

DSAA2011 Machine Learning

L1 Logistics and Introduction

Dr. Zixin Zhong

Data Science and Analytics Thrust
Information Hub
Hong Kong University of Sceience and Technology (Guangzhou)

February 7, 2025

Outline

1 Logistics

2 Introduction to Machine Learning

- Definition and taxonomy

3 Mathematical tools

- Set and function
- Probability and estimation
- Systems of linear equations



Outline

1 Logistics

2 Introduction to Machine Learning

- Definition and taxonomy

3 Mathematical tools

- Set and function
- Probability and estimation
- Systems of linear equations



Course information



Weikai Yang
Assistant Professor

Visual analytics +
Machine learning

W1-316

weikaiyang@hkust-gz.
edu.cn



Zixin Zhong
Assistant Professor

Reinforcement learning
Online learning

W1-308

zixinzhong@hkust-gz.
edu.cn



Zixin Zhong (HKUST-GZ)

- **Lecture time:** 12-3PM Fri
- **Lab time:** 8-9PM Fri
- **Venue:** Rm E1-134
- **Course material:**
Slides, recommended materials
Consultation: Office hour, [email](#), appointment or Canvas, Teams

- **Teaching assistants (3/6):**

1. **Weiwen CHEN**
2. **Guanghua LI**
3. **Yang LUO**
4. **Chunming MA**
5. **Jingyi PAN**
6. **Liangwei WANG**

Office hours (weekly, starting from 11 Feb)

Time	Venue	Instructor/TA	Email
7-8PM Tue (7-8PM Wed)	Rm E2-301	Weiwen CHEN Guanghua LI Yang LUO Chunming MA Jingyi PAN Liangwei WANG	wchen948@connect.hkust-gz.edu.cn gli945@connect.hkust-gz.edu.cn yluo208@connect.hkust-gz.edu.cn cma859@connect.hkust-gz.edu.cn jpan305@connect.hkust-gz.edu.cn lwang344@connect.hkust-gz.edu.cn
3-4PM Fri	Rm W1-316 Rm W1-308	Weikai Yang Zixin Zhong	weikaiyang@hkust-gz.edu.cn zixinzhong@hkust-gz.edu.cn

♥ Please show respect and appreciation to our TAs!



Q&A

- Office hour, email, appointment or Canvas, Microsoft Teams (code: pf0dj1h)

DSAA2011-25sp-Q&A

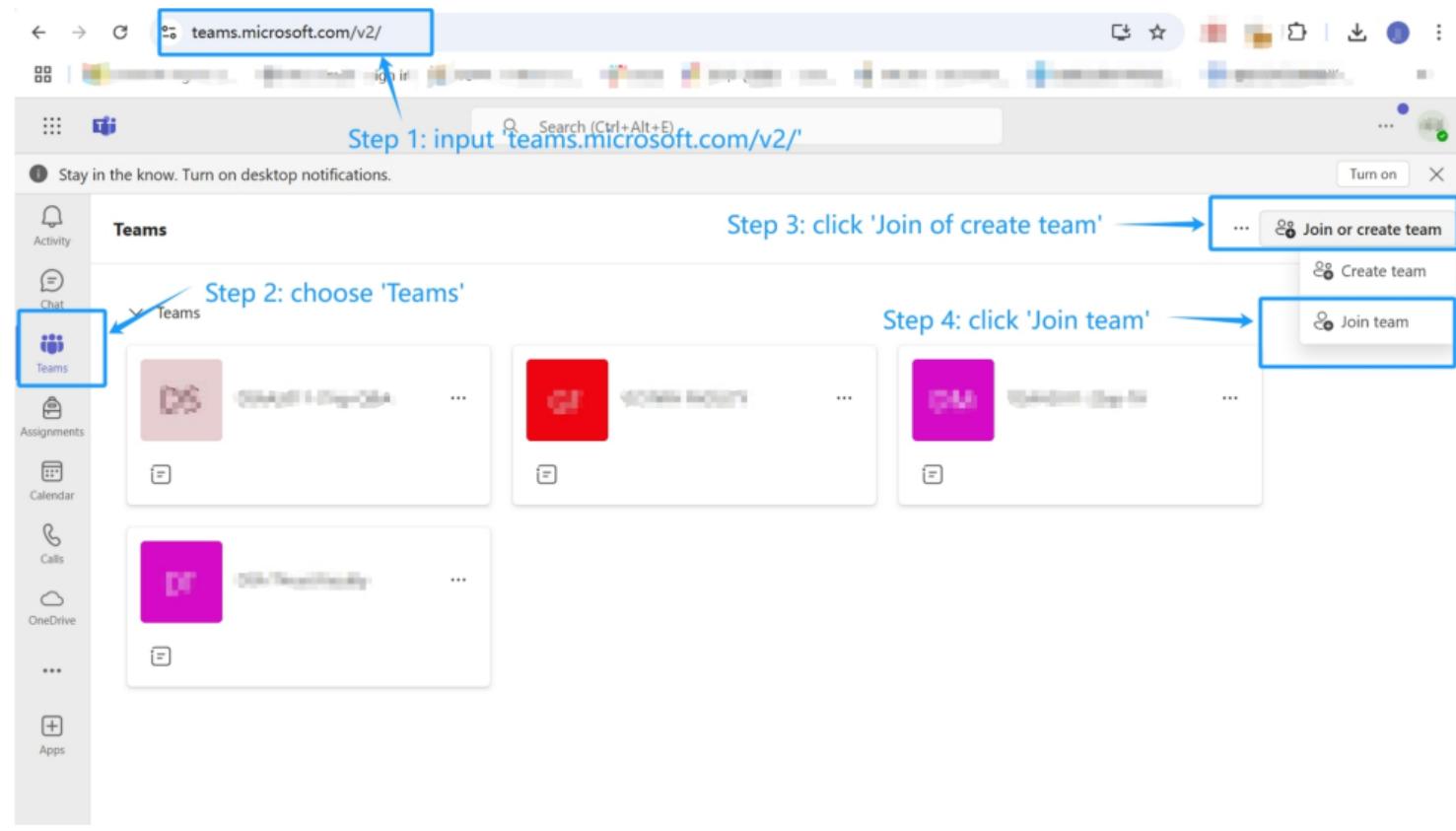
pf0dj1h

Where do I enter the code?

Click [Join or create a team](#) below your teams list and look for the [Join a team with a code](#) card



Q&A via Microsoft Teams: <https://teams.microsoft.com/v2/>



Q&A via Microsoft Teams: join with code pf0dj1h

The screenshot shows the Microsoft Teams interface with the sidebar open. The 'Teams' tab is selected. In the center, there's a search bar with 'Search (Ctrl+Alt+E)' and a 'Turn on' button. On the left, under 'Join a team', there's a search bar with 'Type to search' and a 'Teams for you' section. A blue box highlights the 'Join a team with a code' input field containing 'pf0dj1h' and the 'Add team' button below it. A blue arrow points from the text 'Step 5: type 'pf0dj1h' and click 'Add team'' to the 'Add team' button. The bottom right corner shows a small circular icon with a green checkmark and the number '7'.

Stay in the know. Turn on desktop notifications.

Activity

Chat

Teams

Assignments

Calendar

Calls

OneDrive

...

Apps

Join a team

Type to search

Teams for you

#

Join a team with a code

pf0dj1h

Add team

Step 5: type 'pf0dj1h' and click 'Add team'



Q&A via Microsoft Teams: join with code pf0dj1h

Stay in the know. Turn on desktop notifications.

Turn on X

Activity

Chat

Teams

Assignments

Calendar

Calls

OneDrive

...

Apps

Join or create team

Teams

DSAA2011-25sp-Q&A ...

...

...

...

...

Zixin Zhong (HKUTS-GZ)

February 7, 2025

8 / 81

Q&A via Microsoft Teams: ask questions in the 'Q_A' channel

The screenshot shows the Microsoft Teams application window. On the left is the sidebar with icons for Activity, Chat, Teams, Assignments, Calendar, Cells, OneDrive, and Apps. The main area shows a team named 'DS' with a channel named 'DSAA2011-25sp-Q&A'. Inside this channel, there is a 'General' folder containing a 'Q_A' sub-channel. A blue box highlights the 'Q_A' channel, and a blue arrow points from it to the text 'Step 1: choose the 'Q_A' channel'. Below this, another blue arrow points from the 'Start a post' button to the text 'Step 2: click 'Start a post''. The 'Start a post' button is highlighted with a blue box. At the bottom right of the main area, there is a 'Welcome to the team!' message and a note: 'Try @mentioning the team name or teacher names to begin sharing ideas.' The top right corner of the window has a 'Turn on' button and a close (X) button.

Stay in the know. Turn on desktop notifications.

Activity Chat Teams Assignments Calendar Cells OneDrive Apps

DS < All teams DS Q_A Posts Files Notes +

DSAA2011-25sp-Q&A Main Channels General Q_A

Step 1: choose the 'Q_A' channel

Step 2: click 'Start a post'

Welcome to the team!

Try @mentioning the team name or teacher names to begin sharing ideas.

Start a post



Q&A via Microsoft Teams: ask questions in the 'Q_A' channel

The screenshot shows the Microsoft Teams interface with a blue border around the message input area. On the left is the sidebar with icons for Activity, Chat, Teams, Assignments, Calendar, Calls, OneDrive, and Apps. The main window shows a team named 'DS' with a channel called 'DSAA2011-25sp-Q&A'. The 'Q_A' tab is selected. A welcome message says 'Welcome to the team!' and suggests mentioning team names or teacher names. The message input field has a placeholder 'Add a subject' and a toolbar with various icons. A 'Post' button is at the bottom right.



Grading scheme

- Midterm 20 Points
- Written assignments (Project) 30 Points
- Final exam 50 Points
- More details will be released later

Letter Grade	Raw Mark Range
A	[85, 100]
B	[70, 85)
C	[55, 70)
D	[40, 55)
E	< 40



Syllabus

Week #	Topic	Lecturer
1	Introduction + course Info	Zixin
2-3	Supervised learning: regression and classification	Zixin
4	Model evaluation and choice	Zixin
5	Feature selection	Zixin
6	Boosting methods	Weikai
	Midterm-29 March (Sat): save your day and mark it on calendar!	
7-8	Unsupervised learning: clustering	Weikai
9	Active learning	Weikai
10-11	Markov and graphical models	Weikai
12-13	Online learning	Zixin
14	Final exam	



Programming language policy

- We will do demos in **Python (Jupyter Lab)**.
- You're welcome to use any language you like (that your TAs can read) for homework or project.
- TAs will only support **Python**.



Collaboration policy

Project: yes of course! Clarify the contributions of each member.

- Give credit to the people who have helped you: write on your report the names of the people you worked with.
- Give credit to the other resources that have helped you: please write on your report the GenAI tools, textbooks, notes, or web pages you found useful.

Midterm/Final: Figure out your solutions by yourself.

- All of the text that you submit should be done by you.
- No collaboration/discussion here!



IP policy

Other course note websites:

- do not post any course materials in other websites. this makes the next rendition of the course worse for everyone.
- please report to us any course materials you find online (not on [our websites](#)).



How to do well in this course?

- Pay attention to lecture and lab
- Ask whenever
- Do exercise in lecture notes and lab notes
 - Treat each lab note as a quiz (good practise)
- Find a peer, describe a concept or teach each other
- Practice all you learn in the project



How to do well in this course?

- Pay attention to lecture and lab
- Ask whenever
- Do exercise in lecture notes and lab notes
 - Treat each lab note as a quiz (good practise)
- Find a peer, describe a concept or teach each other
- Practice all you learn in the project
- Be perserverant:

"The brick walls are there for a reason. The brick walls are not there to keep us out. The brick walls are there to give us a chance to show how badly we want something. Because the brick walls are there to stop the people who don't want it badly enough. They're there to stop the other people."

from *The Last Lecture* by Randy Pausch

- Brain power is like muscle power. It can only get better through heavy lifting.
- Everyone can become smarter by training (practicing, analyzing, coding, ...)



Outline

1 Logistics

2 Introduction to Machine Learning

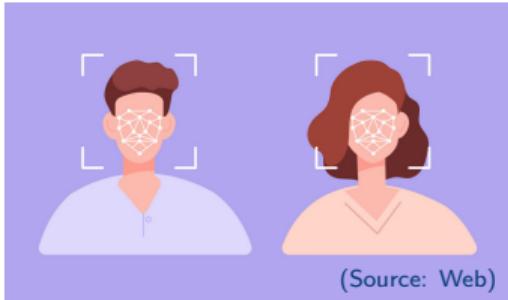
- Definition and taxonomy

3 Mathematical tools

- Set and function
- Probability and estimation
- Systems of linear equations



Machine Learning is everywhere

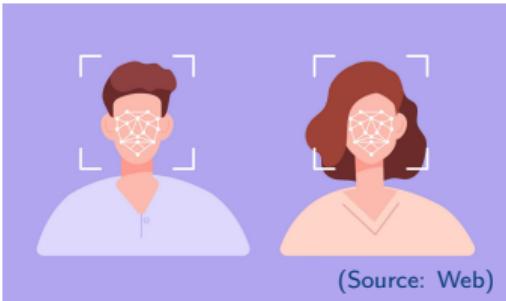


(Source: Web)

Identity verification



Machine Learning is everywhere



(Source: Web)

Identity verification



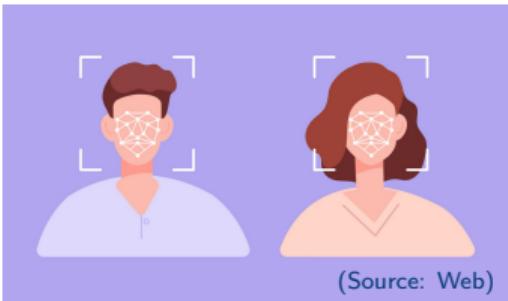
HKUST(GZ) GPT Service System

香港科技大学（广州）GPT服务系统

Large Language Model



Machine Learning is everywhere



(Source: Web)

Identity verification



HKUST(GZ) GPT Service System

香港科技大学（广州）GPT服务系统

Large Language Model

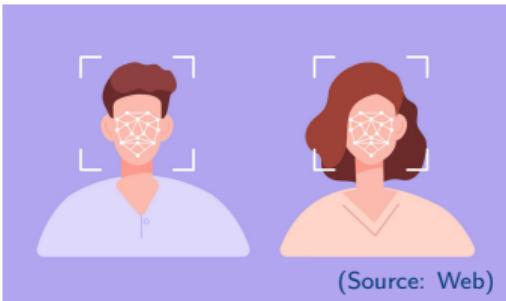


(Source: Web)

Healthcare



Machine Learning is everywhere



(Source: Web)

Identity verification



(Source: Web)

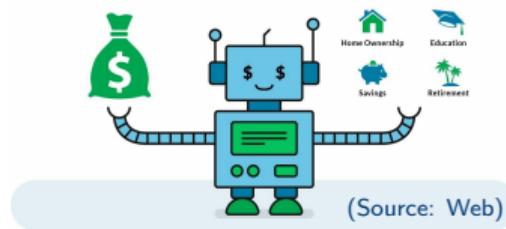
Healthcare



HKUST(GZ) GPT Service System

香港科技大学(广州) GPT服务系统

Large Language Model



(Source: Web)

Wealth management



Machine Learning is everywhere

What is Machine Learning (ML)?



Definition of Machine Learning

Prof. Tom M. Mitchell from CMU (Samuel, 1959):

- 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.



(Source: Web)



Definition of Machine Learning

Prof. Tom M. Mitchell from CMU (Samuel, 1959):

- 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**.



(Source: Web)



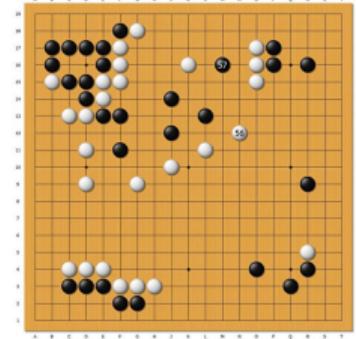
Definition of Machine Learning

Prof. Tom M. Mitchell from CMU (Samuel, 1959):

- 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**.
- **Experience E (data)**: games played by the program or human



(Source: Web)



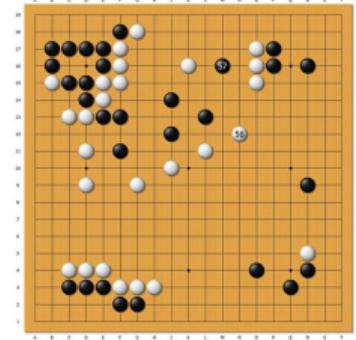
Definition of Machine Learning

Prof. Tom M. Mitchell from CMU (Samuel, 1959):

- 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**.
- **Experience E (data)**: games played by the program or human
- **Performance measure P**: winning rate



(Source: Web)



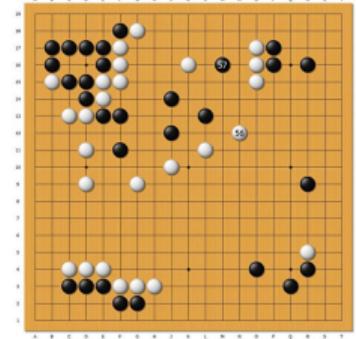
Definition of Machine Learning

Prof. Tom M. Mitchell from CMU (Samuel, 1959):

- 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**.
- **Experience E (data)**: games played by the program or human
- **Performance measure P**: winning rate
- **Task T**: to win

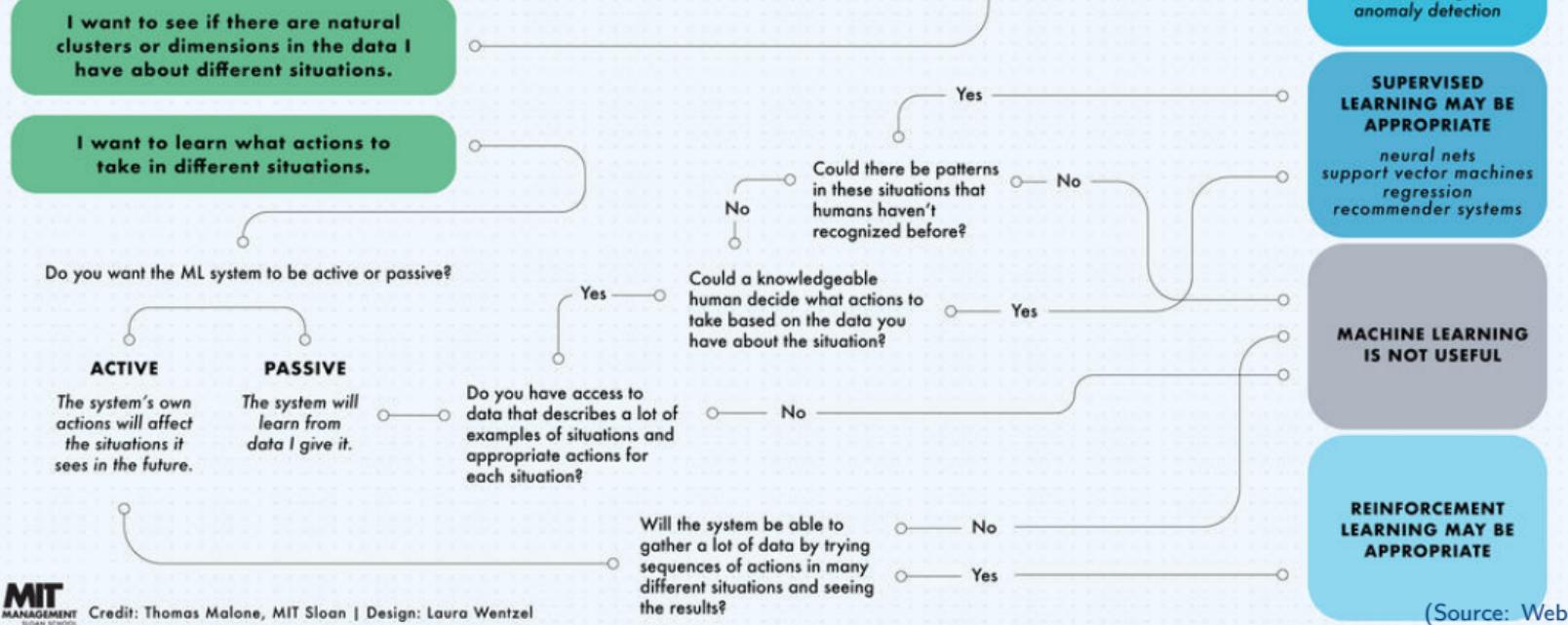


(Source: Web)



Taxonomy of Machine Learning

What do you want the machine learning system to do?



Credit: Thomas Malone, MIT Sloan | Design: Laura Wentzel



Taxonomy of Machine Learning

Supervised Learning

Has outcome information ("labels")

Finds patterns that relate to those outcomes

Uses patterns to predict outcomes not yet known

Unsupervised Learning

No outcome information available

Analyzes or identifies groups without labels or human instruction

Offers insights into characteristics that define groups

Reinforcement Learning

Makes decisions based on trial and error

Decision-making algorithm is constantly refined based on "rewards"

Excels in complex situations



(Source: Web)



Taxonomy of Machine Learning (A Simplistic View)

♠ What type of data?

- **Supervised learning** - labeled data: e.g., Prediction
- **Semi-supervised learning**
- **Unsupervised learning** - unlabeled data: e.g., clustering



Taxonomy of Machine Learning (A Simplistic View)

♠ What type of data?

- **Supervised learning** - labeled data: e.g., Prediction
- **Semi-supervised learning**
- **Unsupervised learning** - unlabeled data: e.g., clustering
- **Reinforcement learning** - environment feedback: e.g., multi-armed bandit



Taxonomy of Machine Learning (A Simplistic View)

♠ What type of data?

- **Supervised learning** - labeled data: e.g., Prediction
 f(x)≈y,
 x is input data;
 y is output data, denoted as 'label';
 → all labelled
- **Semi-supervised learning**
 → partly labelled
- **Unsupervised learning** - unlabeled data: e.g., clustering
 → no labelled
 classification
- **Reinforcement learning** - environment feedback: e.g., multi-armed bandit

♠ When do we collect data?

- **Offline learning**
- **Online learning**



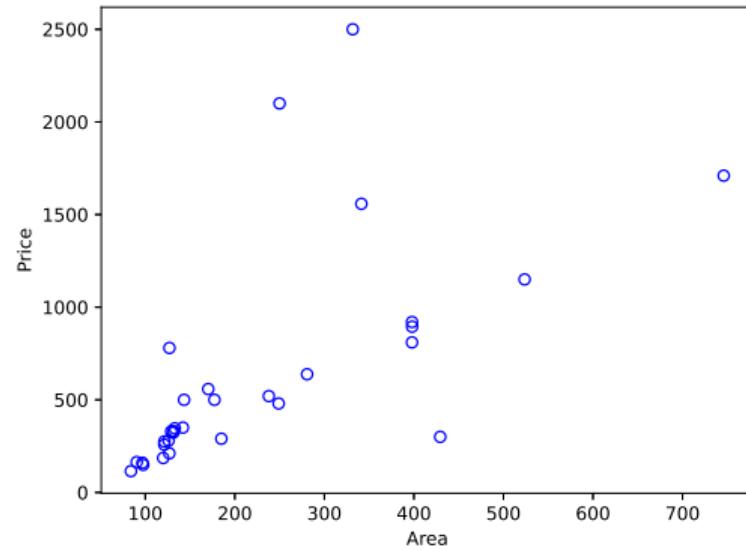
Supervised Learning



Task: housing price prediction

- Given a dataset containing n samples

$$\left(x^{(1)}, y^{(1)} \right), \left(x^{(2)}, y^{(2)} \right), \left(x^{(3)}, y^{(3)} \right), \dots, \left(x^{(n)}, y^{(n)} \right)$$

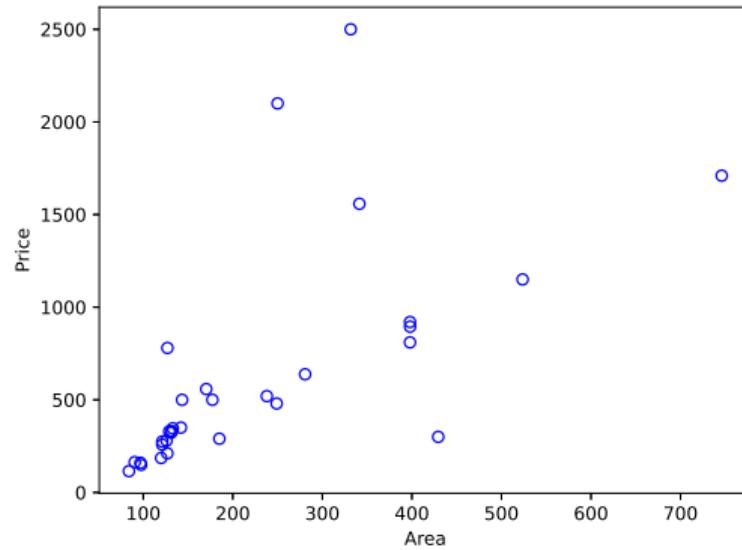


Task: housing price prediction

- Given a dataset containing n samples

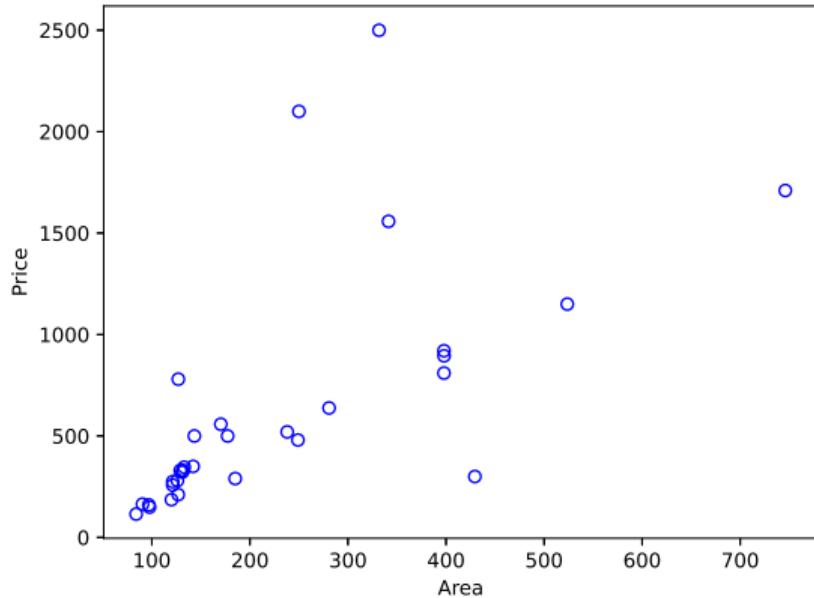
$$\left(x^{(1)}, y^{(1)}\right), \left(x^{(2)}, y^{(2)}\right), \left(x^{(3)}, y^{(3)}\right), \dots, \left(x^{(n)}, y^{(n)}\right)$$

- Task:** if a residence has $x = 320$ square kilometers, can we predict its price?



Task: housing price prediction

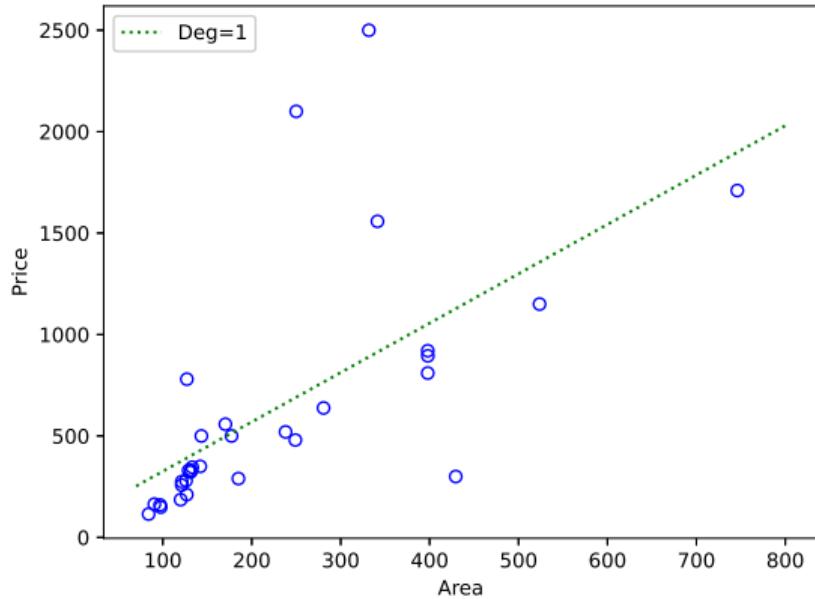
- Linear regression (polynomial regression with degree 1):
Assume $y = p_1(x) = ax + b$ and guess the values of a and b .



Task: housing price prediction

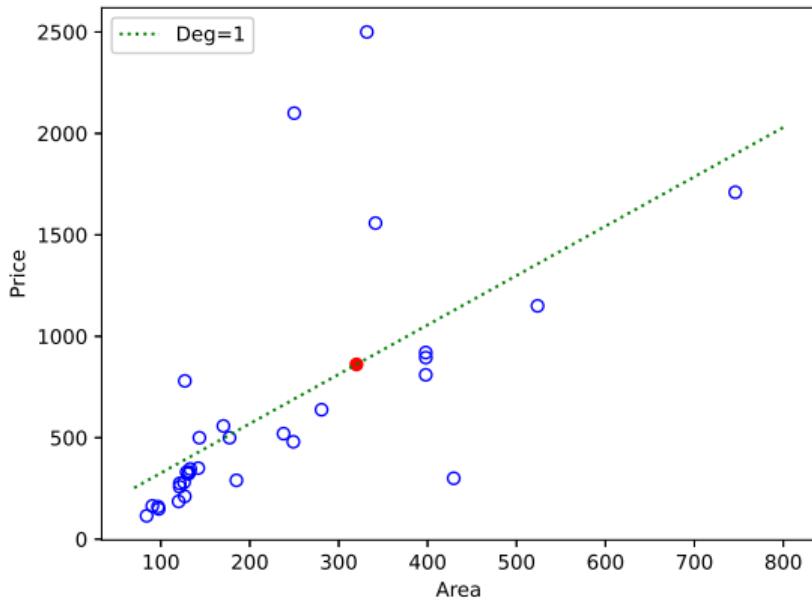
- Linear regression (polynomial regression with degree 1):

Assume $y = p_1(x) = ax + b$ and guess the values of a and b .



Task: housing price prediction

- Linear regression (polynomial regression with degree 1):
Assume $y = p_1(x) = ax + b$ and guess the values of a and b .



$$p_1(320) = 861.29$$

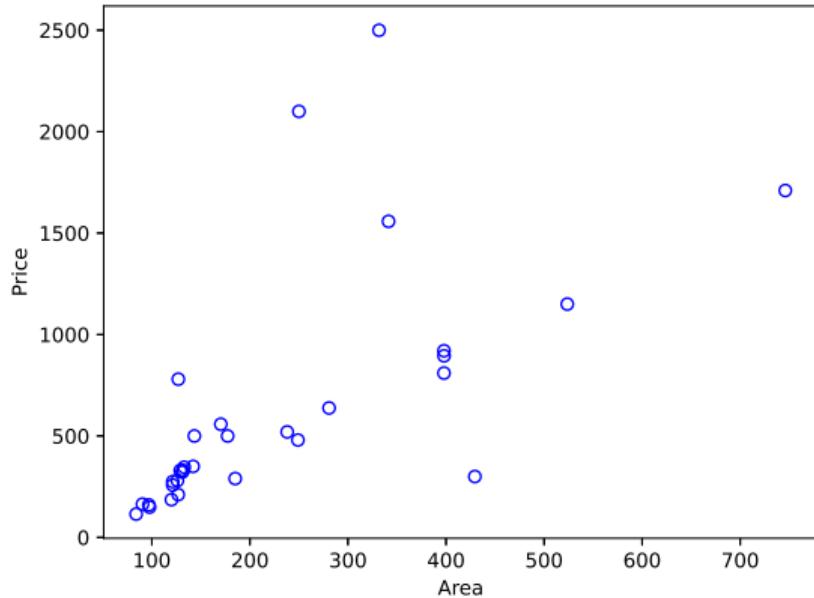
- If a residence has $x = 320$ square kilometers, we predict its price as $y = 861.29$ CNY.



Task: housing price prediction

- Polynomial regression with degree 2:

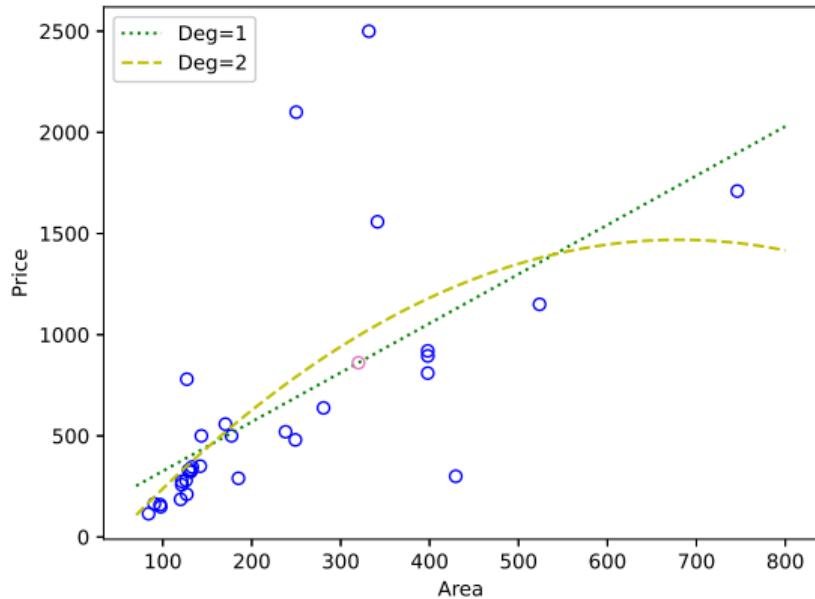
Assume $y = p_2(x) = ax^2 + bx + c$ and guess the values of a , b and c .



Task: housing price prediction

- Polynomial regression with degree 2:

Assume $y = p_2(x) = ax^2 + bx + c$ and guess the values of a , b and c .



$$a = -0.0037, b = 4.97, c = -221.8$$

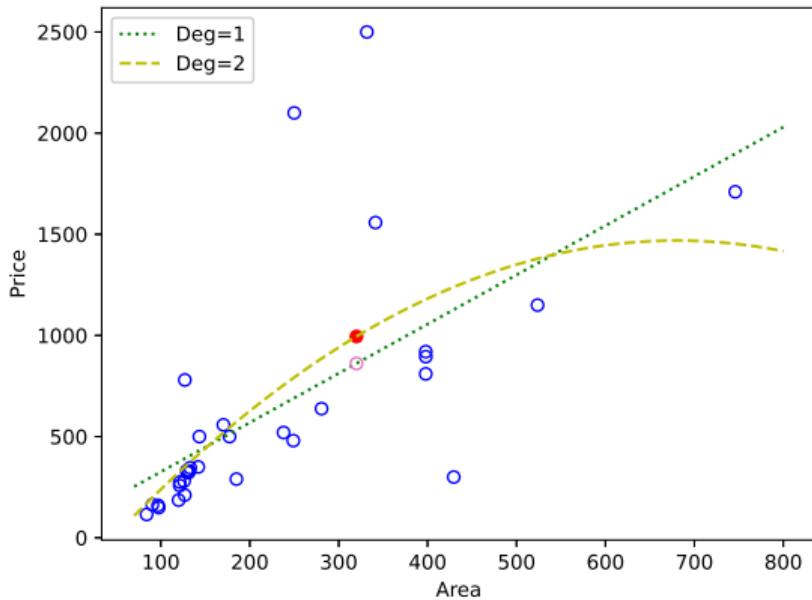
$$p_2(x) = -0.0037x^2 + 4.97x - 221.8$$



Task: housing price prediction

- Polynomial regression with degree 2:

Assume $y = p_2(x) = ax^2 + bx + c$ and guess the values of a , b and c .



$$a = -0.0037, b = 4.97, c = -221.8$$
$$p_2(x) = -0.0037x^2 + 4.97x - 221.8$$

$$p_2(320) = 994.94$$

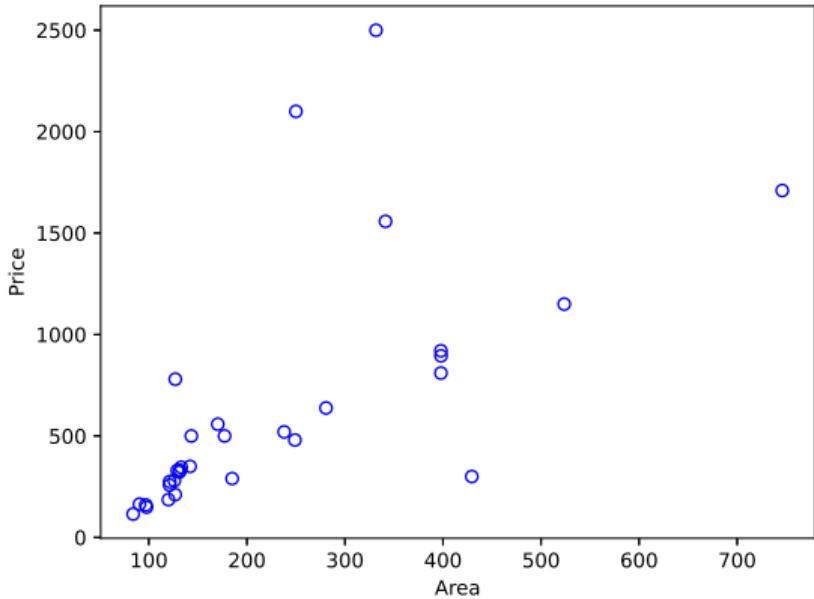
- If a residence has $x = 320$ square kilometers, we predict its price as $y = 994.94$ CNY.



Task: housing price prediction

- Polynomial regression with degree 3:

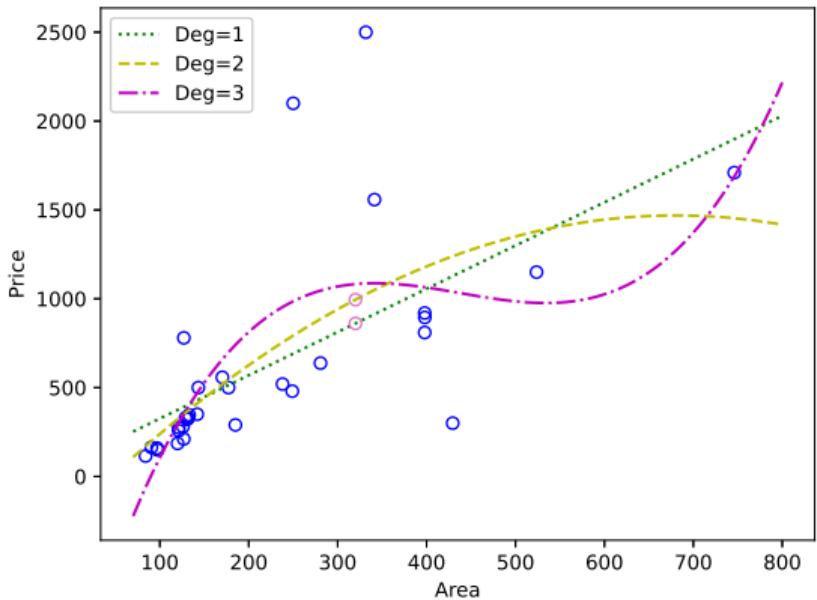
Assume $y = p_3(x) = ax^3 + bx^2 + cx + d$ and guess the values of a, b, c and d .



Task: housing price prediction

- Polynomial regression with degree 3:

Assume $y = p_3(x) = ax^3 + bx^2 + cx + d$ and guess the values of a, b, c and d .



$$a = 3.16 * 10^{-5}, b = -0.042,$$

$$c = 17.33, d = -1243$$

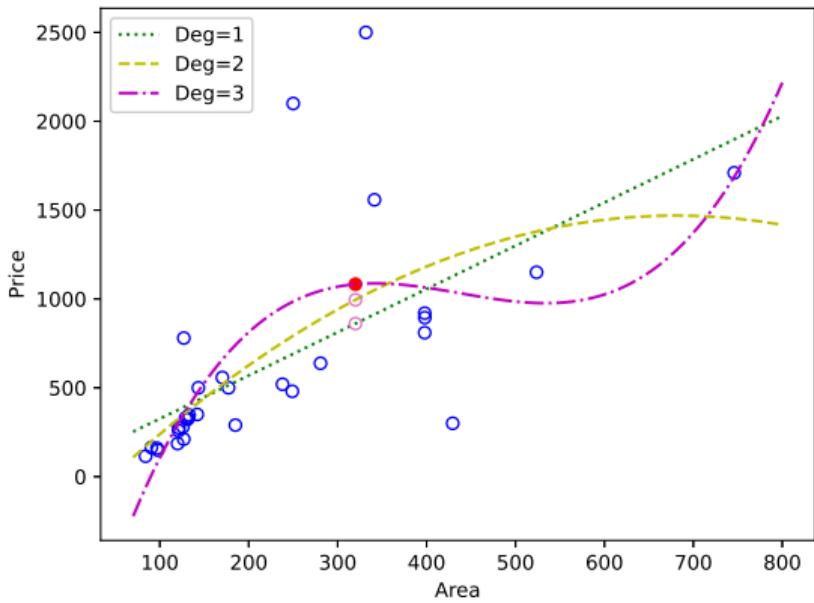
$$p_3(x) = 3.16 * 10^{-5}x - 0.042x + 17.33x - 1243$$



Task: housing price prediction

- Polynomial regression with degree 3:

Assume $y = p_3(x) = ax^3 + bx^2 + cx + d$ and guess the values of a, b, c and d .



$$a = 3.16 * 10^{-5}, b = -0.042,$$

$$c = 17.33, d = -1243$$

$$p_3(x) = 3.16 * 10^{-5}x - 0.042x + 17.33x - 1243$$

$$p_3(320) = 1082.80$$

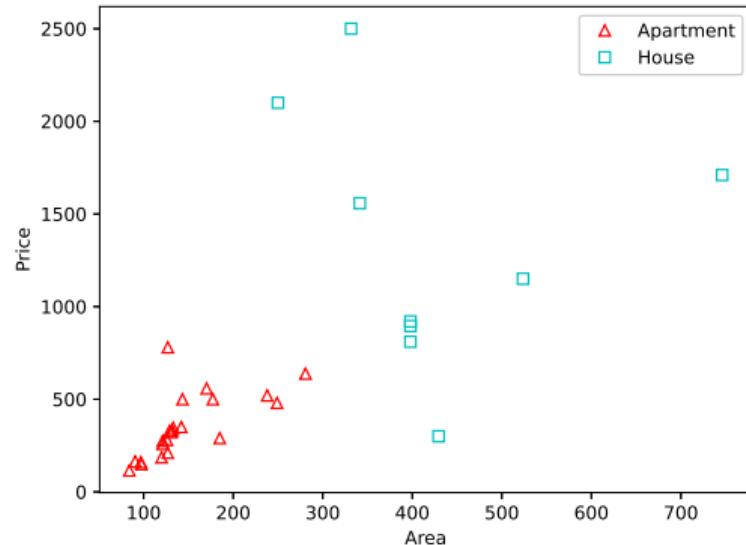
- If a residence has $x = 320$ square kilometers, we predict its price as $y = 1082.80$ CNY.



Task: housing type prediction

- Given a dataset containing n samples ($z \in \{0, 1\}$)

$$\left(x^{(1)}, y^{(1)}, z^{(1)} \right), \left(x^{(2)}, y^{(2)}, z^{(2)} \right), \left(x^{(3)}, y^{(3)}, z^{(3)} \right), \dots, \left(x^{(n)}, y^{(n)}, z^{(n)} \right)$$

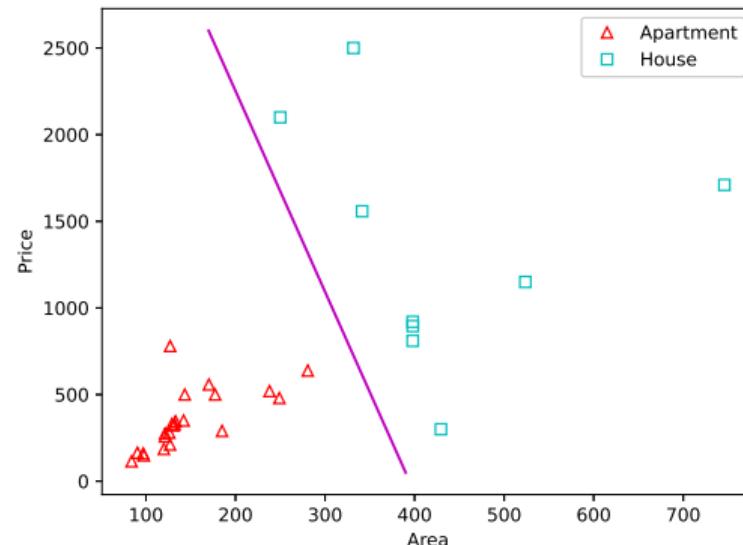


Task: housing type prediction

- Given a dataset containing n samples ($z \in \{0, 1\}$)

$$\left(x^{(1)}, y^{(1)}, z^{(1)}\right), \left(x^{(2)}, y^{(2)}, z^{(2)}\right), \left(x^{(3)}, y^{(3)}, z^{(3)}\right), \dots, \left(x^{(n)}, y^{(n)}, z^{(n)}\right)$$

- Task:** if a residence has $x = 320$ square kilometers and worth $y = 1080.20$ CNY, can we predict its type z ?

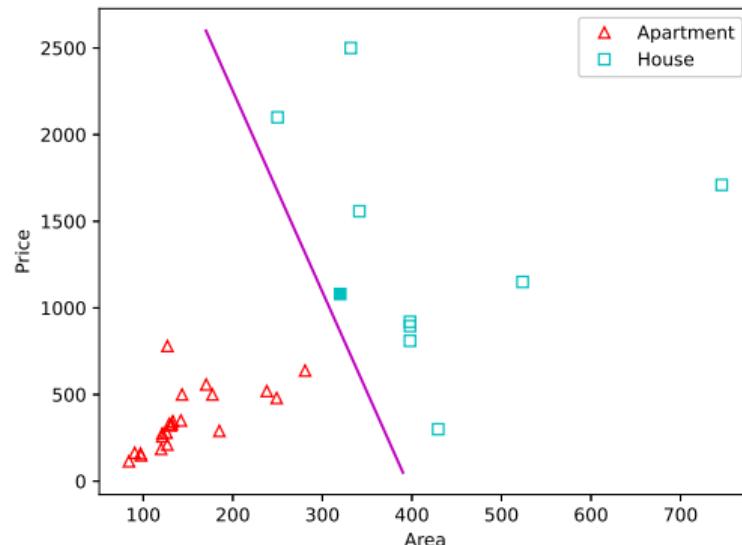


Task: housing type prediction

- Given a dataset containing n samples ($z \in \{0, 1\}$)

$$\left(x^{(1)}, y^{(1)}, z^{(1)}\right), \left(x^{(2)}, y^{(2)}, z^{(2)}\right), \left(x^{(3)}, y^{(3)}, z^{(3)}\right), \dots, \left(x^{(n)}, y^{(n)}, z^{(n)}\right)$$

- Task:** if a residence has $x = 320$ square kilometers and worth $y = 1080.20$ CNY, can we predict its type z ?



Task comparison: regression vs classification

- **Regression:** the prediction result is a **continuous** variable
 - ▶ e.g., price prediction
 $(x = \text{Area}) \rightarrow (y = \text{Price})?$
- **Classification:** the prediction result is a **discrete** variable
 - ▶ e.g., type prediction
 $(x = \text{Area}, y = \text{Price}) \rightarrow (z = \text{Type})?$



Unsupervised Learning

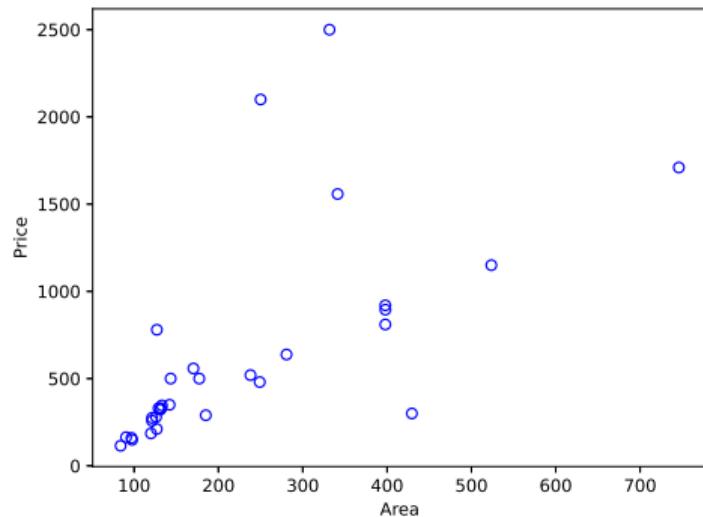


Task: housing type prediction

- Given a dataset containing n samples

$$\left(x^{(1)}, y^{(1)} \right), \left(x^{(2)}, y^{(2)} \right), \left(x^{(3)}, y^{(3)} \right), \dots, \left(x^{(n)}, y^{(n)} \right)$$

- Task:** can we learn the type of each residence?

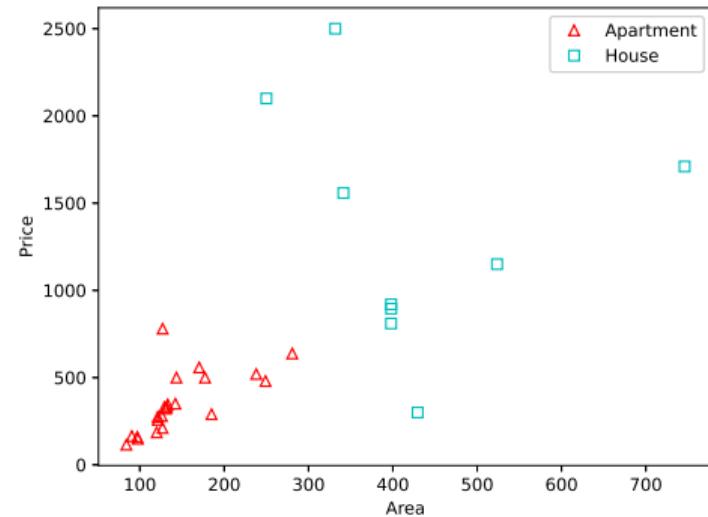
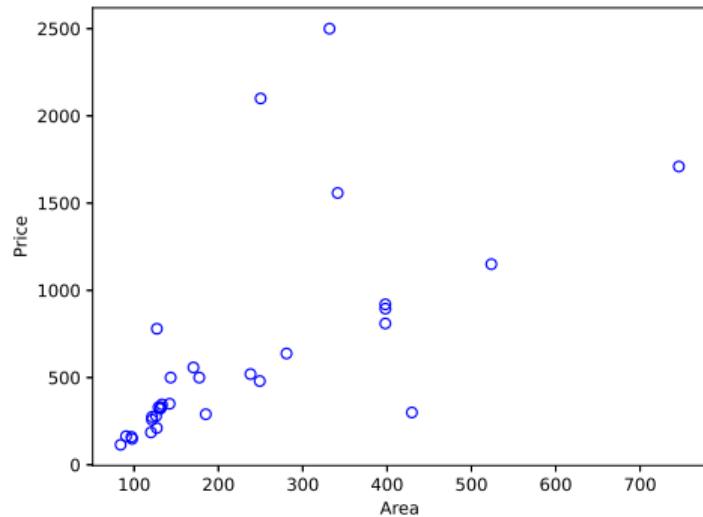


Task: housing type prediction

- Given a dataset containing n samples

$$\left(x^{(1)}, y^{(1)} \right), \left(x^{(2)}, y^{(2)} \right), \left(x^{(3)}, y^{(3)} \right), \dots, \left(x^{(n)}, y^{(n)} \right)$$

- Task:** can we learn the type of each residence?

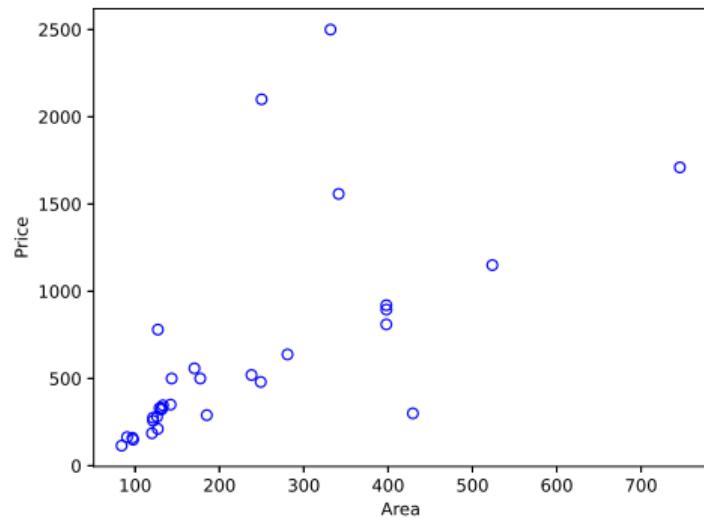


Clustering

- Given a dataset containing n samples

$$\left(x^{(1)}, y^{(1)} \right), \left(x^{(2)}, y^{(2)} \right), \left(x^{(3)}, y^{(3)} \right), \dots, \left(x^{(n)}, y^{(n)} \right)$$

- Task (vague): find interesting structures in the data

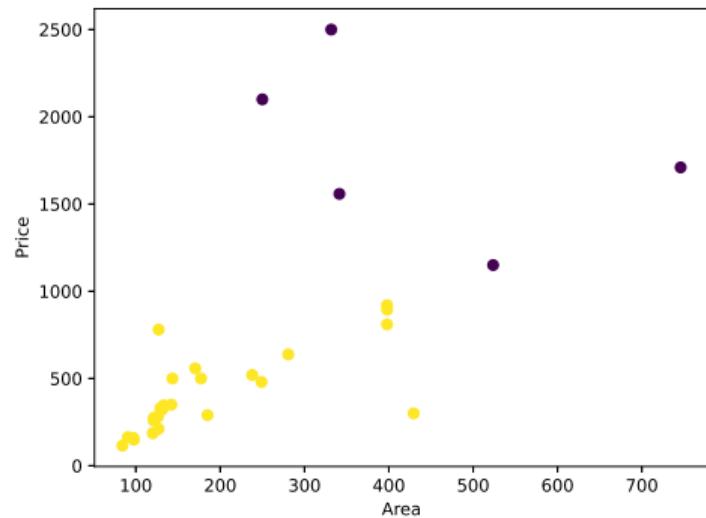
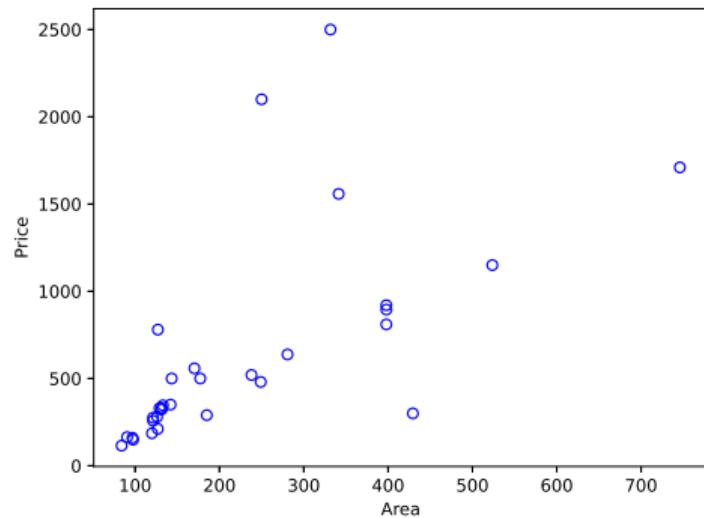


Clustering

- Given a dataset containing n samples

$$\left(x^{(1)}, y^{(1)} \right), \left(x^{(2)}, y^{(2)} \right), \left(x^{(3)}, y^{(3)} \right), \dots, \left(x^{(n)}, y^{(n)} \right)$$

- Task (vague): find interesting structures in the data

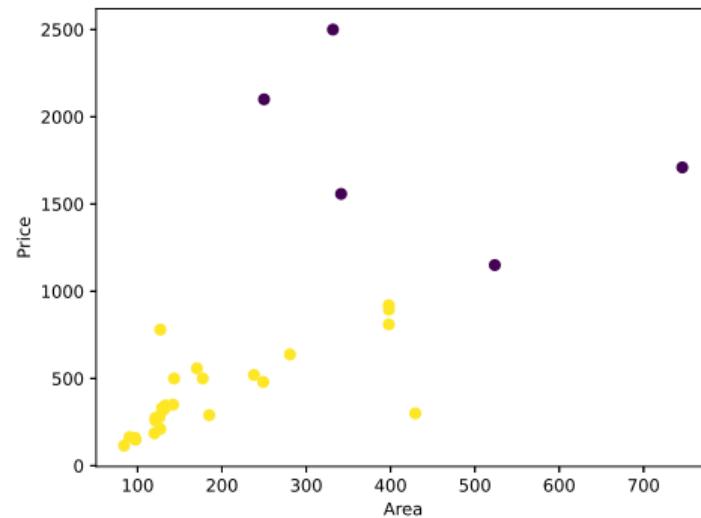
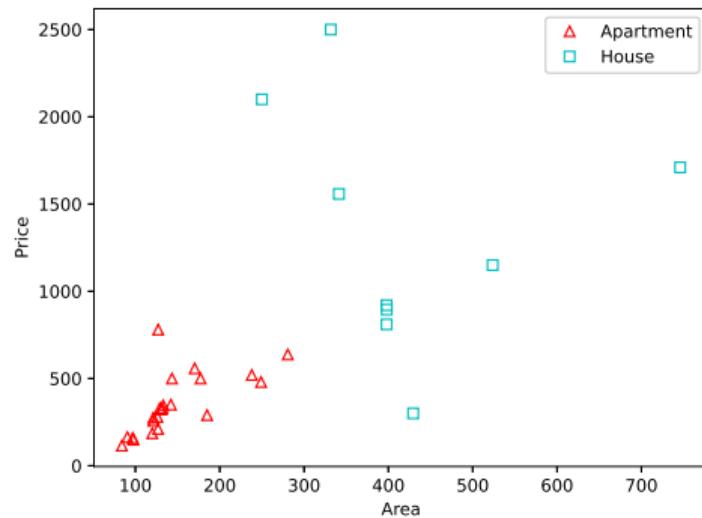


Clustering

- Given a dataset containing n samples

$$\left(x^{(1)}, y^{(1)} \right), \left(x^{(2)}, y^{(2)} \right), \left(x^{(3)}, y^{(3)} \right), \dots, \left(x^{(n)}, y^{(n)} \right)$$

- Task (vague): find interesting structures in the data

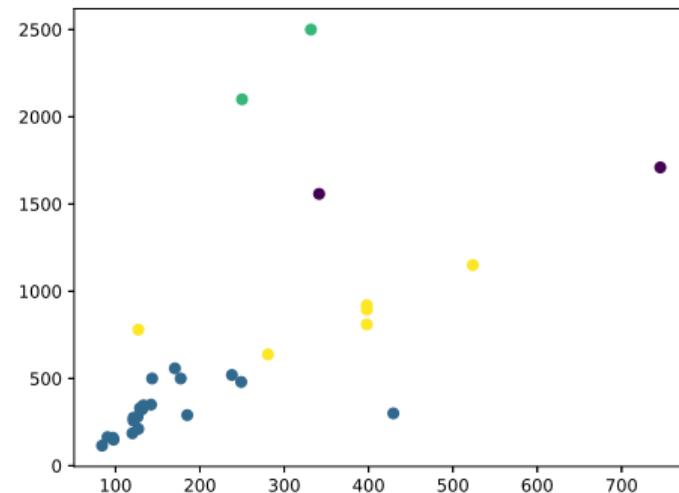
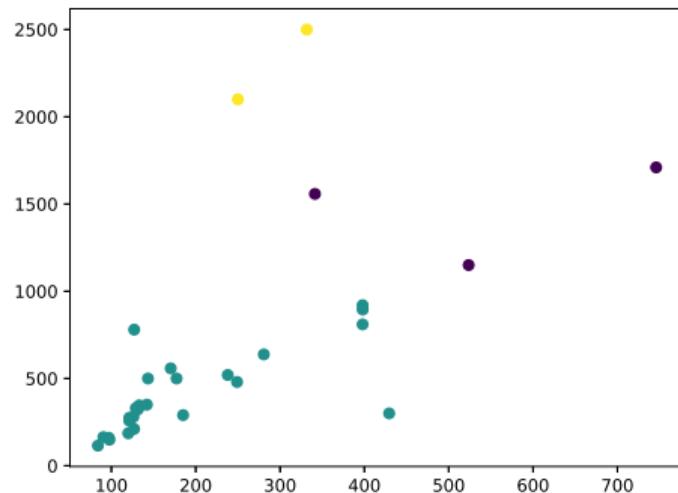


Clustering

- Given a dataset containing n samples

$$\left(x^{(1)}, y^{(1)} \right), \left(x^{(2)}, y^{(2)} \right), \left(x^{(3)}, y^{(3)} \right), \dots, \left(x^{(n)}, y^{(n)} \right)$$

- Task (vague): find interesting structures in the data



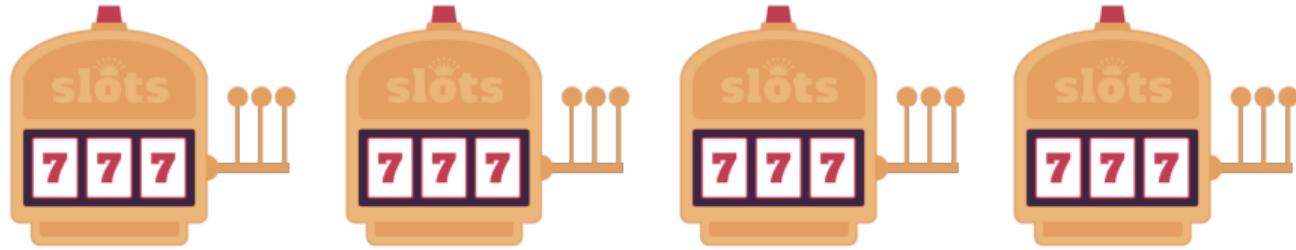
Reinforcement Learning



Task: win more

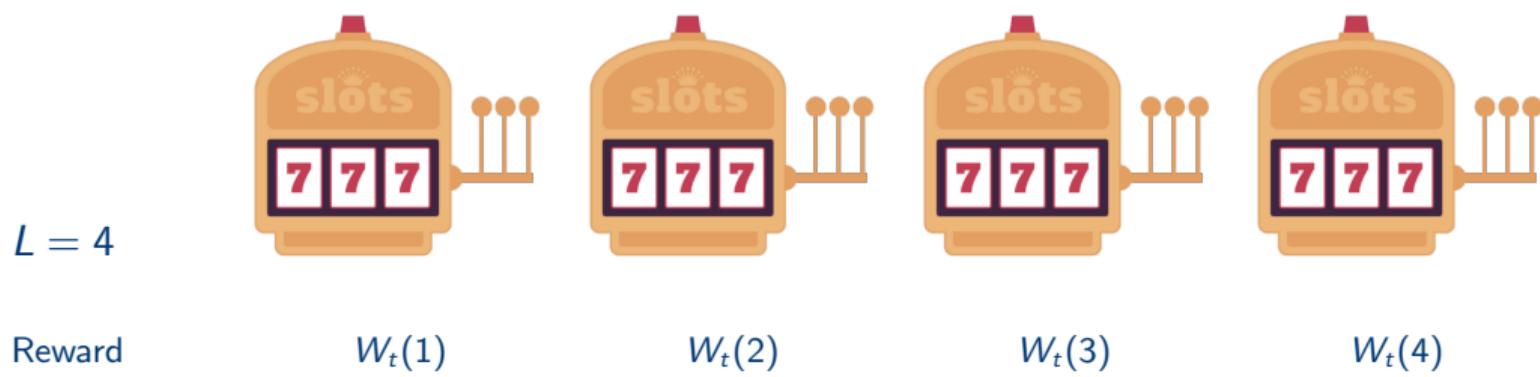
- We pull one of the 4 arms of a slot machine at each round t

$$L = 4$$



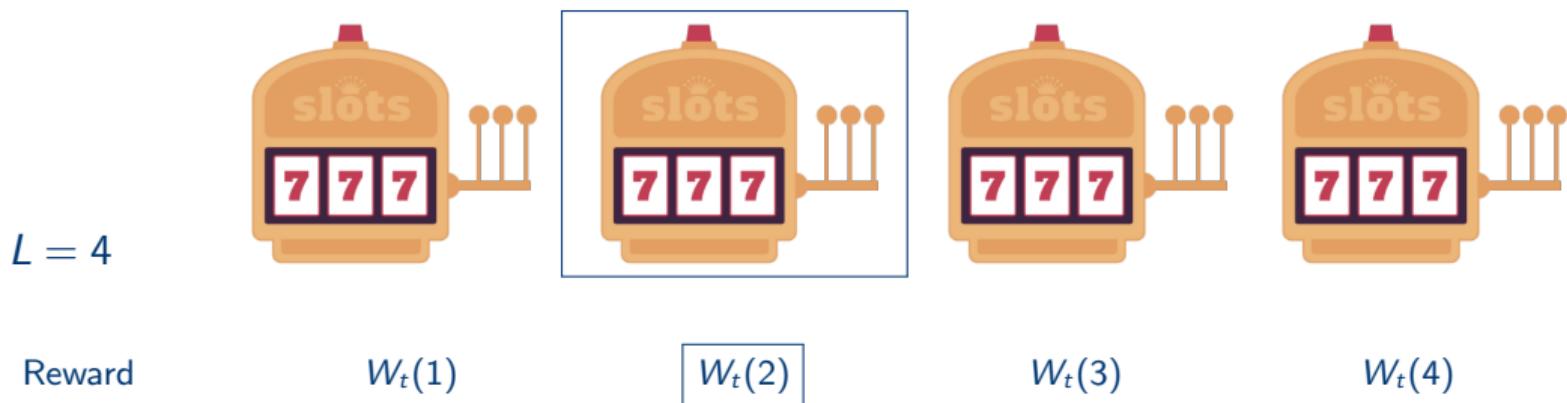
Task: win more

- We pull one of the 4 arms of a slot machine at each round t
- Each arm i yields a random reward $W_t(i)$ at round t



Task: win more

- We pull one of the 4 arms of a slot machine at each round t
- Each arm i yields a random reward $W_t(i)$ at round t
- At time t , if we pull arm 2, we earn $W_t(2)$

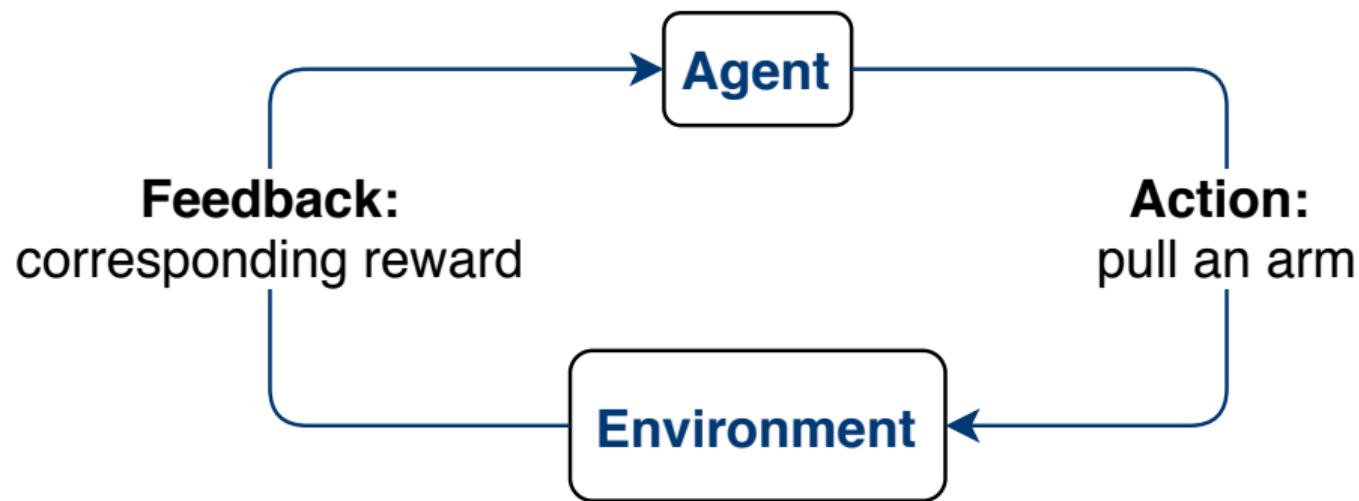


Task: win more

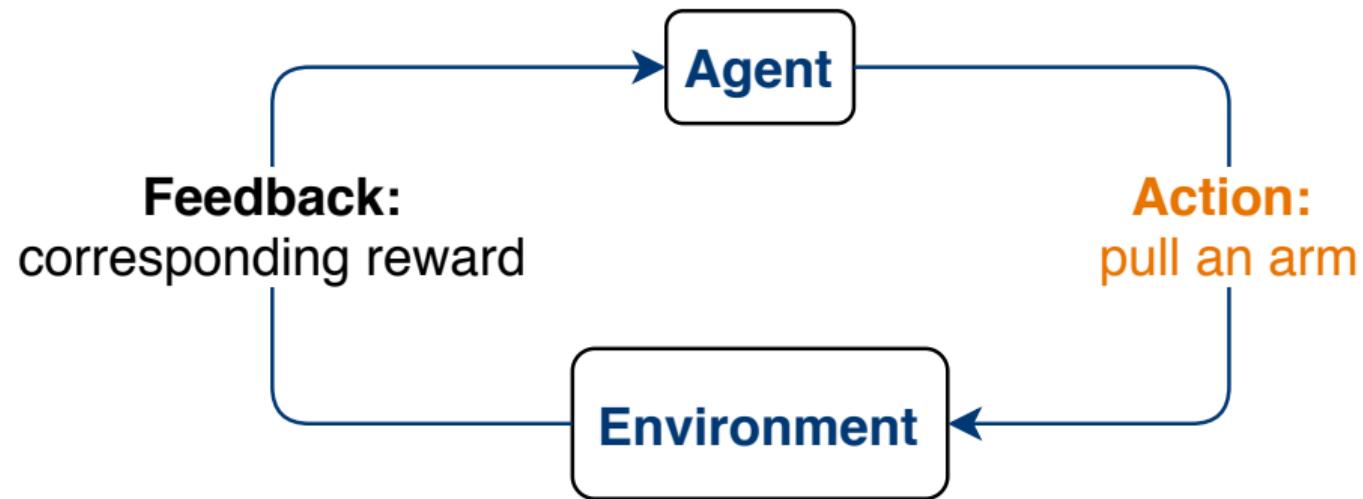
- We pull one of the 4 arms of a slot machine at each round t
- Each arm i yields a random reward $W_t(i)$ at round t
- At time t , if we pull arm 2, we earn $W_t(2)$
- Task: how can we earn as much as possible after 100 rounds?



Multi-armed bandit problem (MAB)



Multi-armed bandit problem (MAB)



- Instead of labeled/unlabeled data, a **reinforcement learning** algorithm learns from environment feedback.



Taxonomy of Machine Learning (A Simplistic View)

♠ What type of data?

- **Supervised learning** - labeled data: e.g., Prediction
- **Semi-supervised learning**
- **Unsupervised learning** - unlabeled data: e.g., clustering
- **Reinforcement learning** - environment feedback: e.g., multi-armed bandit

♠ When do we collect data?

- **Offline learning** provide all data at once
- **Online learning** provide data step by step

both proper for supervised learning, it depends on the actual implementation



Outline

1 Logistics

2 Introduction to Machine Learning

- Definition and taxonomy

3 Mathematical tools

- Set and function
- Probability and estimation
- Systems of linear equations



Sets

- A set S is an [unordered collection of objects](#).
- $S = \{1, 2, 3, 4, 5, 6\}$ is the possible outcomes of the toss of a die.
- $S = [a, b] = \{x \in \mathbb{R} : a \leq x \leq b\}$ is the set of all numbers from a to b inclusive.
- $S = (a, b] = \{x \in \mathbb{R} : a < x \leq b\}$ is the set of all numbers from a to b , excluding a including b .
- \mathbb{R} is the set of all real numbers.
- \mathbb{R}^d is the set of all real vectors of length d .

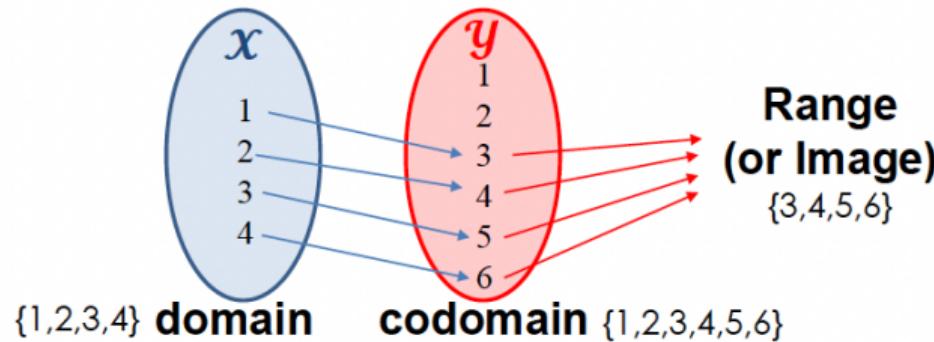


Functions

- A function f is a map from a set X to another set Y . We write this as

$$f : X \rightarrow Y.$$

- For example, the function $f : \mathbb{R} \rightarrow [0, \infty)$ could be given by the recipe $f(x) = x^2$.
- The set of inputs is called the domain; the set of possible outputs is called the codomain; the set $\{f(x) : x \in X\}$ is called the range (or image).
- For example $f : \{1, 2, 3, 4\} \rightarrow \{1, 2, 3, 4, 5, 6\}$ given by the recipe $f(x) = x + 2$ has codomain $\{1, 2, 3, 4, 5, 6\}$ and range $\{3, 4, 5, 6\}$.



Functions

Definition 3.1

Let $f(x)$ be a continuous function defined on domain \mathbb{R} .

- $f(x)$ is **differentiable** if $f'(x) = \lim_{x \rightarrow c} \frac{f(x) - f(c)}{x - c}$ exists for all $x \in \mathbb{R}$.
- $f(x)$ is **convex** if $f(\lambda a + (1 - \lambda)b) \leq \lambda f(a) + (1 - \lambda)f(b)$ for all $a, b \in \mathbb{R}$;
- $f(x)$ is **strictly convex** if $f(\lambda a + (1 - \lambda)b) < \lambda f(a) + (1 - \lambda)f(b)$ for all $a, b \in \mathbb{R}$;
- $f(x)$ is **concave** if $f(\lambda a + (1 - \lambda)b) \geq \lambda f(a) + (1 - \lambda)f(b)$ for all $a, b \in \mathbb{R}$;
- $f(x)$ is **strictly concave** if $f(\lambda a + (1 - \lambda)b) > \lambda f(a) + (1 - \lambda)f(b)$ for all $a, b \in \mathbb{R}$;



Functions

Definition 3.1

Let $f(x)$ be a continuous function defined on domain \mathbb{R} .

- $f(x)$ is **differentiable** if $f'(x) = \lim_{x \rightarrow c} \frac{f(x) - f(c)}{x - c}$ exists for all $x \in \mathbb{R}$.
- $f(x)$ is **convex** if $f(\lambda a + (1 - \lambda)b) \leq \lambda f(a) + (1 - \lambda)f(b)$ for all $a, b \in \mathbb{R}$;
- $f(x)$ is **strictly convex** if $f(\lambda a + (1 - \lambda)b) < \lambda f(a) + (1 - \lambda)f(b)$ for all $a, b \in \mathbb{R}$;
- $f(x)$ is **concave** if $f(\lambda a + (1 - \lambda)b) \geq \lambda f(a) + (1 - \lambda)f(b)$ for all $a, b \in \mathbb{R}$;
- $f(x)$ is **strictly concave** if $f(\lambda a + (1 - \lambda)b) > \lambda f(a) + (1 - \lambda)f(b)$ for all $a, b \in \mathbb{R}$;

- $f(x)$ is differentiable and convex: x_0 such that $f'(x_0) = 0$ minimize function f .
- $f(x)$ is differentiable and concave: x_0 such that $f'(x_0) = 0$ maximize function f .



Linear functions

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is linear if it satisfies

- **(Homogeneity)** For any vector $\mathbf{x} \in \mathbb{R}^d$ and scalar $a \in \mathbb{R}$,

$$f(a \mathbf{x}) = a f(\mathbf{x})$$

- **(Additivity)** For any two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$,

$$f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$$

Note that a linear function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ must pass through the origin, i.e., $f(\mathbf{0}) = 0$ where $\mathbf{0} \in \mathbb{R}^d$ is the zero vector in d dimensions. **Why?**



Affine functions

- An affine function f is a linear function plus possibly a constant.
- More precisely, a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is affine if it can be expressed as

$$f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x} + b$$

for some vector $\mathbf{a} \in \mathbb{R}^d$ and some scalar $b \in \mathbb{R}$.

- The scalar b is called the bias or offset.

Example: The following function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is affine. Why?

$$f(\mathbf{x}) = f(x_1, x_2) = -x_1 + 3x_2 + 7.$$

Exercise: Is a linear function affine? Is an affine function linear?



Linear and Affine Functions

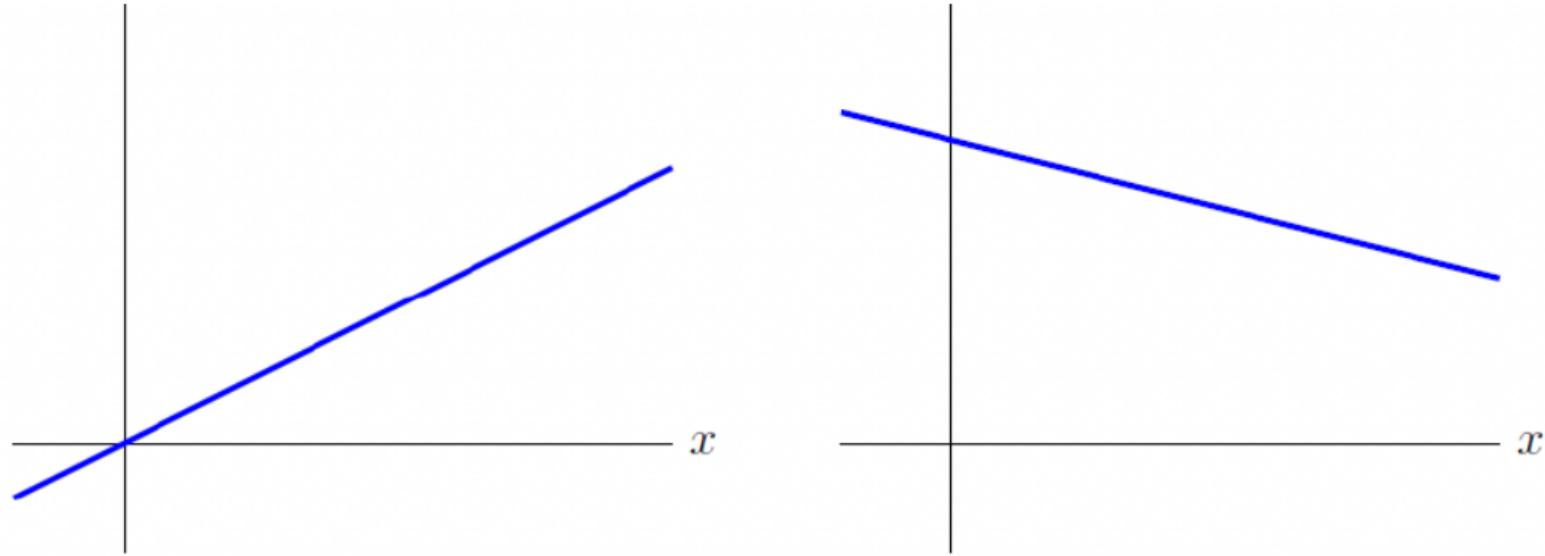


Figure: Left: Linear; Right: Affine but not linear.



Local and global extrema

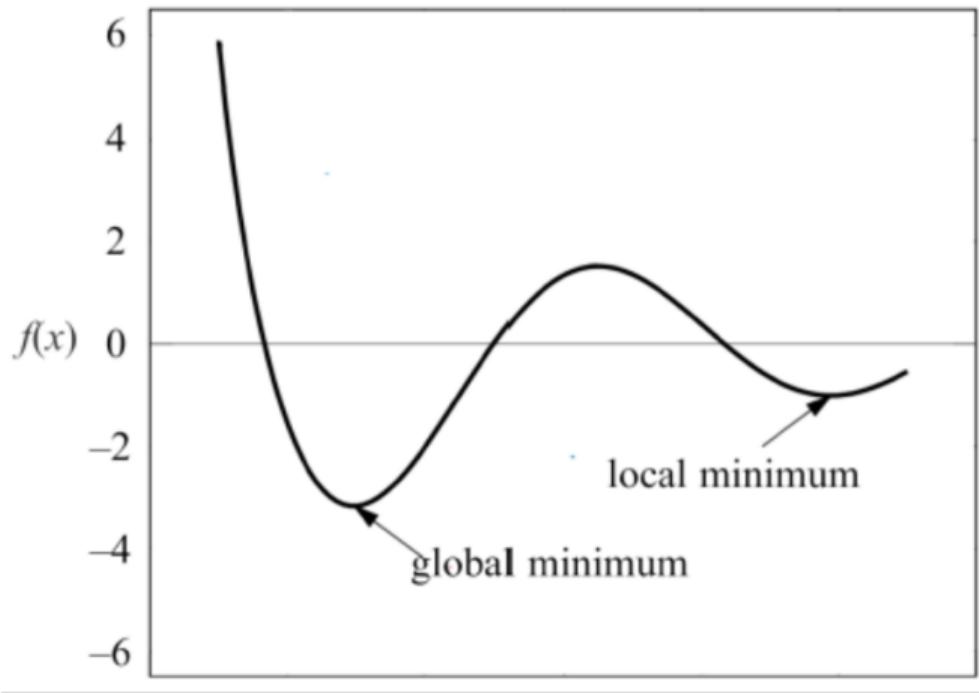
- Consider a function $f : [a, b] \rightarrow \mathbb{R}$.
- The function f has a local minimum at $c \in \mathbb{R}$ if $f(x) \geq f(c)$ for all x in an open neighborhood of c .
- The function f has a global minimum at $c \in \mathbb{R}$ if $f(x) \geq f(c)$ for all $x \in [a, b]$.

Exercise: If c is a local minimum of f , is it a global minimum? If c is a global minimum of f , is it a local minimum?

Exercise: How would you define local maximum and global maximum?



Local and global extrema



min and arg min

- For a function $f : X \rightarrow Y$, the **minimum** $\min_{x \in X} f(x)$ returns the smallest value among all elements in the set $\{f(x) : x \in X\}$.
- For a function $f : X \rightarrow Y$, the **argmin** $x^* = \arg \min_{x \in X} f(x)$ returns the value of $x \in X$ that minimizes $f(x)$, i.e.,

$$f(x^*) = \min_{x \in X} f(x)$$

- $\arg \min$ returns a value from the **domain** of the function X and \min returns from the **range** (codomain) Y of the function.
- Let $X = \{0, 1\}$ and $f(0) = \pi$ and $f(1) = e$. Then

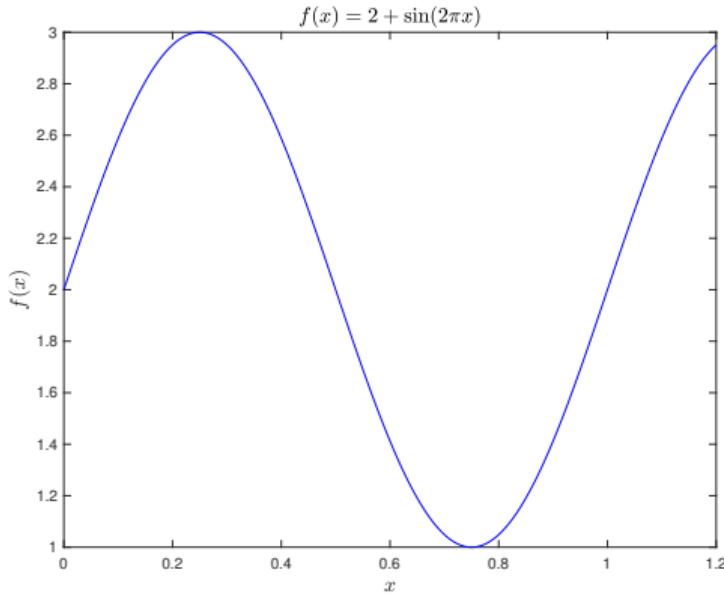
$$\arg \min_{x \in X} f(x) = 1 \quad \min_{x \in X} f(x) = e,$$

and

$$\arg \max_{x \in X} f(x) = 0 \quad \max_{x \in X} f(x) = \pi.$$



min and $\arg \min$



- Let $f : X = [0, 1.2] \rightarrow \mathbb{R}$ be defined as $f(x) = 2 + \sin(2\pi x)$ (see plot above). Then

$$\arg \min_{x \in X} f(x) = 3/4 \quad \min_{x \in X} f(x) = +1$$



Derivatives

- For a multivariable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, its **gradient vector** or **derivative** at $\mathbf{x} \in \mathbb{R}^d$ is the column vector

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \frac{df}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \cdots & \frac{\partial f}{\partial x_d} \end{bmatrix}^\top$$

- Recall that $\frac{\partial f}{\partial x_i}$ is the **partial derivative** of f with respect to the scalar variable x_i .
- For example, if $f(x_1, x_2) = 2x_1^2 + 5x_1x_2 + 3x_2^3$, then

$$\frac{\partial f}{\partial x_1} = 4x_1 + 5x_2 \quad \text{and} \quad \frac{\partial f}{\partial x_2} = 5x_1 + 9x_2^2.$$



Important derivatives

- There are only two derivatives for functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that take vectors to scalars you need to know for now.
- For a fixed vector $\mathbf{a} \in \mathbb{R}^d$, consider $f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x}$ (known as the [inner or dot product](#)). Then

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \mathbf{a}.$$

- For a fixed matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, consider $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$ (known as the [quadratic form](#)). Then

见微知著，先计算矩阵内每一项，找到规律之后再整合

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}.$$

定义入手证明

- In most applications, \mathbf{A} is a [symmetric](#) matrix (i.e., $\mathbf{A} = \mathbf{A}^\top$) so

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = 2\mathbf{A}\mathbf{x}.$$

- This generalizes the basic fact that if $f(x) = ax^2$, then $\frac{df}{dx} = 2ax$.



Important derivatives

Exercise: Show from the definition of $f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x}$ that

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \mathbf{a}.$$

Exercise: Show from the definition of $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A}\mathbf{x}$ that

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = (\mathbf{A} + \mathbf{A}^\top)\mathbf{x}.$$

You'll be forgiven for having to exhibit a substantial amount of meticulous bookkeeping here.

Advice: It is difficult to remember a lot of derivative formulas of complicated multivariate functions. Usually, one consults the Matrix Cookbook

<https://www2.imm.dtu.dk/pubdb/doc/imm3274.pdf>



Basic Probability Theory

Statistical machine learning deals with uncertainty



Definitions of events, probabilities, joint probabilities and conditional probabilities



Notations

- Random variables (r.v.s): upper (e.g., X)
- Values r.v.s take on: lower case (e.g., x, x_i)
- Sets (in which the r.v.s assume their values): calligraphic font (e.g., \mathcal{X})
- Probability, expectation and variance operators: $\Pr(\cdot)$, $\mathbb{E}[\cdot]$, $\text{Var}(\cdot)$



Discrete random variables

Random variables X and Y :

- X can only take values in the finite set $\mathcal{X} = \{x_1, x_2, \dots, x_M\}$.
- Y can only take values in the finite set $\mathcal{Y} = \{y_1, y_2, \dots, y_L\}$.

A random variable is said to be **discrete** if the set of values it can take has a finite number of values.

$$\sum_{i=1}^M \Pr(X = x_i) = 1, \sum_{i=1}^L \Pr(Y = y_i) = 1.$$



Event

Consider m trials (of sampling X and Y) and let the number of trials for which $X = x_i$ and $Y = y_j$ be m_{ij} . Then, if m is large, we can assume that

$$\Pr(X = x_i, Y = y_j) = \frac{m_{ij}}{m}.$$

- The argument of $\Pr(\cdot)$ is known as an **event**.
- **Event:** A set in which a r.v. (or multiples r.v.s) assume some value(s) in some set.
e.g., $\{X = x_i, Y = y_j\}$.



Sum rule and product rule

- **Sum rule:**

$$\Pr(X = x_i) = \sum_{j=1}^L \Pr(X = x_i, Y = y_j), \quad \Pr(Y = y_j) = \sum_{i=1}^M \Pr(X = x_i, Y = y_j).$$

$$\Pr(X = x_i, Y = y_j) = \Pr(X = x_i | Y = y_j) \cdot \Pr(Y = y_j).$$

- **Product rule:**

$$\Pr(X = x) = \sum_{y \in \mathcal{Y}} \Pr(X = x, Y = y).$$

$$\Pr(X = x, Y = y) = \Pr(X = x | Y = y) \cdot \Pr(Y = y).$$



Bayes' rule

- **Probability mass function (pmf):**

$$p(x) := \Pr(X = x), \quad p(x, y) := \Pr(X = x, Y = y), \quad p(x|y) := \Pr(X = x | Y = y).$$

- **Bayes' rule/Bayes' Theorem:**

$$p(y|x) = \frac{p(x,y)}{p(x)} = \frac{p(x|y)p(y)}{p(x)} = \frac{p(x|y)p(y)}{\sum_{y' \in \mathcal{Y}} p(x|y')p(y')}.$$

- ▶ ‘Invert the causal relationship’ between X and Y .
- ▶ $p(x|y) \propto p(y|x)p(x)$.
- ▶ posterior = $\frac{\text{likelihood} \times \text{prior}}{\text{model evidence}} \propto \text{likelihood} \times \text{prior}$.



Continuous random variables

- **Cumulative distribution function (cdf):** $x \in \mathbb{R} \mapsto F_X(x) = \Pr(X \leq x)$ is differentiable on \mathbb{R} .
- **Probability density function (pdf):**

$$f_X(x) = \frac{d}{dx} F_X(x), \quad x \in \mathbb{R}; \quad f_X(x) \geq 0 \quad \forall x \in \mathbb{R} \text{ and } \int_{\mathbb{R}} f_X(x) dx = 1.$$



Continuous random variables

- **Cumulative distribution function (cdf):** $x \in \mathbb{R} \mapsto F_X(x) = \Pr(X \leq x)$ is differentiable on \mathbb{R} .
- **Probability density function (pdf):**

$$f_X(x) = \frac{d}{dx} F_X(x), \quad x \in \mathbb{R}; \quad f_X(x) \geq 0 \quad \forall x \in \mathbb{R} \text{ and } \int_{\mathbb{R}} f_X(x) dx = 1.$$

- **Joint probability density function** of continuous r.v.s.:

$$f_{X,Y}(x,y) = \frac{\partial^2}{\partial x \partial y} F_{X,Y}(x,y), \quad (x,y) \in \mathbb{R}^2.$$

- Conditional probability and independence can be defined analogously to the discrete case:

$$\Pr(a < X \leq b, c < Y \leq d) = \int_c^d \int_a^b f_{X,Y}(x,y) dx dy.$$



Independence

Definition 3.2 (Independence)

Two random variables X and Y are independent if $p_{X,Y}(x,y) = p_X(x)p_Y(y)$ for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.

- X and Y are independent if $p_{X|Y}(x|y) = p_X(x)$ for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$ such that $p_Y(y) > 0$ (**what if $p_Y(y) = 0$?**).



Expectation and variance

- **Expectation:**

- ▶ Discrete r.v.: $\mathbb{E}X = \sum_{x \in \mathcal{X}} xp_X(x).$
- ▶ Continuous r.v.: $\mathbb{E}X = \int_{\mathbb{R}} xf_X(x)dx.$

If g is a function from the domain of X to \mathbb{R} , we can obtain the expectation of $Y = g(X)$ in the same way:

$$\mathbb{E}Y = \mathbb{E}[g(X)] = \int_{\mathbb{R}} yf_Y(y)dy = \int_{\mathbb{R}} g(x)f_X(x)dx.$$

- ▶ **Linearity of expectation:** if $g(X) = aX + b$ (for constants $a, b \in \mathbb{R}$), then $\mathbb{E}[g(X)] = a\mathbb{E}[X] + b = g(\mathbb{E}[X]).$

- **Variance:**

- ▶ Discrete r.v.: $\text{Var}(X) = \sum_{x \in \mathcal{X}} (x - \mathbb{E}X)^2 p_X(x).$
- ▶ Continuous r.v.: $\text{Var}(X) = \int_{\mathbb{R}} (x - \mathbb{E}X)^2 f_X(x)dx.$



Maximum likelihood estimation (MLE)

- In machine learning, we almost always do **NOT** have access to the underlying distributions that the data are generated from.
- Rather, we have access to sample data and we would like to **use it to estimate some parameters**.



Maximum likelihood estimation: Gaussian distribution

- If samples in $\mathcal{D} = \{X_1, \dots, X_m\}$ are independently drawn from the univariate Gaussian distribution

$$f_X(x; \mu, \sigma^2) = \mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) \quad x \in \mathbb{R}, \quad (3.1)$$

we can use the samples in \mathcal{D} to estimate the parameter vector $\theta = (\mu, \sigma^2)$.

- ▶ Note: notation $f_X(x; \mu, \sigma^2)$ or $f_X(x; \theta)$ emphasize that the density is parametrized by the parameters (μ, σ^2) or θ .



Maximum likelihood estimation: Gaussian distribution

- If samples in $\mathcal{D} = \{X_1, \dots, X_m\}$ are independently drawn from the univariate Gaussian distribution

$$f_X(x; \mu, \sigma^2) = \mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) \quad x \in \mathbb{R}, \quad (3.1)$$

we can use the samples in \mathcal{D} to estimate the parameter vector $\theta = (\mu, \sigma^2)$.

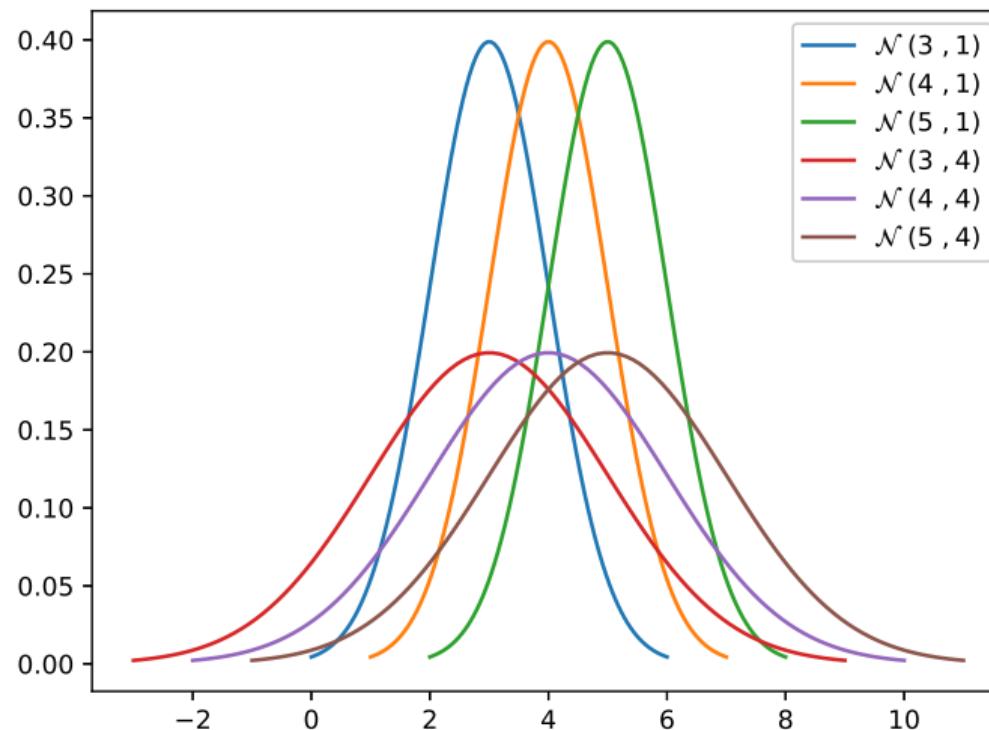
- ▶ Note: notation $f_X(x; \mu, \sigma^2)$ or $f_X(x; \theta)$ emphasize that the density is parametrized by the parameters (μ, σ^2) or θ .

- ④ Maximum likelihood estimates for the mean and variance:

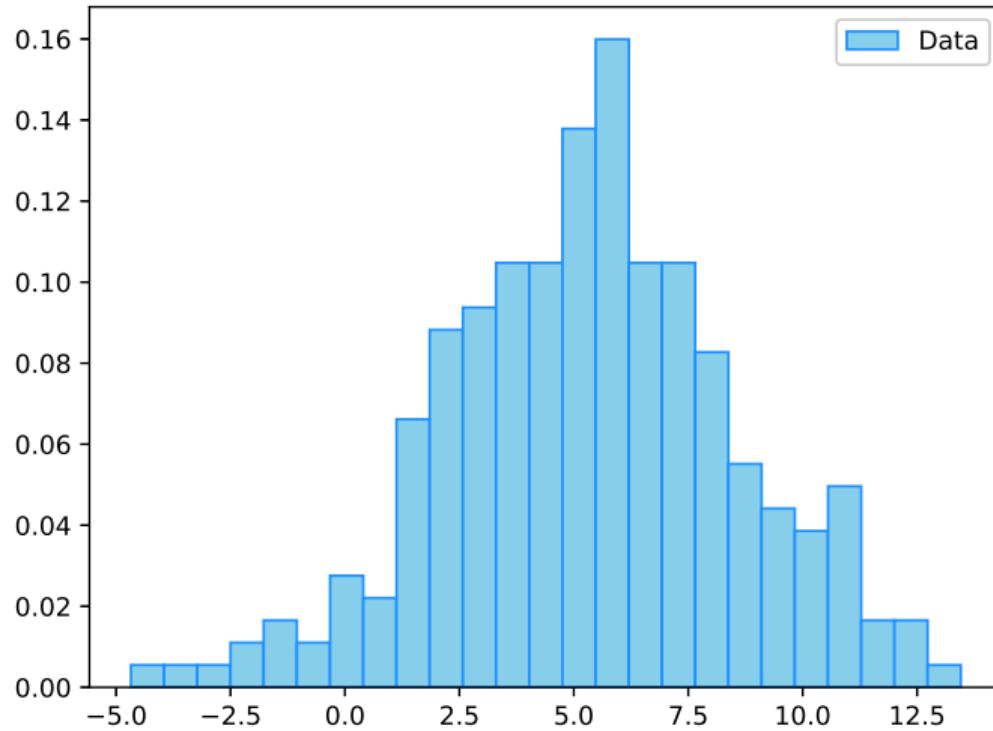
$$\hat{\mu}_{\text{ML}} = \frac{1}{m} \sum_{i=1}^m X_i \quad \text{and} \quad \hat{\sigma}_{\text{ML}}^2 = \frac{1}{m} \sum_{i=1}^m (X_i - \hat{\mu})^2. \quad (3.2)$$



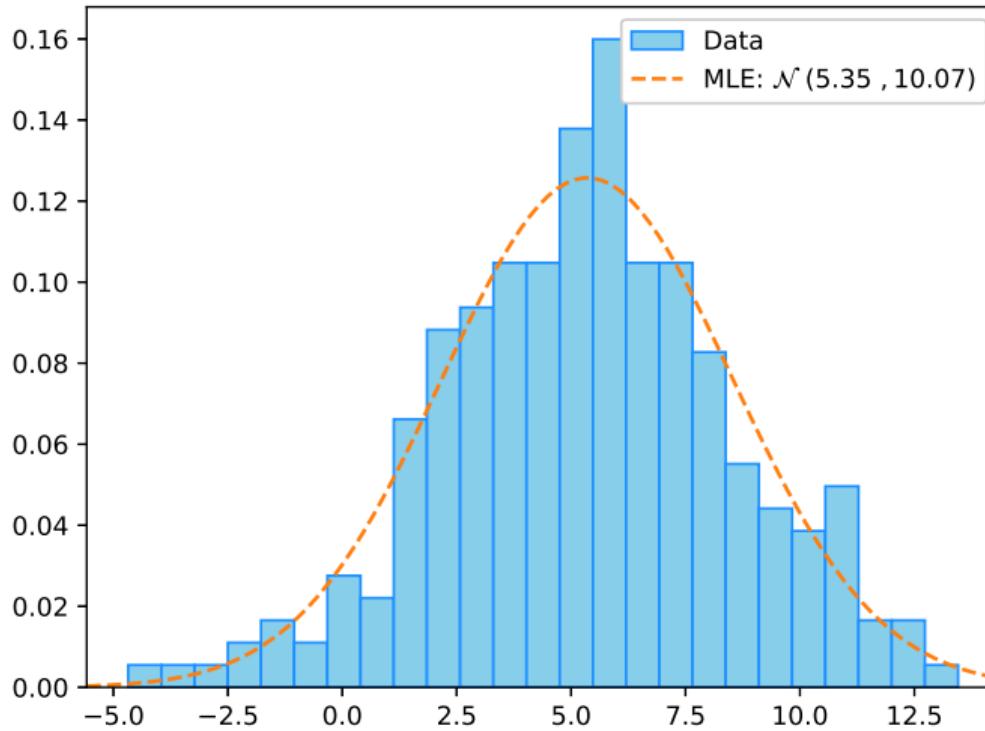
Maximum likelihood estimation: Gaussian distribution



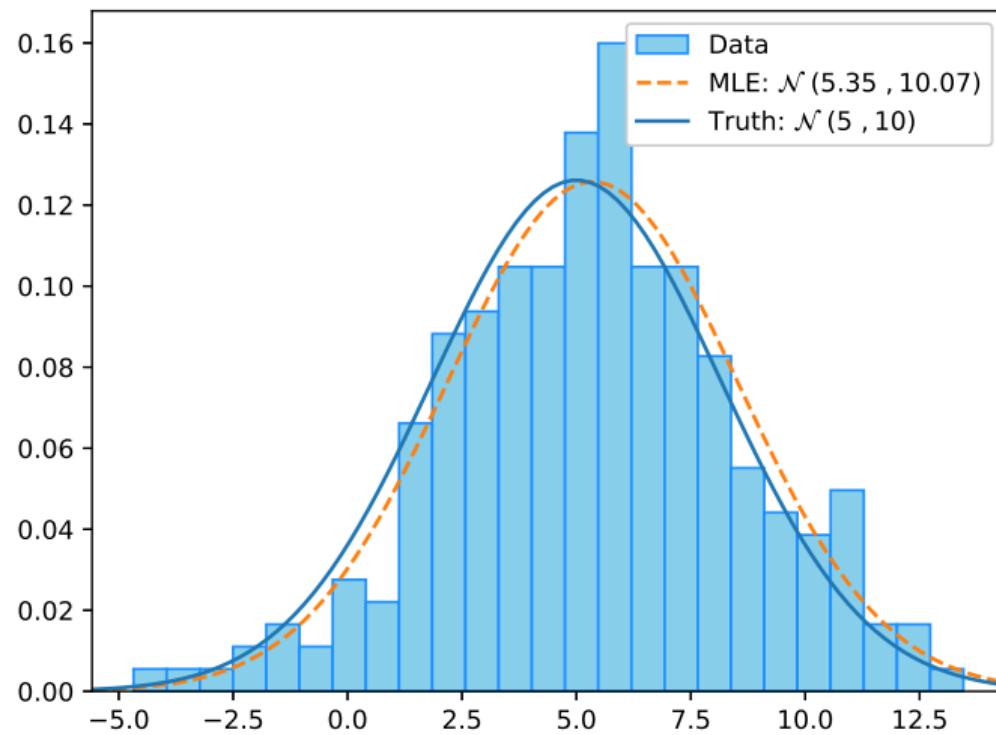
Maximum likelihood estimation: Gaussian distribution



Maximum likelihood estimation: Gaussian distribution



Maximum likelihood estimation: Gaussian distribution



Maximum likelihood estimation: Bernoulli distribution

- Samples in $\mathcal{D} = \{X_1, \dots, X_m\}$ are generated independently from the Bernoulli (coin toss) distribution

$$p_X(x; \theta) = \begin{cases} 1 - \theta & x = 0 \\ \theta & x = 1 \end{cases} \quad (3.3)$$

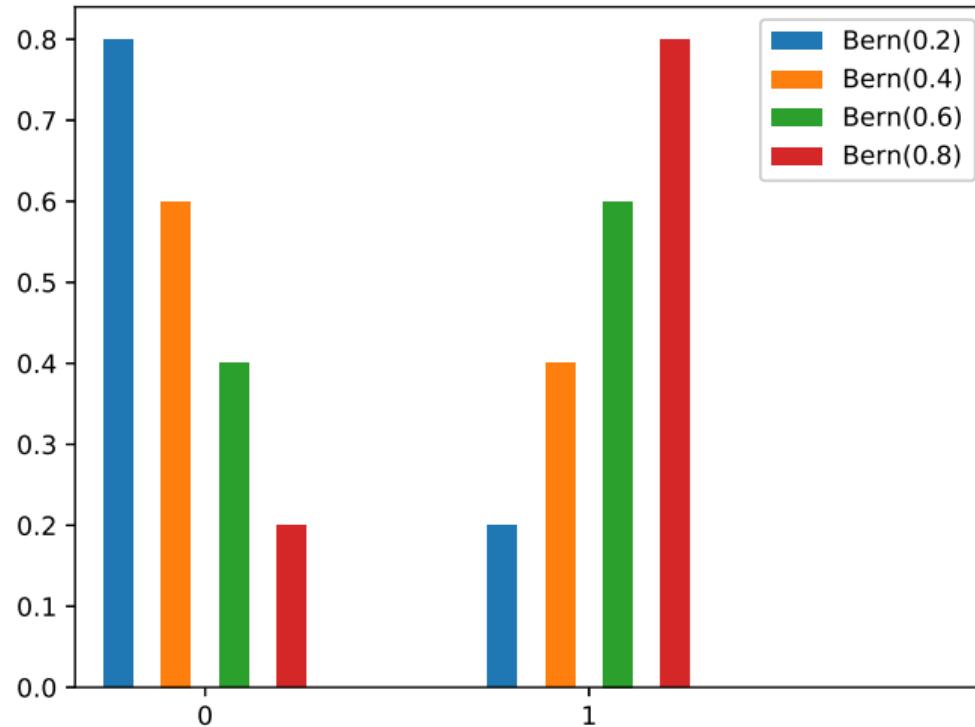
or

$$p_X(x; \theta) = (1 - \theta)^{1-x} \theta^x, \quad x \in \{0, 1\}. \quad (3.4)$$

- How would we estimate θ from samples?



Maximum likelihood estimation: Bernoulli distribution



Maximum likelihood estimation: Bernoulli distribution

- $M_1 := \sum_{i=1}^m X_i$ denote the total number of ones in \mathcal{D} . Consider

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in (0,1)} \prod_{i=1}^m p_X(X_i; \theta) \quad (3.5)$$



Maximum likelihood estimation: Bernoulli distribution

- $M_1 := \sum_{i=1}^m X_i$ denote the total number of ones in \mathcal{D} . Consider

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in (0,1)} \prod_{i=1}^m p_X(X_i; \theta) \quad (3.5)$$

$$= \arg \max_{\theta \in (0,1)} \sum_{i=1}^m \log p_X(X_i; \theta) \quad (3.6)$$



Maximum likelihood estimation: Bernoulli distribution

- $M_1 := \sum_{i=1}^m X_i$ denote the total number of ones in \mathcal{D} . Consider

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in (0,1)} \prod_{i=1}^m p_X(X_i; \theta) \quad (3.5)$$

$$= \arg \max_{\theta \in (0,1)} \sum_{i=1}^m \log p_X(X_i; \theta) \quad (3.6)$$

$$= \arg \max_{\theta \in (0,1)} \sum_{i=1}^m [(1 - X_i) \log(1 - \theta) + X_i \log \theta] \quad (3.7)$$

$$= \arg \max_{\theta \in (0,1)} (m - M_1) \log(1 - \theta) + M_1 \log \theta. \quad (3.8)$$



Maximum likelihood estimation: Bernoulli distribution

- $M_1 := \sum_{i=1}^m X_i$ denote the total number of ones in \mathcal{D} . Consider

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in (0,1)} \underbrace{(m - M_1) \log(1 - \theta) + M_1 \log \theta}_{\text{Objective function } g(\theta)}. \quad (3.9)$$



Maximum likelihood estimation: Bernoulli distribution

- $M_1 := \sum_{i=1}^m X_i$ denote the total number of ones in \mathcal{D} . Consider

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in (0,1)} \underbrace{(m - M_1) \log(1 - \theta) + M_1 \log \theta}_{\text{Objective function } g(\theta)}. \quad (3.9)$$

- $g(\theta)$ is strictly concave \Rightarrow differentiating and setting to zero yields the (unique) maximum:



Maximum likelihood estimation: Bernoulli distribution

- $M_1 := \sum_{i=1}^m X_i$ denote the total number of ones in \mathcal{D} . Consider

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in (0,1)} \underbrace{(m - M_1) \log(1 - \theta) + M_1 \log \theta}_{\text{Objective function } g(\theta)}. \quad (3.9)$$

- $g(\theta)$ is strictly concave \Rightarrow differentiating and setting to zero yields the (unique) maximum:

$$g'(\theta) = -\frac{m - M_1}{1 - \theta} + \frac{M_1}{\theta}, \quad (3.10)$$

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in (0,1)} g(\theta) \implies g'(\hat{\theta}_{\text{ML}}) = 0 \implies \hat{\theta}_{\text{ML}} = \frac{M_1}{m}. \quad (3.11)$$



Maximum likelihood estimation: Bernoulli distribution

- $M_1 := \sum_{i=1}^m X_i$ denote the total number of ones in \mathcal{D} . Consider

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in (0,1)} \underbrace{(m - M_1) \log(1 - \theta) + M_1 \log \theta}_{\text{Objective function } g(\theta)}. \quad (3.9)$$

- $g(\theta)$ is strictly concave \Rightarrow differentiating and setting to zero yields the (unique) maximum:

$$g'(\theta) = -\frac{m - M_1}{1 - \theta} + \frac{M_1}{\theta}, \quad (3.10)$$

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in (0,1)} g(\theta) \implies g'(\hat{\theta}_{\text{ML}}) = 0 \implies \hat{\theta}_{\text{ML}} = \frac{M_1}{m}. \quad (3.11)$$

- Mean θ is estimated by the empirical mean M_1/m : agrees with common sense!



Maximum likelihood estimation: exponential distribution

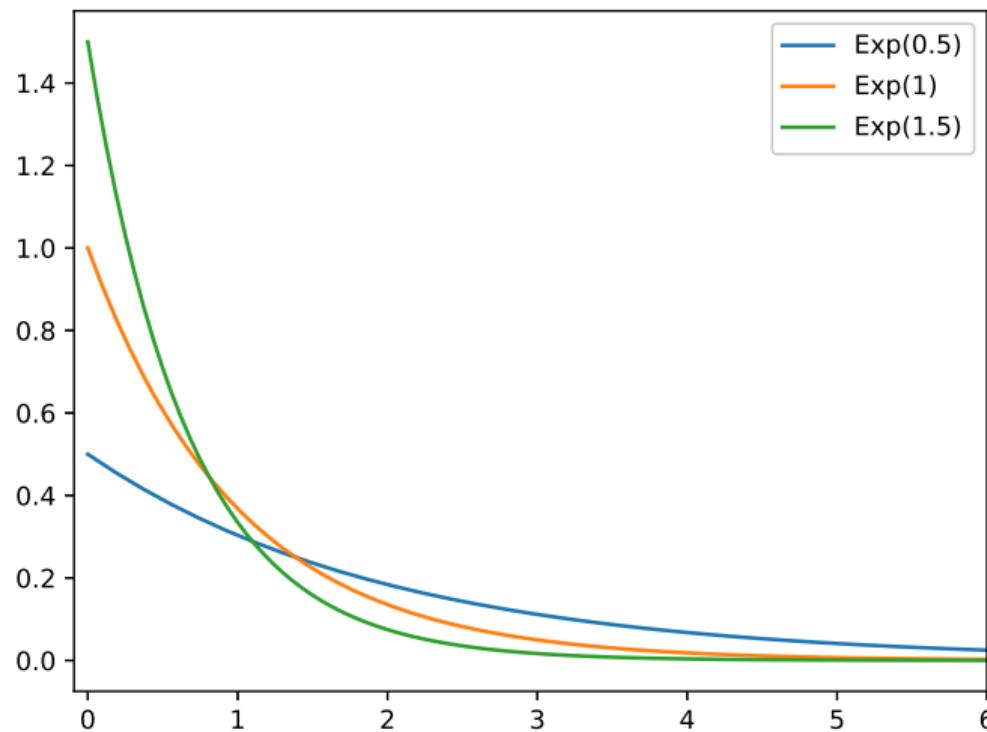
- Samples in $\mathcal{D} = \{X_1, \dots, X_m\}$ are generated independently from the exponential distribution (memory-less)

$$f_X(x; \theta) = \begin{cases} \theta \exp(-\theta x) & x \geq 0 \\ 0 & x < 0 \end{cases}.$$

- $\theta > 0$ is the unknown (rate) parameter. In fact, $\mathbb{E}[X] = 1/\theta$.
- How would we estimate θ from samples?



Maximum likelihood estimation: exponential distribution



Maximum likelihood estimation: exponential distribution

- Consider

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta > 0} \prod_{i=1}^m f_X(X_i; \theta) = \arg \max_{\theta > 0} \sum_{i=1}^m \log f_X(X_i; \theta) \quad (3.12)$$



Maximum likelihood estimation: exponential distribution

- Consider

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta > 0} \prod_{i=1}^m f_X(X_i; \theta) = \arg \max_{\theta > 0} \sum_{i=1}^m \log f_X(X_i; \theta) \quad (3.12)$$

$$= \arg \max_{\theta > 0} \sum_{i=1}^m (\log \theta - \theta X_i) = \arg \max_{\theta > 0} m \log \theta - \sum_{i=1}^m \theta X_i. \quad (3.13)$$



Maximum likelihood estimation: exponential distribution

- Consider

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta > 0} \prod_{i=1}^m f_X(X_i; \theta) = \arg \max_{\theta > 0} \sum_{i=1}^m \log f_X(X_i; \theta) \quad (3.12)$$

$$= \arg \max_{\theta > 0} \sum_{i=1}^m (\log \theta - \theta X_i) = \arg \max_{\theta > 0} m \log \theta - \sum_{i=1}^m \theta X_i. \quad (3.13)$$

- $g(\theta) = m \log \theta - \sum_{i=1}^m \theta X_i$ is strictly concave \Rightarrow differentiating and setting to zero yields the (unique) maximum:



Maximum likelihood estimation: exponential distribution

- Consider

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta > 0} \prod_{i=1}^m f_X(X_i; \theta) = \arg \max_{\theta > 0} \sum_{i=1}^m \log f_X(X_i; \theta) \quad (3.12)$$

$$= \arg \max_{\theta > 0} \sum_{i=1}^m (\log \theta - \theta X_i) = \arg \max_{\theta > 0} m \log \theta - \sum_{i=1}^m \theta X_i. \quad (3.13)$$

- $g(\theta) = m \log \theta - \sum_{i=1}^m \theta X_i$ is strictly concave \Rightarrow differentiating and setting to zero yields the (unique) maximum:

$$g'(\theta) = \frac{m}{\hat{\theta}_{\text{ML}}} - \sum_{i=1}^m X_i, \quad (3.14)$$

$$g'(\hat{\theta}_{\text{ML}}) = 0 \implies \hat{\theta}_{\text{ML}} = \hat{\theta}_{\text{ML}} = \frac{m}{\sum_{i=1}^n X_i} = \left(\frac{1}{m} \sum_{i=1}^m X_i \right)^{-1}. \quad (3.15)$$



Maximum likelihood estimation: exponential distribution

- Exponential distribution (this distribution models waiting times for buses)

$$f_X(x; \theta) = \begin{cases} \theta \exp(-\theta x) & x \geq 0 \\ 0 & x < 0 \end{cases}.$$

- $\mathbb{E}[X] = 1/\theta$, $\hat{\theta}_{\text{ML}} = \left(\frac{1}{m} \sum_{i=1}^m X_i\right)^{-1}$.



Maximum likelihood estimation: exponential distribution

- Exponential distribution (this distribution models waiting times for buses)

$$f_X(x; \theta) = \begin{cases} \theta \exp(-\theta x) & x \geq 0 \\ 0 & x < 0 \end{cases}.$$

- $\mathbb{E}[X] = 1/\theta$, $\hat{\theta}_{\text{ML}} = \left(\frac{1}{m} \sum_{i=1}^m X_i\right)^{-1}$.
- ⑤ $\frac{1}{m} \sum_{i=1}^m X_i$ is the empirical mean:
it seems plausible that $m / (\sum_{i=1}^m X_i)$ is a “good” estimate of θ .
- Why would it not be “good”?



Systems of linear equations

Definition 3.3 (Rules)

A vector space over the reals consists of a set \mathcal{V} , a vector sum operation $+ : \mathcal{V} \times \mathcal{V} \rightarrow \mathcal{V}$ and a scalar multiplication operation $\cdot : \mathbb{R} \times \mathcal{V} \rightarrow \mathcal{V}$ satisfying the following properties.

- Commutativity: $x + y = y + x$ for all $x, y \in \mathcal{V}$;
- Associativity: $(x + y) + z = x + (y + z)$ for all $x, y, z \in \mathcal{V}$;
- Identity element of addition: $x + \mathbf{0} = x$ for all $x \in \mathcal{V}$;
- Inverse element of addition: For every $x \in \mathcal{V}$, there exists an element $-x \in \mathcal{V}$ such that $x + (-x) = \mathbf{0}$;
- Associativity of scalar multiplication: For all $a, b \in \mathbb{R}$ and $x \in \mathcal{V}$, $a(bx) = (ab)x$;
- Identity element of scalar multiplication: $1x = x$ for $x \in \mathcal{V}$;
- Distributivity of scalar multiplication with respect to vector addition: $a(x + y) = ax + ay$ for all $a \in \mathbb{R}$ and $x, y \in \mathcal{V}$;
- Distributivity of scalar multiplication with respect to addition in \mathbb{R} : $(a + b)x = ax + bx$ for all $a, b \in \mathbb{R}$ and $x \in \mathcal{V}$.



Definition 3.4 (Linear independence)

A set of vectors $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$ from a vector space \mathcal{V} is linearly independent if

$$\beta_1 \mathbf{x}_1 + \beta_2 \mathbf{x}_2 + \dots + \beta_k \mathbf{x}_k = \mathbf{0} \quad (3.16)$$

implies that $\beta_1 = \beta_2 = \dots = \beta_k = 0$.

- $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$ are linearly independent \iff no vector \mathbf{x}_i can be expressed as a linear combination of the other vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_k\}$.



Matrix

Definition 3.5 (Basis of vector space)

A set of vectors $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$ is a basis for a vector space \mathcal{V} if

- $\mathcal{V} = \text{span}\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$;
- $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$ is a **linearly independent set** of vectors.

Equivalently, every $\mathbf{x} \in \mathcal{V}$ can be uniquely written as $\sum_{i=1}^k \beta_i \mathbf{x}_i$ for some $\{\beta_i\}_{i=1}^k \subset \mathbb{R}$. The number of vectors in any basis of \mathcal{V} is called the **dimension** of \mathcal{V} , written as $\dim(\mathcal{V})$.

Definition 3.6 (nullspace and range of matrix)

The **nullspace (also called kernel)** of a matrix $\mathbf{A} \in \mathbb{R}^{m \times d}$ is defined as

$$\mathcal{N}(\mathbf{A}) = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{Ax} = 0\}. \quad (3.17)$$

The range or column space of \mathbf{A} is defined as

$$\mathcal{R}(\mathbf{A}) = \{\mathbf{Ax} : \mathbf{x} \in \mathbb{R}^f\} \subset \mathbb{R}^m. \quad (3.18)$$



Definition 3.7 (Rank of matrix)

The **column rank or rank** of a matrix $\mathbf{A} \in \mathbb{R}^{m \times d}$ is defined as

$$\text{rank}(\mathbf{A}) = \dim(\mathcal{R}(\mathbf{A})). \quad (3.19)$$

In other words, the column rank of \mathbf{A} is the dimension of the column space of \mathbf{A} .



Definition 3.7 (Rank of matrix)

The **column rank or rank** of a matrix $\mathbf{A} \in \mathbb{R}^{m \times d}$ is defined as

$$\text{rank}(\mathbf{A}) = \dim(\mathcal{R}(\mathbf{A})). \quad (3.19)$$

In other words, the column rank of \mathbf{A} is the dimension of the column space of \mathbf{A} .

- ◎ **Rank-nullity:** $\text{rank}(\mathbf{A}) + \dim(\mathcal{N}(\mathbf{A})) = d$.
- It is always true that $\text{rank}(\mathbf{A}) \leq \min\{m, d\}$.
- A matrix is **full rank** if $\text{rank}(\mathbf{A}) = \min\{m, d\}$.
- A matrix is **full column rank** (resp. **full row rank**) if the set of columns (resp. rows) of the matrix is linearly independent.
- If the matrix \mathbf{A} is square (i.e., $m = d$) and it is full rank, then the inverse \mathbf{A}^{-1} exists.



Solutions to linear systems

Solve systems of equations of the form (\underline{x}_i , $1 \leq i \leq d$: columns of \mathbf{X})

$$\mathbf{X}\mathbf{w} = \mathbf{y} \text{ or } [\underline{x}_1 \ \underline{x}_2 \ \dots \ \underline{x}_d] \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{bmatrix}, \text{ where } \underline{x}_i = \begin{bmatrix} x_{1,i} \\ x_{2,i} \\ \vdots \\ x_{m,i} \end{bmatrix}. \quad (3.20)$$

Matrix $\mathbf{X} \in \mathbb{R}^{m \times d}$ and vector $\mathbf{y} \in \mathbb{R}^m$ are given and $\mathbf{w} \in \mathbb{R}^d$ is to be found.



Solutions to linear systems

Solve systems of equations of the form (\underline{x}_i , $1 \leq i \leq d$: columns of \mathbf{X})

$$\mathbf{X}\mathbf{w} = \mathbf{y} \text{ or } [\underline{x}_1 \ \underline{x}_2 \ \dots \ \underline{x}_d] \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{bmatrix}, \text{ where } \underline{x}_i = \begin{bmatrix} x_{1,i} \\ x_{2,i} \\ \vdots \\ x_{m,i} \end{bmatrix}. \quad (3.20)$$

Matrix $\mathbf{X} \in \mathbb{R}^{m \times d}$ and vector $\mathbf{y} \in \mathbb{R}^m$ are given and $\mathbf{w} \in \mathbb{R}^d$ is to be found.

- $m = d$: if matrix \mathbf{X} is square and full rank, \mathbf{X}^{-1} exists \implies solve for \mathbf{w} by simple matrix inversion and multiplication $\mathbf{w} = \mathbf{X}^{-1}\mathbf{y}$.



Solutions to linear systems

Solve systems of equations of the form (\underline{x}_i , $1 \leq i \leq d$: columns of \mathbf{X})

$$\mathbf{X}\mathbf{w} = \mathbf{y} \text{ or } [\underline{x}_1 \ \underline{x}_2 \ \dots \ \underline{x}_d] \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{bmatrix}, \text{ where } \underline{x}_i = \begin{bmatrix} x_{1,i} \\ x_{2,i} \\ \vdots \\ x_{m,i} \end{bmatrix}. \quad (3.20)$$

Matrix $\mathbf{X} \in \mathbb{R}^{m \times d}$ and vector $\mathbf{y} \in \mathbb{R}^m$ are given and $\mathbf{w} \in \mathbb{R}^d$ is to be found.

- $m = d$: if matrix \mathbf{X} is square and full rank, \mathbf{X}^{-1} exists \implies solve for \mathbf{w} by simple matrix inversion and multiplication $\mathbf{w} = \mathbf{X}^{-1}\mathbf{y}$.
- $m \neq d$: existence and uniqueness of solutions to the linear system?

Augmented matrix $\tilde{\mathbf{X}} = [\mathbf{X}\mathbf{y}] \in \mathbb{R}^{m \times (d+1)}$: $\text{rank}(\mathbf{X}) \leq \text{rank}(\tilde{\mathbf{X}})$.



Solutions to linear systems

◎ **Method:** apply Gaussian-elimination method to the augmented matrix $\tilde{\mathbf{X}}$.

Theorem 3.8 (Rouché-Capelli Theorem)

- *The system in (3.20) admits a unique solution if and only if $\text{rank}(\mathbf{X}) = \text{rank}(\tilde{\mathbf{X}}) = d$;*
- *The system in (3.20) has no solution if and only if $\text{rank}(\mathbf{X}) < \text{rank}(\tilde{\mathbf{X}})$;*
- *The system in (3.20) has infinitely many solutions if and only if $\text{rank}(\mathbf{X}) = \text{rank}(\tilde{\mathbf{X}}) < d$.*



Thanks for listening.

- Slides credit: some slides are adapted from (alphabetical order)
Haiyun He (Cornell), Tengyu Ma and Chris Re (Stanford) and Vincent Y. F. Tan (NUS).



References

- A. L. Samuel. Some studies in machine learning using the game of checkers. *IBM Journal of Research and Development*, 3(3):210–229, 1959. doi: 10.1147/rd.33.0210.

