

DDA5001 Machine Learning

Introduction

Xiao Li

School of Data Science
The Chinese University of Hong Kong, Shenzhen



About the Course

Machine Learning

Machine Learning

To extract important **patterns**, and use it to **make predictions or decisions**.

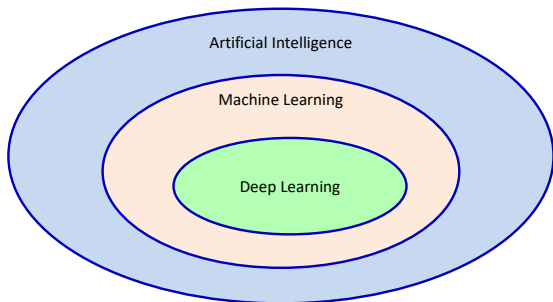
Machine learning is useful in the cases where:

- ▶ **Automate** something a human can do: Face recognition, image classification, autonomous driving, etc.
- ▶ Do things a human **cannot** do, e.g., processing dataset that is large-scale.

Linear algebra, **statistics**, **optimization**, etc are the foundations to machine learning.

Deep Learning, Machine Learning, and AI

- ▶ Nowadays, any field related to images, languages, statistics, algorithms, machine learning, etc is called **Artificial Intelligence (AI)**.
- ▶ AI is to mimic the cognitive capability of humans to reason and to think.

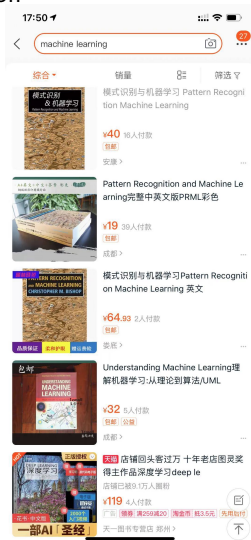


Most of what is labeled AI today, particularly in the public sphere, is highly related to machine learning.

Motivations: Machine Learning is Ubiquitous

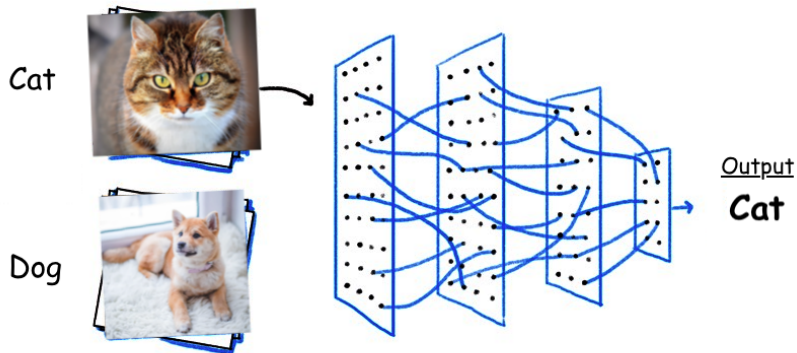
Ubiquity of Machine Learning

► Product recommendation



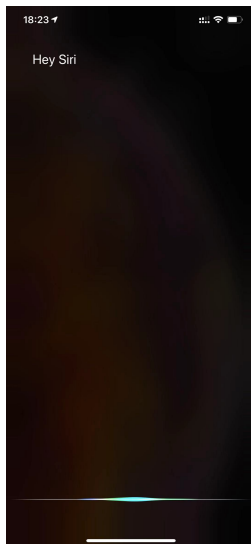
Ubiquity of Machine Learning

► Image classification



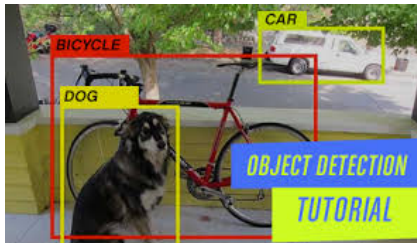
Ubiquity of Machine Learning

► Speech recognition

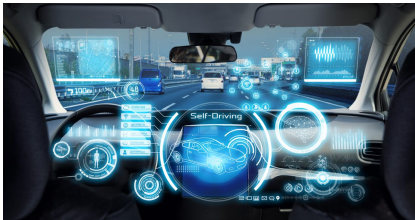


Ubiquity of Machine Learning

► Object detection

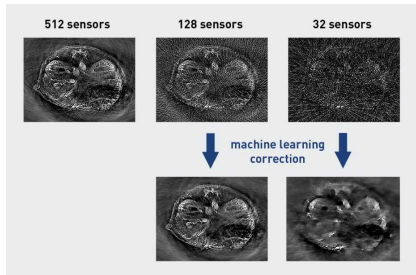


► Autonomous driving



Ubiquity of Machine Learning

► Medical imaging

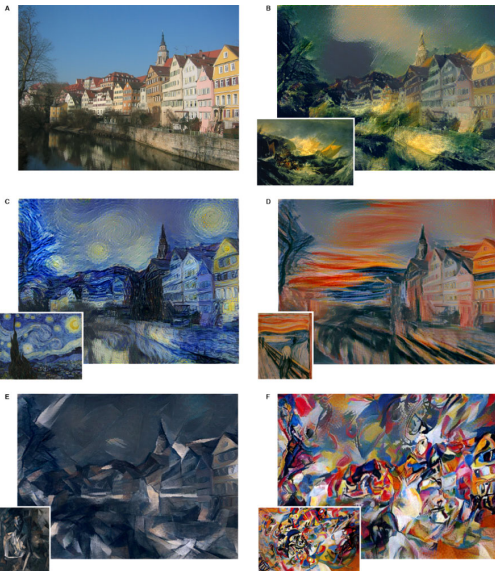


► Health care



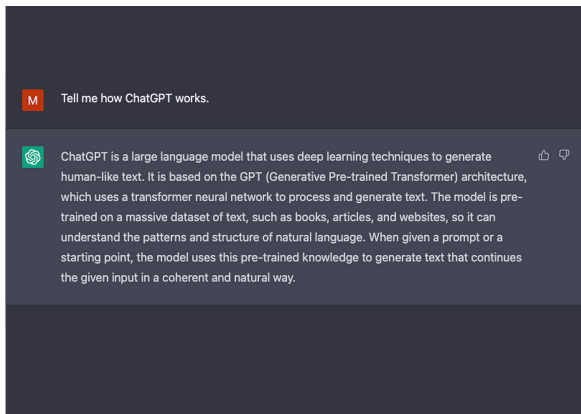
Ubiquity of Machine Learning

- 'Artificial artists'— AIGC (generative AI)



Ubiquity of Machine Learning

► ChatGPT (Large language models)



▪

Course Information

Lecture and Instructors

Lectures

- ▶ Lecture: **TuTh 4:00pm - 5:20pm, TxA402.**

Instructor: Prof. Xiao Li

Email: lixiao@cuhk.edu.cn

Office: Daoyuan 506a.

Office hour: By appointment.

Home page: xiao-li.org

Research interests: Optimization and LLMs.

Course Content

What this course will be?

- ▶ **Methodologies:** The fundamentals and methodologies of machine learning, including MLE, least squares, theory of learning, overfitting, logistic regression, gradient-based algorithms, SVMs, PCA, neural networks, deep learning, transformers, etc.
- ▶ **Applications:** Design and apply machine learning methods (including model and algorithms) to solve practical problems, e.g., language models.

After studying this course, the students are expected to be able to apply machine learning algorithms in practice and be well prepared in conducting advanced machine learning research.

Course Outline and Plan

The main line of this course covers two main types of learning:

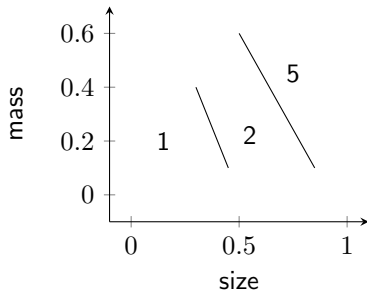
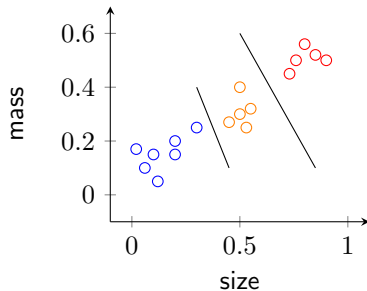
- ▶ **Supervised Learning** (main)
 - ▶ Training versus testing
 - ▶ Linear models
 - ▶ Nonlinear models
 - ▶ Learning algorithms
 - ▶ Advanced topics such as Adam, transformers, and language models
- ▶ **Unsupervised learning**
 - ▶ Dimensionality reduction
 - ▶ Clustering

We have prepared the “**DDA5001_25Fall_plan**” excel file. You can find it on the “Information” page on Blackboard for more detailed information.

Supervised Learning

Example from vending machine — coin recognition:

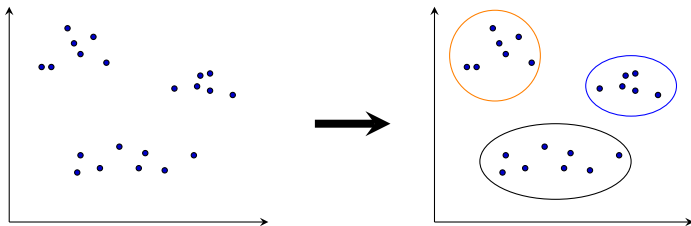
input data + label



Unsupervised Learning

Instead of input + label, we only have:

input data



Prerequisites

Machine learning lies in the intersection of several disciplines and it is a fundamental course.

To study machine learning, we need:

- ▶ Probability and statistics;
- ▶ Multi-variate calculus (mainly calculation of the gradient);
- ▶ Linear algebra (vector, matrix, norm, subspaces);
- ▶ Optimization (algorithms).

We will go through basic mathematics next lecture. These math notions are important. For instance, linear algebra is the basic “language” of this course, we have to represent all the data and operations in vectors and matrices.

Tutors

We have 3 tutors in this semester:

- ▶ Chaolong Ying
Email: 222043012@link.cuhk.edu.cn
- ▶ Jie Mao
121090413@link.cuhk.edu.cn
- ▶ Mengqi Li
223040115@link.cuhk.edu.cn

Their office hours are by appointment.

Tutorials

We will have about **6 tutorials**. Note that these tutorials **do not** take place every week. The tentative plan for tutorials is:

- ▶ Tutorial 1 is in the 2nd teaching Week (Sep 9). Main contents: Install Python and introduce useful grammar for machine learning.
- ▶ Tutorial 2 is in the 3rd teaching Week (Sep 16). Main contents: Implement perceptron, including how to load data, implement algorithm, and plot figures.
- ▶ Tutorial 3 is in the 6th teaching week (Oct 14). Main contents: Implement gradient-based learning algorithms for logistic regression.
- ▶ Tutorials 4-6 are about our course project. They will take place in the 9th teaching week (Nov 4), 11th teaching week (Nov 18), and 13th teaching week (Dec 2), respectively.

The tutorial venue and time are:

- ▶ **6:30pm - 7:30pm, online**, the zoom link will be provided.

Grades

The evaluation of this course contains three parts:

- ▶ Homework (30%)
- ▶ Course project (30%)
- ▶ Final exam (40%)

Homework (30%)

We will have about 3 homework.

- ▶ In principle, the homework is not very easy. You will have two weeks for completing each assignment.
- ▶ Through completing the homework, you are supposed to learn how to apply machine learning methods to practical problems. It includes programming (Python) and algorithm implementation (more than call software package).
- ▶ Tutorials 1-3 are designed to help programming in homework.
- ▶ Late submission will not be graded.

I encourage discussion. However, all the homework must be written by yourself. You should understand your answer. Two very similar homework submissions will result in plagiarism.

Course Project (30%) I

We will have a course project.

Format:

- ▶ It is about Large Language Models (LLMs).
- ▶ The project is cut into 3 consecutive parts (part I, part II, and part III). For each part, you have 2 weeks to complete. Thus, the project will last 6 weeks (1.5 months) in total.
- ▶ The project will be released around the 8th week (Nov 2).
- ▶ Group work (6-8 students). The first part is individual work, while the remaining two parts are group work. For group work, each member in a group should clearly state which part you are taking in charge. Find a suitable group by yourself in this semester (the best way we figured out in the past). If you have difficulty to find group members, please email TA for help.

Course Project (30%) II

Submission materials: Assessment will be conducted according to the following materials:

- ▶ A PDF report for part I (up to one-page), a PDF report for part II (up to two-page), and a PDF report for part III (up to two-page). These page limits do not contain references pages.
- ▶ Code.

Project grade:

- ▶ The project contains 100 points. The distribution is: part I: 30, part II: 40, part III: 30.
- ▶ Part III will have an exploration direction. If you complete the exploration direction, you can earn another 10 bonus points. Hence, one can earn at most 110 points from the course project.

Project tutorial: We will have about 3 tutorials (Tutorials 4-6) to explain large language models and provide guidelines on the project. The specific time of these tutorials will be announced later, roughly around the 9-12th weeks.

Final Exam (40%)

- ▶ Final exam covers all parts of this course.
- ▶ It exams basic knowledge of machine learning.
- ▶ We will provide sample final exam for review.
- ▶ Time: 3:00pm - 5:00pm, December 13 (Saturday), 2025.
- ▶ Venue: TBD

Course Materials

- **Main material:** Lecture slides.
- **Recommended Reading books (not required):**
 - ▶ Abu-Mostafa, Y. S., Magdon-Ismail, M., & Lin, H. T. (2012). Learning from data. New York, NY, USA: AMLBook.
 - ▶ Bishop C., & Bishop H. (2024). Deep Learning: Foundations and Concepts. Springer.
- A useful website: <https://scikit-learn.org/stable/>.

Course Syllabus and Policies

We have prepared the “**DDA5001_25Fall_Syllabus**” PDF file. You can find it on the “Information” page on Blackboard. It contains:

- ▶ Recommended reference books.
- ▶ Course summary.
- ▶ Homework information.
- ▶ Tutor and tutorial information.
- ▶ Project information.
- ▶ Grade distribution and final exam time.
- ▶ Course policies.

The first six parts are included in this lecture slides, please refer to the course syllabus for the **course policies**, e.g., policies on exam, homework submission, academic honesty, etc.

How to Study this Course

Lots of students worry about if they can pass this course. It should not be hard by doing the following things:

- ▶ **Math:** Try to understand basic math knowledge in the lecture.
- ▶ **Tutorials:** Try to attend the tutorial. It will include knowledge about programming in homework and our course project.
- ▶ **Homework:** Discuss with classmates (but remember to write down everything using your own understanding), and **submit every homework**.
- ▶ **Project:** Start early, explore more possible ideas, and spend some time on writing the report.
- ▶ **Exam:** Understand most of the basic concepts included in the lecture, and **be sure to attend the exam** (exam is easier than homework).

⇒ Next lecture: Basic math and concepts of machine learning.