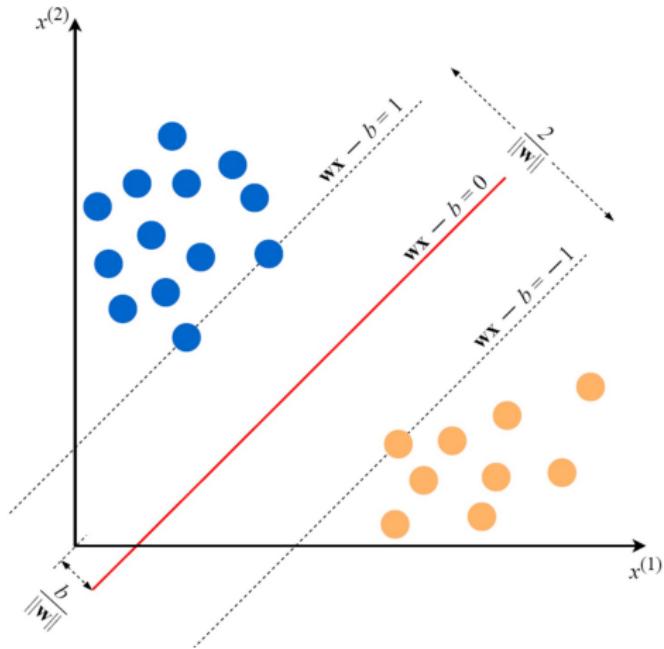
- **Data collection:**  $\{(x_i, y_i)\}_{i=1}^N$
- **Training:** conducting model training on the training data to learn the model parameters
- **Inference (test):** Using the trained model to predict the output of unseen data  $x$

# Regression example



- **Task:** To predict the price of diamond as a function of mass
- **Performance:** Accuracy of predicted values
- **Experience:** Historical data

# Classification example



- **Task:** To classify the input data into two categories
- **Performance:** Classification accuracy
- **Experience:** Historical data

# Regression vs. classification

Suppose that you are running a company, and you want to develop learning algorithms to address the following two problems.

**Problem 1:** You have a large inventory of identical items. You want to predict how many of these items will sell over the next 3 months.

**Problem 2:** You'd like a program to examine individual customer accounts, and for each account, decide if it has been hacked or compromised.

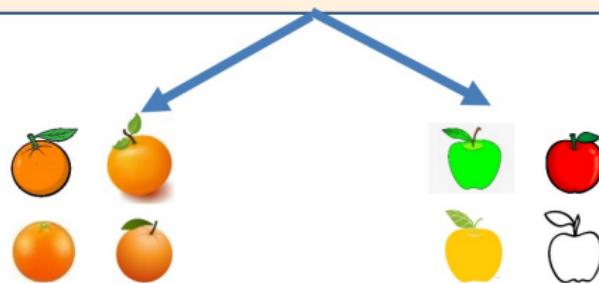
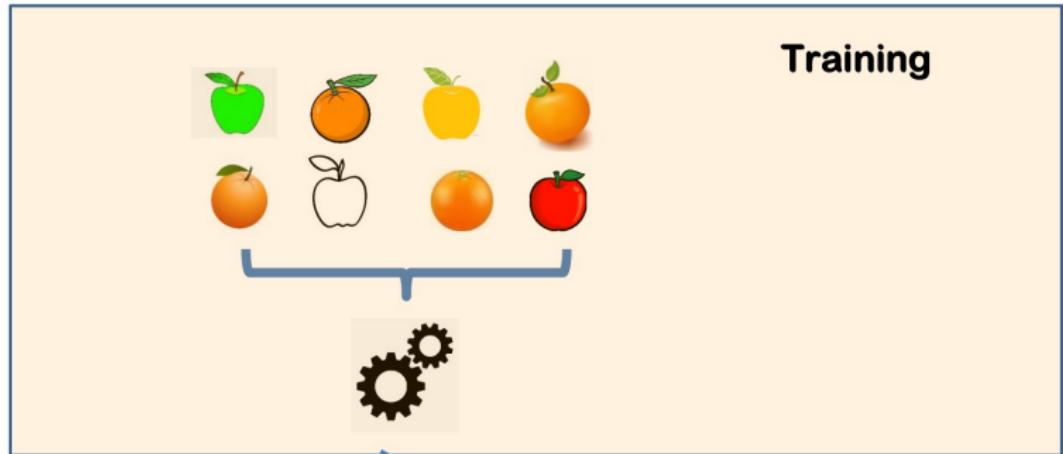
Should you treat these as classification or as regression problems?

- Treat both as classification problems.
- Treat problem 1 as classification, problem 2 as regression.
- Treat problem 1 as regression, problem 2 as classification.
- Treat both as regression problems.

Regression with **continuous output** vs. Classification with **finite and discrete or categorical outputs**

- 1 About this course
- 2 What is machine learning
- 3 Supervised learning
- 4 Unsupervised learning
- 5 Some basic concepts in machine learning

# Unsupervised learning

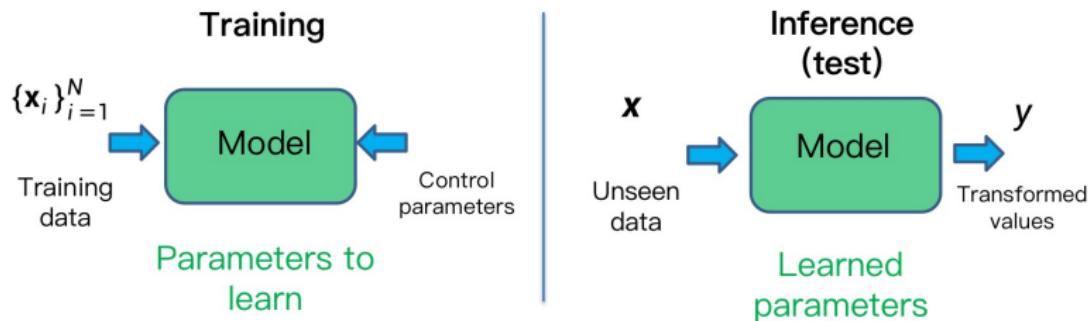


I found two  
types of fruits!

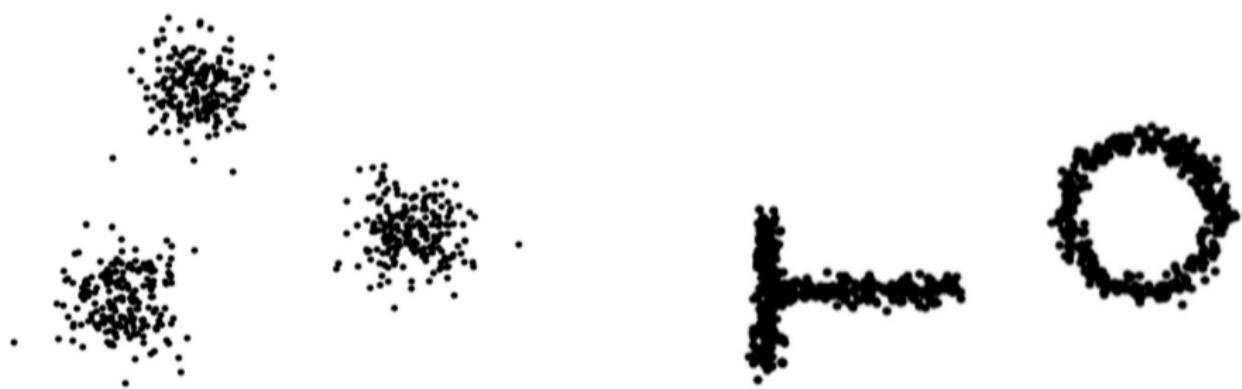
# Unsupervised learning

- In **unsupervised learning**, the **dataset** is a collection of **unlabeled examples**, *i.e.*,  $\{\mathbf{x}_i\}_{i=1}^N$ .
- Again,  $\mathbf{x}$  is a feature vector, and the goal of **an unsupervised learning algorithm** is to create a **model** that takes a feature vector  $\mathbf{x}$  as input and either **transforms it into another vector or into a value** that can be used to solve a practical problem.
- Its **main task** is to **analyze the structure of data** for future inference.

# How unsupervised learning works



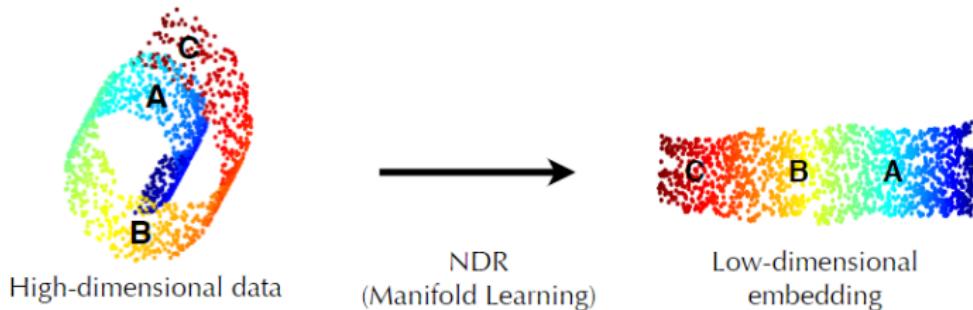
# Clustering example



- **Task:** to partition a set of unlabeled points to clusters.
- **Performance (for test data):**
  - points within the same cluster are **close** to each other
  - points from different clusters are **far** from each other
  - the clusters have an appropriate coverage of all data
- **Experience:** available data

The key question is that how to define and measure **close** or **far**, depending on what features are adopted (*e.g.*, global Euclidean distance, local distance, shape ...)

# Dimensionality reduction example



The purposes of dimensionality reduction:

- **Data simplification:** non-linear → linear
- **Data visualization:** high dimensional → low-dimensional
- **Reduce noise:** some dimensions of the input data may be noises
- **Variable selection for prediction:** learn a sparse model, if there are redundancies among different dimensions

- 1 About this course
- 2 What is machine learning
- 3 Supervised learning
- 4 Unsupervised learning
- 5 Some basic concepts in machine learning

# Basic concepts in supervised learning

- **Data for supervised learning:**

- **Training set:**  $D_{tr} = \{(\mathbf{x}_i, y_i)\}_i^n$ , with  $\mathbf{x}_i \in \mathcal{X}$  being the feature presentation (could be scalar, image/video, text sequence, graph, cloud point, *etc.*),  $y_i \in \mathcal{Y}$  being the supervision/ground-truth value (could be discrete label, continuous value, vector, sequence, tensor, *etc.*).
- **Testing set:**  $D_{test} = \{(\mathbf{x}_i, y_i)\}_i^m$  is used to evaluate the performance of the trained model.

- **Data for unsupervised learning**

- **Training set:**  $D_{tr} = \{\mathbf{x}_i\}_i^n$ , with  $\mathbf{x}_i \in \mathcal{X}$  being the feature presentation (could be scalar, image/video, text sequence, graph, cloud point, *etc.*).
- **Testing set:**  $D_{test} = \{\mathbf{x}_i\}_i^m$  is used to evaluate the performance of the trained model.
- **Independent and identically distributed (i.i.d.) assumption:** In standard machine learning, it assumes that the all samples are observations/realizations of **independent and identical random variables**, and training and testing sets follow **the same distribution**.

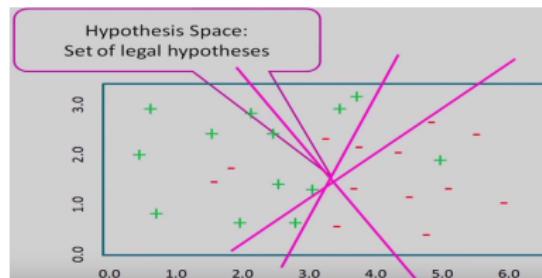
# Basic concepts in supervised learning

- **Target function**  $t : \mathcal{X} \rightarrow \mathcal{Y}$ : the ground-truth mapping function from the input to the output behind the training/testing data. It is unknown, and our goal is to find it.
- **Hypothesis**  $h$ : A hypothesis is a candidate function that describes the unknown target function, for example

$$h(\mathbf{x}) = 1 \times x_1 + 2 \times x_2 = [1, 2] \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

- **Hypothesis space**  $\mathcal{H}$ : Hypothesis space is the set of all the possible legal hypothesis, for example

$$\mathcal{H}_{\mathbf{w}}(\mathbf{x}) = w_1 \times x_1 + w_2 \times x_2 = \mathbf{w}^\top \mathbf{x}$$



# Basic concepts in supervised learning

- **Cost function** is a measure of how good/bad the hypothesis is in terms of its ability to estimate the relationship between  $\mathbf{x}$  and  $y$ , such as **the square loss**  $(h(\mathbf{x}) - y)^2$ .
- **Objective function** is the function that we want to optimize (minimize, maximize or minimax). When we are **minimizing** it, we may also call it the **cost function**, loss function, or error function. In this course, we use these terms interchangeably.
- **Training/learning** is the process of searching a good hypothesis  $h$  in the hypothesis space  $\mathcal{H}$ , through optimizing the objective function on a training set  $D_{tr}$  using the optimization method, such as

$$h^* = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{(\mathbf{x}_i, y_i) \in D_{tr}} (h(\mathbf{x}_i) - y_i)^2$$

- **Testing/evaluation:** evaluating the performance of the learned model  $h^*$  on a test set  $D_{te}$ , such as

$$\frac{1}{m} \sum_{(\mathbf{x}_i, y_i) \in D_{te}} (h^*(\mathbf{x}_i) - y_i)^2$$

# A general machine learning workflow

- ① **Collecting data:** Be it the raw data from excel, access, text files etc., this step (gathering past data) forms the foundation of the future learning. The better the variety, density and volume of relevant data, better the learning prospects for the machine becomes.
- ② **Preprocessing data:** One needs to spend time determining the quality of data and then taking steps for fixing issues such as missing data and treatment of outliers.
- ③ **Determining the hypothesis space, objective function, optimization method.**
- ④ **Training:** learning the parameters of the hypothesis function through optimizing the objective function.
- ⑤ **Testing:** evaluating the learned model on testing data.
- ⑥ **Improving the performance:** This step might involve choosing a different model altogether or introducing more variables to augment the efficiency. That's why significant amount of time needs to be spent in data collection and preparation.

# Review of this week

- **Machine learning definition:** “A computer program is said to learn from **experience E** with respect to some class of **tasks T** and **performance measure P**, if its performance at tasks in T, as measured by P, improves with experience E.”
- **Three main paradigms of machine learning:**
  - **Supervised learning:** learning from teacher
  - **Unsupervised learning:** learning by oneself
  - Reinforcement learning: learning from rewards/punishments
  - Many other learning paradigms ...
- **Supervised learning:** regression and classification.
- **Unsupervised learning:** clustering and dimensionality reduction.
- **Some basic concepts**

# Further reading

## Other learning paradigms:

- **Reinforcement learning:** an area of machine learning concerned with how intelligent agents ought to take actions in an environment in order to maximize the notion of cumulative reward. ([Link](#))
- **Semi-supervised learning:** labeled data + unlabeled data ([Link](#))
- **Ensemble learning:** learning with multiple ML models, and the final prediction is obtained by combining these models. **Improving the performance of individual models** 三个臭皮匠, 顶个诸葛亮 ([Link](#))
- **Transfer learning:** source domain data + target domain data, useful especially when the target data is insufficient ([Link](#))
- **Federated learning:** learning the model at local servers using local data, then uploading locally updated parameters to the central server to obtain the unified parameters. **Protecting users' privacy** ([Link](#))
- **Machine unlearning:** erasing the effect of some particular training samples from a trained model. **Protecting users' privacy** ([Link](#))