

Probability Theory

SUZUKI, Atsushi

The whole contents

1. Random variables
2. Multiple Random Variables
3. Continuous Random Variables
4. Sample Statistics
5. Statistical Test

Outline

1. Random variables

1.1 Introduction: why do we learn random variables?

1.2 Univariate discrete random variable

1.3 Visualization of a distribution

1.4 Summary statistics for a univariate random variable

1.5 Expectation

1.6 Median

1.7 Variance and a function of a random variable

1.8 Exercises

Outline

1. Random variables

1.1 Introduction: why do we learn random variables?

1.2 Univariate discrete random variable

1.3 Visualization of a distribution

1.4 Summary statistics for a univariate random variable

1.5 Expectation

1.6 Median

1.7 Variance and a function of a random variable

1.8 Exercises

Probability theory handles random events

Probability theory handles random **events**, where the probability $\Pr(A) \in [0, 1]$ is defined for each event A . Here, an event is a subset of the **sample space**, the set of all the possible outcomes. Each element in the sample space is called an **elementary event**.

Example (Weather forecast)

Consider a weather forecast of 24 hours later. The sample space $S = \{(\text{It will be}) \text{ sunny, cloudy, rainy, snowy}\}$. Suppose that the probability for each elementary event is given by

Event A	{sunny}	{cloudy}	{rainy}	{snowy}
The probability $\Pr(A)$	0.4	0.2	0.3	0.1

If an event A includes multiple elements in the sample space S , the probability $\Pr(A)$ is given by the sum of the probabilities of those elements. For example, $\Pr(\{\text{sunny, rainy}\}) = \Pr(\{\text{sunny}\}) + \Pr(\{\text{rainy}\}) = 0.4 + 0.3 = 0.7$.

Probability theory handles random events

Probability theory handles random **events**, where the probability $\Pr(A) \in [0, 1]$ is defined for each event A . Here, an event is a subset of the **sample space**, the set of all the possible outcomes. Each element in the sample space is called an **elementary event**.

Example (Weather forecast)

Consider a weather forecast of 24 hours later. The sample space $S = \{(\text{It will be}) \text{ sunny, cloudy, rainy, snowy}\}$. Suppose that the probability for each elementary event is given by

Event A	{sunny}	{cloudy}	{rainy}	{snowy}
The probability $\Pr(A)$	0.4	0.2	0.3	0.1

In the above example, each event is a set of real phenomena, which we do not regard as a numeric value directly. However, in the following, **we always assume that each event is a set of numeric values**.

Random variable

When each elementary event is associated with a real value, then the set of those random events is called a ***random variable (RV)***.

Example (RVs in real life)

- A stock price in finance
- The remainder of one's life in medicine
- The intensity of the acoustic signal in speech recognition

Random variable

When each elementary event is associated with a real value, then the set of those random events is called a **random variable (RV)**.

Example (RVs in real life)

- A stock price in finance
- The remainder of one's life in medicine
- The intensity of the acoustic signal in speech recognition

But **WHY** do we limit the discussion to RVs only, instead of considering general random events? The reasons are the following:

- RVs, i.e., numeric random events, are **all the random events we need to handle in computer science**, including AI, since a computer can only handle numeric values.
- If random events are RVs, i.e., numeric, we can discuss their random behaviors **quantitatively** based on the RVs' numeric values.

Learning outcomes

By the end of this section, you should be able to:

- Explain the difference between random events and random variables,
- Represent the probability distribution of a random variable using the probability mass function and cumulative distribution function, and
- Describe a probability distribution using summary statistics.

Outline

1. Random variables

1.1 Introduction: why do we learn random variables?

1.2 Univariate discrete random variable

1.3 Visualization of a distribution

1.4 Summary statistics for a univariate random variable

1.5 Expectation

1.6 Median

1.7 Variance and a function of a random variable

1.8 Exercises

Discrete random variable: motivation

In general, a random variable may take all the real values.

Still, when considering applications in computer science, including artificial intelligence, we do not need to handle all the real values. Specifically, we can assume that a random variable always takes a value in a finite subset of \mathbb{R} (the set of real numbers).

¹ Nevertheless, we need to learn more general cases later even if we are interested in finite value cases only.

Discrete random variable: motivation

In general, a random variable may take all the real values.

Still, when considering applications in computer science, including artificial intelligence, we do not need to handle all the real values. Specifically, we can assume that a random variable always takes a value in a finite subset of \mathbb{R} (the set of real numbers).

This is because a computer can handle a finite number of real numbers. For example, a computer usually uses 64 bits to represent a real value. In this case, the computer can represent only $2^{64} \approx 1.84 \times 10^{19}$ real numbers.

Hence, it is good to begin with such finite cases¹.

¹ Nevertheless, we need to learn more general cases later even if we are interested in finite value cases only.

Discrete random variables

Definition

A random variable taking a value randomly in a discrete subset² of \mathbb{R} (the set of real numbers) is called a ***discrete random variable***.

The subset of \mathbb{R} in which a discrete random variable X takes a value is called the ***support*** or ***target space*** of X .

²Strictly speaking, “discrete” stands for “at most countable.” Here, we say a set is at most countable if and only if there exists a surjective map from the set of integers to the set.

Discrete random variables

Definition

A random variable taking a value randomly in a discrete subset² of \mathbb{R} (the set of real numbers) is called a **discrete random variable**.

The subset of \mathbb{R} in which a discrete random variable X takes a value is called the **support** or **target space** of X .

Example (Rolling an ideal six-sided dice)

Let X be the number that lands face-up when we roll an ideal six-sided dice. The support of X is $\{1, 2, 3, 4, 5, 6\}$. The probability of each event is given by:

x	1	2	3	4	5	6
$\Pr(X = x)$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

Rolling an ideal six-sided dice

²Strictly speaking, “discrete” stands for “at most countable.” Here, we say a set is at most countable if and only if there exists a surjective map from the set of integers to the set.

Probability mass function (PMF)

When we consider a univariate discrete random variable taking a value in a discrete set $\mathcal{X} = \{x_1, x_2, \dots\} \subset \mathbb{R}$, we can completely understand the behaviour of X by knowing the probability of X taking a value x , where $x \in \mathcal{X}$. Hence, we define a function describing those probabilities.

Definition (probability mass function (PMF))

Let X be a discrete random variable taking a value in a discrete set $\mathcal{X} \subset \mathbb{R}$. We define the **probability mass function (PMF)** $P_X : \mathcal{X} \rightarrow [0, 1]$ of the random variable X by

$$P_X(x) := \Pr(X = x). \quad (1)$$

The relation between the value that a RV takes and its probability is called the **distribution** of the RV. The PMF is the most fundamental way to represent the distribution of a discrete RV.

Properties of a PMF

A PMF must satisfy the following:

- **(Nonnegativity)** $P_X(x) \geq 0$ for all $x \in \mathcal{X}$.
- **(The sum)** $\sum_{x \in \mathcal{X}} P_X(x) = 1$.

PMF tells us all we want to know.

If we want to know, for example, $\Pr(a \leq X \leq b)$, we can find it by the PMF:

$$\Pr(a \leq X \leq b) = \sum_{a \leq x \leq b} P_X(x). \quad (2)$$

Example (Rolling an ideal six-sided dice)

Let X be the number that lands face-up when we roll an ideal six-sided dice. The support of X is $\{1, 2, 3, 4, 5, 6\}$. The PMF of each event is given by:

x	1	2	3	4	5	6
$P_X(x) := \Pr(X = x)$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

Rolling an ideal six-sided dice

Here, $\Pr(2 \leq x \leq 4)$ is given by

$$\sum_{2 \leq x \leq 4} P_X(x) = P_X(2) + P_X(3) + P_X(4) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{1}{2}.$$

A frequency is a discrete random variable

The probability theory can handle data points by considering its **frequency**. This is the first step of **data science**.

Suppose that we have m data points taking values in \mathbb{R} . For the probability theory to handle the data points, we need to construct a random variable.

Specifically, we sample a data point uniform-randomly. Then, the value of the sampled data point is a discrete random variable.

The probability distribution of the random variable constructed from the data points this way is called the **frequency** or **empirical distribution**.

Example of frequency

Example (Exam results)

Suppose that we have $m = 20$ students and consider their results in an exam. For $x \in \mathcal{X} = \{0, 1, 2, 3, 4, 5\}$, we denote the number of the students who got a score x by m_x . Let X be the score of the student sampled uniform-randomly from the 20 students. The probability $\Pr(X = x)$ equals to $\frac{m_x}{m}$. For example,

Score x	0	1	2	3	4	5
# students m_x	3	2	3	5	6	1
$P_X(x) := \Pr(X = x) = \frac{m_x}{m}$	0.15	0.10	0.15	0.25	0.30	0.05

Exam result data points and the frequency.

Outline

1. Random variables

1.1 Introduction: why do we learn random variables?

1.2 Univariate discrete random variable

1.3 Visualization of a distribution

1.4 Summary statistics for a univariate random variable

1.5 Expectation

1.6 Median

1.7 Variance and a function of a random variable

1.8 Exercises

How to visualize a distribution?

If a distribution is complicated, then you might want to understand it from a figure, not from a long table.

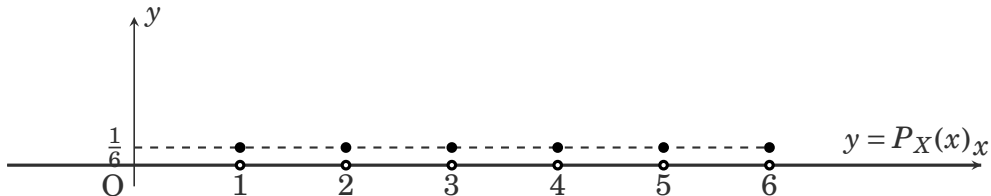
One way is to draw a graph of the PMF.

Example of a PMF graph: rolling an ideal dice

Suppose that we roll an ideal six-sided dice. The PMF is given as follows.

x	1	2	3	4	5	6
$P_X(x) := \Pr(X = x)$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

The PMF of rolling an ideal six-sided dice

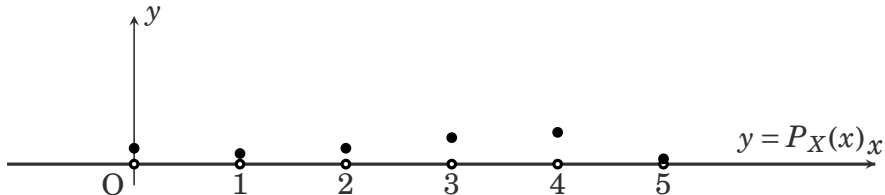


Example of a PMF graph: rolling an ideal dice

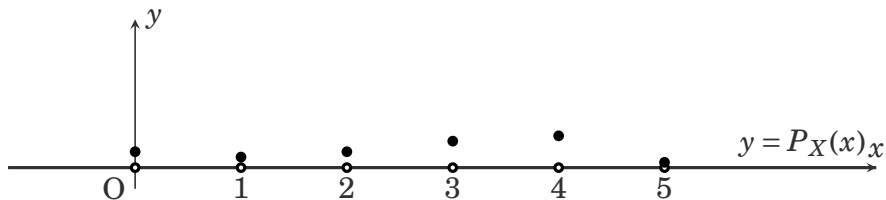
Suppose that we roll an ideal six-sided dice. The PMF is given as follows.

Score x	0	1	2	3	4	5
$P_X(x) := \Pr(X = x)$	0.15	0.10	0.15	0.25	0.30	0.05

The PMF of the frequency of exam results



Pros and cons of the PMF graph



Pros: From the PMF graph, we can easily see which value the RV takes more and less frequently.

Cons: A PMF is not suitable to calculate the probability of a RV taking a value in a certain range, e.g., $\Pr(1.5 \leq X \leq 3.8)$.

Cumulative distribution function (CDF)

Any random variable has a ***cumulative distribution function (CDF)*** defined as follows.

Definition

Let X be a random variable. The ***cumulative distribution function (CDF)*** $F_X : \mathbb{R} \rightarrow [0, 1]$ of X is defined by

$$F_X(x) := \Pr(X \leq x). \quad (3)$$

The CDF gives formulae to evaluate a section's probability

In the following, let $a, b \in \mathbb{R}$ and $a < b$.

We have that $\Pr(X < a) = \lim_{x \nearrow a} F_X(x)$, where the right hand side is the left limit of F_X at a , given by evaluating $F_X(x - \epsilon)$ while diminishing ϵ to a positive value infinitely close to zero.

Using the above fact, we can calculate the probability of a random variable taking a value in a section using the CDF as follows.

Theorem

- $\Pr(a \leq X \leq b) = \Pr(X \leq b) - \Pr(X < a) = F_X(b) - \lim_{x \nearrow a} F_X(x).$
- $\Pr(a < X < b) = \Pr(X < b) - \Pr(X \leq a) = \lim_{x \nearrow b} F_X(x) - F_X(a).$
- $\Pr(a < X \leq b) = \Pr(X \leq b) - \Pr(X \leq a) = F_X(b) - F_X(a).$
- $\Pr(a \leq X < b) = \Pr(X < b) - \Pr(X < a) = \lim_{x \nearrow b} F_X(x) - \lim_{x \nearrow a} F_X(x).$

Example of CDF: rolling an ideal dice

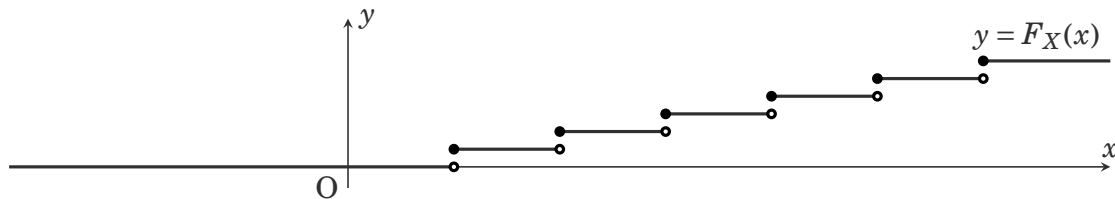
Suppose that we roll an ideal six-sided dice. The PMF is given as follows.

x	1	2	3	4	5	6
$P_X(x) := \Pr(X = x)$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

The PMF of rolling an ideal six-sided dice

Example of CDF: rolling an ideal dice

Suppose that we roll an ideal six-sided dice. The CDF is given as follows.

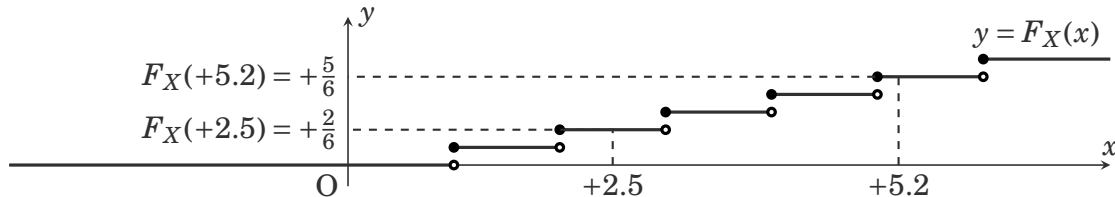


x	$(-\infty, 1)$	$[1, 2)$	$[2, 3)$	$[3, 4)$	$[4, 5)$	$[5, 6)$	$[6, +\infty)$
$F_X(x) := \Pr(X = x)$	0	$\frac{1}{6}$	$\frac{2}{6}$	$\frac{3}{6}$	$\frac{4}{6}$	$\frac{5}{6}$	1

The CDF of rolling an ideal six-sided dice

Example of CDF: rolling an ideal dice

Suppose that we roll an ideal six-sided dice. The CDF is given as follows.

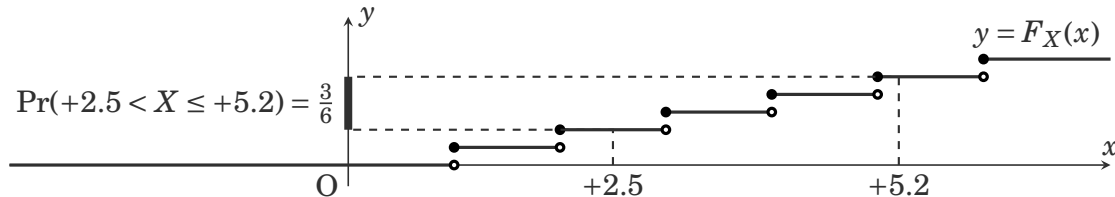


Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(+2.5 < X \leq +5.2) &= \Pr(X \leq +5.2) - \Pr(X \leq +2.5) \\ &= F_X(+5.2) - F_X(+2.5) \\ &= \frac{5}{6} - \frac{2}{6} = \frac{3}{6}.\end{aligned}\tag{4}$$

Example of CDF: rolling an ideal dice

Suppose that we roll an ideal six-sided dice. The CDF is given as follows.

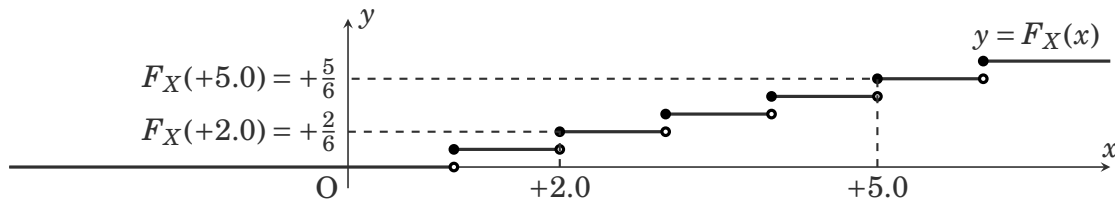


Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(+2.5 < X \leq +5.2) &= \Pr(X \leq +5.2) - \Pr(X \leq +2.5) \\ &= F_X(+5.2) - F_X(+2.5) \\ &= \frac{5}{6} - \frac{2}{6} = \frac{3}{6}.\end{aligned}\tag{4}$$

Example of CDF: rolling an ideal dice

Suppose that we roll an ideal six-sided dice. The CDF is given as follows.

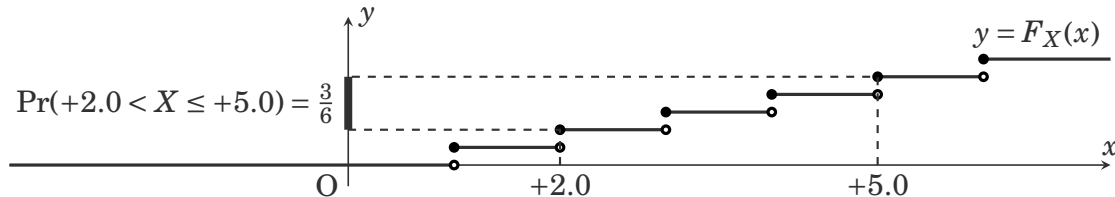


Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(+2.0 < X \leq +5.0) &= \Pr(X \leq +5.0) - \Pr(X \leq +2.0) \\ &= F_X(+5.0) - F_X(+2.0) \\ &= \frac{5}{6} - \frac{2}{6} = \frac{3}{6}.\end{aligned}\tag{4}$$

Example of CDF: rolling an ideal dice

Suppose that we roll an ideal six-sided dice. The CDF is given as follows.

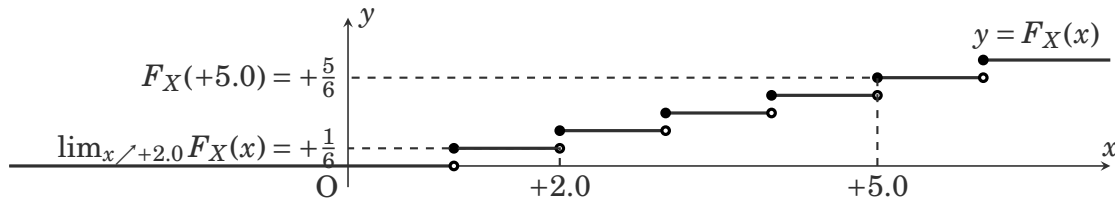


Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(+2.0 < X \leq +5.0) &= \Pr(X \leq +5.0) - \Pr(X \leq +2.0) \\ &= F_X(+5.0) - F_X(+2.0) \\ &= \frac{5}{6} - \frac{2}{6} = \frac{3}{6}.\end{aligned}\tag{4}$$

Example of CDF: rolling an ideal dice

Suppose that we roll an ideal six-sided dice. The CDF is given as follows.

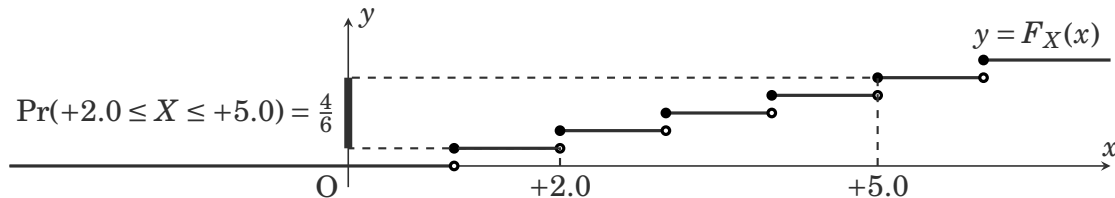


Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(+2.0 \leq X \leq +5.0) &= \Pr(X \leq +5.0) - \Pr(X < +2.0) \\ &= F_X(+5.0) - \lim_{x \nearrow +2.0} F_X(x) \\ &= \frac{5}{6} - \frac{1}{6} = \frac{4}{6}.\end{aligned}\tag{4}$$

Example of CDF: rolling an ideal dice

Suppose that we roll an ideal six-sided dice. The CDF is given as follows.

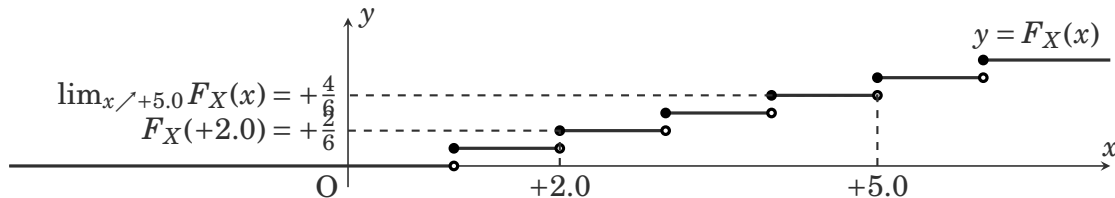


Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(+2.0 \leq X \leq +5.0) &= \Pr(X \leq +5.0) - \Pr(X < +2.0) \\ &= F_X(+5.0) - \lim_{x \nearrow +2.0} F_X(x) \\ &= \frac{5}{6} - \frac{1}{6} = \frac{4}{6}.\end{aligned}\tag{4}$$

Example of CDF: rolling an ideal dice

Suppose that we roll an ideal six-sided dice. The CDF is given as follows.

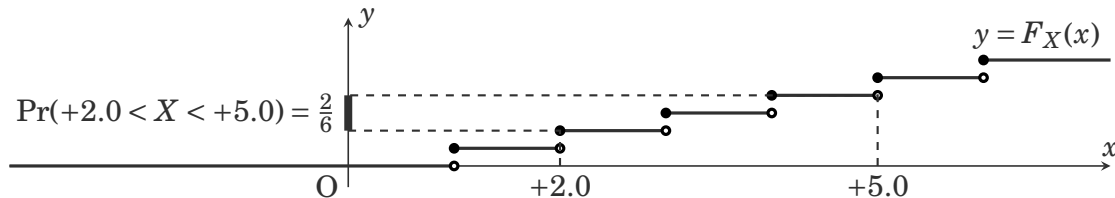


Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(+2.0 < X < +5.0) &= \Pr(X < +5.0) - \Pr(X \leq +2.0) \\ &= \lim_{x \nearrow +5.0} F_X(x) - F_X(+2.0) \\ &= \frac{4}{6} - \frac{2}{6} = \frac{2}{6}.\end{aligned}\tag{4}$$

Example of CDF: rolling an ideal dice

Suppose that we roll an ideal six-sided dice. The CDF is given as follows.

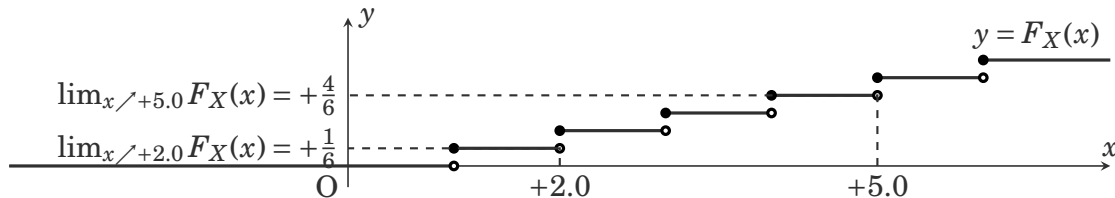


Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(+2.0 < X < +5.0) &= \Pr(X < +5.0) - \Pr(X \leq +2.0) \\ &= \lim_{x \nearrow +5.0} F_X(x) - F_X(+2.0) \\ &= \frac{4}{6} - \frac{2}{6} = \frac{2}{6}.\end{aligned}\tag{4}$$

Example of CDF: rolling an ideal dice

Suppose that we roll an ideal six-sided dice. The CDF is given as follows.

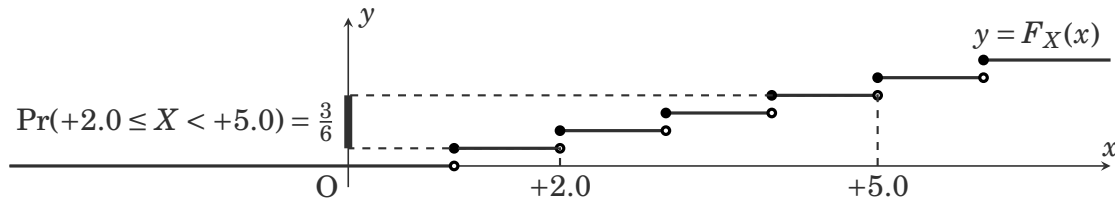


Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(+2.0 \leq X < +5.0) &= \Pr(X < +5.0) - \Pr(X < +2.0) \\ &= \lim_{x \nearrow +5.0} F_X(x) - \lim_{x \nearrow +2.0} F_X(x) \\ &= \frac{4}{6} - \frac{1}{6} = \frac{3}{6}.\end{aligned}\tag{4}$$

Example of CDF: rolling an ideal dice

Suppose that we roll an ideal six-sided dice. The CDF is given as follows.



Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(+2.0 \leq X < +5.0) &= \Pr(X < +5.0) - \Pr(X < +2.0) \\ &= \lim_{x \nearrow +5.0} F_X(x) - \lim_{x \nearrow +2.0} F_X(x) \\ &= \frac{4}{6} - \frac{1}{6} = \frac{3}{6}.\end{aligned}\tag{4}$$

Example of CDF: student score frequency

Suppose that X is a random variable whose PMF is given as follows.

x	0	1	2	3	4	5
$P_X(x) := \Pr(X = x)$	0.15	0.10	0.15	0.25	0.30	0.05

The PMF of a student exam result frequency

The CDF is given as the cumulative sum of the PMF, as follows.

x	$(-\infty, 0)$	$[0, 1)$	$[1, 2)$	$[2, 3)$	$[3, 4)$	$[4, 5)$	$[5, +\infty)$
$F_X(x) := \Pr(X \leq x)$	0.00	0.15	0.25	0.40	0.65	0.95	1.00

The CDF of a student exam result frequency

Example of CDF: student score frequency

The CDF is given as the cumulative sum of the PMF, as follows.

x	$(-\infty, 0)$	$[0, 1)$	$[1, 2)$	$[2, 3)$	$[3, 4)$	$[4, 5)$	$[5, +\infty)$
$F_X(x) := \Pr(X = x)$	0.00	0.15	0.25	0.40	0.65	0.95	1.00

The CDF of a student exam result frequency

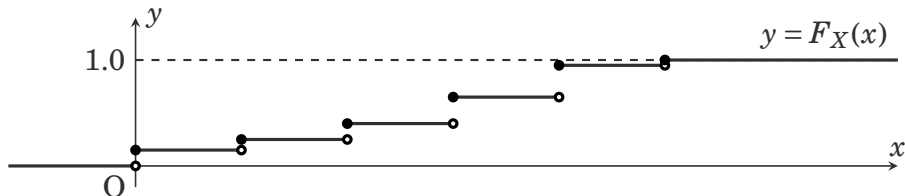
Example of CDF: student score frequency

The CDF is given as the cumulative sum of the PMF, as follows.

x	$(-\infty, 0)$	$[0, 1)$	$[1, 2)$	$[2, 3)$	$[3, 4)$	$[4, 5)$	$[5, +\infty)$
$F_X(x) := \Pr(X = x)$	0.00	0.15	0.25	0.40	0.65	0.95	1.00

The CDF of a student exam result frequency

The graph of the CDF is as follows.



Properties of CDF

For any random variable X , its CDF F_X satisfies

- $\lim_{x \rightarrow -\infty} F_X(x) = 0$.
- $\lim_{x \rightarrow +\infty} F_X(x) = 1$.
- The CDF is everywhere right-continuous, i.e., $\lim_{x \searrow x_0} F_X(x) = F_X(x_0)$ for all $x_0 \in \mathbb{R}$.
- The CDF has its left-limit $\lim_{x \nearrow x_0} F_X(x)$ for all $x_0 \in \mathbb{R}$.

Appendix: the definition of the left limit

Definition (left/right limit/continuous)

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a real function and a be a real value.

Suppose that for all $\delta > 0$ there exists $\epsilon > 0$ such that $|f(a - \epsilon') - c| < \delta$ for all ϵ' that satisfies $0 < \epsilon' < \epsilon$.

Then the value c is called the **left limit** of a function f at $a \in \mathbb{R}$ and denoted by $\lim_{x \nearrow a} f(x)$.

We have the definition of the **right limit** by replacing $(a - \epsilon')$ with $(a + \epsilon')$ in the definition of the left limit.

The right limit is denoted by $\lim_{x \searrow a} f(x)$.

A function f is called **left continuous** at $a \in \mathbb{R}$ if $\lim_{x \nearrow a} f(x) = f(a)$ and **right continuous** at $a \in \mathbb{R}$ if $\lim_{x \searrow a} f(x) = f(a)$.

If a function is left/right continuous at every value in its domain, then we simply call the function left/right continuous.

Appendix: relation between the limit and the left and right limits

Theorem

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a real function and a be a real value.

- $\lim_{x \rightarrow a} f(x) = c$ if and only if $\lim_{x \nearrow a} f(x) = \lim_{x \searrow a} f(x) = c$.
- f is continuous at a if and only if f is left continuous and right continuous at a .

Outline

1. Random variables

1.1 Introduction: why do we learn random variables?

1.2 Univariate discrete random variable

1.3 Visualization of a distribution

1.4 Summary statistics for a univariate random variable

1.5 Expectation

1.6 Median

1.7 Variance and a function of a random variable

1.8 Exercises

Summary statistics

Motivation: A probability mass function might have too much information to understand the behaviour of a random variable intuitively.

Hence, we often want to calculate a single value (or a few values) that describes a distribution, called a ***descriptive statistic*** or ***summary statistic***³.

³These words are often used to distinguish them from inferential statistics.

Summary statistics: examples

Central tendency measures give a representative value of the values that the random variable takes, e.g., ***expectation, median, mode***, etc.

Variability measures show how spread values the random variable takes, e.g., ***range, variance, standard deviation, quartile deviation***.

Other measures e.g., kurtosis, skewness.

Outline

1. Random variables

1.1 Introduction: why do we learn random variables?

1.2 Univariate discrete random variable

1.3 Visualization of a distribution

1.4 Summary statistics for a univariate random variable

1.5 Expectation

1.6 Median

1.7 Variance and a function of a random variable

1.8 Exercises

Definition of expectation (mean)

The most fundamental central tendency measure of a distribution is the **expectation**.

Definition (Expectation of a discrete RV)

The **expectation** of a discrete random variable X , denoted by $\mathbb{E}X$, $\mathbf{E}X$, $\langle X \rangle$, or \overline{X} , is the weighted mean of the values with the probability masses as weights. That is

$$\mathbb{E}X := \sum_{x \in \mathcal{X}} x P_X(x). \quad (5)$$

The expectation is also called the **mean**. Indeed, if the probability distribution is a frequency of data points, the expectation is nothing but the mean of the data points.

Example of expectation calculation

Suppose that X is a random variable whose PMF P_X is given by the following table.

x	-2	-1	0	+1	+2
$P_X(x)$	0.05	0.10	0.20	0.10	0.55

Example random function and its PMF.

We can calculate the expectation $\mathbb{E}X$ by the following procedure.

- **Step 1:**
- **Step 2:**

Example of expectation calculation

Suppose that X is a random variable whose PMF P_X is given by the following table.

x	-2	-1	0	+1	+2
$P_X(x)$	0.05	0.10	0.20	0.10	0.55
$xP_X(x)$					

Example random function and its PMF.

We can calculate the expectation $\mathbb{E}X$ by the following procedure.

- **Step 1:** Calculate $xP_X(x)$ for each $x \in \mathcal{X}$.
- **Step 2:**

Example of expectation calculation

Suppose that X is a random variable whose PMF P_X is given by the following table.

x	-2	-1	0	+1	+2
$P_X(x)$	0.05	0.10	0.20	0.10	0.55
$xP_X(x)$	-0.10				

Example random function and its PMF.

We can calculate the expectation $\mathbb{E}X$ by the following procedure.

- **Step 1:** Calculate $xP_X(x)$ for each $x \in \mathcal{X}$.
- **Step 2:**

Example of expectation calculation

Suppose that X is a random variable whose PMF P_X is given by the following table.

x	-2	-1	0	+1	+2
$P_X(x)$	0.05	0.10	0.20	0.10	0.55
$xP_X(x)$	-0.10	-0.10			

Example random function and its PMF.

We can calculate the expectation $\mathbb{E}X$ by the following procedure.

- **Step 1:** Calculate $xP_X(x)$ for each $x \in \mathcal{X}$.
- **Step 2:**

Example of expectation calculation

Suppose that X is a random variable whose PMF P_X is given by the following table.

x	-2	-1	0	+1	+2
$P_X(x)$	0.05	0.10	0.20	0.10	0.55
$xP_X(x)$	-0.10	-0.10	0.00		

Example random function and its PMF.

We can calculate the expectation $\mathbb{E}X$ by the following procedure.

- **Step 1:** Calculate $xP_X(x)$ for each $x \in \mathcal{X}$.
- **Step 2:**

Example of expectation calculation

Suppose that X is a random variable whose PMF P_X is given by the following table.

x	-2	-1	0	+1	+2
$P_X(x)$	0.05	0.10	0.20	0.10	0.55
$xP_X(x)$	-0.10	-0.10	0.00	+0.10	+1.10

Example random function and its PMF.

We can calculate the expectation $\mathbb{E}X$ by the following procedure.

- **Step 1:** Calculate $xP_X(x)$ for each $x \in \mathcal{X}$.
- **Step 2:**

Example of expectation calculation

Suppose that X is a random variable whose PMF P_X is given by the following table.

x	-2	-1	0	+1	+2
$P_X(x)$	0.05	0.10	0.20	0.10	0.55
$xP_X(x)$	-0.10	-0.10	0.00	+0.10	+1.10

Example random function and its PMF.

We can calculate the expectation $\mathbb{E}X$ by the following procedure.

- **Step 1:** Calculate $xP_X(x)$ for each $x \in \mathcal{X}$.
- **Step 2:** Evaluate the sum $\sum_{x \in \mathcal{X}} xP_X(x)$, which equals the expectation $\mathbb{E}X$.

Example of expectation calculation

Suppose that X is a random variable whose PMF P_X is given by the following table.

x	-2	-1	0	+1	+2
$P_X(x)$	0.05	0.10	0.20	0.10	0.55
$xP_X(x)$	-0.10	-0.10	0.00	+0.10	+1.10

Example random function and its PMF.

We can calculate the expectation $\mathbb{E}X$ by the following procedure.

- **Step 1:** Calculate $xP_X(x)$ for each $x \in \mathcal{X}$.
- **Step 2:** Evaluate the sum $\sum_{x \in \mathcal{X}} xP_X(x)$, which equals the expectation $\mathbb{E}X$.
In the above case, the expectation $\mathbb{E}X$ is given by
 $\mathbb{E}X = (-0.10) + (-0.10) + 0.00 + 0.10 + 1.10 = 1.00$.

Expectation of a function

If X is a random variable and f is a function, $f(X)$ is again a random variable. Hence, we can define the expectation of $f(X)$.

The expectation $\mathbb{E} f(X)$ often gives us important information as well as the original expectation $\mathbb{E} X$. The most important example is the **variance** of a random variable, which is the most frequently used variability measure.

The expectation is easily “warped” by outliers.

If a distribution takes some extremely large or small values, called **outliers**, the expectation is significantly influenced by the probability of the random variable taking such values.

Example (Imbalanced score distribution)

Suppose you got a score of 99 in an exam where 100 students participated and the expectation was 98, you might feel you did very well.

However, it might be just that one student who got a score of 1 decreased the expectation significantly, as follows.

Score x	1	99	100
# students m_x	1	1	98
$P_X := \Pr(X = x) = \frac{m_x}{m}$	0.01	0.01	0.98

Exam result data points and the frequency.

Outline

1. Random variables

1.1 Introduction: why do we learn random variables?

1.2 Univariate discrete random variable

1.3 Visualization of a distribution

1.4 Summary statistics for a univariate random variable

1.5 Expectation

1.6 Median

1.7 Variance and a function of a random variable

1.8 Exercises

Median's idea

If a random value takes an extremely large or small value in a small probability, some might want to use the ***median*** as a summary statistic.

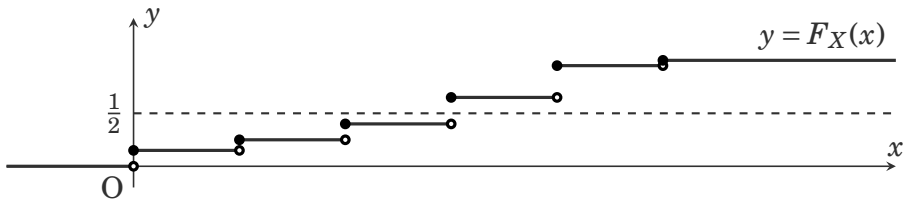
Roughly speaking, the median is defined so that the random variable is larger than the median in 50% probability and smaller than the median in 50% probability.

Median's idea

If a random value takes an extremely large or small value in a small probability, some might want to use the **median** as a summary statistic.

Roughly speaking, the median is defined so that the random variable is larger than the median in 50% probability and smaller than the median in 50% probability.

In other words, the median is the value x such that the graph $y = F_X(x)$ of the CDF crosses the horizontal line $y = \frac{1}{2}$.

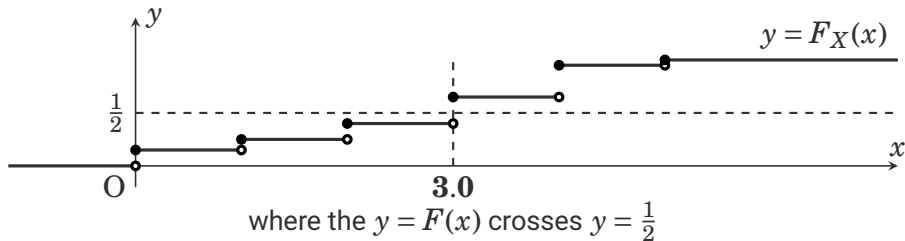


Median's idea

If a random value takes an extremely large or small value in a small probability, some might want to use the **median** as a summary statistic.

Roughly speaking, the median is defined so that the random variable is larger than the median in 50% probability and smaller than the median in 50% probability.

In other words, the median is the value x such that the graph $y = F_X(x)$ of the CDF crosses the horizontal line $y = \frac{1}{2}$.

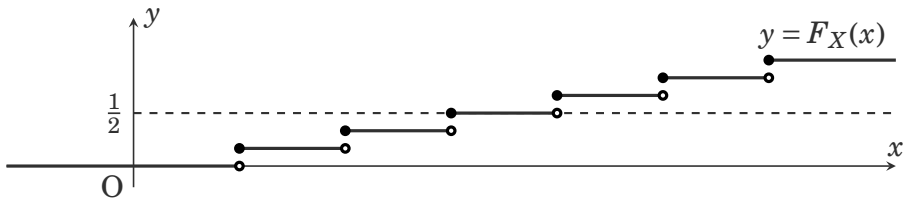


Median's idea

If a random value takes an extremely large or small value in a small probability, some might want to use the **median** as a summary statistic.

Roughly speaking, the median is defined so that the random variable is larger than the median in 50% probability and smaller than the median in 50% probability.

If the CDF graph has a horizontal segment on $y = \frac{1}{2}$, the median is the middle point of the segment.

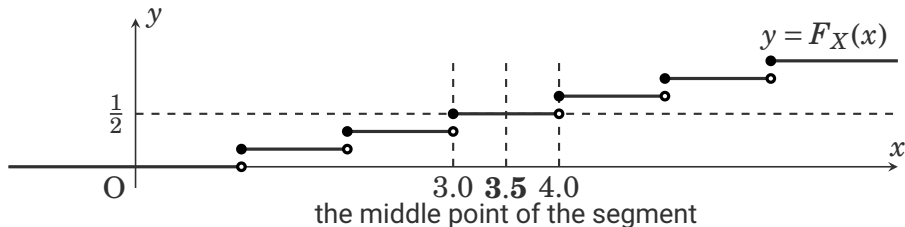


Median's idea

If a random value takes an extremely large or small value in a small probability, some might want to use the **median** as a summary statistic.

Roughly speaking, the median is defined so that the random variable is larger than the median in 50% probability and smaller than the median in 50% probability.

If the CDF graph has a horizontal segment on $y = \frac{1}{2}$, the median is the middle point of the segment.



Definition of median

Definition (The definition of the median)

Let $P : \mathbb{R} \rightarrow [0, 1]$ be the probability mass function of a univariate discrete random variable X . If a real value $M \in \mathbb{R}$ satisfies the following equation, then M is called a **median** of the distribution of X :

$$\Pr(X \leq M) \geq \frac{1}{2} \text{ and } \Pr(X \geq M) \geq \frac{1}{2}. \quad (6)$$

We can often see the above definition in the context of probability theory.

The definition of the median

Definition (The definition of the median)

Let $P : \mathbb{R} \rightarrow [0, 1]$ be the probability mass function of a univariate discrete random variable X . Define the values \underline{M} and \overline{M} by If a real value $M \in \mathbb{R}$ satisfies the following equation, then M is called a **median** of the distribution of X :

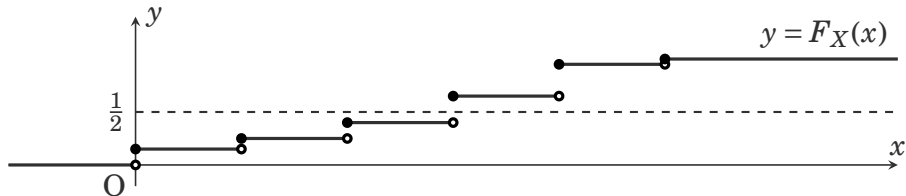
$$\begin{aligned}\underline{M} &:= \min \left\{ M \in \mathbb{R} \left| \Pr(X \leq M) \geq \frac{1}{2} \text{ and } \Pr(X \geq M) \geq \frac{1}{2} \right. \right\}, \\ \overline{M} &:= \max \left\{ M \in \mathbb{R} \left| \Pr(X \leq M) \geq \frac{1}{2} \text{ and } \Pr(X \geq M) \geq \frac{1}{2} \right. \right\}.\end{aligned}\tag{7}$$

The **median** M is defined as the midpoint of \underline{M} and \overline{M} , i.e., $M := \frac{\underline{M} + \overline{M}}{2}$.

The above definition looks complicated, but it is in fact easy if we see the CDF graph.

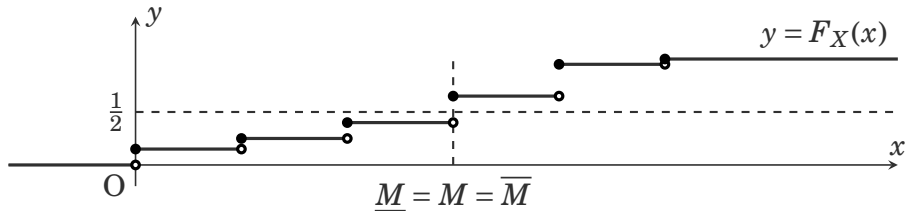
The definition of the median by the CDF graph

If the CDF graph “crosses” the graph of $y = \frac{1}{2}$,



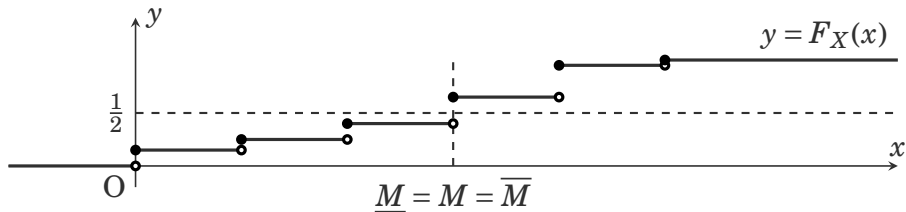
The definition of the median by the CDF graph

If the CDF graph “crosses” the graph of $y = \frac{1}{2}$,

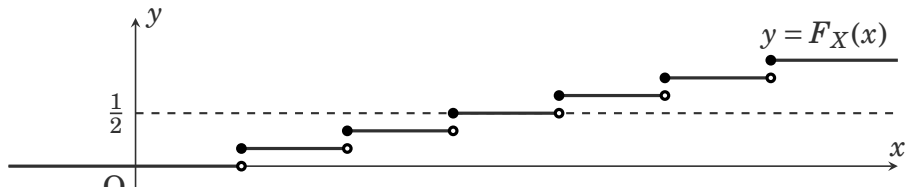


The definition of the median by the CDF graph

If the CDF graph “crosses” the graph of $y = \frac{1}{2}$,

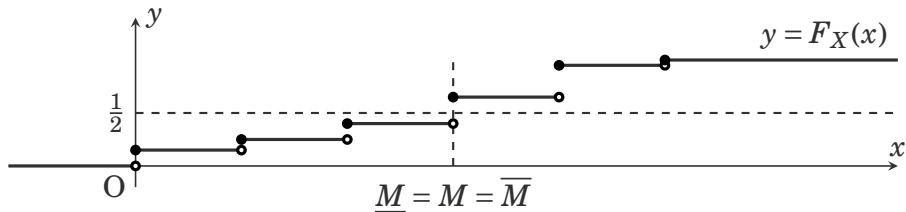


If the CDF graph has a horizontal segment on $y = \frac{1}{2}$,

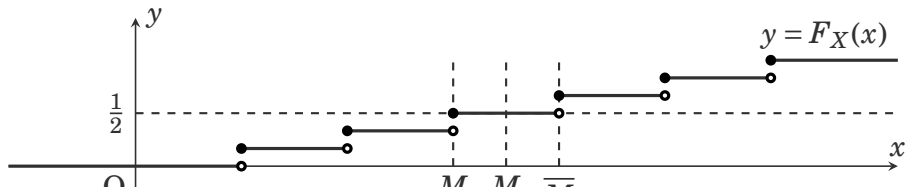


The definition of the median by the CDF graph

If the CDF graph “crosses” the graph of $y = \frac{1}{2}$,



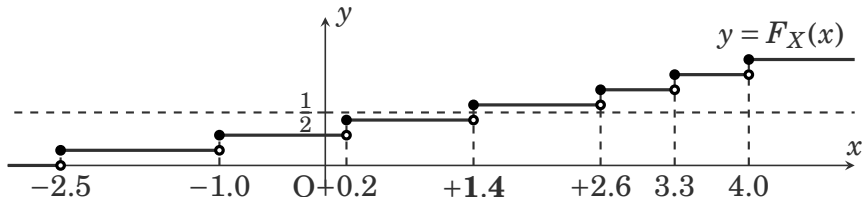
If the CDF graph has a horizontal segment on $y = \frac{1}{2}$,



Median of frequency for an odd data point case

By definition, the median of the frequency of $(2k + 1)$ data points is the value of the $(k + 1)$ th largest data point. This is equivalent to the $(k + 1)$ th smallest data point. In this sense, the definition is symmetric. The value is simply called the median of the data points.

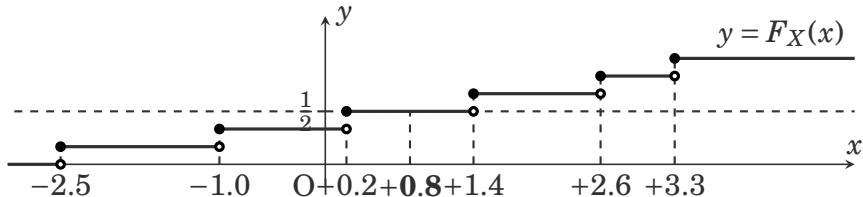
For example, if we have 7 sorted data points $(-2.5, -1.0, +0.2, +1.4, +2.6, +3.3, +4.0)$, then the median is the value of the 4th largest (or equivalently, the 4th smallest) data point, which is $+1.4$.



Median of frequency for an even data point case

By definition, the median of the frequency of $2k$ datapoints is the middle point of the values of the k th and $k + 1$ th largest data points. This is equivalent to the middle point of the values of the k th and $k + 1$ th largest data points. In this sense, the definition is symmetric. The value is simply called the median of the data points.

For example, if we have 6 sorted data points $(-2.5, -1.0, 0.2, +1.4, +2.6, +3.5)$, then the median is the middle point of the values of the 3rd and 4th largest (or equivalently, the 3rd and 4th smallest) data points, which is $\frac{0.2+1.4}{2} = 0.8$.



Median of imbalanced data

Example

Consider the following exam results of 100 participants given by the following table and the frequency of the data points.

Score x	1	99	100
# students m_x	1	1	98
$P_X := \Pr(X = x) = \frac{m_x}{m}$	0.01	0.01	0.98

Exam result data points and the frequency.

Since we have 100 students, which is an even number, the median is the middle point of the 50th-best student's score and the 51th-best student's score, which is 100.

Median tends to ignore “minor” data points

It is not that the median is a perfect statistic. Indeed, the median tends to ignore a relatively minor cohort even though the size of the cohort is not ignorable.

Example

Consider the following exam results of 100 participants given by the following table and the frequency of the data points.

Score x	0	100
# students m_x	49	51
$P_X := \Pr(X = x) = \frac{m_x}{m}$	0.49	0.51

Exam result data points and the frequency.

Then, the median is the middle point of the 50th-best student's score and the 51th-best student's score, which is 100. However, this median ignores the 49%, who received zero scores.

Outline

1. Random variables

1.1 Introduction: why do we learn random variables?

1.2 Univariate discrete random variable

1.3 Visualization of a distribution

1.4 Summary statistics for a univariate random variable

1.5 Expectation

1.6 Median

1.7 Variance and a function of a random variable

1.8 Exercises

The basic idea of variance as a variability measure

Variability measures show how much the random variable deviates from the “center”.

The most representative one is the **variance**, defined based on the **square deviation**.

Let X be a random variable and μ be its expectation. The **square deviation** of X is defined as $(X - \mu)^2$. If X is far (whether large or not) from μ , the square deviation $(X - \mu)^2$ is large.

Hence, we expect to create a variability measure using $(X - \mu)^2$.

But, what is $(X - \mu)^2$?

The basic idea of variance as a variability measure

Variability measures show how much the random variable deviates from the “center”.

The most representative one is the **variance**, defined based on the **square deviation**.

Let X be a random variable and μ be its expectation. The **square deviation** of X is defined as $(X - \mu)^2$. If X is far (whether large or not) from μ , the square deviation $(X - \mu)^2$ is large.

Hence, we expect to create a variability measure using $(X - \mu)^2$.

But, what is $(X - \mu)^2$? Since it depends on the value of X , $(X - \mu)^2$ is (the output value of) a function of X , and since X is a random variable, $(X - \mu)^2$ is **also a random variable**!

The **variance** is nothing but the expectation of the RV $(X - \mu)^2$. To understand this amount, let's discuss the function of random variables in general.

A function of a random variable

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a function and X be a random variable.

If we input X to f , the return value $f(X)$ is also a random variable.

In particular, if X is a discrete RV, then $f(X)$ is also a discrete RV. Specifically, if the support of X is \mathcal{X} , then the support of $f(X)$ is $\{f(x) | x \in \mathcal{X}\}$.

Let's find its PMF $P_{f(X)}$.

Example of a function of a RV

Define f by $f(x) = x^2$, and suppose the PMF P_X is given by the following table.

x	-1	0	+1
$P_X(x)$	0.2	0.3	0.5

Example random function and its PMF.

Suppose that we are interested in the behavior of $f(X)$. The variable $f(X)$ is also a random variable since it depends on the random behavior of the RV X .

Now, what are the support, the PMF, and the expectation of $f(X) = X^2$?

Example of a function of a RV

Define f by $f(x) = x^2$, and suppose the PMF P_X is given by the following table.

x	-1	0	+1
$P_X(x)$	0.2	0.3	0.5

Example random function and its PMF.

Let's find the **support** of $f(X) = X^2$. The RV X takes a value in $\mathcal{X} = \{-1, 0, +1\}$. Since $f(-1) = (-1)^2 = +1$, $f(0) = (0)^2 = 0$, and $f(+1) = (+1)^2 = +1$, The RV $f(X) = X^2$ only takes a value 0 or +1 only. Hence the support of $f(X) = X^2$ is $\{0, +1\}$. In particular, $f(X)$ is also a discrete RV.

Example of a function of a RV

Define f by $f(x) = x^2$, and suppose the PMF P_X is given by the following table.

x	-1	0	+1
$P_X(x)$	0.2	0.3	0.5

Example random function and its PMF.

Let's find the **PMF** P_{X^2} .

By definition $P_{X^2}(0) = \Pr(X^2 = 0)$.

Since $X^2 = 0 \Leftrightarrow X = 0$ holds,⁴ we have that $\Pr(X^2 = 0) = \Pr(X = 0) = 0.3$.

This case is easy since only one value of X corresponds to $X^2 = 0$.

⁴The symbol \Leftrightarrow indicates a necessary and sufficient condition, or equivalence.

Example of a function of a RV

Define f by $f(x) = x^2$, and suppose the PMF P_X is given by the following table.

x	-1	0	+1
$P_X(x)$	0.2	0.3	0.5

Example random function and its PMF.

Let's find the **PMF** P_{X^2} .

By definition $P_{X^2}(1) = \Pr(X^2 = 1)$.

Since " $X^2 = 1$ " \Leftrightarrow " $X = -1$ or $X = +1$ " holds, we have that

$$\begin{aligned}\Pr(X^2 = 1) &= \Pr("X = -1 \text{ or } X = +1") \\ &= \Pr(X = -1) + \Pr(X = +1) = 0.2 + 0.5 = 0.7.\end{aligned}\tag{8}$$

Here, the second equation comes from the sum law since " $X = -1$ and $X = +1$ " do not happen at the same time.

Example of a function of a RV

Define f by $f(x) = x^2$, and suppose the PMF P_X is given by the following table.

x	-1	0	+1
$P_X(x)$	0.2	0.3	0.5

Example random function and its PMF.

To wrap up,

y	0	+1
$P_{X^2}(y)$	0.3	0.7

The PMF of X^2 .

Example of a function of a RV

Define f by $f(x) = x^2$, and suppose the PMF P_X is given by the following table.

x	-1	0	+1
$P_X(x)$	0.2	0.3	0.5

Example random function and its PMF.

Let's evaluate the **expectation** $\mathbb{E} f(X) = \mathbb{E} X^2$.

Example of a function of a RV

Define f by $f(x) = x^2$, and suppose the PMF P_X is given by the following table.

x	-1	0	+1
$P_X(x)$	0.2	0.3	0.5

Example random function and its PMF.

Let's evaluate the **expectation** $\mathbb{E} f(X) = \mathbb{E} X^2$. If we use the PMF of X^2 , it looks like

$$\begin{aligned}\mathbb{E} X^2 &= 0 \cdot P_{X^2}(0) + (+1) \cdot P_{X^2}(+1) \\ &= 0 \cdot 0.3 + (+1) \cdot 0.7.\end{aligned}\tag{8}$$

Example of a function of a RV

Define f by $f(x) = x^2$, and suppose the PMF P_X is given by the following table.

x	-1	0	+1
$P_X(x)$	0.2	0.3	0.5

Example random function and its PMF.

Let's evaluate the **expectation** $\mathbb{E} f(X) = \mathbb{E} X^2$. Since $P_{X^2}(0)$ equals $P_X(0)$ and $P_{X^2}(+1)$ equals the sum $P_{X^2}(-1) + P_{X^2}(+1)$, we can rewrite it using P_X only.

$$\begin{aligned}\mathbb{E} X^2 &= 0 \cdot P_{X^2}(0) + (+1) \cdot P_{X^2}(+1) \\ &= 0^2 \cdot P_X(0) + [(-1)^2 \cdot P_X(-1) + (+1)^2 \cdot P_X(+1)] \\ &= \sum_{x \in \{-1, 0, +1\}} f(x) P_X(x)\end{aligned}\tag{8}$$

Behaviors of A function of a RV

If we generalize the previous discussion, we have the following theorem.

Theorem

Suppose that X is a RV and $f : \mathbb{R} \rightarrow \mathbb{R}$ are a real-valued function taking a real variable as an input. Then, $f(X)$ is also a RV.

In particular, if X is a discrete RV, $f(X)$ is also a discrete RV. Furthermore, if the support and PMF of X are denoted by \mathcal{X} and P_X , respectively,

- *The support of $f(X)$ is $\{f(x)|x \in \mathcal{X}\}$,*
- *The PMF $P_{f(X)}$ is given by $P_{f(X)}(y) = \sum_{x \in \{x' | f(x')=y\}} P_X(x)$,*
- *The expectation $\mathbb{E} f(X)$ is given by $\mathbb{E} f(X) = \sum_{x \in \mathcal{X}} f(x)P_X(x)$.*

The linearity of the expectation

The expectation operator \mathbb{E} has the property called **linearity**, which often makes the expectation calculation of a complicated function easier.

Theorem (The linearity of the expectation)

Let X be a random variable, $a, b \in \mathbb{R}$ be real numbers, and $f, g : \mathbb{R} \rightarrow \mathbb{R}$ be real-valued functions taking a real variable. Then, we have that

$$\mathbb{E}[af(X) + bg(X)] = a\mathbb{E}f(X) + b\mathbb{E}g(X). \quad (9)$$

The above theorem provides us with the formula for the expectation calculation of a linear function of a RV.

Corollary

Let X be a random variable and $a, b \in \mathbb{R}$ be real numbers. Then, we have that

$$\mathbb{E}[aX + b] = a\mathbb{E}X + b. \quad (10)$$

Example of a linear function's expectation

Example

Suppose that X is a random variable whose PMF P_X is given by the following table.

x	-2	-1	0	+1	+2
$P_X(x)$	0.05	0.10	0.20	0.10	0.55

Example random function and its PMF.

The expectation is given by $\mathbb{E}X = 1.00$.

Let's consider the random function given by $-3X + 5$ and its expectation.

According to the formula, $\mathbb{E}[-3X + 5] = -3\mathbb{E}X + 5 = (-3) \cdot 1.00 + 5 = 2.00$.

Note that the PMF P_{-3X+5} is given by the following, which we did not use to calculate $\mathbb{E}[-3X + 5]$.

x	+11	+8	+5	+2	-1
$P_{-3X+5}(x)$	0.05	0.10	0.20	0.10	0.55

The PMF of $-3X + 5$.

Proof: the linearity of the expectation

Proof.

$$\begin{aligned}\mathbb{E}[af(X) + bg(X)] &= \sum_{x \in \mathcal{X}} [af(x) + bg(x)]P_X(x) \\ &= a \sum_{x \in \mathcal{X}} f(x)P_X(x) + b \sum_{x \in \mathcal{X}} g(x)P_X(x) \\ &= a\mathbb{E}f(X) + b\mathbb{E}g(X).\end{aligned}\tag{11}$$



Definition of variance

Recall the basic idea of the variance.

Let X be a random variable and μ be its expectation. The **square deviation** of X is defined as $(X - \mu)^2$. If X is far (whether large or not) from μ , the square deviation $(X - \mu)^2$ is large. Hence, we can regard its expectation as a variability measure. This is the idea of the variance.

Definition (Variance)

Let X be a random variable and assume that the expectation $\mu := \mathbb{E}X$ exists. Then, the **variance** $\mathbb{V}[X] \in \mathbb{R}_{\geq 0}$ is defined as the expectation of the squared deviation ⁴ $(X - \mu)^2$, that is,

$$\mathbb{V}[X] := \mathbb{E}(X - \mu)^2. \quad (12)$$

⁴One reason for considering the square is to ignore the sign. For the same reason, the expectation of the absolute deviation is also used. However, the variance, the expectation of the squared deviation, is much more often used owing to the central limit theorem.

Calculating the variance

Recall that the variance is defined by $\mathbb{V}[X] := \mathbb{E}(X - \mu)^2$. Using the formula to calculate the expectation of the discrete RV, we get the following formula to calculate the variance of a discrete random variable.

Theorem

Let X be a discrete random variable taking values in $\mathcal{X} \subset \mathbb{R}$. Also, suppose that $\mu := \mathbb{E}X$ and $P_X : \mathcal{X} \rightarrow [0, 1]$ are its expectation and PMF, respectively.

The variance $\mathbb{V}[X]$ is given by

$$\mathbb{V}[X] = \sum_{x \in \mathcal{X}} (x - \mu)^2 P_X(x). \quad (13)$$

Example of variance calculation

x	-2	-1	0	+1	+2
Probability mass $P_X(x)$	0.05	0.10	0.20	0.10	0.55
Deviation $x - \mu_x$					
Square deviation $(x - \mu_X)^2$					
Weighted sq. dev. $(x - \mu_X)^2 P_X(x)$					

Example random function and its PMF.

- **Step 1:** Calculate the expectation $\mu_X = \mathbb{E}X$ of X .
In the above example, we have $\mu_X = \mathbb{E}X = +1.00$.
- **Step 2:**
- **Step 3:**

Example of variance calculation

x	-2	-1	0	+1	+2
Probability mass $P_X(x)$	0.05	0.10	0.20	0.10	0.55
Deviation $x - \mu_x$					
Square deviation $(x - \mu_X)^2$					
Weighted sq. dev. $(x - \mu_X)^2 P_X(x)$					

Example random function and its PMF.

- **Step 1:** Calculate the expectation $\mu_X = \mathbb{E}X$ of X .
In the above example, we have $\mu_X = \mathbb{E}X = +1.00$.
- **Step 2:** Calculate the deviation $x - \mu_X$, the square deviation $(x - \mu_X)^2$, and the weighted square deviation $(x - \mu_X)^2 P_X(x)$ for every $x \in \mathcal{X}$.
- **Step 3:**

Example of variance calculation

x	-2	-1	0	+1	+2
Probability mass $P_X(x)$	0.05	0.10	0.20	0.10	0.55
Deviation $x - \mu_x$	-3.00				
Square deviation $(x - \mu_X)^2$					
Weighted sq. dev. $(x - \mu_X)^2 P_X(x)$					

Example random function and its PMF.

- **Step 1:** Calculate the expectation $\mu_X = \mathbb{E}X$ of X .
In the above example, we have $\mu_X = \mathbb{E}X = +1.00$.
- **Step 2:** Calculate the deviation $x - \mu_X$, the square deviation $(x - \mu_X)^2$, and the weighted square deviation $(x - \mu_X)^2 P_X(x)$ for every $x \in \mathcal{X}$.
- **Step 3:**

Example of variance calculation

x	-2	-1	0	+1	+2
Probability mass $P_X(x)$	0.05	0.10	0.20	0.10	0.55
Deviation $x - \mu_x$	-3.00				
Square deviation $(x - \mu_X)^2$	9.00				
Weighted sq. dev. $(x - \mu_X)^2 P_X(x)$					

Example random function and its PMF.

- **Step 1:** Calculate the expectation $\mu_X = \mathbb{E}X$ of X .
In the above example, we have $\mu_X = \mathbb{E}X = +1.00$.
- **Step 2:** Calculate the deviation $x - \mu_X$, the square deviation $(x - \mu_X)^2$, and the weighted square deviation $(x - \mu_X)^2 P_X(x)$ for every $x \in \mathcal{X}$.
- **Step 3:**

Example of variance calculation

x	-2	-1	0	+1	+2
Probability mass $P_X(x)$	0.05	0.10	0.20	0.10	0.55
Deviation $x - \mu_x$	-3.00				
Square deviation $(x - \mu_X)^2$	9.00				
Weighted sq. dev. $(x - \mu_X)^2 P_X(x)$	0.45				

Example random function and its PMF.

- **Step 1:** Calculate the expectation $\mu_X = \mathbb{E}X$ of X .
In the above example, we have $\mu_X = \mathbb{E}X = +1.00$.
- **Step 2:** Calculate the deviation $x - \mu_X$, the square deviation $(x - \mu_X)^2$, and the weighted square deviation $(x - \mu_X)^2 P_X(x)$ for every $x \in \mathcal{X}$.
- **Step 3:**

Example of variance calculation

x	-2	-1	0	+1	+2
Probability mass $P_X(x)$	0.05	0.10	0.20	0.10	0.55
Deviation $x - \mu_x$	-3.00	-2.00			
Square deviation $(x - \mu_X)^2$	9.00				
Weighted sq. dev. $(x - \mu_X)^2 P_X(x)$	0.45				

Example random function and its PMF.

- **Step 1:** Calculate the expectation $\mu_X = \mathbb{E}X$ of X .
In the above example, we have $\mu_X = \mathbb{E}X = +1.00$.
- **Step 2:** Calculate the deviation $x - \mu_X$, the square deviation $(x - \mu_X)^2$, and the weighted square deviation $(x - \mu_X)^2 P_X(x)$ for every $x \in \mathcal{X}$.
- **Step 3:**

Example of variance calculation

x	-2	-1	0	+1	+2
Probability mass $P_X(x)$	0.05	0.10	0.20	0.10	0.55
Deviation $x - \mu_x$	-3.00	-2.00			
Square deviation $(x - \mu_X)^2$	9.00	4.00			
Weighted sq. dev. $(x - \mu_X)^2 P_X(x)$	0.45				

Example random function and its PMF.

- **Step 1:** Calculate the expectation $\mu_X = \mathbb{E}X$ of X .
In the above example, we have $\mu_X = \mathbb{E}X = +1.00$.
- **Step 2:** Calculate the deviation $x - \mu_X$, the square deviation $(x - \mu_X)^2$, and the weighted square deviation $(x - \mu_X)^2 P_X(x)$ for every $x \in \mathcal{X}$.
- **Step 3:**

Example of variance calculation

x	-2	-1	0	+1	+2
Probability mass $P_X(x)$	0.05	0.10	0.20	0.10	0.55
Deviation $x - \mu_x$	-3.00	-2.00			
Square deviation $(x - \mu_X)^2$	9.00	4.00			
Weighted sq. dev. $(x - \mu_X)^2 P_X(x)$	0.45	0.40			

Example random function and its PMF.

- **Step 1:** Calculate the expectation $\mu_X = \mathbb{E}X$ of X .
In the above example, we have $\mu_X = \mathbb{E}X = +1.00$.
- **Step 2:** Calculate the deviation $x - \mu_X$, the square deviation $(x - \mu_X)^2$, and the weighted square deviation $(x - \mu_X)^2 P_X(x)$ for every $x \in \mathcal{X}$.
- **Step 3:**

Example of variance calculation

x	-2	-1	0	+1	+2
Probability mass $P_X(x)$	0.05	0.10	0.20	0.10	0.55
Deviation $x - \mu_x$	-3.00	-2.00	-1.00	± 0.00	+1.00
Square deviation $(x - \mu_X)^2$	9.00	4.00	1.00	0.00	1.00
Weighted sq. dev. $(x - \mu_X)^2 P_X(x)$	0.45	0.40	0.20	0.00	0.55

Example random function and its PMF.

- **Step 1:** Calculate the expectation $\mu_X = \mathbb{E}X$ of X .
In the above example, we have $\mu_X = \mathbb{E}X = +1.00$.
- **Step 2:** Calculate the deviation $x - \mu_X$, the square deviation $(x - \mu_X)^2$, and the weighted square deviation $(x - \mu_X)^2 P_X(x)$ for every $x \in \mathcal{X}$.
- **Step 3:**

Example of variance calculation

x	-2	-1	0	+1	+2
Probability mass $P_X(x)$	0.05	0.10	0.20	0.10	0.55
Deviation $x - \mu_x$	-3.00	-2.00	-1.00	± 0.00	+1.00
Square deviation $(x - \mu_X)^2$	9.00	4.00	1.00	0.00	1.00
Weighted sq. dev. $(x - \mu_X)^2 P_X(x)$	0.45	0.40	0.20	0.00	0.55

Example random function and its PMF.

- **Step 1:** Calculate the expectation $\mu_X = \mathbb{E}X$ of X .
In the above example, we have $\mu_X = \mathbb{E}X = +1.00$.
- **Step 2:** Calculate the deviation $x - \mu_X$, the square deviation $(x - \mu_X)^2$, and the weighted square deviation $(x - \mu_X)^2 P_X(x)$ for every $x \in \mathcal{X}$.
- **Step 3:** Take the sum $\sum_{x \in \mathcal{X}} (x - \mu_X)^2 P_X(x)$.
In the above example, we have
$$\mathbb{V}[X] = \sum_{x \in \mathcal{X}} (x - \mu_X)^2 P_X(x) = 0.45 + 0.40 + 0.20 + 0.00 + 0.55$$

Example of variance calculation

x	-2	-1	0	+1	+2
Probability mass $P_X(x)$	0.05	0.10	0.20	0.10	0.55
Deviation $x - \mu_x$	-3.00	-2.00	-1.00	± 0.00	+1.00
Square deviation $(x - \mu_X)^2$	9.00	4.00	1.00	0.00	1.00
Weighted sq. dev. $(x - \mu_X)^2 P_X(x)$	0.45	0.40	0.20	0.00	0.55

Example random function and its PMF.

- **Step 1:** Calculate the expectation $\mu_X = \mathbb{E}X$ of X .
In the above example, we have $\mu_X = \mathbb{E}X = +1.00$.
- **Step 2:** Calculate the deviation $x - \mu_X$, the square deviation $(x - \mu_X)^2$, and the weighted square deviation $(x - \mu_X)^2 P_X(x)$ for every $x \in \mathcal{X}$.
- **Step 3:** Take the sum $\sum_{x \in \mathcal{X}} (x - \mu_X)^2 P_X(x)$.
In the above example, we have
$$\mathbb{V}[X] = \sum_{x \in \mathcal{X}} (x - \mu_X)^2 P_X(x) = 0.45 + 0.40 + 0.20 + 0.00 + 0.55 = 1.60.$$

Another formula of the variance

The following formula is also useful.

Theorem

Let X be a discrete random variable taking values in $\mathcal{X} \subset \mathbb{R}$. Also, suppose that $\mu := \mathbb{E}X$ and $P_X : \mathcal{X} \rightarrow [0, 1]$ are its expectation and PMF, respectively.

The variance $\mathbb{V}[X]$ is given by

$$\mathbb{V}[X] = \mathbb{E}X^2 - \mu^2 = \sum_{x \in \mathcal{X}} x^2 P_X(x) - \left(\sum_{x \in \mathcal{X}} x P_X(x) \right)^2. \quad (14)$$

Proof.

$$\mathbb{V}[X] = \mathbb{E} \left[(X - \mu)^2 \right] = \mathbb{E} [X^2 - 2\mu X + \mu^2] = \mathbb{E}X^2 - 2\mu \cdot \mu + \mu^2 = \mathbb{E}X^2 - \mu^2. \quad (15)$$



The variance of a linear function

Theorem

Let X be a random variable and $\alpha, b \in \mathbb{R}$ be real numbers. Then we have that

$$\mathbb{V}[\alpha X + b] = \alpha^2 \mathbb{V}[X]. \quad (16)$$

In particular, the variance does not depend on b .

Example of calculating the variance of a linear function

Example

Suppose that X is a random variable whose PMF P_X is given by the following table.

x	-2	-1	0	+1	+2
$P_X(x)$	0.05	0.10	0.20	0.10	0.55

Example random function and its PMF.

The variance is given by $\mathbb{V}[X] = 1.60$.

Let's consider the random function given by $-3X + 5$ and its variance.

According to the formula, $\mathbb{V}[-3X + 5] = (-3)^2 \mathbb{V}[X] = (-3)^2 \cdot 1.60 = 14.40$.

Note that we did not use the PMF of $-3X + 5$ to calculate $\mathbb{V}[-3X + 5]$.

Standard deviation

Variance's interpretation is somewhat tricky since its effect against scaling is not “linear.” Specifically, the variance of $10X$ is 100 times as large as that of X .

To make it “linear”, we consider the square root of the variance, called the **standard deviation** of the random variable.

Definition (Standard deviation)

The **standard deviation** $\sigma[X] \in \mathbb{R}$ of the random variable X is defined as

$$\sigma[X] := \sqrt{\mathbb{V}[X]}. \quad (17)$$

Example of the standard deviation calculation

Example

Suppose that X is a random variable whose PMF P_X is given by the following table.

x	-2	-1	0	+1	+2
$P_X(x)$	0.05	0.10	0.20	0.10	0.55

Example random function and its PMF.

The variance is given by $\mathbb{V}[X] = 1.60$.

Hence, the standard deviation $\sigma[X]$ is given by $\sigma[X] = \sqrt{\mathbb{V}[X]} = \sqrt{1.60} = 1.2649\dots$

The standard deviation of a linear function

Theorem

If f is a linear function, i.e., if $f(x) = ax + b$, where $a, b \in \mathbb{R}$, then we have that

$$\sigma[f(X)] = \sigma[aX + b] = |a|\sigma[X]. \quad (18)$$

In particular, the standard deviation does not depend on b .

Hence, as we expected, the standard deviation of $10X$ is 10 times as large as that of X . In this sense, the standard deviation is “linear.”

Note that the standard deviation is always non-negative. In particular, $\sigma[-10X]$ equals $10\sigma[X]$, but not $-10\sigma[X]$. This is an expected behavior since we originally wanted to measure the variability, which does not change even if we flip the sign.

Outline

1. Random variables

1.1 Introduction: why do we learn random variables?

1.2 Univariate discrete random variable

1.3 Visualization of a distribution

1.4 Summary statistics for a univariate random variable

1.5 Expectation

1.6 Median

1.7 Variance and a function of a random variable

1.8 Exercises

Exercise (Frequency)

Suppose that we have $m = 20$ students and consider their results in an exam. For $x \in \mathcal{X} = \{0, 1, 2, 3, 4, 5\}$, we denote the number of the students who got a score x by m_x . Suppose that m_x is given by the following table.

Score x	0	1	2	3	4	5
# students m_x	3	2	3	5	6	1

Exam results.

Let X be the score of the student sampled uniform-randomly from the 20 students. Find the PMF of X .

Exercise (Expectation and median)

Suppose that X is a random variable whose PMF P_X is given by the following table.

x	-2	-1	0	+1	+2
$P_X(x)$	0.05	0.10	0.20	0.10	0.55

Example random function and its PMF.

- (i) Find the value of the expectation $\mathbb{E}X$.
- (ii) Find the value of the median of X .

Exercise (Cumulative distribution function)

Suppose that X is a random variable whose PMF is given as follows.

x	0	1	2	3	4	5
$P_X(x) := \Pr(X = x)$	0.15	0.10	0.15	0.25	0.30	0.05

The PMF of a student exam result frequency

Find the cumulative distribution function of X and plot the graph.

Exercise (Variance)

Let X be a discrete random variable, whose PMF is given by the following table.

x	-2	-1	0	+1	+2
Probability mass $P_X(x)$	0.05	0.10	0.20	0.10	0.55

Example random function and its PMF.

1. Calculate the variance $\mathbb{V}[X]$.
2. Write down the PMF table of $X - 2$.
3. Calculate the variance $\mathbb{V}[X - 2]$.
4. Write down the PMF table of $5X$.
5. Calculate the variance $\mathbb{V}[5X]$.

Outline

2. Multiple Random Variables

2.1 Introduction: why are multiple random variables less trivial?

2.2 Joint distribution

2.3 Marginal distribution

2.4 Conditional distribution

2.5 Independence of random variables

2.6 Summary statistics for multiple RVs and covariance

2.7 Correlation

2.8 Exercises

Outline

2. Multiple Random Variables

2.1 Introduction: why are multiple random variables less trivial?

2.2 Joint distribution

2.3 Marginal distribution

2.4 Conditional distribution

2.5 Independence of random variables

2.6 Summary statistics for multiple RVs and covariance

2.7 Correlation

2.8 Exercises

Multiple random variables

Example

- The prices of multiple stocks.
- The pixels of an image.
- The values at each time frame in a wave file.

When we consider multiple random variables, knowing each probability mass function is not sufficient to know their stochastic behavior completely.

Knowing multiple random variables \neq knowing multiple PMFs

If we have two discrete random variables X and Y , then just knowing each probability mass function is not sufficient.

Rather, what we need to know is the distribution of the **pair** (X, Y) , which is called the ***joint distribution*** of the random variables X and Y .

Learning outcomes

By the end of this section, you should be able to:

- Explain why two probability mass functions are not sufficient to describe multiple random variables,
- Describe multiple random variables using the joint probability mass function and conditional probability mass function,
- Describe the relation between multiple random variables using covariance, correlation, and independence, and
- Explain the difference between covariance, correlation, independence, and causality.

Outline

2. Multiple Random Variables

2.1 Introduction: why are multiple random variables less trivial?

2.2 Joint distribution

2.3 Marginal distribution

2.4 Conditional distribution

2.5 Independence of random variables

2.6 Summary statistics for multiple RVs and covariance

2.7 Correlation

2.8 Exercises

Joint distribution and marginal distribution

In general, the **joint distribution** refers to the distribution of the tuple of multiple random variables. For example, if we have two random variables X and Y , the joint distribution refers to the distribution of the pair (X, Y) .

In contrast, when we consider multiple random variables, the distribution of a single random variable is called the **marginal distribution** of the random variable to distinguish it from the joint distribution.

Joint probability mass function (two variable cases)

If we have two discrete random variables X and Y , then just knowing each probability mass function is not sufficient. Rather, what we need to know is the probability of the pair (X, Y) taking every pair of values $(x, y) \in \mathcal{X} \times \mathcal{Y}$. That is, the following **joint probability mass function (joint PMF)** has all the information that we need.

Definition (two-variable Joint PMF)

Let X and Y be discrete random variables taking a value in discrete sets \mathcal{X} and \mathcal{Y} , respectively, where $\mathcal{X}, \mathcal{Y} \subset \mathbb{R}$. We define the **joint probability mass function (joint PMF)** $P_{X,Y} : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$ of the pair of random variables X, Y by

$$P_{X,Y}(x, y) := \Pr(X = x, Y = y). \quad (19)$$

Joint PMF example

The random variables X and Y are the scores of a math test and a history test, respectively, where we uniform-randomly sample a student.

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.16	0.04	0.02	0.06	0.28
	1	0.18	0.04	0.04	0.16	0.42
	2	0.06	0.02	0.08	0.14	0.30
$P_X(x)$		0.40	0.10	0.14	0.36	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

Joint probability mass function (general cases)

If we have m discrete random variables X_1, X_2, \dots, X_m , then all we need to know is the following joint PMF.

Definition (Joint PMF (general cases))

Let X_1, X_2, \dots, X_m be discrete random variables taking a value in discrete sets $\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_m \subset \mathbb{R}$, respectively. We define the **joint probability mass function (joint PMF)** $P_{X_1, X_2, \dots, X_m} : \mathcal{X}_1 \times \mathcal{X}_2 \times \dots \times \mathcal{X}_m \rightarrow [0, 1]$ of random variables X_1, X_2, \dots, X_m by

$$P_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) := \Pr(X_1 = x_1, X_2 = x_2, \dots, X_m = x_m). \quad (20)$$

Outline

2. Multiple Random Variables

2.1 Introduction: why are multiple random variables less trivial?

2.2 Joint distribution

2.3 Marginal distribution

2.4 Conditional distribution

2.5 Independence of random variables

2.6 Summary statistics for multiple RVs and covariance

2.7 Correlation

2.8 Exercises

Marginal PMF (two variable cases)

The joint PMF can tell us the PMFs of each discrete random variable, called **marginal PMF**. For two discrete random variables X and Y that takes a value in \mathcal{X} and \mathcal{Y} , respectively, suppose that the joint PMF is $P_{X,Y} : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$. Then, the marginal PMFs P_X and P_Y are given by

$$P_X(x) = \sum_{y \in \mathcal{Y}} P_{X,Y}(x, y), \quad P_Y(y) = \sum_{x \in \mathcal{X}} P_{X,Y}(x, y), \quad (21)$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	
	1	0.24	0.08	0.12	0.06	
	2	0.06	0.00	0.06	0.18	
$P_X(x)$						

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_X(0) &= P_{X,Y}(0,0) + P_{X,Y}(0,1) + P_{X,Y}(0,2) \\ &= 0.10 + 0.24 + 0.06 \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	
	1	0.24	0.08	0.12	0.06	
	2	0.06	0.00	0.06	0.18	
$P_X(x)$		0.40				

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_X(0) &= P_{X,Y}(0,0) + P_{X,Y}(0,1) + P_{X,Y}(0,2) \\ &= 0.10 + 0.24 + 0.06 = \mathbf{0.40}. \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	
	1	0.24	0.08	0.12	0.06	
	2	0.06	0.00	0.06	0.18	
$P_X(x)$		0.40				

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_X(1) &= P_{X,Y}(1,0) + P_{X,Y}(1,1) + P_{X,Y}(1,2) \\ &= 0.02 + 0.08 + 0.00 \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	
	1	0.24	0.08	0.12	0.06	
	2	0.06	0.00	0.06	0.18	
$P_X(x)$		0.40	0.10			

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_X(1) &= P_{X,Y}(1,0) + P_{X,Y}(1,1) + P_{X,Y}(1,2) \\ &= 0.02 + 0.08 + 0.00 = \mathbf{0.10}. \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	
	1	0.24	0.08	0.12	0.06	
	2	0.06	0.00	0.06	0.18	
$P_X(x)$		0.40	0.10			

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_X(2) &= P_{X,Y}(2,0) + P_{X,Y}(2,1) + P_{X,Y}(2,2) \\ &= 0.02 + 0.12 + 0.06 \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	
	1	0.24	0.08	0.12	0.06	
	2	0.06	0.00	0.06	0.18	
$P_X(x)$		0.40	0.10	0.20		

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_X(2) &= P_{X,Y}(2,0) + P_{X,Y}(2,1) + P_{X,Y}(2,2) \\ &= 0.02 + 0.12 + 0.06 = \mathbf{0.20}. \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	
	1	0.24	0.08	0.12	0.06	
	2	0.06	0.00	0.06	0.18	
$P_X(x)$		0.40	0.10	0.20		

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_X(3) &= P_{X,Y}(3,0) + P_{X,Y}(3,1) + P_{X,Y}(3,2) \\ &= 0.06 + 0.06 + 0.18 \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	
	1	0.24	0.08	0.12	0.06	
	2	0.06	0.00	0.06	0.18	
$P_X(x)$		0.40	0.10	0.20	0.30	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_X(3) &= P_{X,Y}(3,0) + P_{X,Y}(3,1) + P_{X,Y}(3,2) \\ &= 0.06 + 0.06 + 0.18 = \mathbf{0.30}. \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	
	1	0.24	0.08	0.12	0.06	
	2	0.06	0.00	0.06	0.18	
$P_X(x)$		0.40	0.10	0.20	0.30	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_Y(0) &= P_{X,Y}(0,0) + P_{X,Y}(1,0) + P_{X,Y}(2,0) + P_{X,Y}(3,0) \\ &= 0.10 + 0.02 + 0.02 + 0.06 \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	0.20
	1	0.24	0.08	0.12	0.06	
	2	0.06	0.00	0.06	0.18	
$P_X(x)$		0.40	0.10	0.20	0.30	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_Y(0) &= P_{X,Y}(0,0) + P_{X,Y}(1,0) + P_{X,Y}(2,0) + P_{X,Y}(3,0) \\ &= 0.10 + 0.02 + 0.02 + 0.06 = \mathbf{0.20}. \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	0.20
	1	0.24	0.08	0.12	0.06	
	2	0.06	0.00	0.06	0.18	
$P_X(x)$		0.40	0.10	0.20	0.30	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_Y(1) &= P_{X,Y}(0,1) + P_{X,Y}(1,1) + P_{X,Y}(2,1) + P_{X,Y}(3,1) \\ &= 0.24 + 0.08 + 0.12 + 0.06 \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	0.20
	1	0.24	0.08	0.12	0.06	0.50
	2	0.06	0.00	0.06	0.18	
$P_X(x)$		0.40	0.10	0.20	0.30	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_Y(1) &= P_{X,Y}(0,1) + P_{X,Y}(1,1) + P_{X,Y}(2,1) + P_{X,Y}(3,1) \\ &= 0.24 + 0.08 + 0.12 + 0.06 = \mathbf{0.50}. \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	0.20
	1	0.24	0.08	0.12	0.06	0.50
	2	0.06	0.00	0.06	0.18	
$P_X(x)$		0.40	0.10	0.20	0.30	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_Y(2) &= P_{X,Y}(0,2) + P_{X,Y}(1,2) + P_{X,Y}(2,2) + P_{X,Y}(3,2) \\ &= 0.06 + 0.00 + 0.06 + 0.18 \end{aligned} \tag{22}$$

Marginal distribution example

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	0.20
	1	0.24	0.08	0.12	0.06	0.50
	2	0.06	0.00	0.06	0.18	0.30
$P_X(x)$		0.40	0.10	0.20	0.30	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

In the above example, we can calculate the marginal PMF from the joint PMF as follows.

$$\begin{aligned} P_Y(2) &= P_{X,Y}(0,2) + P_{X,Y}(1,2) + P_{X,Y}(2,2) + P_{X,Y}(3,2) \\ &= 0.06 + 0.00 + 0.06 + 0.18 = \mathbf{0.30}. \end{aligned} \tag{22}$$

Outline

2. Multiple Random Variables

2.1 Introduction: why are multiple random variables less trivial?

2.2 Joint distribution

2.3 Marginal distribution

2.4 Conditional distribution

2.5 Independence of random variables

2.6 Summary statistics for multiple RVs and covariance

2.7 Correlation

2.8 Exercises

Conditional distribution

If two RVs are “related,” then we get more precise information about a RV’s distribution by knowing the value of the other RV.

The ***conditional distribution*** is a piece of such information.

The conditional distribution is the distribution of one RV when we know the value of the other RV.

The probability mass function (PMF) of the conditional distribution is called the ***conditional PMF***.

Conditional distribution example

Let X and Y be discrete RVs, and suppose that their joint PMF $P_{X,Y}$ and marginal PMFs P_X and P_Y are given by the following table.

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	0.20
	1	0.24	0.08	0.12	0.06	0.50
	2	0.06	0.00	0.06	0.18	0.30
$P_X(x)$		0.40	0.10	0.20	0.30	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

Conditional distribution example

Let X and Y be discrete RVs, and suppose that their joint PMF $P_{X,Y}$ and marginal PMFs P_X and P_Y are given by the following table.

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	0.20
	1	0.24	0.08	0.12	0.06	0.50
	2	0.06	0.00	0.06	0.18	0.30
$P_X(x)$		0.40	0.10	0.20	0.30	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

Suppose that we know that $Y = 2$. This information changes the distribution of X . For example, $X = 1$ no longer happens, so the probability of the event $X = 1$ is now zero.

So, for $x = 0, 1, 2, 3$, what is the probability of " $X = x$ " when we know $Y = 2$? It is called the **conditional probability** of $X = x$ given $Y = 2$ and denoted by $P_{X|Y}(x|2)$.

Conditional probability calculation

	x				$P_Y(y)$
	0	1	2	3	
$P_{X,Y}(x,2)$	0.06	0.00	0.06	0.18	0.30

Joint PMF and conditional PMF

- If we know $Y = 2$, then the probability masses of $X = 0, 1, 2, 3$ are proportional to the joint masses $P_{X,Y}(0,2), P_{X,Y}(1,2), P_{X,Y}(2,2), P_{X,Y}(3,2)$, shown above.
- The sum $P_{X|Y}(0|2) + P_{X|Y}(1|2) + P_{X|Y}(2|2) + P_{X|Y}(3|2)$ of the conditional probabilities must be 1 for them to be probabilities.

Hence, the conditional probability $P_{X|Y}(x|2)$ is each joint probability over the sum, i.e.,

$$P_{X|Y}(x|2) = \frac{P_{X,Y}(x,2)}{P_{X,Y}(0,2) + P_{X,Y}(1,2) + P_{X,Y}(2,2) + P_{X,Y}(3,2)} = \frac{P_{X,Y}(x,2)}{P_Y(2)}. \quad (23)$$

Conditional probability calculation example

	x				$P_Y(y)$
	0	1	2	3	
$P_{X,Y}(x,2)$	0.06	0.00	0.06	0.18	0.30
$P_{X Y}(x y)$					

Joint PMF and conditional PMF

For example,

Conditional probability calculation example

	x				$P_Y(y)$
	0	1	2	3	
$P_{X,Y}(x,2)$	0.06	0.00	0.06	0.18	0.30
$P_{X Y}(x y)$?				

Joint PMF and conditional PMF

For example,

$$P_{X|Y}(0|2) = \frac{P_{X,Y}(0,2)}{P_Y(2)} \quad (23)$$

Conditional probability calculation example

	x				$P_Y(y)$
	0	1	2	3	
$P_{X,Y}(x,2)$	0.06	0.00	0.06	0.18	0.30
$P_{X Y}(x y)$	0.20				

Joint PMF and conditional PMF

For example,

$$P_{X|Y}(0|2) = \frac{P_{X,Y}(0,2)}{P_Y(2)} = \frac{0.06}{0.30} = 0.20 \quad (23)$$

Conditional probability calculation example

	x				$P_Y(y)$
	0	1	2	3	
$P_{X,Y}(x,2)$	0.06	0.00	0.06	0.18	0.30
$P_{X Y}(x y)$	0.20	?			

Joint PMF and conditional PMF

For example,

$$P_{X|Y}(1|2) = \frac{P_{X,Y}(1,2)}{P_Y(2)} \quad (23)$$

Conditional probability calculation example

	x				$P_Y(y)$
	0	1	2	3	
$P_{X,Y}(x,2)$	0.06	0.00	0.06	0.18	0.30
$P_{X Y}(x y)$	0.20	0.00			

Joint PMF and conditional PMF

For example,

$$P_{X|Y}(1|2) = \frac{P_{X,Y}(1,2)}{P_Y(2)} = \frac{0.00}{0.30} = 0.00 \quad (23)$$

Conditional probability calculation example

	x				$P_Y(y)$
	0	1	2	3	
$P_{X,Y}(x,2)$	0.06	0.00	0.06	0.18	0.30
$P_{X Y}(x y)$	0.20	0.00	?		

Joint PMF and conditional PMF

For example,

$$P_{X|Y}(2|2) = \frac{P_{X,Y}(2,2)}{P_Y(2)} \quad (23)$$

Conditional probability calculation example

	x				$P_Y(y)$
	0	1	2	3	
$P_{X,Y}(x,2)$	0.06	0.00	0.06	0.18	0.30
$P_{X Y}(x y)$	0.20	0.00	0.20		

Joint PMF and conditional PMF

For example,

$$P_{X|Y}(2|2) = \frac{P_{X,Y}(2,2)}{P_Y(2)} = \frac{0.06}{0.30} = 0.20 \quad (23)$$

Conditional probability calculation example

	x				$P_Y(y)$
	0	1	2	3	
$P_{X,Y}(x,2)$	0.06	0.00	0.06	0.18	0.30
$P_{X Y}(x y)$	0.20	0.00	0.20	?	

Joint PMF and conditional PMF

For example,

$$P_{X|Y}(3|2) = \frac{P_{X,Y}(3,2)}{P_Y(2)} \quad (23)$$

Conditional probability calculation example

	x				$P_Y(y)$
	0	1	2	3	
$P_{X,Y}(x,2)$	0.06	0.00	0.06	0.18	0.30
$P_{X Y}(x y)$	0.20	0.00	0.20	0.60	

Joint PMF and conditional PMF

For example,

$$P_{X|Y}(3|2) = \frac{P_{X,Y}(3,2)}{P_Y(2)} = \frac{0.18}{0.30} = 0.60 \quad (23)$$

Conditional probability calculation example

	x				$P_Y(y)$
	0	1	2	3	
$P_{X,Y}(x,2)$	0.06	0.00	0.06	0.18	0.30
$P_{X Y}(x y)$	0.20	0.00	0.20	0.60	
$P_X(x)$	0.40	0.10	0.20	0.30	

Joint PMF and conditional PMF

You can see that

- The conditional probabilities are different from the marginal probabilities.
- The sum $P_{X|Y}(0|2) + P_{X|Y}(1|2) + P_{X|Y}(2|2) + P_{X|Y}(3|2)$ of the conditional probabilities is one.

We call the function $P_{X|Y}$ the **conditional PMF** of X given Y .

Definition of the conditional PMF

Definition

Let X and Y be discrete random variables, whose supports are \mathcal{X} and \mathcal{Y} , respectively. In other words, for any $x \in \mathcal{X}$ and $y \in \mathcal{Y}$, $P_X(x) > 0$ and $P_Y(y) > 0$ holds, where P_X and P_Y are the marginal PMFs of X and Y , respectively.

Let $P_{X,Y}$ be the joint PMF of X and Y .

We define the conditional PMF $P_{X|Y}$ by

$$P_{X|Y}(x|y) := \frac{P_{X,Y}(x,y)}{P_Y(y)}. \quad (23)$$

Likewise, we define the conditional PMF $P_{Y|X}$ by

$$P_{Y|X}(y|x) := \frac{P_{X,Y}(x,y)}{P_X(x)}. \quad (24)$$

Note: The conditional probability is not commutable.

Note that $P_{X|Y}(x|y) \neq P_{Y|X}(y|x)$ in general.

In this sense, the conditional probability is **NOT commutable**.

Conditional probability calculation from joint PMF

In general, we can calculate the conditional PMF from the joint PMF and the marginal PMF as follows:

$$P_{X|Y}(x|y) = \frac{P_{X,Y}(x,y)}{P_Y(y)}. \quad (25)$$

Since we can calculate the marginal probability $P_Y(y)$ by $P_Y(y) = \sum_{x \in \mathcal{X}} P_{X,Y}(x,y)$ using the joint PMF $P_{X,Y}$, we can calculate the conditional PMF only from the joint PMF in theory.

Outline

2. Multiple Random Variables

2.1 Introduction: why are multiple random variables less trivial?

2.2 Joint distribution

2.3 Marginal distribution

2.4 Conditional distribution

2.5 Independence of random variables

2.6 Summary statistics for multiple RVs and covariance

2.7 Correlation

2.8 Exercises

Independence of random variables

Suppose that the conditional PMF always equals the marginal PMF, i.e., $P_{X|Y}(x|y) = P_X(x)$ for all x and y .

It means that Y has no relation to X . In this case, we say that X and Y are ***independent***.

Definition

Let X and Y be discrete random variables. If one of the following equivalent conditions⁵ holds, we say that X and Y are independent.

- $P_{X|Y}(x|y) = P_X(x)$ for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.
- $P_{Y|X}(y|x) = P_Y(y)$ for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.
- $P_{X,Y}(x,y) = P_X(x)P_Y(y)$ for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.

⁵Specifically, if one condition holds, then the other two conditions also hold.

Example of independent random variables

Suppose that the joint PMF of random variables X and Y is given by:

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.08	0.02	0.04	0.06	0.20
	1	0.20	0.05	0.10	0.15	0.50
	2	0.12	0.03	0.06	0.09	0.30
$P_X(x)$		0.40	0.10	0.20	0.30	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

We can confirm that X and Y are mutually independent by checking that $P_{X,Y}(x,y) = P_X(x)P_Y(y)$ holds for every $x \in \mathcal{X} = \{0, 1, 2, 3\}$ and $y \in \mathcal{Y} = \{0, 1, 2\}$.

Example of independent random variables

Suppose that the joint PMF of random variables X and Y is given by:

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.08	0.02	0.04	0.06	0.20
	1	0.20	0.05	0.10	0.15	0.50
	2	0.12	0.03	0.06	0.09	0.30
$P_X(x)$		0.40	0.10	0.20	0.30	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

For example, $P_{X,Y}(0,0) = 0.08$, which equals to $P_X(0)P_Y(0) = 0.40 \times 0.20$.

Example of independent random variables

Suppose that the joint PMF of random variables X and Y is given by:

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.08	0.02	0.04	0.06	0.20
	1	0.20	0.05	0.10	0.15	0.50
	2	0.12	0.03	0.06	0.09	0.30
$P_X(x)$		0.40	0.10	0.20	0.30	

An example of $P_{X,Y}(x, y) := \Pr(X = x \wedge Y = y)$

For example, $P_{X,Y}(2, 1) = 0.10$, which equals to $P_X(2)P_Y(1) = 0.20 \times 0.50$.

Outline

2. Multiple Random Variables

2.1 Introduction: why are multiple random variables less trivial?

2.2 Joint distribution

2.3 Marginal distribution

2.4 Conditional distribution

2.5 Independence of random variables

2.6 Summary statistics for multiple RVs and covariance

2.7 Correlation

2.8 Exercises

Summary statistics for multiple RVs to show the relation

When we have multiple variables, we can calculate summary statistics for each of the variables. However, they do not give us information about the relation between multiple variables.

There are some statistics to show the relation between two RVs.

One principal question about the relation between two random variables X and Y is: “Do the RVs tend to take (relatively) large values simultaneously?”

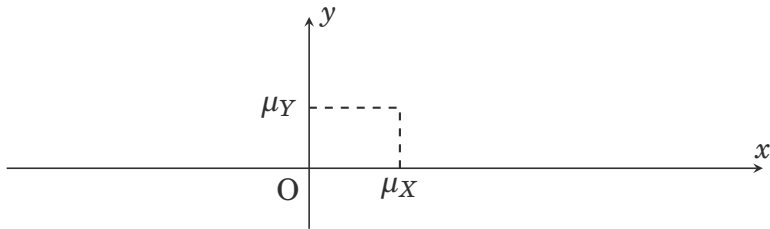
If X is easily observable and Y is the value of some product in the near future, then the information about the above relation financially benefits us.

The idea of covariance

The question is “Do the RVs tend to take (relatively) large values simultaneously?”

To answer the question, we consider the product of $X - \mu_X$ and $Y - \mu_Y$, where $\mu_X := \mathbb{E}X$ and $\mu_Y := \mathbb{E}Y$ are the expectations of X and Y , respectively.

The value $X - \mu_X$ is positive if X takes a relatively large value and negative if X takes a relatively small value.

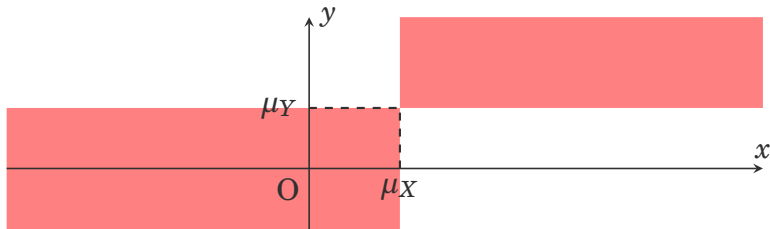


The idea of covariance

The question is “Do the RVs tend to take (relatively) large values simultaneously?”

To answer the question, we consider the product of $X - \mu_X$ and $Y - \mu_Y$, where $\mu_X := \mathbb{E}X$ and $\mu_Y := \mathbb{E}Y$ are the expectations of X and Y , respectively.

If X and Y tend to take large values simultaneously and small values simultaneously as well, then the product $(X - \mu_X)(Y - \mu_Y)$ tends to be positive.



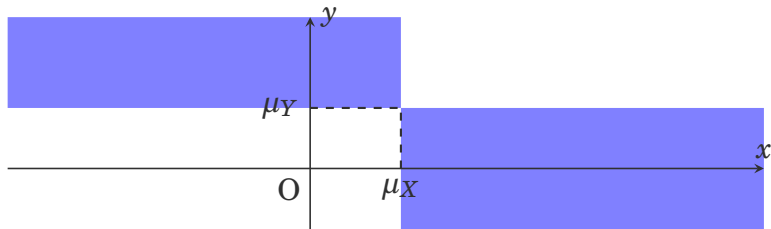
The area where $(X - \mu_X)(Y - \mu_Y)$ takes a positive value.

The idea of covariance

The question is “Do the RVs tend to take (relatively) large values simultaneously?”

To answer the question, we consider the product of $X - \mu_X$ and $Y - \mu_Y$, where $\mu_X := \mathbb{E}X$ and $\mu_Y := \mathbb{E}Y$ are the expectations of X and Y , respectively.

Conversely, if one tends to be small when the other is large, then the product $(X - \mu_X)(Y - \mu_Y)$ tends to be negative.



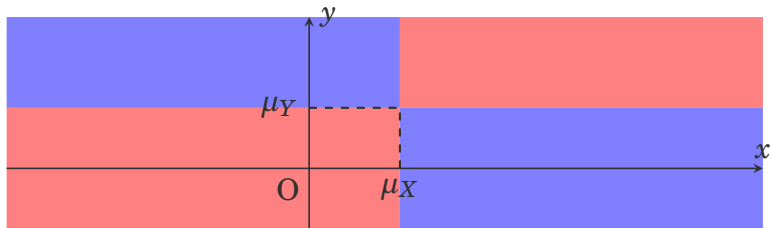
The area where $(X - \mu_X)(Y - \mu_Y)$ takes a negative value.

The idea of covariance

The question is “Do the RVs tend to take (relatively) large values simultaneously?”

To answer the question, we consider the product of $X - \mu_X$ and $Y - \mu_Y$, where $\mu_X := \mathbb{E}X$ and $\mu_Y := \mathbb{E}Y$ are the expectations of X and Y , respectively.

Hence, we are interested in the value of $(X - \mu_X)(Y - \mu_Y)$. This is the basic idea of **covariance**. But what is $(X - \mu_X)(Y - \mu_Y)$?



A function of multiple RVs

We say that the variable $(X - \mu_X)(Y - \mu_Y)$ is a function of RVs X and Y since it depends on the RVs X and Y .

We remark that $(X - \mu_X)(Y - \mu_Y)$ is a random variable. In particular, it is a discrete RV since X and Y are discrete RVs. Since it is a random variable, we can define its expectation $\mathbb{E}(X - \mu_X)(Y - \mu_Y)$.

Let's discuss the general function of multiple RVs and define its expectations.

A function of multiple RVs and its expectation

Theorem

Suppose that X and Y are random variables and $f : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ are a real-valued function taking two real values as an input. Then, $f(X, Y)$ is a random variable. In particular, suppose that X and Y are discrete RVs, their supports are \mathcal{X} and \mathcal{Y} , respectively, and their joint PMF is $P_{X,Y}$. Then, $f(X, Y)$ is also a discrete RV and

- The support of $f(X, Y)$ is $\{f(x, y) | x \in \mathcal{X}, y \in \mathcal{Y}\}$,
- The PMF $P_{f(X,Y)}$ is given by

$$P_{f(X,Y)}(z) = \sum_{(x,y) \in \{(x',y') | f(x',y')=z\}} P_{X,Y}(x, y), \quad (26)$$

- The expectation $\mathbb{E} f(X, Y)$ is given by

$$\mathbb{E} f(X, Y) = \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} f(x, y) P_{X,Y}(x, y). \quad (27)$$

The linearity of the expectation: the multi-variable case

From the linearity of the expectation operator \mathbb{E} , the following holds.

Theorem (The linearity of the expectation)

Let X, Y be random variables, $a, b \in \mathbb{R}$ be real numbers, and $f, g : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be real-valued functions taking two real variables as an input. Then, we have that

$$\mathbb{E}[af(X, Y) + bg(X, Y)] = a\mathbb{E}f(X, Y) + b\mathbb{E}g(X, Y). \quad (28)$$

The above theorem provides us with the formula for the expectation calculation of a linear function of multiple variables.

Corollary

Let X, Y be random variables and $a, b, c \in \mathbb{R}$ be real numbers. Then, we have that

$$\mathbb{E}[aX + bY + c] = a\mathbb{E}X + b\mathbb{E}Y + c. \quad (29)$$

Definition of the covariance

Now, we are ready to define the **covariance**. Recall that the idea of covariance is to evaluate the behavior of $(X - \mu_X)(Y - \mu_Y)$. In fact, the covariance is nothing but the expectation of $(X - \mu_X)(Y - \mu_Y)$.

Definition (Covariance)

Let X and Y be RVs and $\mu_X := \mathbb{E}X$ and $\mu_Y := \mathbb{E}Y$ be their expectations. We define the **covariance** $\text{Cov}(X, Y) \in \mathbb{R}$ between the two random variables X and Y by

$$\text{Cov}(X, Y) := \mathbb{E}[(X - \mu_X)(Y - \mu_Y)]. \quad (30)$$

Note that the covariance is symmetric, i.e., $\text{Cov}(X, Y) = \text{Cov}(Y, X)$.

A positive covariance indicates that the two random variables tend to take relatively large values simultaneously. A negative covariance indicates that when one of the two takes a relatively large value, then the other tends to take a relatively small value.

Formulae to calculate the covariance

We provide the explicit calculation formula of the covariance.

Theorem

Suppose that X and Y are discrete RVs, their supports are \mathcal{X} and \mathcal{Y} , respectively, and their joint PMF is $P_{X,Y}$.

Then, the covariance $\text{Cov}(X, Y) \in \mathbb{R}$ between the two random variables X and Y is given by

$$\text{Cov}(X, Y) = \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} (x - \mu_X)(y - \mu_Y) P_{X,Y}(x, y). \quad (31)$$

Example of covariance calculation

Example

		x		$P_Y(y)$
		0	+1	
y	0	0.25	0.00	0.25
	+1	0.25	0.25	0.50
	+2	0.00	0.25	0.25
$P_X(x)$		0.50	0.50	

The joint PMF $P_{X,Y}$. The RVs X and Y have a positive covariance.

We can calculate the covariance $\text{Cov}(X, Y)$ of RVs X, Y from its joint PMF $P_{X,Y}$.

-

Example of covariance calculation

Example

		x		$P_Y(y)$
		0	+1	
y	0	0.25	0.00	0.25
	+1	0.25	0.25	0.50
	+2	0.00	0.25	0.25
$P_X(x)$		0.50	0.50	

The joint PMF $P_{X,Y}$. The RVs X and Y have a positive covariance.

We can calculate the covariance $\text{Cov}(X, Y)$ of RVs X, Y from its joint PMF $P_{X,Y}$.

- **Step 1:** Calculate the expectations $\mu_X = \mathbb{E}X$ and $\mu_Y = \mathbb{E}Y$. Then memorize the value $x - \mu_X$ for all $x \in \mathcal{X}$ and $y - \mu_Y$ for all $y \in \mathcal{Y}$.

Example of covariance calculation

Example

		x		$P_Y(y)$
		0	+1	
y	0	0.25	0.00	0.25
	+1	0.25	0.25	0.50
	+2	0.00	0.25	0.25
$P_X(x)$		0.50	0.50	

The joint PMF $P_{X,Y}$. The RVs X and Y have a positive covariance.

We can calculate the covariance $\text{Cov}(X, Y)$ of RVs X, Y from its joint PMF $P_{X,Y}$.

- **Step 1:** Calculate the expectations $\mu_X = \mathbb{E}X$ and $\mu_Y = \mathbb{E}Y$. Then memorize the value $x - \mu_X$ for all $x \in \mathcal{X}$ and $y - \mu_Y$ for all $y \in \mathcal{Y}$.

In the above example, we have $\mu_X = \mathbb{E}X = +0.50$ and $\mu_Y = \mathbb{E}Y = +1.00$.

Example of covariance calculation

Example

		$x - \mu_X$		$P_Y(y)$
		-0.5	+0.5	
$y - \mu_Y$	-1	0.25	0.00	0.25
	0	0.25	0.25	0.50
	+1	0.00	0.25	0.25
$P_X(x)$		0.50	0.50	

The joint PMF $P_{X,Y}$. The RVs X and Y have a positive covariance.

We can calculate the covariance $\text{Cov}(X, Y)$ of RVs X, Y from its joint PMF $P_{X,Y}$.

- **Step 1:** Calculate the expectations $\mu_X = \mathbb{E}X$ and $\mu_Y = \mathbb{E}Y$. Then memorize the value $x - \mu_X$ for all $x \in \mathcal{X}$ and $y - \mu_Y$ for all $y \in \mathcal{Y}$.

In the above example, we have $\mu_X = \mathbb{E}X = +0.50$ and $\mu_Y = \mathbb{E}Y = +1.00$.

Example of covariance calculation

Example

		$x - \mu_X$		$P_Y(y)$
		-0.5	+0.5	
$y - \mu_Y$	-1	0.25	0.00	0.25
	0	0.25	0.25	0.50
	+1	0.00	0.25	0.25
$P_X(x)$		0.50	0.50	

The joint PMF $P_{X,Y}$. The RVs X and Y have a positive covariance.

We can calculate the covariance $\text{Cov}(X, Y)$ of RVs X, Y from its joint PMF $P_{X,Y}$.

- Step 2:** Calculate the weighted product of the deviations $(x - \mu_X)(y - \mu_Y)P_{X,Y}(x, y)$ for every $(x, y) \in \mathcal{X} \times \mathcal{Y}$ and take the sum.

In the above example, we have

$$\begin{aligned}\text{Cov}(X, Y) &= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} (x - \mu_X)(y - \mu_Y)P_{X,Y}(x, y) \\ &= (-0.5) \cdot (-1) \cdot 0.25\end{aligned}$$

Example of covariance calculation

Example

		$x - \mu_X$		$P_Y(y)$
		-0.5	+0.5	
$y - \mu_Y$	-1	0.25	0.00	0.25
	0	0.25	0.25	0.50
	+1	0.00	0.25	0.25
$P_X(x)$		0.50	0.50	

The joint PMF $P_{X,Y}$. The RVs X and Y have a positive covariance.

We can calculate the covariance $\text{Cov}(X, Y)$ of RVs X, Y from its joint PMF $P_{X,Y}$.

- Step 2:** Calculate the weighted product of the deviations $(x - \mu_X)(y - \mu_Y)P_{X,Y}(x, y)$ for every $(x, y) \in \mathcal{X} \times \mathcal{Y}$ and take the sum.

In the above example, we have

$$\begin{aligned}\text{Cov}(X, Y) &= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} (x - \mu_X)(y - \mu_Y)P_{X,Y}(x, y) \\ &= (-0.5) \cdot (-1) \cdot 0.25 + (+0.5) \cdot (-1) \cdot 0.00\end{aligned}$$

Example of covariance calculation

Example

		$x - \mu_X$		$P_Y(y)$
		-0.5	+0.5	
$y - \mu_Y$	-1	0.25	0.00	0.25
	0	0.25	0.25	0.50
	+1	0.00	0.25	0.25
$P_X(x)$		0.50	0.50	

The joint PMF $P_{X,Y}$. The RVs X and Y have a positive covariance.

We can calculate the covariance $\text{Cov}(X, Y)$ of RVs X, Y from its joint PMF $P_{X,Y}$.

- Step 2:** Calculate the weighted product of the deviations $(x - \mu_X)(y - \mu_Y)P_{X,Y}(x, y)$ for every $(x, y) \in \mathcal{X} \times \mathcal{Y}$ and take the sum.

In the above example, we have

$$\begin{aligned}\text{Cov}(X, Y) &= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} (x - \mu_X)(y - \mu_Y)P_{X,Y}(x, y) \\ &= (-0.5) \cdot (-1) \cdot 0.25 + (+0.5) \cdot (-1) \cdot 0.00 + (-0.5) \cdot 0 \cdot 0.25\end{aligned}$$

Example of covariance calculation

Example

		$x - \mu_X$		$P_Y(y)$
		-0.5	+0.5	
$y - \mu_Y$	-1	0.25	0.00	0.25
	0	0.25	0.25	0.50
	+1	0.00	0.25	0.25
$P_X(x)$		0.50	0.50	

The joint PMF $P_{X,Y}$. The RVs X and Y have a positive covariance.

We can calculate the covariance $\text{Cov}(X, Y)$ of RVs X, Y from its joint PMF $P_{X,Y}$.

- **Step 2:** Calculate the weighted product of the deviations $(x - \mu_X)(y - \mu_Y)P_{X,Y}(x, y)$ for every $(x, y) \in \mathcal{X} \times \mathcal{Y}$ and take the sum.

In the above example, we have

$$\begin{aligned}\text{Cov}(X, Y) &= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} (x - \mu_X)(y - \mu_Y)P_{X,Y}(x, y) \\ &= (-0.5) \cdot (-1) \cdot 0.25 + (+0.5) \cdot (-1) \cdot 0.00 + (-0.5) \cdot 0 \cdot 0.25 + \cdots + 0.5 \cdot 1 \cdot 0.25\end{aligned}$$

Example of covariance calculation

Example

		$x - \mu_X$		$P_Y(y)$
		-0.5	+0.5	
$y - \mu_Y$	-1	0.25	0.00	0.25
	0	0.25	0.25	0.50
	+1	0.00	0.25	0.25
$P_X(x)$		0.50	0.50	

The joint PMF $P_{X,Y}$. The RVs X and Y have a positive covariance.

We can calculate the covariance $\text{Cov}(X, Y)$ of RVs X, Y from its joint PMF $P_{X,Y}$.

- **Step 2:** Calculate the weighted product of the deviations $(x - \mu_X)(y - \mu_Y)P_{X,Y}(x, y)$ for every $(x, y) \in \mathcal{X} \times \mathcal{Y}$ and take the sum.

In the above example, we have

$$\begin{aligned}\text{Cov}(X, Y) &= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} (x - \mu_X)(y - \mu_Y)P_{X,Y}(x, y) \\ &= (-0.5) \cdot (-1) \cdot 0.25 + (+0.5) \cdot (-1) \cdot 0.00 + (-0.5) \cdot 0 \cdot 0.25 + \dots + 0.5 \cdot 1 \cdot 0.25 = 0.25.\end{aligned}$$

The variance is a special case of the covariance

The covariance between a random variable and itself is the variance of the random variable. In other words:

Theorem

$$\text{Cov}(X, X) = \mathbb{V}[X]. \quad (32)$$

Covariance matrix

Definition

Let X_1, X_2, \dots, X_m be RVs. The $m \times m$ real matrix

$$\begin{bmatrix} \text{Cov}(X_1, X_1) & \text{Cov}(X_1, X_2) & \cdots & \text{Cov}(X_1, X_m) \\ \text{Cov}(X_2, X_1) & \text{Cov}(X_2, X_2) & \cdots & \text{Cov}(X_2, X_m) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_m, X_1) & \text{Cov}(X_m, X_2) & \cdots & \text{Cov}(X_m, X_m) \end{bmatrix} \quad (33)$$

is called the **covariance matrix** of RVs X_1, X_2, \dots, X_m .

Example of the covariance matrix

Let X and Y be random variables whose joint PMF $P_{X,Y}$ are given by the following table.

		x		$P_Y(y)$
		0	+1	
y	0	0.25	0.00	0.25
	+1	0.25	0.25	0.50
	+2	0.00	0.25	0.25
$P_X(x)$		0.50	0.50	

The joint PMF $P_{X,Y}$. The RVs X and Y have a positive covariance.

In the above example, $\text{Cov}(X,X) = \mathbb{V}[X] = 0.25$, $\text{Cov}(Y,Y) = \mathbb{V}[Y] = 0.5$, and $\text{Cov}(X,Y) = \text{Cov}(Y,X) = 0.25$.

Hence, the covariance matrix is
$$\begin{bmatrix} \text{Cov}(X,X) & \text{Cov}(X,Y) \\ \text{Cov}(Y,X) & \text{Cov}(Y,Y) \end{bmatrix} = \begin{bmatrix} 0.25 & 0.25 \\ 0.25 & 0.5 \end{bmatrix}.$$

Outline

2. Multiple Random Variables

2.1 Introduction: why are multiple random variables less trivial?

2.2 Joint distribution

2.3 Marginal distribution

2.4 Conditional distribution

2.5 Independence of random variables

2.6 Summary statistics for multiple RVs and covariance

2.7 Correlation

2.8 Exercises

Correlation

The covariance considers the scale of each random variable, not only the relation between them. Specifically, for $a, b \in \mathbb{R}$, we have that

$$\text{Cov}(aX, bY) = ab \text{Cov}(X, Y). \quad (34)$$

This implies that just multiplying the random variables by some factors changes the value of the correlation although the relation between aX and bY would be “qualitatively” the same as that of X and Y .

To see the “qualitative” relation between X and Y , we normalize it by dividing it by the covariance by the sum of the standard deviations of X and Y . The normalized covariance is called the **correlation coefficient** of X and Y .

Definition of the correlation coefficient

Definition (Correlation coefficient)

Let X and Y be random variables. The **correlation coefficient** $\text{corr}[X, Y]$ between X and Y is given by

$$\text{corr}[X, Y] := \frac{\text{Cov}[X, Y]}{\sigma[X]\sigma[Y]}. \quad (35)$$

The correlation coefficient is often denoted by ρ .

As expected, for positive real numbers a and b , we have that

$$\text{corr}[aX, bY] = \text{corr}[X, Y]. \quad (36)$$

Example of the correlation coefficient

Let X and Y be random variables whose joint PMF $P_{X,Y}$ are given by the following table.

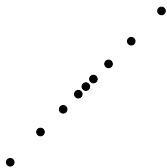
		x		$P_Y(y)$
		0	+1	
y	0	0.25	0.00	0.25
	+1	0.25	0.25	0.50
	+2	0.00	0.25	0.25
$P_X(x)$		0.50	0.50	

The joint PMF $P_{X,Y}$. The RVs X and Y have a positive covariance.

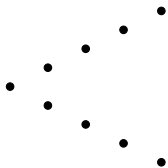
In the above example, $\text{Cov}(X, X) = \mathbb{V}[X] = 0.25$, $\text{Cov}(Y, Y) = \mathbb{V}[Y] = 0.5$, and $\text{Cov}(X, Y) = \text{Cov}(Y, X) = 0.25$.

Hence, the correlation coefficient between X and Y is $\text{corr}(X, Y) = \frac{0.25}{\sqrt{0.25}\sqrt{0.5}} = \frac{1}{\sqrt{2}}$.

Correlation coefficient examples



$$\rho = 1.0$$



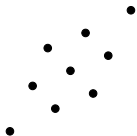
$$\rho = 0.0$$



$$\rho = 0.0$$



$$\rho = -1.0$$



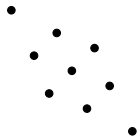
$$\rho = 0.735$$



$$\rho = 0.385$$



$$\rho = -0.385$$



$$\rho = -0.735$$

Independence implies no-correlation

Theorem

Let random variables X and Y be mutually independent. Then the covariance $\text{Cov}(X, Y)$ and the correlation $\text{corr}[X, Y]$ are zero.

Note: the converse of the above theorem is FALSE (see the next slide).

No correlation does NOT imply independence!

Example

Let X and Y be random variables whose joint PMF $P_{X,Y}$ are given by the following table.

		x		$P_Y(y)$
		-1	+1	
y	-1	0.0	0.25	0.25
	0	0.5	0.0	0.5
	+1	0.0	0.25	0.25
$P_X(x)$		0.5	0.5	

The joint PMF $P_{X,Y}$. The RVs X and Y are uncorrelated but mutually independent.

Then, the covariance $\text{Cov}(X, Y)$ and the correlation $\text{corr}[X, Y]$ are zero. However, X and Y are not independent. For example, $P_{X,Y}(-1, -1) \neq P_X(-1)P_Y(-1)$. The LHS is 0.0, while the RHS is $0.5 \times 0.25 = 0.125$.

No correlation does NOT imply independence!

Example

Let X and Y be random variables whose joint PMF $P_{X,Y}$ are given by the following table.

		x		$P_Y(y)$
		-1	+1	
y	-1	0.0	0.25	0.25
	0	0.5	0.0	0.5
	+1	0.0	0.25	0.25
$P_X(x)$		0.5	0.5	

The joint PMF $P_{X,Y}$. The RVs X and Y are uncorrelated but mutually independent.

Indeed, we cannot say Y increases as X increases since the expectation of Y is invariant when X . Hence, the correlation is zero. On the other hand, the variance of Y is 0 when $X = -1$ but it is non-zero if $X = +1$, hence X has some information about Y . These are intuitive explanations of zero correlation and non-independence of X and Y .

Correlation \neq Causality

If two random variables X and Y have a correlation, i.e., $\text{corr}[X, Y] \neq 0$, you might expect that X is the cause of Y .

However, there are many possibilities behind the correlation, e.g.,

1. X is a cause of Y .
2. Y is a cause of X .
3. There exists a random variable Z that causes the both X and Y .
4. (When we estimate the correlation coefficient) There is no relation between X and Y but our estimation of the correlation coefficient is non-zero by estimation errors.

Hence, we cannot conclude that X is a cause of Y just by $\text{corr}[X, Y] \neq 0$.

Outline

2. Multiple Random Variables

2.1 Introduction: why are multiple random variables less trivial?

2.2 Joint distribution

2.3 Marginal distribution

2.4 Conditional distribution

2.5 Independence of random variables

2.6 Summary statistics for multiple RVs and covariance

2.7 Correlation

2.8 Exercises

Exercise (Independent RVs (1))

Let X and Y be discrete RVs, whose supports are $\mathcal{X} = \{0, 1, 2, 3\}$ and $\mathcal{Y} = \{0, 1, 2\}$, respectively. Suppose that the joint PMF $P_{X,Y}$ is given by the following table.

		x			
		0	1	2	3
y	0	0.08	0.02	0.04	0.06
	1	0.20	0.05	0.10	0.15
	2	0.12	0.03	0.06	0.09

An example of $P_{X,Y}(x,y) := \Pr(X=x \wedge Y=y)$

- (i) Find the values of the marginal PMFs $P_X(x)$ and $P_Y(y)$ for every $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.
- (ii) Judge whether X and Y are mutually independent or not.

Exercise (Independent RVs (2))

Let X and Y be discrete RVs, whose supports are $\mathcal{X} = \{0, 1, 2, 3\}$ and $\mathcal{Y} = \{0, 1, 2\}$, respectively. Suppose that the joint PMF $P_{X,Y}$ is given by the following table.

		x			
		0	1	2	3
y	0	0.10	0.02	0.02	0.06
	1	0.24	0.08	0.12	0.06
	2	0.06	0.00	0.06	0.18

An example of $P_{X,Y}(x,y) := \Pr(X=x \wedge Y=y)$

- (i) Find the values of the marginal PMFs $P_X(x)$ and $P_Y(y)$ for every $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.
- (ii) Judge whether X and Y are mutually independent or not.

Exercise (Joint PMF of independent RVs)

Let X and Y be mutually independent discrete RVs, whose supports are $\mathcal{X} = \{0, 1, 2, 3\}$ and $\mathcal{Y} = \{0, 1, 2\}$, respectively. Suppose that those marginal PMFs, P_X and P_Y , are given by the following tables.

	x			
	0	1	2	3
$P_X(x)$	0.40	0.10	0.20	0.30

	y		
	0	1	2
$P_Y(y)$	0.20	0.50	0.30

The marginal PMFs $P_X(x)$ and $P_Y(y)$.

Find the values of the Joint PMF $P_{X,Y}(x,y)$ for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.

Exercise (Exercise: conditional PMF calculation)

Let X and Y be discrete RVs, whose supports are $\mathcal{X} = \{0, 1, 2, 3\}$ and $\mathcal{Y} = \{0, 1, 2\}$, respectively. Suppose that the joint PMF $P_{X,Y}$ is given by the following table.

		x			
		0	1	2	3
y	0	0.08	0.02	0.04	0.06
	1	0.20	0.05	0.10	0.15
	2	0.12	0.03	0.06	0.09

An example of $P_{X,Y}(x,y) := \Pr(X=x \wedge Y=y)$

- (i) Find the values of the marginal PMFs $P_X(x)$ and $P_Y(y)$ for every $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.
- (ii) Find the values of the conditional PMFs $P_{X|Y}(x|y)$ and $P_{Y|X}(y|x)$ for every $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.

Exercise (Exercise: conditional PMF calculation)

Let X and Y be discrete RVs, whose supports are $\mathcal{X} = \{0, 1, 2, 3\}$ and $\mathcal{Y} = \{0, 1, 2\}$, respectively. Suppose that the joint PMF $P_{X,Y}$ is given by the following table.

		x			
		0	1	2	3
y	0	0.10	0.02	0.02	0.06
	1	0.24	0.08	0.12	0.06
	2	0.06	0.00	0.06	0.18

An example of $P_{X,Y}(x,y) := \Pr(X=x \wedge Y=y)$

- (i) Find the values of the marginal PMFs $P_X(x)$ and $P_Y(y)$ for every $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.
- (ii) Find the values of the conditional PMFs $P_{X|Y}(x|y)$ and $P_{Y|X}(y|x)$ for every $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.

Exercise (Marginal distribution)

Suppose that the joint PMF of random variables X and Y is given by the following table. Find the marginal PMFs of X and Y .

		x				$P_Y(y)$
		0	1	2	3	
y	0	0.10	0.02	0.02	0.06	$P_Y(0) = ?$
	1	0.24	0.08	0.12	0.06	$P_Y(1) = ?$
	2	0.06	0.00	0.06	0.18	$P_Y(2) = ?$
$P_X(x)$		$P_X(0) = ?$	$P_X(1) = ?$	$P_X(2) = ?$	$P_X(3) = ?$	

An example of $P_{X,Y}(x,y) := \Pr(X = x \wedge Y = y)$

Outline

3. Continuous Random Variables

3.1 Introduction: why are continuous random variables less trivial?

3.2 Probability density function

3.3 Area, integration, and properties of PDF.

3.4 Calculating integral

3.5 Summary statistics of continuous RV and integral

3.6 Multivariate random variables and multiple integral

3.7 Relation among RVs in a continuous multivariate RV

3.8 Exercises

Outline

3. Continuous Random Variables

3.1 Introduction: why are continuous random variables less trivial?

3.2 Probability density function

3.3 Area, integration, and properties of PDF.

3.4 Calculating integral

3.5 Summary statistics of continuous RV and integral

3.6 Multivariate random variables and multiple integral

3.7 Relation among RVs in a continuous multivariate RV

3.8 Exercises

Continuous random variables in real AI applications

A discrete RV can take only limited values. However, many real-world phenomena are represented as random variables which can take any real value in a continuous section.

- Inflation rate (economics),
- Position of a vehicle,
- The brightness of scenery,
- The intensity of an acoustic signal,
- Density of air pollution.

Hence, when we want to analyze those phenomena using probability theory, we cannot always use mathematical tools to handle discrete RVs.

For example, those random variables typically have **no probability mass function (PMF)**.

A random variable may not have a PMF.

Consider a simple random variable uniformly distributed in $[0, 1]$. Here $\Pr(0 \leq X \leq 1) = 1$.

This random variable have nowhere probability mass, i.e., $\Pr(X = x) = 0$. for any $X \in \mathbb{R}$.

Proof.

Since its support is $[0, 1]$, it is trivial that $\Pr(X = x) = 0$ for $x \notin [0, 1]$. For $x \in [0, 1]$, assume, for the sake of contradiction, that $\Pr(X = x) = \epsilon$, where $\epsilon > 0$. From its uniformity, if $\Pr(X = x) = \epsilon$ holds for one value $x \in [0, 1]$, then it holds for all $x \in [0, 1]$. Hence, if $A \subset [0, 1]$ and A has at least N elements, $\Pr(X \in A) \leq N\epsilon$. However, there are an infinite number of real numbers in $[0, 1]$, so $\Pr(X \in [0, 1])$ is infinity. It contradicts $\Pr(X \in [0, 1]) = 1$. □

A random variable may not have a PMF.

Consider a simple random variable uniformly distributed in $[0, 1]$. Here $\Pr(0 \leq X \leq 1) = 1$.

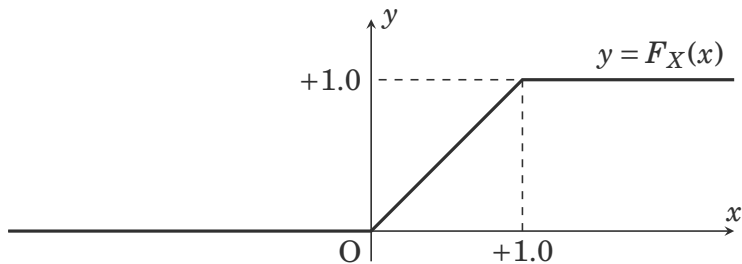
This random variable have nowhere probability mass, i.e., $\Pr(X = x) = 0$. for any $X \in \mathbb{R}$.

Other random variables whose support is a section in the real line have the same problem. Hence, we need another way to represent a random variable.

Fortunately, any univariate random variable has a cumulative distribution function (CDF)

Example of CDF for a non-discrete RV

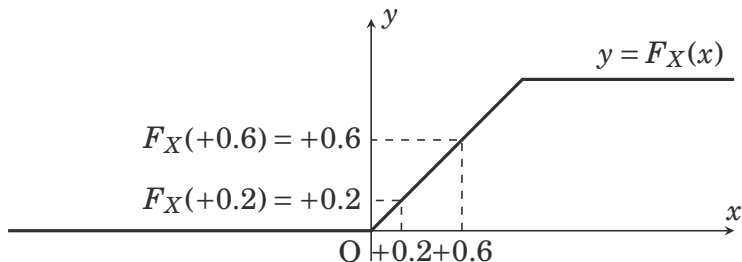
The CDF of a random variable X uniformly distributed in $[0, 1]$ is:



$$F_X(x) = \begin{cases} 0 & \text{if } x \leq 0, \\ x & \text{if } 0 \leq x \leq 1, \\ 1 & \text{if } x \geq 1. \end{cases} \quad (37)$$

Example of CDF for a non-discrete RV

The CDF of a random variable X uniformly distributed in $[0, 1]$ is:

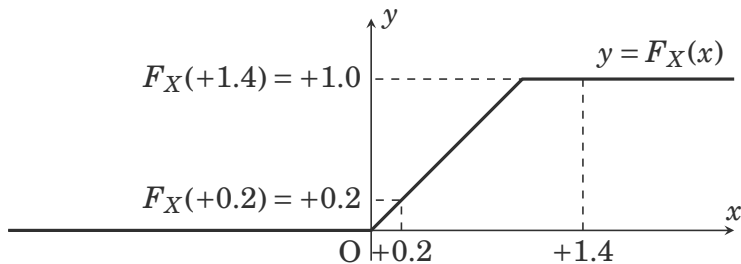


Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(0.2 < X \leq 0.6) &= \Pr(X \leq 0.6) - \Pr(X \leq 0.2) \\ &= F_X(0.6) - F_X(0.2) \\ &= 0.6 - 0.2 = 0.4.\end{aligned}\tag{37}$$

Example of CDF for a non-discrete RV

The CDF of a random variable X uniformly distributed in $[0, 1]$ is:

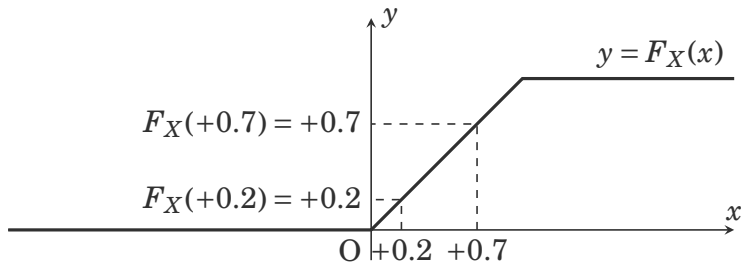


Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(0.2 < X \leq 1.4) &= \Pr(X \leq 1.4) - \Pr(X \leq 0.2) \\ &= F_X(1.4) - F_X(0.2) \\ &= 1.0 - 0.2 = 0.8.\end{aligned}\tag{37}$$

Example of CDF for a non-discrete RV

The CDF of a random variable X uniformly distributed in $[0, 1]$ is:



Using the CDF, we can calculate the probability of various events. For example,

$$\begin{aligned}\Pr(0.2 \leq X \leq 0.7) &= \Pr(X \leq 0.7) - \lim_{x \nearrow 0.2} \Pr(x) \\ &= F_X(0.7) - \lim_{x \nearrow 0.2} F_X(x) \\ &= 0.7 - 0.2 = 0.5.\end{aligned}\tag{37}$$

Why are we not satisfied with the CDF?

However, the CDF is not always welcomed. It is because

- The CDF is not intuitive. At one glance, we do not know around which value the random variable tends to take a value.
- The CDF can be extremely complex even for a practically important distribution.

Although there exists no PMF for a continuous RV in general, we want to indicate which values the RV tends to take frequently as the PMF does for a discrete RV.

The ***probability density function (PDF)*** achieves this objective.

Learning outcomes

By the end of this section, you should be able to:

- Explain what a probability density function represents,
- Explain the relation between the probability density function and cumulative distribution function,
- Calculate the probability of an event using the integral and the probability density function, and
- Calculate summary statistics of continuous random variables.

Notation: sections

In the following, \mathbb{R} , $\mathbb{R}_{\geq 0}$, and $\mathbb{R}_{>0}$ are the sets of real numbers, nonnegative real numbers, and positive real numbers, respectively.

Let a and b be real values. By $[a, b]$, (a, b) , we denote the closed and open sections defined by

- $[a, b] = \{x \in \mathbb{R} | a \leq x \leq b\},$
- $(a, b) = \{x \in \mathbb{R} | a < x < b\},$

respectively. Likewise, by $(a, b]$ and $[a, b)$, we denote the semi-open sets defined by

- $(a, b] = \{x \in \mathbb{R} | a < x \leq b\},$
- $[a, b) = \{x \in \mathbb{R} | a \leq x < b\},$

Notation: Napier's constant and the exponential function

The real number constant e , called **Napier's constant** or **Euler's number**, is defined by
$$e := \lim_{n \rightarrow +\infty} \left(1 + \frac{1}{n}\right)^n.$$
 Note that $e = 2.718281828\dots$ and is the only real value that satisfies $\frac{d}{dx}e^x = e^x$.

We define the **(natural) exponential function** $\exp : \mathbb{R} \rightarrow \mathbb{R}_{>0}$ by $\exp(x) = e^x$.

Outline

3. Continuous Random Variables

3.1 Introduction: why are continuous random variables less trivial?

3.2 Probability density function

3.3 Area, integration, and properties of PDF.

3.4 Calculating integral

3.5 Summary statistics of continuous RV and integral

3.6 Multivariate random variables and multiple integral

3.7 Relation among RVs in a continuous multivariate RV

3.8 Exercises

Idea of the probability density function

As we have seen in the case of the uniform distribution in the section $[0, 1]$, the probability $\Pr(X = c)$ might be zero for a real value c in many cases. In this case, we cannot say which values the RV tend to take more frequently than others.

Idea of the probability density function

As we have seen in the case of the uniform distribution in the section $[0, 1]$, the probability $\Pr(X = c)$ might be zero for a real value c in many cases. In this case, we cannot say which values the RV tend to take more frequently than others.

Hence, we evaluate the probability of the RV taking a value **in a section**. For example, instead of evaluating $\Pr(X = c)$, we evaluate the probability $\Pr(a < X \leq b)$ for real values a, b around c such that $a < b$. If the probability is high and the section length $b - a$ is short, we can say that the RV X takes a value around c frequently.

Idea of the probability density function

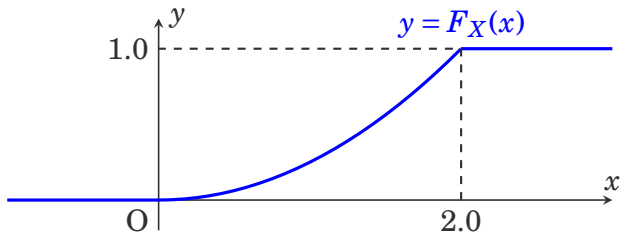
As we have seen in the case of the uniform distribution in the section $[0, 1]$, the probability $\Pr(X = c)$ might be zero for a real value c in many cases. In this case, we cannot say which values the RV tend to take more frequently than others.

Hence, we evaluate the probability of the RV taking a value **in a section**. For example, instead of evaluating $\Pr(X = c)$, we evaluate the probability $\Pr(a < X \leq b)$ for real values a, b around c such that $a < b$. If the probability is high and the section length $b - a$ is short, we can say that the RV X takes a value around c frequently.

So, we can regard the probability per the section length as the **density** of the probability distribution of the RV X around the section. A high density around a value c indicates that the X tends to take a value around c .

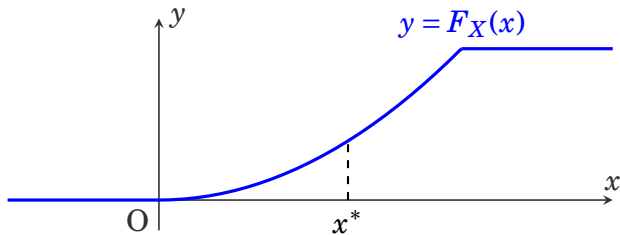
Based on the above idea, we can formulate the **probability density function (PDF)** from the cumulative distribution function (CDF) as follows.

CDF to the density.



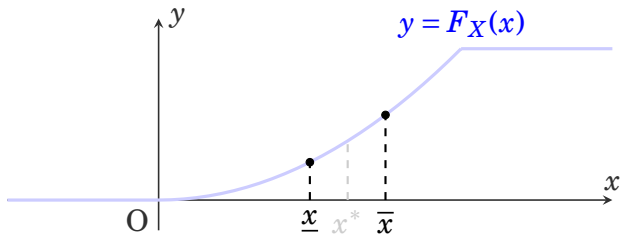
The CDF of a RV X .

CDF to the density.



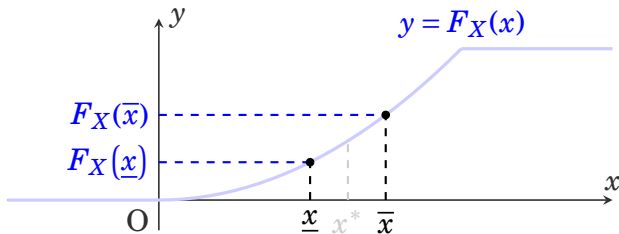
Suppose we want to know how frequently the RV X takes a value “around” x^* .

CDF to the density.



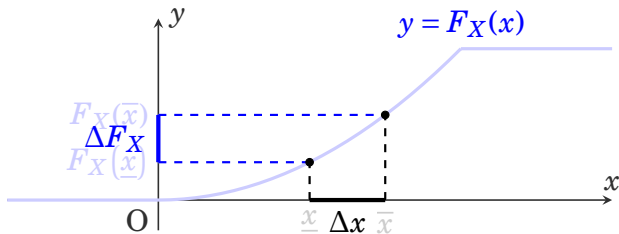
We consider an interval $[\underline{x}, \bar{x}]$ including x^* .

CDF to the density.



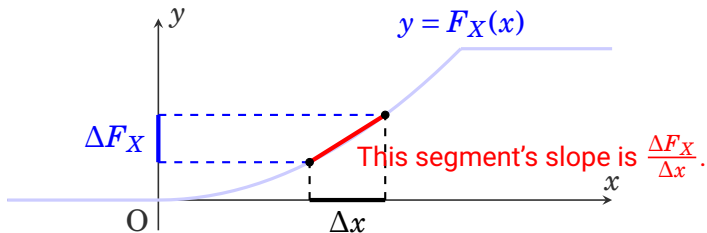
We find the probability $\Pr(X \in [\underline{x}, \bar{x}])$, given by $F_X(\bar{x}) - F_X(\underline{x})$.

CDF to the density.



Define $\Delta x := \bar{x} - \underline{x}$ and $\Delta F_X := F_X(\bar{x}) - F_X(\underline{x}) = \Pr(X \in [\underline{x}, \bar{x}])$.

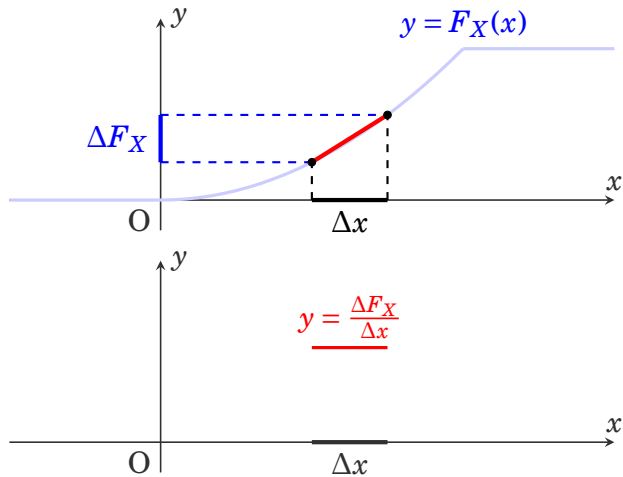
CDF to the density.



Define $\Delta x := \bar{x} - \underline{x}$ and $\Delta F_X := F_X(\bar{x}) - F_X(\underline{x}) = \Pr(X \in [\underline{x}, \bar{x}])$.

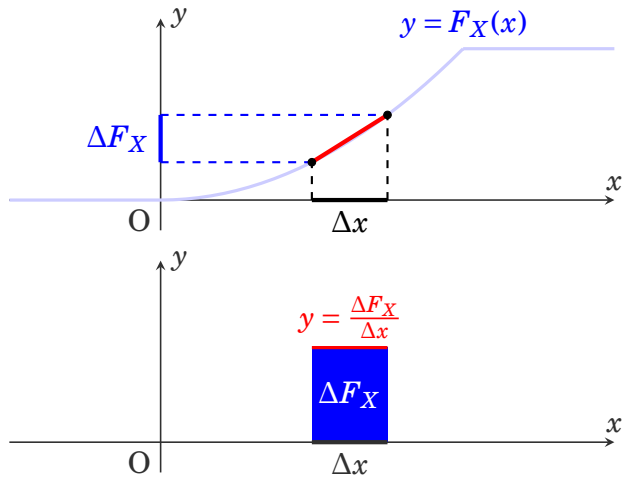
The RV X tends to take a value around x^*
if the probability per length $\frac{\Delta F_X}{\Delta x}$, or the "density" is large.

CDF to the density.



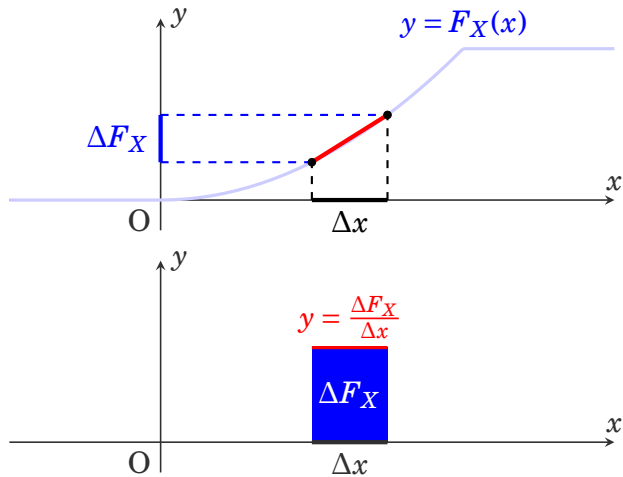
The graph plot of the probability per length $\frac{\Delta F_X}{\Delta x}$, or the “density”

CDF to the density.



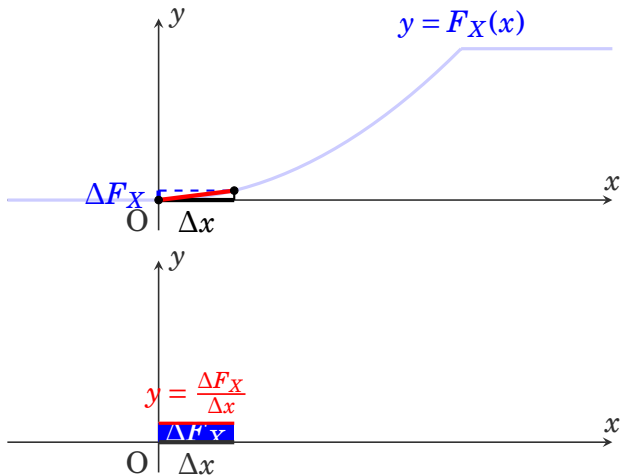
The area under the graph of $\frac{\Delta F_X}{\Delta x}$ is given by $\Delta x \cdot \frac{\Delta F_X}{\Delta x} = \Delta F_X$.

CDF to the density.



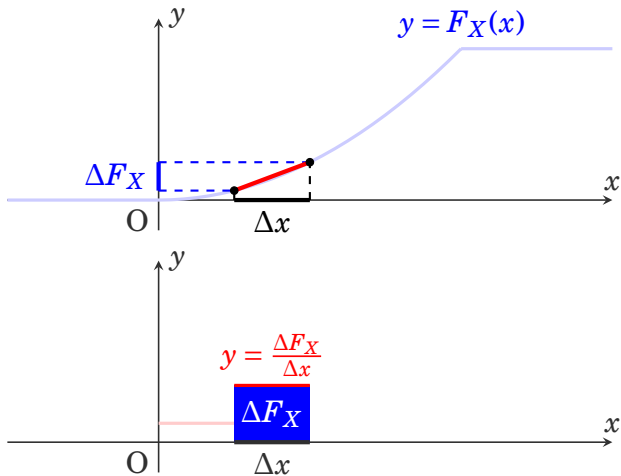
Let's see how the density varies in other sections.

CDF to the density.



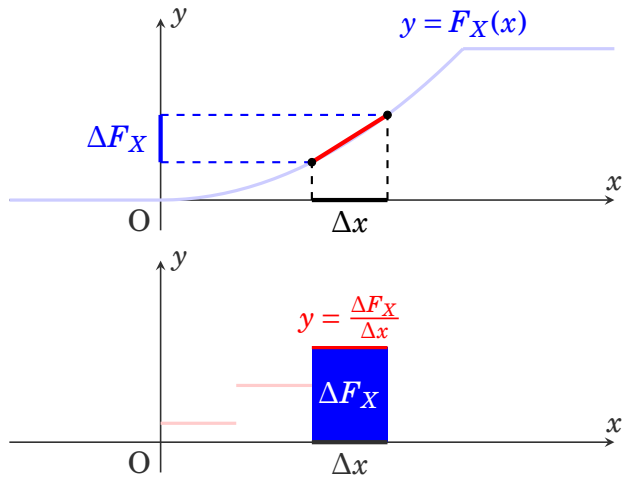
The density at the section $[0.00, 0.50]$ is small, reflecting the gentle slope of the graph of F_X .

CDF to the density.



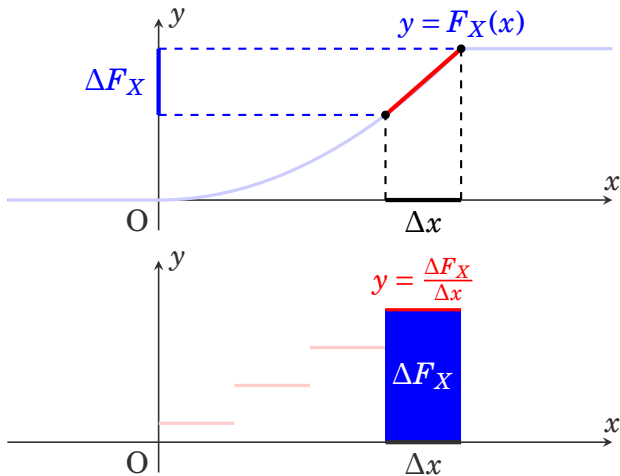
The density at the section $[0.50, 1.00]$ is larger, reflecting the steeper slope of the graph of F_X .

CDF to the density.



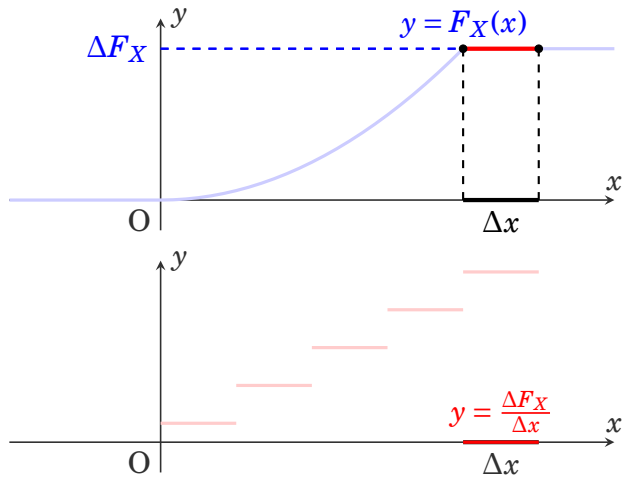
The density at the section $[1.00, 1.50]$ is even larger, reflecting the steeper slope of the graph of F_X .

CDF to the density.



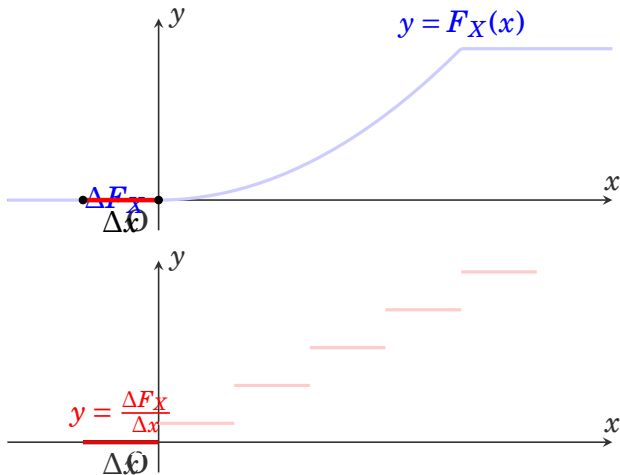
The density at the section $[1.50, 2.00]$ is even larger, reflecting the steeper slope of the graph of F_X .

CDF to the density.



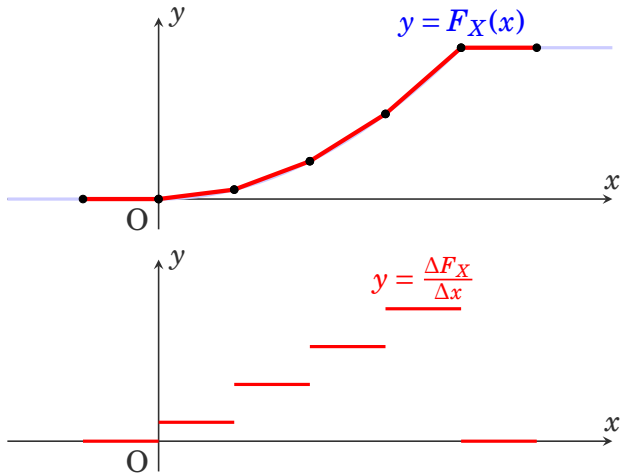
The density at the section $[2.00, 2.50]$ is zero since the graph of F_X is horizontal.

CDF to the density.



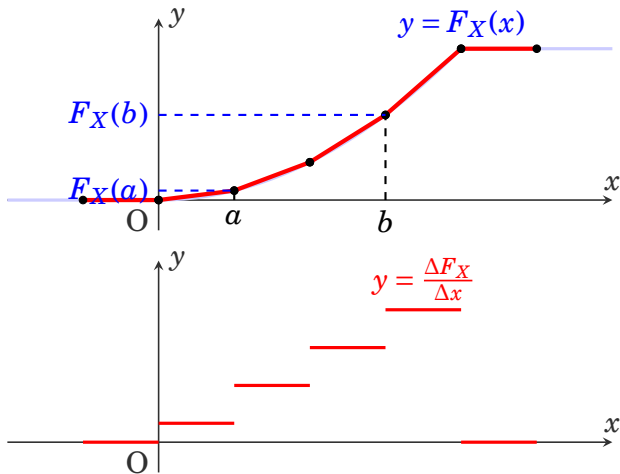
The density at the section $[-0.50, 0.90]$ is zero since the graph of F_X is horizontal.

CDF to the density.



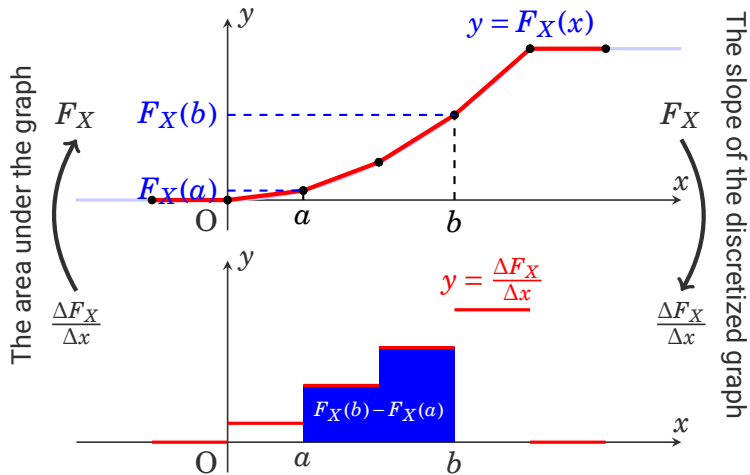
We have the graph of the piece-wise density of the distribution at each section.

CDF to the density.



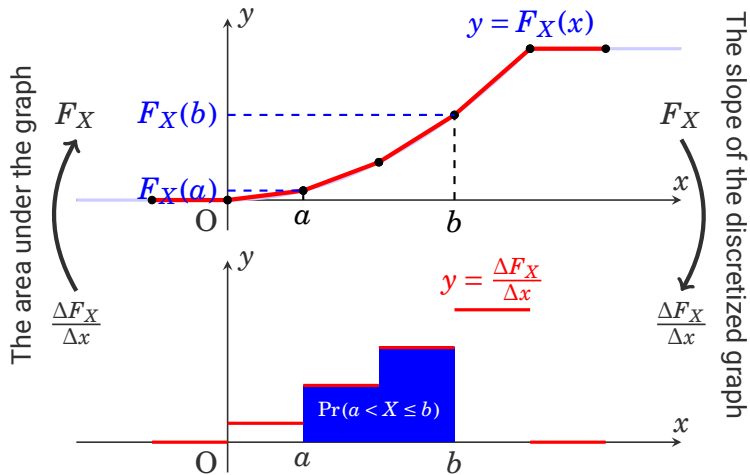
Using the piece-wise density, we can calculate $\Pr(a < X \leq b)$ if each of a and b is an end point of a section.

CDF to the density.



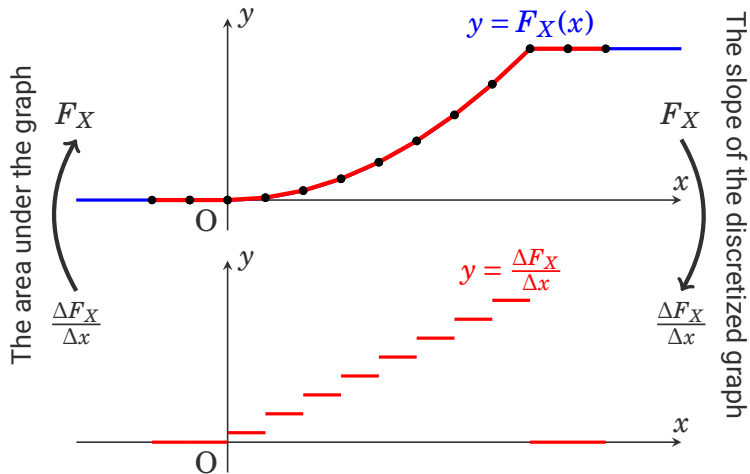
Actually, the probability $\Pr(a < X \leq b) = F_X(b) - F_X(a)$ is the area under the piece-wise density graph.

CDF to the density.



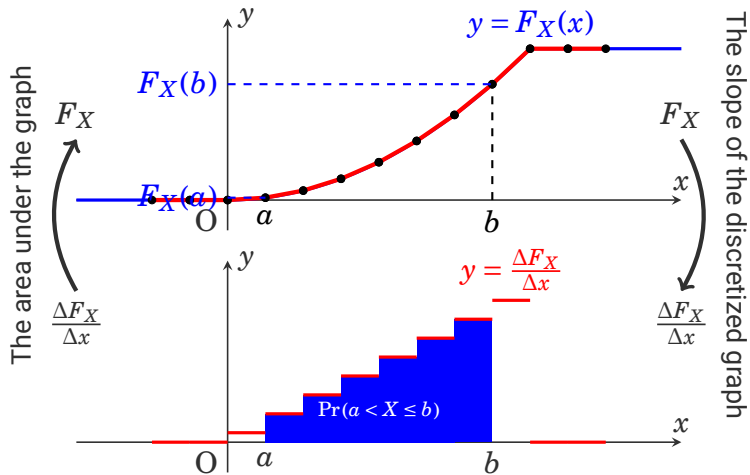
Actually, the probability $\Pr(a < X \leq b) = F_X(b) - F_X(a)$ is the area under the piece-wise density graph.

CDF to the density.



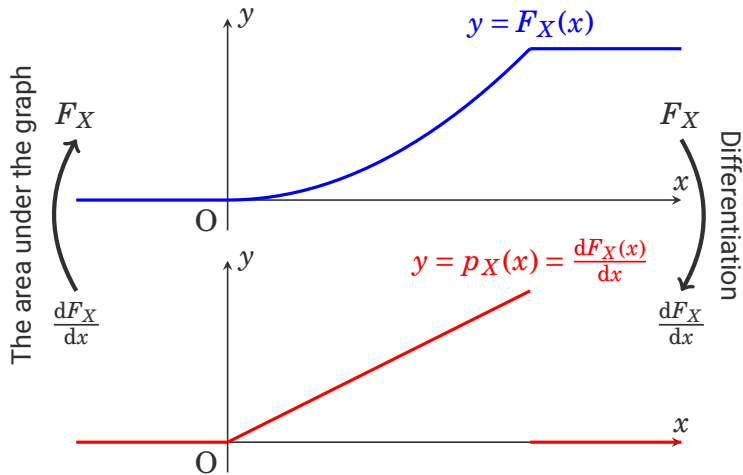
If the sections are shorter, we can evaluate $\Pr(a < X \leq b)$ for more pairs of a and b .

CDF to the density.



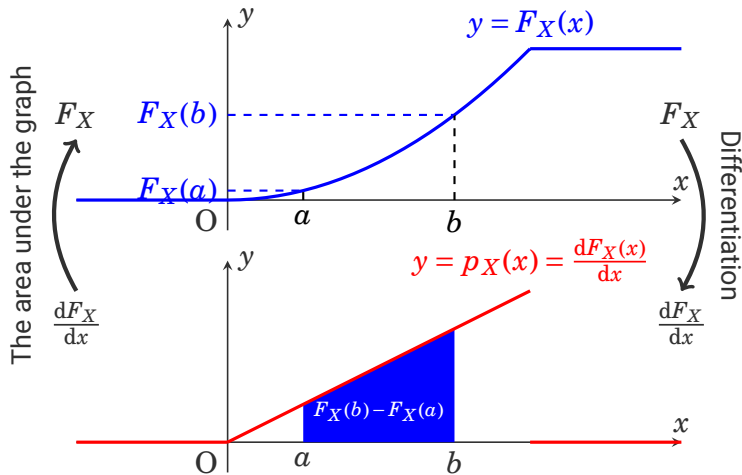
If the sections are shorter, we can evaluate $\Pr(a < X \leq b)$ for more pairs of a and b .

CDF to the density.



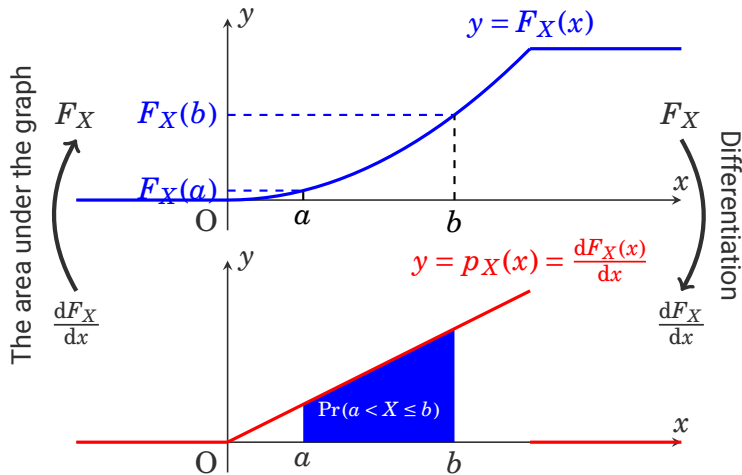
If we consider infinitely short sections, the density $\frac{\Delta F_X}{\Delta x}$ converges to the derivative $\frac{dF_X}{dx}$.

CDF to the density.



Actually, the probability $\Pr(a < X \leq b) = F_X(b) - F_X(a)$ is the area under the graph of $\frac{dF_X}{dx}$.

CDF to the density.



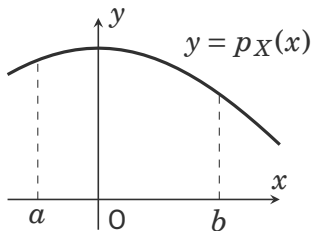
Actually, the probability $\Pr(a < X \leq b) = F_X(b) - F_X(a)$ is the area under the graph of $\frac{dF_X}{dx}$.

Probability density function (PDF)

Definition (Probability density function and continuous random variable)

Let X be a RV. A function $p_X : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$ is called a **probability density function (PDF)** of X if the probability $\Pr(a < X \leq b)$ equals to the area bounded by the graph of $y = p_X(x)$ and $y = 0$ between $x = a$ and $x = b$ for all a and b such that $a \leq b$.

If a RV has at least one PDF, the RV is called a **continuous random variable**.

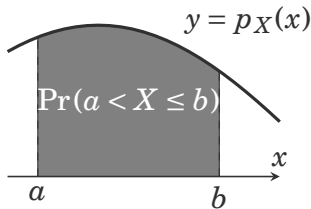


If p_X is a PDF of X , the probability $\Pr(a < X \leq b)$ is given by the area under the PDF in the domain $(a, b]$.

Probability density function (PDF)

Definition (Probability density function and continuous random variable)

Let X be a RV. A function $p_X : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$ is called a **probability density function (PDF)** of X if the probability $\Pr(a < X \leq b)$ equals to the area bounded by the graph of $y = p_X(x)$ and $y = 0$ between $x = a$ and $x = b$ for all a and b such that $a \leq b$. If a RV has at least one PDF, the RV is called a **continuous random variable**.



If p_X is a PDF of X , the probability $\Pr(a < X \leq b)$ is given by the area under the PDF in the domain $(a, b]$.

A continuous RV has nowhere a “mass.”

The area under a curve in a zero-length section is zero. Hence, if a RV is continuous, it has no probability mass anywhere. That is,

Theorem

If X is a continuous RV, the probability $\Pr(X = c)$ is zero for any $c \in \mathbb{R}$.

Hence, when we discuss a continuous RV, we do not need to discuss whether or not a section includes the endpoints. That is,

Corollary

Let X be a continuous RV and a and b be real values such that $a < b$. Then we have,

$$\Pr(a \leq X \leq b) = \Pr(a < X \leq b) = \Pr(a \leq X < b) = \Pr(a < X < b). \quad (38)$$

Hence, we can replace $a < X \leq b$ with $a \leq X \leq b$ or another in the definition of the PDF⁶.

⁶To see this strictly, we need Carathéodory's extension theorem on semiring of sets

Note: the end-points are not ignorable for a discrete RV.

A discrete RV has a probability mass on any value in its support. Hence, for example, $\Pr(a \leq X \leq b) \neq \Pr(a < X \leq b)$ in general.

For example, if X is the value when we roll an ideal six-sided dice,
 $\Pr(3 \leq X \leq 6) = \frac{4}{6} \neq \Pr(3 < X \leq 6) = \frac{3}{6}.$

CDF and PDF

Assume that the CDF is differentiable at all the points on the real number line except for finite points. As we can see in the construction of the PDF from the CDF, we can get the PDF by differentiating the CDF.

In practice, we usually know the PDF in advance but the CDF is unknown. Hence, we need to understand how to evaluate the area bounded by the graph of a general PDF.

Outline

3. Continuous Random Variables

3.1 Introduction: why are continuous random variables less trivial?

3.2 Probability density function

3.3 Area, integration, and properties of PDF.

3.4 Calculating integral

3.5 Summary statistics of continuous RV and integral

3.6 Multivariate random variables and multiple integral

3.7 Relation among RVs in a continuous multivariate RV

3.8 Exercises

How to mathematically calculate the area under the curve?

Let X be a continuous RV and p_X be its PDF. Recall that the probability $\Pr(a \leq X \leq b)$ is given by the area under the graph of PDF p_X in the section $[a, b]$.

Hence, we need a mathematical tool to evaluate the area under the curve of a function in general.

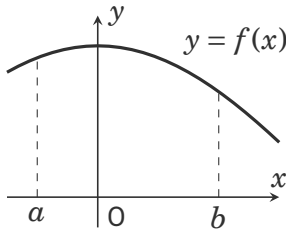
Integration is the area to discuss the area under the graph of a function, (or the volume under the graph of a function in higher-dimensional space). We will learn it in the following.

Definite Integral

Suppose that $a \leq b$.

The (signed) area bounded by the graph of $y = f(x)$ and $y = 0$ between $x = a$ and $x = b$ is called the **definite integral** of f between a and b , which is denoted by $\int_a^b f(x) \mathrm{d}x$.

We also define $\int_b^a f(x) \mathrm{d}x := -\int_a^b f(x) \mathrm{d}x$.



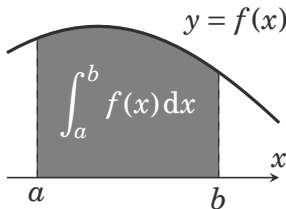
The definite integral is the area bounded by the graph of the function.

Definite Integral

Suppose that $a \leq b$.

The (signed) area bounded by the graph of $y = f(x)$ and $y = 0$ between $x = a$ and $x = b$ is called the **definite integral** of f between a and b , which is denoted by $\int_a^b f(x) \, dx$.

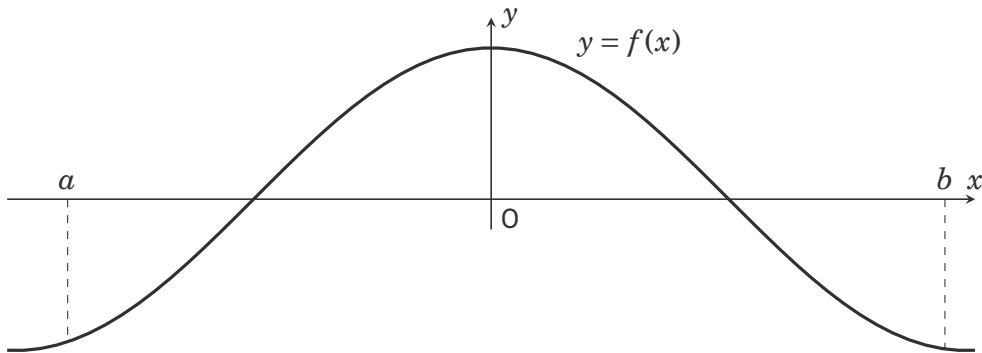
We also define $\int_b^a f(x) \, dx := -\int_a^b f(x) \, dx$.



The definite integral is the area bounded by the graph of the function.

Definite Integral: When the function takes negative values

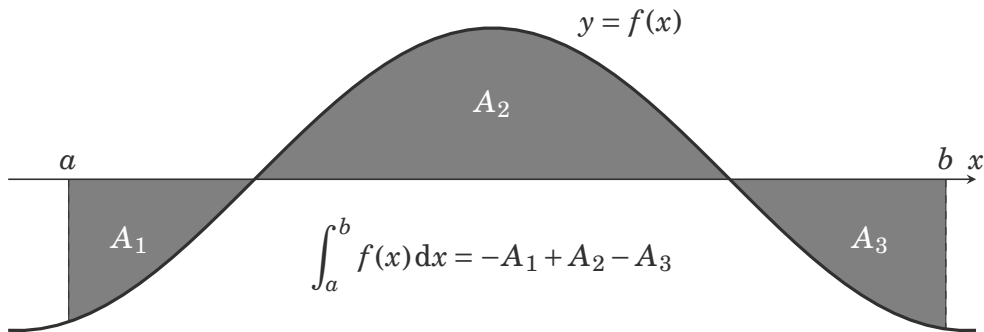
Areas bounded by the graph taking negative values are counted as negative values.



Areas bounded by the graph taking negative values are counted as negative values.

Definite Integral: When the function takes negative values

Areas bounded by the graph taking negative values are counted as negative values.



Areas bounded by the graph taking negative values are counted as negative values.

Any continuous RV has a nonnegative PDF

Assume that X is a continuous RV let p_X be a PDF of X . The probability

$\Pr(a < X \leq b) = \int_a^b p_X(x) dx$ is always nonnegative, so we expect the PDF p_X to be a nonnegative function.

Strictly speaking, a PDF of a RV is not unique, since the area bounded by the graph does not change even if we change the value of the function at finite or countable points⁷.

Nevertheless, if a RV has a PDF, we can assume that it is a nonnegative function without loss of generality.

Theorem

*Let X be a continuous RV, i.e., there is a PDF of X . Then, there exists a **nonnegative** PDF of X , i.e, a PDF p_X such that $p_X(x) \geq 0$ at any $x \in \mathbb{R}$.*

In the following, we always assume that a PDF is nonnegative.

⁷Strictly speaking, at points in a set with measurement zero.

Probability of a RV being in a complicated shape

If we want to calculate the probability $\Pr(X \in A)$, where $A \subset \mathbb{R}$ has a complicated shape, we can calculate it using the sum rule.

Specifically, suppose that we have a decomposition $A = \bigcup_{i=1}^n (a_i, b_i]$, where $(a_i, b_i] \cap (a_j, b_j] = \emptyset$. Then, we have that

$$\Pr(X \in A) = \sum_{i=1}^n \Pr(a_i < X \leq b_i). \quad (39)$$

If X is a continuous RV and p_X is its PDF, the above value equals $\sum_{i=1}^n \int_{a_i}^{b_i} p_X(x) dx$.

The same discussion holds even if the decomposition includes open sections like (a_i, b_i) or closed sections like $[a_j, b_j]$.

Note that the above calculation is not always correct if the decomposition includes an uncountably infinite number of sections.

Probability of a RV being in an infinite length section

If we need to evaluate the probability $\Pr(a < X)$, what we do is consider $\Pr(a < X \leq b)$ for an infinitely large b . Hence, we have that $\Pr(a < X) = \lim_{b \rightarrow +\infty} \Pr(a < X \leq b)$. The reverse holds for $\Pr(X \leq b)$. In other words, we can evaluate those probabilities by taking the limit of a definite integral as follows.

Theorem

Let X be a continuous RV, whose PDF is p_X , and a and b be real values. Then,

- $\Pr(a < X) = \Pr(a \leq X) = \lim_{b \rightarrow +\infty} \int_a^b p_X(x) dx,$
- $\Pr(X < b) = \Pr(X \leq b) = \lim_{a \rightarrow -\infty} \int_a^b p_X(x) dx.$

The “sum” of the PDF is one.

The section $(a, 0]$ includes all the nonpositive numbers if a is infinitely small and the section $(0, b]$ includes all the positive numbers b is infinitely large. Since a continuous RV X always takes a real value, the sum of the probabilities $\Pr(a < X \leq 0) + \Pr(0 < X \leq b)$ is 1 if a is infinitely small and b is infinitely large. Hence, the following always hold.

Theorem

Let X be a continuous RV whose PDF is p_X . We have that

$$\lim_{a \rightarrow -\infty} \int_a^0 p_X(x) dx + \lim_{b \rightarrow +\infty} \int_0^b p_X(x) dx = 1 \quad (40)$$

The above property is similar to a property of the probability mass function (PMF) of a discrete RV. To see that, we will introduce the ***improper integral***.

Improper integral

As we have seen, we often want to calculate limits of the definite integral. We call them ***improper integrals***, and use special notations as follows.

Definition (Improper integrals)

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a function and a and b be real values. We define the value

$\int_a^{+\infty} f(x) dx$, $\int_{-\infty}^b f(x) dx$, and $\int_{-\infty}^{+\infty} f(x) dx$ by the following.

- $\int_a^{+\infty} f(x) dx := \lim_{b \rightarrow +\infty} \int_a^b f(x) dx,$
- $\int_{-\infty}^b f(x) dx := \lim_{a \rightarrow -\infty} \int_a^b f(x) dx,$
- $\int_{-\infty}^{+\infty} f(x) dx := \int_{-\infty}^0 f(x) dx + \int_0^{+\infty} f(x) dx.$

Interpretation of improper integrals

We can regard improper integrals $\int_a^{+\infty} f(x) dx$, $\int_{-\infty}^b f(x) dx$, and $\int_{-\infty}^{+\infty} f(x) dx$ as the signed areas bounded by a graph of f in the section $(a, +\infty)$, $(-\infty, b)$, and $(-\infty, +\infty)$, respectively.

Rewriting properties of the PDF using improper integrals

We can rewrite the properties of the PDF in previous slides as follows.

Theorem

Let X be a continuous RV, whose PDF is p_X , and a and b be real values. Then,

- $\Pr(a < X) = \Pr(a \leq X) = \int_a^{+\infty} p_X(x) dx,$
- $\Pr(X < b) = \Pr(X \leq b) = \int_{-\infty}^b p_X(x) dx,$
- $\int_{-\infty}^{+\infty} p_X(x) dx = 1.$

The third property is similar to a property of the PMF: $\sum_{x \in \mathcal{X}} P_X(x) = 1$, where X is a discrete RV, \mathcal{X} is its support and P_X is its PMF.

Properties of the definite integral

Let a, b, c be real numbers and f and g be functions of a real value.

- $\int_b^a f(x) \mathrm{d}x := -\int_a^b f(x) \mathrm{d}x$ (by definition),
- $\int_a^a f(x) \mathrm{d}x = 0$ (The area is zero in a zero length section.).
- $\int_a^c f(x) \mathrm{d}x + \int_c^b f(x) \mathrm{d}x = \int_a^b f(x) \mathrm{d}x$ (horizontal concatenation).

Outline

3. Continuous Random Variables

3.1 Introduction: why are continuous random variables less trivial?

3.2 Probability density function

3.3 Area, integration, and properties of PDF.

3.4 Calculating integral

3.5 Summary statistics of continuous RV and integral

3.6 Multivariate random variables and multiple integral

3.7 Relation among RVs in a continuous multivariate RV

3.8 Exercises

Calculating definite integrals

Let X be a continuous RV and p_X be its PDF. Since the probability $\Pr(a < X \leq b)$ is given by the definite integral $\int_a^b p(x) \mathrm{d}x$, we need to know **how to calculate definite integrals** to understand the behavior of the continuous RV X .

⁸e.g., the trapezoidal rule, the Gauss-Legendre quadrature rule, the double exponential formula

Calculating definite integrals

Let X be a continuous RV and p_X be its PDF. Since the probability $\Pr(a < X \leq b)$ is given by the definite integral $\int_a^b p(x) dx$, we need to know **how to calculate definite integrals** to understand the behavior of the continuous RV X .

There are two directions to calculate definite integrals.

- **Numerical integration** by approximating the area by shapes of which we can calculate the area easier.
- **Analytical integration** by conducting integration as the inverse operation of differentiation.

⁸e.g., the trapezoidal rule, the Gauss-Legendre quadrature rule, the double exponential formula

Calculating definite integrals

Let X be a continuous RV and p_X be its PDF. Since the probability $\Pr(a < X \leq b)$ is given by the definite integral $\int_a^b p(x) dx$, we need to know **how to calculate definite integrals** to understand the behavior of the continuous RV X .

There are two directions to calculate definite integrals.

- **Numerical integration** by approximating the area by shapes of which we can calculate the area easier.
- **Analytical integration** by conducting integration as the inverse operation of differentiation.

In general, numerical integration methods⁸ can apply to a variety of cases but cause an approximation error. The analytical integration methods can give us the exact value but have limited applications. In practice, we combine them depending on the situation. In this lecture, **we focus on analytical methods**. It also helps us learn numerical integration.

⁸e.g., the trapezoidal rule, the Gauss-Legendre quadrature rule, the double exponential formula

Basic idea of calculating an integral

When we constructed the probability density function (PDF) p_X , we differentiated the cumulative distribution function (CDF) F_X . Specifically, $p_X(x) = \frac{d}{dx}F_X(x)$. Conversely, we observed that the area $\int_a^b p_X(x) dx$ under the graph of the PDF corresponds to the difference $F_X(b) - F_X(a)$.

Basic idea of calculating an integral

When we constructed the probability density function (PDF) p_X , we differentiated the cumulative distribution function (CDF) F_X . Specifically, $p_X(x) = \frac{d}{dx}F_X(x)$. Conversely,

we observed that the area $\int_a^b p_X(x) dx$ under the graph of the PDF corresponds to the difference $F_X(b) - F_X(a)$.

To wrap up, to calculate the definite integral $\int_a^b p_X(x) dx$, we can use a function whose derivative is p_X .

Basic idea of calculating an integral

When we constructed the probability density function (PDF) p_X , we differentiated the cumulative distribution function (CDF) F_X . Specifically, $p_X(x) = \frac{d}{dx}F_X(x)$. Conversely,

we observed that the area $\int_a^b p_X(x) dx$ under the graph of the PDF corresponds to the difference $F_X(b) - F_X(a)$.

To wrap up, to calculate the definite integral $\int_a^b p_X(x) dx$, we can use a function whose derivative is p_X .

According to the ***fundamental theorem of calculus (FTC)***, this relation between the derivative and the definite integral applies to a general function. We can use this relation to calculate a definite integral.

Integral is the “inverse” of differentiation

Definition (Primitive function)

Let a and b be real numbers such that $a < b$ and $f : [a, b] \rightarrow \mathbb{R}$. If $F : [a, b] \rightarrow \mathbb{R}$ satisfies $F' = f$, i.e., $\frac{d}{dx}F(x) = f(x)$ for all $x \in [a, b]$, then F is called a **primitive function** or an **antiderivative function** of f .

Theorem (The fundamental theorem of calculus (FTC))

Let a and b be real numbers such that $a < b$ and $f : [a, b] \rightarrow \mathbb{R}$ be integrable. Suppose that there exists a primitive function $F : [a, b] \rightarrow \mathbb{R}$ of f , then we have that

$$\int_a^b f(t) dt = F(b) - F(a). \quad (41)$$

We often denote $F(b) - F(a)$ by $[F(x)]_a^b$.

According to the FTC, we can **calculate an integral using a primitive function!**

Calculating the definite integral

To calculate the definite integral

$$\int_a^b f(x) \mathrm{d}x, \quad (42)$$

the following steps suffice.

- **Step 1:** Find a primitive (antiderivative) function $F : [\alpha, b] \rightarrow \mathbb{R}$, which satisfies $F' = f$.
- **Step 2:** Evaluate the value of $[F(x)]_a^b := F(b) - F(a)$.

Examples of calculating a definite integral

Example

Let $f(x) = x$.

We can calculate the definite integral $\int_{-4}^5 f(x) \mathrm{d}x = \int_{-4}^5 x \mathrm{d}x$ as follows.

- **Step 1:**

- **Step 2:**

Examples of calculating a definite integral

Example

Let $f(x) = x$.

We can calculate the definite integral $\int_{-4}^5 f(x) dx = \int_{-4}^5 x dx$ as follows.

- **Step 1:** Find a primitive (antiderivative) function F , which satisfies $F' = f$. In this example case, we can use a function $F(x) = \frac{1}{2}x^2$ as a primitive function since $\frac{d}{dx} \frac{1}{2}x^2 = x$.
- **Step 2:**

Examples of calculating a definite integral

Example

Let $f(x) = x$.

We can calculate the definite integral $\int_{-4}^5 f(x) dx = \int_{-4}^5 x dx$ as follows.

- **Step 1:** Find a primitive (antiderivative) function F , which satisfies $F' = f$. In this example case, we can use a function $F(x) = \frac{1}{2}x^2$ as a primitive function since $\frac{d}{dx} \frac{1}{2}x^2 = x$.
- **Step 2:** Evaluate the value of $[F(x)]_{-4}^5 := F(5) - F(-4)$. In this example case, $F(5) - F(-4) = \frac{1}{2}(5)^2 - \frac{1}{2}(-4)^2 = \frac{25}{2} - 8 = \frac{9}{2}$.

Examples of calculating a definite integral

Example

Let $f(x) = x$.

We can calculate the definite integral $\int_{-4}^5 f(x) dx = \int_{-4}^5 x dx$ as follows.

- **Step 1:** Find a primitive (antiderivative) function F , which satisfies $F' = f$. In this example case, we can use a function $F(x) = \frac{1}{2}x^2$ as a primitive function since $\frac{d}{dx} \frac{1}{2}x^2 = x$.
- **Step 2:** Evaluate the value of $[F(x)]_{-4}^5 := F(5) - F(-4)$. In this example case, $F(5) - F(-4) = \frac{1}{2}(5)^2 - \frac{1}{2}(-4)^2 = \frac{25}{2} - 8 = \frac{9}{2}$.

Hence, we have that

$$\int_{-4}^5 f(x) dx = \frac{9}{2}. \quad (43)$$

A primitive function is not unique.

As we have seen, finding a primitive function is essential to calculate the definite integral. Here, we must note that a primitive function is not unique.

If a function $F_1 : [a, b] \rightarrow \mathbb{R}$ is a primitive function of $f : [a, b] \rightarrow \mathbb{R}$, then $F_2 : [a, b] \rightarrow \mathbb{R}$ defined by $F_2(x) = F_1(x) + C$ is also a primitive function, where $C \in \mathbb{R}$ is a constant.

Example

Both $F_1(x) = \frac{1}{2}x^2$ and $F_2(x) = \frac{1}{2}x^2 + 5$ are primitive functions of $f(x) = x$.

A primitive function is not unique but unique up to an additive constant.

A primitive function is not unique. **However**, it is unique **up to an additive constant** in the following sense.

A primitive function is not unique but unique up to an additive constant.

A primitive function is not unique. **However**, it is unique **up to an additive constant** in the following sense.

Theorem (The primitive function is unique up to an additive constant)

Let α and b be real values such that $\alpha < b$. If both $F_1 : [\alpha, b] \rightarrow \mathbb{R}$ and $F_2 : [\alpha, b] \rightarrow \mathbb{R}$ are primitive functions of f , the difference between F_1 and F_2 is a constant function. In other words, there exists a constant $C \in \mathbb{R}$ such that $F_2(x) - F_1(x) = C$.

A primitive function is not unique but unique up to an additive constant.

A primitive function is not unique. **However**, it is unique **up to an additive constant** in the following sense.

Theorem (The primitive function is unique up to an additive constant)

Let a and b be real values such that $a < b$. If both $F_1 : [a, b] \rightarrow \mathbb{R}$ and $F_2 : [a, b] \rightarrow \mathbb{R}$ are primitive functions of f , the difference between F_1 and F_2 is a constant function. In other words, there exists a constant $C \in \mathbb{R}$ such that $F_2(x) - F_1(x) = C$.

To wrap up, if F is a primitive function of f , then, for any constant C , the function given by $F(x) + C$ is also a primitive function of f , and conversely, all the primitive functions are written in this form. We write this fact as follows.

$$\int f(x) dx = F(x) + C, \tag{44}$$

Here, the symbol $\int f(x) dx$ in the LHS denotes all the primitive functions of f . Here, the constant C in the RHS is called the **constant of integration**.

Examples of primitive functions

Example

The function $F(x) = \frac{1}{2}x^2$ is a primitive function of $f(x) = x$ since $F'(x) = f(x)$. Hence,
$$\int f(x) dx = \frac{1}{2}x^2 + C.$$
 Here, C is the constant of integration.

Note about the proof

We can prove the uniqueness of the primitive function up to an additive constant by the *mean value theorem*.

Indefinite integral

Let f be a function and a be a real value. The function defined by the following form is called an ***indefinite integral*** of f .

$$\int_a^x f(t) dt. \tag{45}$$

It is known that if f be continuous, then an indefinite integral is a primitive function of f . Note that some literature use the term “indefinite integral” to refer to a primitive function for this reason, while not all primitive functions are written in the above form.

Linearity of the antidifferentiation and integral

Since the derivation is a linear operator, the antidifferentiation, the operation to find a primitive function, is linear as well in the following sense.

Theorem (Linearity of antidifferentiation)

Let $f, g : \mathbb{R} \rightarrow \mathbb{R}$ be functions and F and G be the primitive functions of f and g , respectively. Also, let α and β be real values.

Then, $\alpha F + \beta G$ is a primitive function of $\alpha f + \beta g$.

In other words,

$$\int (\alpha f(x) + \beta g(x)) \, dx = \alpha \int f(x) \, dx + \beta \int g(x) \, dx. \quad (46)$$

We can easily prove the above by taking the derivatives of both sides⁹.

⁹Strictly speaking, we should consider the uniqueness of the primitive function up to an additive constant

Linearity of the definite integral

By combining the linearity of the antidifferentiation and the FTC, we can immediately get the linearity of the definite integral, which is a useful formula.

Corollary (Linearity of the definite integral)

Let $f, g : \mathbb{R} \rightarrow \mathbb{R}$ be functions and a, b, α and β be real values. Then,

$$\int_a^b (\alpha f(x) + \beta g(x)) \, dx = \alpha \int_a^b f(x) \, dx + \beta \int_a^b g(x) \, dx. \quad (47)$$

Example of the linearity of antidifferentiation

Example

In the following, C is the constant of integration.

- $\int (\cos x + x^2) dx = \int \cos x dx + \int x^2 dx = \sin x + \frac{1}{3}x^3 + C$. Hence,
$$\int_0^\pi (\cos x + x^2) dx = \left(\sin \pi + \frac{1}{3}\pi^3 \right) - \left(\sin 0 + \frac{1}{3} \cdot 0^3 \right) = \frac{1}{3}\pi^3.$$
- $\int 5 \exp(x) dx = 5 \int \exp(x) dx = 5 \exp(x) + C$. Hence,
$$\int_1^3 5 \exp(x) dx = (5 \exp(3)) - (5 \exp(1)) = 5e(e^2 - 1).$$

Finding the primitive function is not always easy.

To calculate the derivative, we had many useful formulae. Let $f, g : \mathbb{R} \rightarrow \mathbb{R}$ be differentiable functions, then, e.g.,

- $(fg)' = f'g + fg'$ for the product,
- $(g \circ f)' = (g' \circ f)f'$ for the composition.

Recall that the composition $g \circ f$ is defined by $(g \circ f)(x) = g(f(x))$.

However, generally speaking, antidifferentiation is more difficult than differentiation. Specifically, we have no formulae to find a primitive function of a general product or composition like in differentiation. Nevertheless, we have some techniques to make such calculation more feasible for some cases, called ***integration by parts*** and ***integration by substitution***.

Integration by parts

Let f and g be real functions and F and G be those primitive functions. While we cannot generally write the primitive function of the product fg only by F and G , the technique, called **integration by parts**, based on the following equation might help.

$$\int f(x)g(x)dx = f(x)G(x) - \int f'(x)G(x)dx. \quad (48)$$

Note that we assume that f is differentiable in the above.

By the above equation, we can find the primitive function of fg as long as we know that of $f'G$.

The proof of the above equation is easy if we differentiate the RHS.

Integration by parts

Let f and g be real functions and F and G be those primitive functions. While we cannot generally write the primitive function of the product $f g$ only by F and G , the technique, called **integration by parts**, based on the following equation might help.

$$\int f(x)g(x) \mathrm{d}x = f(x)G(x) - \int f'(x)G(x) \mathrm{d}x. \quad (48)$$

Example

$$\begin{aligned} \int x \cos(x) \mathrm{d}x &= x \sin(x) - \int (x)' \sin(x) \mathrm{d}x \\ &= x \sin(x) - \int 1 \cdot \sin(x) \mathrm{d}x \\ &= x \sin(x) - (-\cos x) + C. \end{aligned} \quad (49)$$

Integration by parts

Let f and g be real functions and F and G be those primitive functions. While we cannot generally write the primitive function of the product $f g$ only by F and G , the technique, called **integration by parts**, based on the following equation might help.

$$\int f(x)g(x) \mathrm{d}x = f(x)G(x) - \int f'(x)G(x) \mathrm{d}x. \quad (48)$$

Example

$$\begin{aligned} \int \log(x) \mathrm{d}x &= \int \log(x) \cdot 1 \mathrm{d}x = \log(x) \cdot x - \int (\log(x))' \cdot x \mathrm{d}x \\ &= \log(x) \cdot x - \int \frac{1}{x} \cdot x \mathrm{d}x \\ &= x \log(x) - x + C. \end{aligned} \quad (49)$$

Integration by substitution

Let f and g be real functions and assume f be differentiable. If the integrand includes the composition $g \circ f$, we cannot generally write the primitive function only by the primitive functions of f and g . However, we may find it by the following technique, called ***integration by substitution***.

Theorem (Integration by substitution for indefinite integral)

$$\int g(f(t))f'(t)dt = \int g(x)dx \Big|_{x=f(t)}, \quad (50)$$

where the RHS means the function we obtain by substituting $x = f(t)$ to a primitive function of g .

Both directions of the above equation are useful.

Integration by substitution

Let f and g be real functions and assume f be differentiable. If the integrand includes the composition $g \circ f$, we cannot generally write the primitive function only by the primitive functions of f and g . However, we may find it by the following technique, called ***integration by substitution***.

Theorem (Integration by substitution for definite integral)

$$\int_a^b g(f(t))f'(t)dt = \int_{f(a)}^{f(b)} g(x)dx, \quad (50)$$

where the RHS means the function we obtain by substituting $x = f(t)$ for a primitive function of g .

Both directions of the above equation are useful.

Why do we call it integration by substitution?

The previous page's formula is called integration by substitution because the formula is informally given by substituting $x = f(t)$ as follows.

$$\begin{aligned}\int_a^b g(f(t))f'(t) dt &= \int_{t=a}^{t=b} g(f(t)) \frac{df(t)}{dt} dt \\ &= \int_{t=a}^{t=b} g(x) \frac{dx}{dt} dt \\ &= \int_{t=a}^{t=b} g(x) dx \\ &= \int_{x=f(a)}^{x=f(b)} g(x) dx.\end{aligned}\tag{51}$$

Note that the above discussion is mathematically inaccurate (especially where we used $\frac{dx}{dt} dt = dx$). If we want to formally prove the formula, we should simply differentiate both sides of the formula for indefinite integral.

Examples of integration by substitution.

Recall the formula.

$$\int_a^b g(f(t))f'(t)dt = \int_{f(a)}^{f(b)} g(x)dx, \quad (52)$$

Example (integration by substitution: from left to right)

$$\begin{aligned} \int_0^{+2} t \exp(-t^2) dt &= -\frac{1}{2} \int_0^{+2} \exp(-t^2) \cdot (-2t) dt \\ &= -\frac{1}{2} \int_0^{+2} \exp(-t^2) \cdot (-t^2)' dt \\ &= -\frac{1}{2} \int_{-0^2}^{-2^2} \exp(x) dx \\ &= -\frac{1}{2} [\exp(x)]_{-0^2}^{-2^2} = -\frac{1}{2} [\exp(-4) - \exp(0)] = \frac{1}{2} [1 - \exp(-4)]. \end{aligned} \quad (53)$$

Examples of integration by substitution.

Recall the formula.

$$\int_a^b g(f(t))f'(t)dt = \int_{f(a)}^{f(b)} g(x)dx, \quad (52)$$

Example (integration by substitution: from right to left)

$$\begin{aligned} \int_0^1 \sqrt{1-x^2} dx &= \int_{\pi}^0 \sqrt{1-\cos^2(t)}(\cos(t))' dt \quad \text{since } \cos(\pi) = 0, \cos(0) = 1, \\ &= \int_{\pi}^0 \sqrt{1-\cos^2(t)}(-\sin(t)) dt \\ &= \int_0^{\pi} \sin^2(t) dt = \int_0^{\pi} \frac{1-\cos(2t)}{2} dt = \left[\frac{1}{2}t - \frac{1}{4}\sin(2t) \right]_0^{\pi} = \frac{1}{2}\pi. \end{aligned} \quad (53)$$

Example of definite integral calculation in probability theory

Example (Exponential distribution)

The distribution of a RV X is called the **exponential distribution** with mean μ if it has a PDF p_X given by

$$p_X(x) := \begin{cases} 0 & \text{if } x < 0, \\ \frac{1}{\mu} \exp\left(-\frac{x}{\mu}\right) & \text{if } x \geq 0. \end{cases} \quad (54)$$

For nonnegative numbers a and b , the probability $\Pr(a < X \leq b)$ is given by

$$\begin{aligned} \Pr(a < X \leq b) &= \int_a^b p_X(x) dx = \int_a^b \frac{1}{\mu} \exp\left(-\frac{x}{\mu}\right) dx \\ &= \left[-\exp\left(-\frac{x}{\mu}\right) \right]_a^b = \exp\left(-\frac{a}{\mu}\right) - \exp\left(-\frac{b}{\mu}\right). \end{aligned} \quad (55)$$

Finding a primitive function of the product and composition is not easy.

We know that the primitive functions of $\frac{1}{x}$ and \sin , or \exp and $-x^2$. Indeed,

$$\int \frac{1}{x} dx = \log|x| + C, \int \sin x dx = -\cos x + C, \int (-x^2) dx = -\frac{1}{3}x^3 + C, \int \exp(x) dx = \exp(x) + C \quad (56)$$

However, it is known that the primitive functions of $\frac{1}{x} \sin x$ and $\exp(-x^2)$ are not **elementary**, although $\frac{1}{x} \sin x$ and $\exp(-x^2)$ themselves are elementary.

Here, we call a function **elementary** if we can write the function as a composition of finitely many

- algebraic functions, functions represented as a root of polynomial-function-coefficient polynomial equations, including polynomial, rational functions and fractional powers, e.g., $5x^2 + x - 3$, $\sqrt{3}x + 5$, $\frac{3x+1}{-2x^2+x+5}$, etc.
- trigonometric functions, e.g., $\sin x$, $\cos x$ etc.,
- exponential function $\exp x$,
- logarithmic function $\log x$.

Finding a primitive function of the product and composition is not easy.

We know that the primitive functions of $\frac{1}{x}$ and \sin , or \exp and $-x^2$. Indeed,

$$\int \frac{1}{x} dx = \log|x| + C, \int \sin x dx = -\cos x + C, \int (-x^2) dx = -\frac{1}{3}x^3 + C, \int \exp(x) dx = \exp(x) + C \quad (56)$$

However, it is known that the primitive functions of $\frac{1}{x} \sin x$ and $\exp(-x^2)$ are not **elementary**, although $\frac{1}{x} \sin x$ and $\exp(-x^2)$ themselves are elementary.

Roughly speaking, most functions we can imagine without the inverse function and the primitive function are elementary.

The fact that the primitive functions of $\frac{1}{x} \sin x$ and $\exp(-x^2)$ are not elementary means we have no way to write those primitive functions.

From the computer science viewpoint, the above fact means that we cannot easily find the exact value of the integrals of those functions. Some non-elementary primitive functions might be implemented by some libraries if they are famous. If they are not

Finding a primitive function of the product and composition is not easy.

We know that the primitive functions of $\frac{1}{x}$ and \sin , or \exp and $-x^2$. Indeed,

$$\int \frac{1}{x} dx = \log|x| + C, \int \sin x dx = -\cos x + C, \int (-x^2) dx = -\frac{1}{3}x^3 + C, \int \exp(x) dx = \exp(x) + C \quad (56)$$

However, it is known that the primitive functions of $\frac{1}{x} \sin x$ and $\exp(-x^2)$ are not **elementary**, although $\frac{1}{x} \sin x$ and $\exp(-x^2)$ themselves are elementary.

In fact, these functions are important in many areas.

- The PDF of the normal distribution is proportional to $\exp(-x^2)$. The normal distribution is the most important distribution in probability theory, owing to the central limit theorem.
- The sine cardinal function $\frac{\sin x}{x}$ appears in many application areas, including physics, probability theory, signal processing, optics, etc., because it is the Fourier transform of the rectangle function.

Outline

3. Continuous Random Variables

3.1 Introduction: why are continuous random variables less trivial?

3.2 Probability density function

3.3 Area, integration, and properties of PDF.

3.4 Calculating integral

3.5 Summary statistics of continuous RV and integral

3.6 Multivariate random variables and multiple integral

3.7 Relation among RVs in a continuous multivariate RV

3.8 Exercises

Expectation (mean) of a continuous random variable

The expectation of a continuous RV is defined similarly to that of a discrete RV. Specifically, we get the definition for a continuous RV by replacing the PMF and the sum with the PDF and the integration in the definition for a discrete RV.

Expectation (mean) of a continuous random variable

The expectation of a continuous RV is defined similarly to that of a discrete RV. Specifically, we get the definition for a continuous RV by replacing the PMF and the sum with the PDF and the integration in the definition for a discrete RV.

Definition (Expectation of a continuous RV)

Let X be a continuous RV and p_X be its probability density function (PDF). Then, the expectation $\mathbb{E}X$ of X is defined by

$$\mathbb{E}X := \int_{-\infty}^{+\infty} xp(x)dx. \quad (57)$$

Cf.) The expectation of a discrete RV X is given by $\sum_{x \in \mathcal{X}} xP_X(x)$, where P_X is the probability mass function.

Example: expectation of exponential distribution

Example (Expectation of the exponential distribution)

The PDF p_X of a RV X following the exponential distribution with mean μ is given by

$$p_X(x) := \begin{cases} 0 & \text{if } x < 0, \\ \frac{1}{\mu} \exp\left(-\frac{x}{\mu}\right) & \text{if } x \geq 0. \end{cases} \quad (58)$$

Noting that the density is zero for the negative domain, we can calculate the expectation $\mathbb{E}X$ using integration by parts as follows.

$$\begin{aligned} \mathbb{E}X &= \int_{-\infty}^{+\infty} x p_X(x) dx = \int_0^{+\infty} x \cdot \frac{1}{\mu} \exp\left(-\frac{x}{\mu}\right) dx = \int_0^{+\infty} x \left(-\exp\left(-\frac{x}{\mu}\right)\right)' dx \\ &= \left[x \cdot \left(-\exp\left(-\frac{x}{\mu}\right)\right) \right]_0^{+\infty} - \int_0^{+\infty} (x)' \cdot \left(-\exp\left(-\frac{x}{\mu}\right)\right) dx = - \int_0^{+\infty} \left(-\exp\left(-\frac{x}{\mu}\right)\right) dx \\ &= - \left[\mu \exp\left(-\frac{x}{\mu}\right) \right]_0^{+\infty} = \mu. \end{aligned}$$

The expectation of a function of a continuous RV

A function of a discrete RV is always a discrete RV. However, a function of a continuous RV is not always a continuous RV.

The expectation of a function of a continuous RV

A function of a discrete RV is always a discrete RV. However, a function of a continuous RV is not always a continuous RV.

For example, if f is the sign function defined by

$$f(x) := \begin{cases} -1 & \text{if } x < 0, \\ 0 & \text{if } x = 0, \\ +1 & \text{if } x > 0, \end{cases} \quad (60)$$

and X is a continuous RV whose PDF p_X is given by

$$p_X(x) = \begin{cases} +1 & \text{if } -\frac{1}{2} \leq x \leq +\frac{1}{2}, \\ 0 & \text{otherwise.} \end{cases} \quad (61)$$

Then, the RV $f(X)$ takes values -1 and $+1$ with equal probability. In particular, it is a discrete RV, whose support is $\{-1, +1\}$.

The expectation of a function of a continuous RV

A function of a discrete RV is always a discrete RV. However, a function of a continuous RV is not always a continuous RV.

Even though a function of a continuous RV may not be a continuous RV, its expectation can always be calculated by the following formula, which is similar to the formula for a discrete RV.

Theorem

Let X be a continuous RV and its PDF be p_X . Also, let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a real-valued function taking a real value as an input. The expectation $\mathbb{E} f(X)$ of the random variable X is given as follows.

$$\mathbb{E} f(X) = \int_{-\infty}^{+\infty} f(x)p_X(x)dx. \quad (60)$$

Cf.) For a discrete RV whose support and PMF are \mathcal{X} and P_X , respectively, we have that $\mathbb{E} f(X) = \sum_{x \in \mathcal{X}} f(x)P_X(x)$.

The linearity of the expectation operation also applies to a continuous RV

The following theorem, which holds for a discrete RV, also holds for a continuous RV.

Theorem (The linearity of the expectation)

Let X be a random variable, $\alpha, b \in \mathbb{R}$ be real numbers, and $f, g : \mathbb{R} \rightarrow \mathbb{R}$ be real-valued functions taking a real variable. Then, we have that

$$\mathbb{E}[\alpha f(X) + b g(X)] = \alpha \mathbb{E} f(X) + b \mathbb{E} g(X). \quad (61)$$

Variance and standard deviation of a continuous random variable

The definitions of the variance and standard deviation are the same for a continuous RV. Specifically, for a continuous RV X , whose expectation is μ_X , its variance $\mathbb{V}(X)$ is defined by $\mathbb{V}(X) := \mathbb{E}(X - \mu_X)^2$. The standard deviation is defined by $\sigma_X := \sqrt{\mathbb{V}(X)}$.

When we know the explicit form of the PDF, we can use the following formulae.

Theorem

Let X be a continuous RV and its PDF be p_X . Suppose that the expectation $\mathbb{E}X = \int_{-\infty}^{+\infty} x p_X(x) dx$ exists and denote it by μ_X . The variance $\mathbb{V}(X)$ is given by the following formula.

$$\mathbb{V}(X) = \int_{-\infty}^{+\infty} (x - \mu_X)^2 p_X(x) dx = \int_{-\infty}^{+\infty} x^2 p_X(x) dx - (\mu_X)^2. \quad (62)$$

Outline

3. Continuous Random Variables

3.1 Introduction: why are continuous random variables less trivial?

3.2 Probability density function

3.3 Area, integration, and properties of PDF.

3.4 Calculating integral

3.5 Summary statistics of continuous RV and integral

3.6 Multivariate random variables and multiple integral

3.7 Relation among RVs in a continuous multivariate RV

3.8 Exercises

Handling a continuous multivariate random variable

Similar to the univariate random variable case, we can define the cumulative distribution function (CDF) of the distribution.

Definition (The CDF of bivariate RV)

Let X and Y be random variables. The **cumulative distribution function (CDF)** $F_{X,Y} : \mathbb{R}^2 \rightarrow [0, 1]$ of X and Y is defined by

$$F_{X,Y}(x, y) := \Pr(X \leq x \wedge Y \leq y), \quad (63)$$

where \wedge indicates the logical “and” statement.

Using the CDF, we can calculate the probability $\Pr(a_1 < X \leq b_1 \wedge a_2 < Y \leq b_2)$ by

$$\Pr(a_1 < X \leq b_1 \wedge a_2 < Y \leq b_2) = F_{X,Y}(b_1, b_2) - F_{X,Y}(a_1, b_2) - F_{X,Y}(b_1, a_2) + F_{X,Y}(a_1, a_2) \quad (64)$$

To define the probability density function for a multivariate random variable

Let (X_1, X_2, \dots, X_m) be a multivariate random variable

As in the univariate random variable case, the CDF may not be easy to interpret or not be elementary even in practical cases. Hence, we want to define the probability density function (PDF) for a multivariate random variable.

Similar to multivariate discrete RV cases, it is not sufficient to evaluate the PDF for each RV to see the behavior of a multivariate discrete RV. In discrete RV cases, we evaluated the Joint PMF, which returns the probability mass of the event $(X_1, X_2, \dots, X_m) = (x_1, x_2, \dots, x_m)$. Similarly, we want to define the function that returns the probability density at $(X_1, X_2, \dots, X_m) = (x_1, x_2, \dots, x_m)$. Since the PDF for a univariate continuous RV was defined using the area under the graph, let us define the graph of a multivariate function and the high-dimensional area (volume) in the following.

The graph of a multivariate function and multiple integral

Let $D = [a_1, b_1] \times [a_2, b_2] \times \cdots \times [a_m, b_m]$ be a m -dimensional hyper-rectangle. Let $f : D \rightarrow \mathbb{R}$ be a function of a m -dimensional variable. Similar to one-dimensional function cases, we call the set of points

$$\{(x_1, x_2, \dots, x_m, f(x_1, x_2, \dots, x_m)) | (x_1, x_2, \dots, x_m) \in D\} \quad (65)$$

the **graph** of a function f . The (signed) volume in the domain D bounded by the graph of $y = f(\mathbf{x})$ and $y = 0$ is called the **multiple integral** of f on D , denoted by $\int_D f(\mathbf{x}) \, d\mathbf{x}$.

Joint PDF

Based on the definition of multiple integration, we can define the joint probability density function (joint PDF) of a multivariate random variable.

Definition

Let (X_1, X_2, \dots, X_m) be a multivariate random variable. If $p_{X_1, X_2, \dots, X_m} : \mathbb{R}^m \rightarrow \mathbb{R}_{\geq 0}$ satisfies

$$\Pr((X_1, X_2, \dots, X_m) \in D) = \int_D p_{X_1, X_2, \dots, X_m}(\mathbf{x}) d\mathbf{x} \quad (66)$$

for any m -dimensional hyper-rectangle D , then the function p_{X_1, X_2, \dots, X_m} is called the **joint probability density function (joint PDF)** of X_1, X_2, \dots, X_m .

If (X_1, X_2, \dots, X_m) have a joint PDF, we call it a **continuous multivariate random variable**.

Multiple continuous random variables are not always a continuous multivariate random variable

Let X_1, X_2, \dots, X_m be continuous random variables. In other words, suppose that there exist PDFs $p_{X_1}, p_{X_2}, \dots, p_{X_m}$ for X_1, X_2, \dots, X_m , respectively.

Even under this assumption, **it is possible that (X_1, X_2, \dots, X_m) has no joint PDF.**

Multiple continuous random variables are not always a continuous multivariate random variable

Let X_1, X_2, \dots, X_m be continuous random variables. In other words, suppose that there exist PDFs $p_{X_1}, p_{X_2}, \dots, p_{X_m}$ for X_1, X_2, \dots, X_m , respectively.

Even under this assumption, **it is possible that (X_1, X_2, \dots, X_m) has no joint PDF.**

For example, if X is a continuous RV and $Y = X$, then the probability concentrate on the line $y = x$, so (X, Y) .

In this case, (X, Y) is not a continuous multivariate random variable.

Multiple integral on a complicated shape

In high-dimensional space, we might want to consider the volume bounded by a function in a complicated shape, say A , that cannot be represented as a union of hyper-rectangles.

Multiple integral on a complicated shape

In high-dimensional space, we might want to consider the volume bounded by a function in a complicated shape, say A , that cannot be represented as a union of hyper-rectangles.

For example, we might want to consider the probability $\Pr((X, Y) \in A)$, where $A := \{(x, y) \in \mathbb{R}^2 \mid x^2 + y^2 \leq 1\}$ is defined as the unit disk centered at the origin. We cannot decompose the disk into rectangles, so we cannot evaluate the probability by the sum rule if we can only define the probability of the multivariate RV being in a rectangle.

Multiple integral on a complicated shape

In high-dimensional space, we might want to consider the volume bounded by a function in a complicated shape, say A , that cannot be represented as a union of hyper-rectangles.

Hence, we want to define the volume bounded by a function on a general set A . We can do it by multiplying the value of the function by zero everywhere outside of A as follows.

Definition

Let $D = [a_1, b_1] \times [a_2, b_2] \times \cdots \times [a_m, b_m]$ be a m -dimensional hyper-rectangle.

For a general subset $A \subset D$, we define the multiple integral of f on A by

$$\int_A f(\mathbf{x}) d\mathbf{x} := \int_D 1_A(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}, \quad (67)$$

where the indicator function 1_A is defined by $1_A(\mathbf{x}) := \begin{cases} 1 & \text{if } \mathbf{x} \in A, \\ 0 & \text{if } \mathbf{x} \notin A. \end{cases}$

Probability on a complicated shape

The probability density function can be applied to a complicated shape.

Theorem

Let (X_1, X_2, \dots, X_m) be a multivariate continuous random variable, and let p_{X_1, X_2, \dots, X_m} be its joint PDF. Then, for $A \in \mathbb{R}^m$, we have that

$$\Pr((X_1, X_2, \dots, X_m) \in A) = \int_A p_{X_1, X_2, \dots, X_m}(\mathbf{x}) d\mathbf{x} \quad (68)$$

Calculating multi integral by iterated integral

We can calculate a multi-integral by an *iterated integral*.

Theorem

Under some loose conditions¹⁰, we have that

$$\begin{aligned} & \iint_A p(x, y) \, dx \, dy \\ &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} 1_A(x, y) p(x, y) \, dx \right] dy \\ &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} 1_A(x, y) p(x, y) \, dy \right] dx. \end{aligned} \tag{69}$$

¹⁰We refer the readers wanting to know the exact conditions to the Fubini-Tonelli theorem.

Joint PDF example 1: Uniform distribution

Example (Uniform distribution)

Let X and Y be RVs following the bivariate uniform distribution with the support $[0, 3] \times [-1, +1]$. The RVs X and Y have has the joint PDF

$$p_{X,Y}(x,y) = \begin{cases} \frac{1}{6} & \text{if } (x,y) \in [0,3] \times [-1,+1], \\ 0 & \text{if } (x,y) \notin [0,3] \times [-1,+1]. \end{cases} \quad (70)$$

For example, the probability $\Pr((X,Y) \in [0, \frac{1}{2}] \times [0, \frac{1}{4}])$ is given by

$$\int_0^{\frac{1}{4}} \int_0^{\frac{1}{2}} \frac{1}{6} dx dy = \int_0^{\frac{1}{4}} \left[\frac{1}{6}x \right]_0^{\frac{1}{2}} dy = \int_0^{\frac{1}{4}} \frac{1}{12} dy = \left[\frac{1}{12}y \right]_0^{\frac{1}{4}} = \frac{1}{48} \quad (71)$$

Joint PDF example 2: Multivariable normal distribution

Example (Multivariable normal distribution)

Let $\boldsymbol{\mu}$ be a real m -dimensional vector and $\boldsymbol{\Sigma}$ be a real $m \times m$ positive definite matrix, i.e., a $m \times m$ matrix such that $\mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x} > 0$ for any non-zero m -dimensional vector \mathbf{x} . We call the distribution of a m -tuple (X_1, X_2, \dots, X_m) of RVs a ***multivariable normal distribution*** if it has the following joint PDF $p_{X,Y}$.

$$p_{X_1, X_2, \dots, X_m}(\mathbf{x}) = \frac{1}{\sqrt{2^m \pi \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (72)$$

Joint PDF example 2: Multivariable normal distribution

Example (Multivariable normal distribution)

Recall that the joint PDF of a multivariable normal distribution is given as follows.

$$p_{X_1, X_2, \dots, X_m}(\mathbf{x}) = \frac{1}{\sqrt{2^m \pi \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (72)$$

For a bivariable case $m = 2$, the joint PDF is given by

$$p_{X_1, X_2}(x_1, x_2) = \frac{1}{\sqrt{2^2 \pi \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}([x_1 \ x_2] - [\mu_1 \ \mu_2]) \begin{bmatrix} s_{11} & s_{12} \\ s_{21} & s_{22} \end{bmatrix}^{-1} \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \right)\right) \quad (73)$$

where $\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ is a 2-dimensional vector and $\boldsymbol{\Sigma} = \begin{bmatrix} s_{11} & s_{12} \\ s_{21} & s_{22} \end{bmatrix}$ is a real 2×2 positive definite matrix.

Joint PDF example 2: Multivariable normal distribution

Example (Multivariable normal distribution)

Recall that the joint PDF of a multivariable normal distribution is given as follows.

$$p_{X_1, X_2, \dots, X_m}(\mathbf{x}) = \frac{1}{\sqrt{2^m \pi \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (72)$$

We can see that $p_{X,Y}(x, y)$ takes its maximum if $s = \boldsymbol{\mu}$ since $(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})$ is zero if $s = \boldsymbol{\mu}$ and positive otherwise, according to the positive definite assumption on $\boldsymbol{\Sigma}$.

Joint PDF example 2: Multivariable normal distribution

Example (Multivariable normal distribution)

Recall that the joint PDF of a multivariable normal distribution is given as follows.

$$p_{X_1, X_2, \dots, X_m}(\mathbf{x}) = \frac{1}{\sqrt{2^m \pi \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (72)$$

Unfortunately, we cannot calculate the probability $\Pr((X, Y) \in A)$ analytically for general A .

Marginal PDF (bivariable cases)

Suppose that the joint PDF $p_{X,Y}$ is given. The **marginal probability density functions (marginal PDFs)** p_X and p_Y are given by

$$\begin{aligned} p_X(x) &= \int_{-\infty}^{+\infty} p_{X,Y}(x,y)dy, \\ p_Y(y) &= \int_{-\infty}^{+\infty} p_{X,Y}(x,y)dx. \end{aligned} \tag{73}$$

Outline

3. Continuous Random Variables

3.1 Introduction: why are continuous random variables less trivial?

3.2 Probability density function

3.3 Area, integration, and properties of PDF.

3.4 Calculating integral

3.5 Summary statistics of continuous RV and integral

3.6 Multivariate random variables and multiple integral

3.7 Relation among RVs in a continuous multivariate RV

3.8 Exercises

Conditional PDF (bivariate cases)

The ***conditional probability distribution function (conditional PDF)*** is defined by

$$p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)}, \quad p_{Y|X}(y|x) = \frac{p_{X,Y}(x,y)}{p_X(x)}. \quad (74)$$

Independency

Let X and Y be RVs and assume that they have a joint PDF $p_{X,Y}$ and let their marginal PDFs be p_X and p_Y . Also, denote the conditional PDF of X given Y and that of Y given X by $p_{X|Y}$ and $p_{Y|X}$, respectively.

We say that the RVs X and Y are (mutually) independent if one of the following equivalent conditions holds

- $p_{X,Y}(x,y) = p_X(x)p_Y(y)$ for all (x,y) .
- $p_{X|Y}(x|y) = p_X(x)$ for all (x,y) such that $p_Y(y) \neq 0$.
- $p_{Y|X}(y|x) = p_Y(y)$ for all (x,y) such that $p_X(x) \neq 0$.

Calculating the expectation of a function from joint PDF

Let X_1, X_2, \dots, X_m be random variables and p_{X_1, X_2, \dots, X_m} be the joint PDF. Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$ be a function. The expectation of the random variable $f(X_1, X_2, \dots, X_m)$ is given by

$$\int_{\mathbb{R}^m} f(\mathbf{x}) p_{X_1, X_2, \dots, X_m}(\mathbf{x}) d\mathbf{x}. \quad (75)$$

Covariance

Let X and Y are random variables and μ_X and μ_Y be the expectation of X and Y , respectively. Suppose that $p_{X,Y}$ is a joint PDF of X and Y . Then, the covariance $\text{Cov}(X, Y)$ is given by

$$\text{Cov}(X, Y) = \int_{\mathbb{R}^2} (x - \mu_X)(y - \mu_Y) p_{X,Y}(x, y) dx dy \quad (76)$$

Integration by substitution for a double integral

Theorem (Integration by substitution for a double integral)

$$\begin{aligned} & \int_U f(\varphi_1(u_1, u_2), \varphi_2(u_1, u_2)) \left| \det \left(\begin{bmatrix} \frac{\partial \varphi_1}{\partial u_1}(u_1, u_2) & \frac{\partial \varphi_1}{\partial u_2}(u_1, u_2) \\ \frac{\partial \varphi_2}{\partial u_1}(u_1, u_2) & \frac{\partial \varphi_2}{\partial u_2}(u_1, u_2) \end{bmatrix} \right) \right| du_1 du_2 \\ &= \int_{\varphi(U)} f(x_1, x_2) dx_1 dx_2 \end{aligned} \quad (77)$$

Here, recall that

$$\det \left(\begin{bmatrix} a & b \\ c & d \end{bmatrix} \right) = ad - bc. \quad (78)$$

An example of integration by substitution

Most practical substitutions are given by the polar coordinate: $x = r \cos \theta, y = r \sin \theta$.

By this substitution, we have that $\sqrt{x^2 + y^2} = r$.

Also, the determinant of the Jacobian of the coordinate transform is given by

$$\det \left(\begin{bmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{bmatrix} \right) = \det \left(\begin{bmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{bmatrix} \right) = r \cos^2 \theta - (-r \sin^2 \theta) = r. \quad (79)$$

Using the above results, we can calculate, for example,

$$\begin{aligned} \iint_{x^2+y^2 \leq 1} \left(1 - \sqrt{x^2 + y^2}\right) dx dy &= \int_0^{2\pi} \int_0^1 (1-r) \left| \det \left(\begin{bmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{bmatrix} \right) \right| dr d\theta \\ &= \int_0^{2\pi} \int_0^1 (1-r) |r| dr d\theta \\ &= \int_0^{2\pi} \left[\int_0^1 (r - r^2) dr \right] d\theta = \int_0^{2\pi} \frac{1}{6} d\theta = \frac{1}{3} \pi. \end{aligned} \quad (80)$$

Integration by substitution for a multiple integral

Theorem (Integration by substitution for a multiple integral)

$$\int_U f(\boldsymbol{\varphi}(\boldsymbol{u})) \left| \det \left(\frac{\partial \boldsymbol{\varphi}}{\partial \boldsymbol{u}}(\boldsymbol{u}) \right) \right| d\boldsymbol{u} = \int_{\boldsymbol{\varphi}(U)} f(\boldsymbol{x}) d\boldsymbol{x} \quad (81)$$

Outline

3. Continuous Random Variables

3.1 Introduction: why are continuous random variables less trivial?

3.2 Probability density function

3.3 Area, integration, and properties of PDF.

3.4 Calculating integral

3.5 Summary statistics of continuous RV and integral

3.6 Multivariate random variables and multiple integral

3.7 Relation among RVs in a continuous multivariate RV

3.8 Exercises

Exercise (CDF of a continuous RV)

Let X be a random variable whose CDF is given by

$$F_X(x) = \begin{cases} 0 & \text{if } x \leq 0, \\ x & \text{if } 0 \leq x \leq 1, \\ 1 & \text{if } x \geq 1. \end{cases} \quad (82)$$

Find the probability of $0.2 \leq X \leq 0.7$.

Exercise (CDF to PDF)

Let X be a random variable whose CDF is given by

$$F_X(x) = \begin{cases} 0 & \text{if } x \leq 0, \\ \frac{1}{4}x^2 & \text{if } 0 \leq x \leq 2, \\ 1 & \text{if } x \geq 2. \end{cases} \quad (83)$$

Find the PDF p_X .

Exercise (PDF to probability)

Let X be a random variable whose PDF is given by

$$p_X(x) = \begin{cases} 0 & \text{if } x \leq 0, \\ \exp(-x) & \text{if } x \geq 0. \end{cases} \quad (84)$$

Find the probability of the event $1 \leq X \leq 2$.

Exercise (Indefinite integral)

Find the following indefinite integrals

- $\int (\cos x + x^2) \, dx$
- $\int 5 \exp(x) \, dx$
- $\int x \cos(x) \, dx$
- $\int \log(x) \, dx$

Exercise (Definite integral)

Find the following definite integrals

- $\int_{\pi}^{2\pi} (\cos x + x^2) \, dx$
- $\int_0^1 5 \exp(x) \, dx$
- $\int_{\pi}^0 x \cos(x) \, dx$
- $\int_1^e \log(x) \, dx$
- $\int_0^{+2} t \exp(-t^2) \, dt$
- $\int_0^1 \sqrt{1-x^2} \, dx$

Exercise (Improper integral)

Find the following improper integrals

- $\int_0^{+\infty} \exp(-x) dx.$
- $\int_{-\infty}^{+\infty} \frac{1}{1+x^2} dx.$

You can use the fact that $\frac{d}{dx} \arctan(x) = \frac{1}{1+x^2}$, where \arctan is the inverse function of \tan with the domain limited to $(-\pi, +\pi)$.

Exercise (Double integral)

- Evaluate the definite integral $\iint_D 3x^2 y \, dx \, dy$ where D is defined by $0 \leq x \leq 1$ and $0 \leq y \leq 2$.
- Evaluate the definite integral $\iint_D (x^2 + y^2) \, dx \, dy$ where D is the region bounded by the curves $y = x$, $y = 2x$, $x = 1$, and $x = 2$.

Exercise (Integration by substitution for a double integral)

Evaluate the following values

- $\iint_{x^2+y^2 \leq 1} \left(1 - \sqrt{x^2 + y^2}\right) dx dy,$
- $\iint_{x^2+y^2 \leq 1} \exp\left(-\sqrt{x^2 + y^2}\right) dx dy,$
- $\int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \exp(-(x^2 + y^2)) dx dy.$

Outline

4. Sample Statistics

4.1 Introduction: why do we learn sample statistics?

4.2 Terminology

4.3 Sample mean, law of large numbers, and central limit theorem

4.4 Estimation of distribution and parametric model

4.5 Likelihood

4.6 Maximum likelihood estimator

4.7 Exercises

Outline

4. Sample Statistics

4.1 Introduction: why do we learn sample statistics?

4.2 Terminology

4.3 Sample mean, law of large numbers, and central limit theorem

4.4 Estimation of distribution and parametric model

4.5 Likelihood

4.6 Maximum likelihood estimator

4.7 Exercises

Sample and sample statistics

In real applications, we **rarely know the true distribution**, behind the data.

On the other hand, we often **have many data points** that we can assume follow the same distribution (often independently). Such a series of data points is called **sample** of the distribution.

Statistics, data science, machine learning, etc., aim to **extract information about the true distribution from available data points**. **Sample statistics are the basis of those pieces of technology**.

Learning outcomes

By the end of this section, you should be able to:

- Explain the difference between summary statistics and sample statistics,
- Estimate the true mean of an unknown distribution by finite size sample,
- Explain why many random variables in the real world follow a normal distribution, and
- Estimate an unknown distribution using a parametric model and maximum likelihood estimator.

Outline

4. Sample Statistics

4.1 Introduction: why do we learn sample statistics?

4.2 Terminology

4.3 Sample mean, law of large numbers, and central limit theorem

4.4 Estimation of distribution and parametric model

4.5 Likelihood

4.6 Maximum likelihood estimator

4.7 Exercises

Population and sample

In the context of statistics,

- The true distribution is often called the ***population***.
- A series of data points that we can assume follow the same distribution is called ***sample***. If it has many data points, we say that the sample is large, and if it has few data, we say that the sample is small.

Summary statistics and sample statistics

- **Summary statistics** aims to describe characteristics of a (known or true) distribution by a few values.
- **Sample statistics** aims to estimate some information about the true distribution from finite sample data.

We only have **finite** data points in real applications, so sample statistics are practically necessary to handle probability.

Outline

4. Sample Statistics

4.1 Introduction: why do we learn sample statistics?

4.2 Terminology

4.3 Sample mean, law of large numbers, and central limit theorem

4.4 Estimation of distribution and parametric model

4.5 Likelihood

4.6 Maximum likelihood estimator

4.7 Exercises

Sample mean

One principal summary statistic is the expectation.

For data points X_1, X_2, \dots, X_m , we can easily calculate the **sample mean**

$$\overline{X}_m = \frac{1}{m}(X_1 + X_2 + \dots + X_m), \quad (85)$$

the mean of the data points.

If we can assume that those data points are the values of random variables following the same distribution with a true mean μ , we expect \overline{X}_m to approximate the true mean μ , which is unknown.

Is it correct? The answer is YES, according to the **law of large numbers**.

Law of large numbers

Theorem ((Strong) law of large numbers)

Let X_1, X_2, \dots be an infinite sequence of independently and identically distributed (i.i.d.) random variables and assume that the mean of the distribution is $\mu \in \mathbb{R}$.

Let \overline{X}_m be the sample mean

$$\overline{X}_m := \frac{1}{m}(X_1 + X_2 + \dots + X_m). \quad (86)$$

Then \overline{X}_m converges to μ in probability 1.

Thus, the sample mean tells us some information about the unknown true distribution!

How the sample mean behaves?

The sample mean converges to the expectation. Now,

- How close to the expectation will the sample mean get as we increase the data points?
- What does the distribution of the sample mean look like?

The answer is

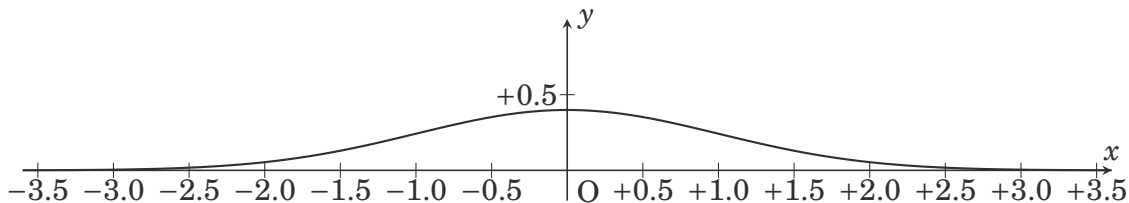
- The difference between the sample mean and the true expectation is proportional to the standard deviation σ of the true distribution and $\frac{1}{\sqrt{m}}$,
- With appropriate scaling, the distribution of the sample mean converges to a ***normal distribution (Gaussian distribution)***,

according to the ***central limit theorem***.

What is the standard normal distribution?

The ***standard normal distribution***, also known as the ***standard Gaussian distribution*** is the distribution with the following PDF:

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right). \quad (87)$$



The standard normal distribution's PDF.

Mean: 0, Variance: 1. The PDF is symmetric about $x = 0$ and it is dense around $x = 0$.

Central limit theorem (CLT)

Theorem (Central limit theorem (CLT))

Let X_1, X_2, \dots be an infinite sequence of independently and identically distributed (i.i.d.) random variables and assume that the mean and variance of the distribution are $\mu \in \mathbb{R}$ and $\sigma^2 \in \mathbb{R}_{\geq 0}$, respectively.

Let \bar{X}_m be the sample mean

$$\bar{X}_m := \frac{1}{m}(X_1 + X_2 + \dots + X_m). \quad (88)$$

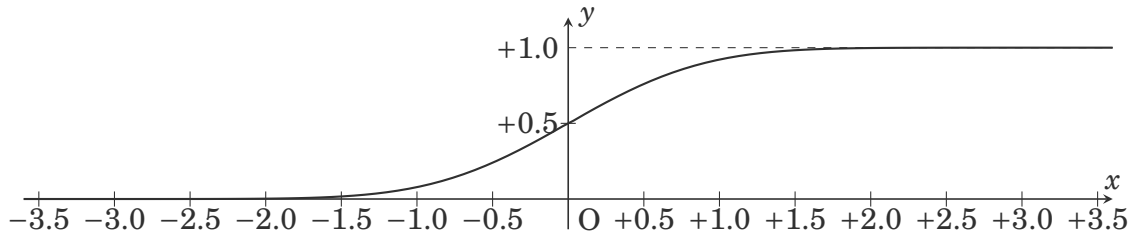
Then, the CDF of $\sqrt{m} \frac{\bar{X}_m - \mu}{\sigma}$ converges to the CDF of the standard normal distribution at any point in \mathbb{R} .

The standard normal distribution's CDF

By definition, the CDF $F : \mathbb{R} \rightarrow [0, 1]$ of the standard normal distribution is given by

$$F(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x'^2}{2}\right) dx'. \quad (89)$$

It is known that this function is not elementary.



The standard normal distribution's CDF.

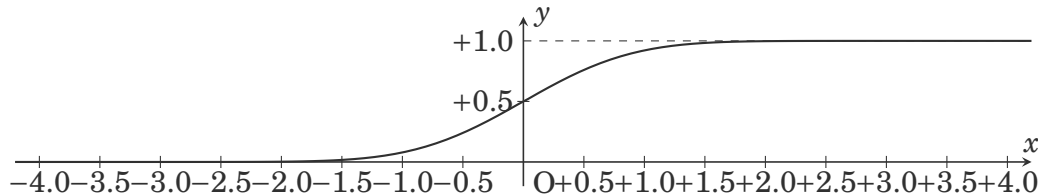
Example of the convergence by the CLT

Example

Let X_1, X_2, \dots be an infinite sequence of independently identically distributed RVs, where X_i takes 0 or +1 with probability $\frac{1}{2}$ for each.

Then the mean and the variance of X_i are $\frac{1}{2}$ and $\frac{1}{4}$, respectively.

According to the CLT, the CDF of $2\sqrt{m}\left(\overline{X_m} - \frac{1}{2}\right)$ converges to that of the standard normal distribution $\mathcal{N}(0, 1)$.



Dashed: the standard normal distribution's CDF.

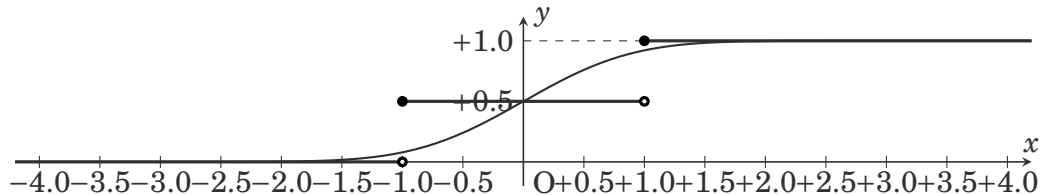
Example of the convergence by the CLT

Example

Let X_1, X_2, \dots be an infinite sequence of independently identically distributed RVs, where X_i takes 0 or +1 with probability $\frac{1}{2}$ for each.

Then the mean and the variance of X_i are $\frac{1}{2}$ and $\frac{1}{4}$, respectively.

According to the CLT, the CDF of $2\sqrt{m}\left(\overline{X_m} - \frac{1}{2}\right)$ converges to that of the standard normal distribution $\mathcal{N}(0, 1)$.



Dashed: the standard normal distribution's CDF. Solid: the CDF of $2\sqrt{m}\left(\overline{X_m} - \frac{1}{2}\right)$, where $m = 1$.

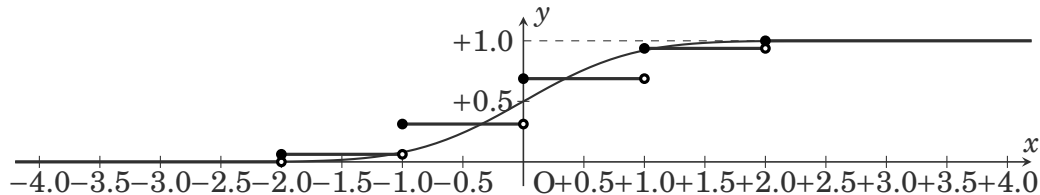
Example of the convergence by the CLT

Example

Let X_1, X_2, \dots be an infinite sequence of independently identically distributed RVs, where X_i takes 0 or +1 with probability $\frac{1}{2}$ for each.

Then the mean and the variance of X_i are $\frac{1}{2}$ and $\frac{1}{4}$, respectively.

According to the CLT, the CDF of $2\sqrt{m}\left(\overline{X_m} - \frac{1}{2}\right)$ converges to that of the standard normal distribution $\mathcal{N}(0, 1)$.



Dashed: the standard normal distribution's CDF. Solid: the CDF of $2\sqrt{m}\left(\overline{X_m} - \frac{1}{2}\right)$, where $m = 4$.

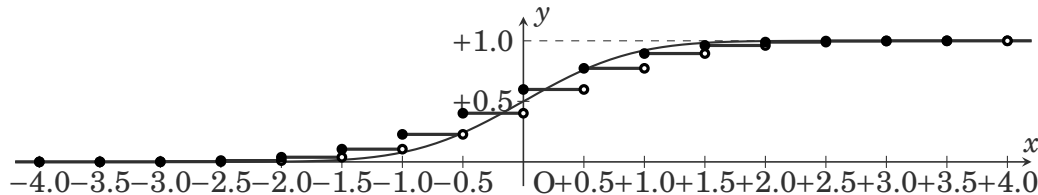
Example of the convergence by the CLT

Example

Let X_1, X_2, \dots be an infinite sequence of independently identically distributed RVs, where X_i takes 0 or +1 with probability $\frac{1}{2}$ for each.

Then the mean and the variance of X_i are $\frac{1}{2}$ and $\frac{1}{4}$, respectively.

According to the CLT, the CDF of $2\sqrt{m}\left(\overline{X}_m - \frac{1}{2}\right)$ converges to that of the standard normal distribution $\mathcal{N}(0, 1)$.



Dashed: the standard normal distribution's CDF. Solid: the CDF of $2\sqrt{m}\left(\overline{X}_m - \frac{1}{2}\right)$, where $m = 16$.

Example of the convergence by the CLT

Example

Let X_1, X_2, \dots be an infinite sequence of independently identically distributed RVs, where X_i takes 0 or +1 with probability $\frac{1}{2}$ for each.

Then the mean and the variance of X_i are $\frac{1}{2}$ and $\frac{1}{4}$, respectively.

According to the CLT, the CDF of $2\sqrt{m}\left(\overline{X_m} - \frac{1}{2}\right)$ converges to that of the standard normal distribution $\mathcal{N}(0, 1)$.

Note that the CLT is about the CDF, but **NOT about the PDF**. The convergence of the PDF does not always hold. Specifically, in the above case, $\overline{X_m}$ is a discrete random variable since each X_i is. Hence, the random variable $2\sqrt{m}\left(\overline{X_m} - \frac{1}{2}\right)$ does not have a PDF.

Therefore, we **CANNOT** say that the PDF of $2\sqrt{m}\left(\overline{X_m} - \frac{1}{2}\right)$ converges to that of the standard normal distribution.

The implications of the CLT

- The error $\bar{X}_m - \mu$ in estimating the true mean μ is almost proportional to $\frac{1}{\sqrt{m}}$. In particular, the more data points, the more accurate the estimate is.
- The sum of sufficiently many independent random variables approximately follows a normal distribution. In particular, various types of random variables decomposable to many independent factors follow a normal distribution. This is why **the normal distribution appears everywhere in the real world.**

Outline

4. Sample Statistics

4.1 Introduction: why do we learn sample statistics?

4.2 Terminology

4.3 Sample mean, law of large numbers, and central limit theorem

4.4 Estimation of distribution and parametric model

4.5 Likelihood

4.6 Maximum likelihood estimator

4.7 Exercises

Estimation of a distribution

We have estimated the expectation only. In real applications, we might want to estimate the distribution itself. However, if the support of the distribution is an infinite set¹¹, it is not practical to determine a PMF or PDF from finite data points with no assumptions.

We often assume that the distribution is in a parametric model, which is a set of distributions parametrized by a few values.

¹¹This is almost always the case if we consider a continuous RV

Parametric model

Definition (A parametric model)

- **A discrete parametric model** on support $\mathcal{X} \subset \mathbb{R}^n$ is a pair of a parameter set $\Theta \subset \mathbb{R}^k$ and a parametrized PMF $P : \mathcal{X} \times \Theta \rightarrow [0, 1]$ such that $P(\mathbf{x}; \boldsymbol{\theta})$ is a PMF on \mathcal{X} as a function of \mathbf{x} for all $\boldsymbol{\theta} \in \Theta$.
- **A continuous parametric model** on support \mathbb{R}^n is a pair of a parameter set $\Theta \subset \mathbb{R}^k$ and a parametrized PDF $p : \mathbb{R}^n \times \Theta \rightarrow \mathbb{R}_{\geq 0}$ such that $p(\mathbf{x}; \boldsymbol{\theta})$ is a PDF on \mathbb{R}^n as a function of \mathbf{x} for all $\boldsymbol{\theta} \in \Theta$.

Here, the nonnegative integer k is the dimension of the parameter.

When we have a parametric model, estimating a parameter corresponds to estimating a distribution.

Parametric model example 1: Bernoulli distribution

Example (Bernoulli distribution)

The Bernoulli distribution¹² is a discrete parametric model with a sole parameter, which is usually denoted by θ . The support and the parameter set are $\mathcal{X} = \{0, 1\}$ and $\Theta = [0, 1]$, respectively. The parametrized PMF $P(x; \theta)$ is given by $P(1; \theta) = \theta$. Thus, we have $P(0; \theta) = 1 - \theta$.

Theorem

The mean and the variance of a RV following the Bernoulli distribution with the parameter θ are θ and $\theta(1 - \theta)$, respectively.

¹²A parametric model is often called like the XXX distribution, but it is, indeed, a parametrized **set** of distributions.

Parametric model example 2: normal distribution

Example (Normal distribution)

The normal distribution, also known as the **Gaussian distribution**, is a continuous parametric model, which has mean parameter $\mu \in \mathbb{R}$ and variance parameter $\sigma^2 \in \mathbb{R}_{>0}$. That is, the parameter set is $\Theta = \mathbb{R} \times \mathbb{R}_{>0}$. The parametrized PDF $p(x; \mu, \sigma^2)$ is given by

$$p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

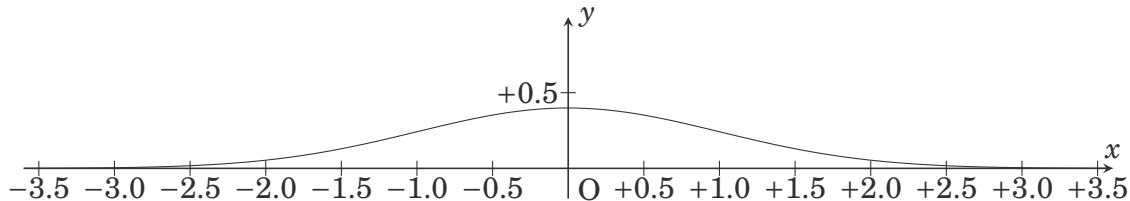
Theorem

The mean and the variance of a RV following the normal distribution with mean parameter μ and variance parameter σ^2 are μ and σ^2 , respectively.

PDF of the normal distribution

The **normal distribution**, also known as the **Gaussian distribution** with a mean parameter $\mu \in \mathbb{R}$ and a variance parameter $\sigma^2 \in \mathbb{R}_{>0}$ is a distribution with the following PDF:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right). \quad (90)$$



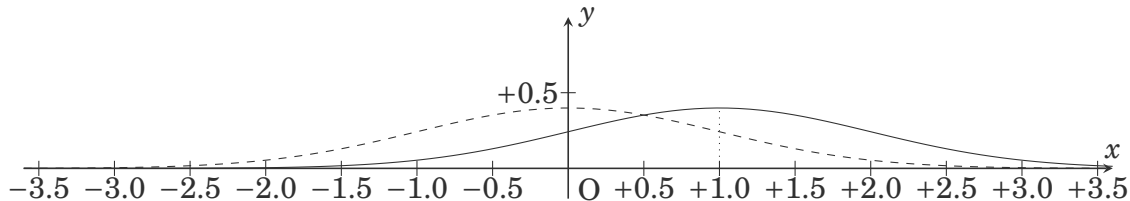
Normal distributions' PDF ($\mu = 0, \sigma = 1$).

The mean, the variance, and the standard deviation are μ , σ^2 , and $\sigma := \sqrt{\sigma^2}$, respectively.

PDF of the normal distribution

The **normal distribution**, also known as the **Gaussian distribution** with a mean parameter $\mu \in \mathbb{R}$ and a variance parameter $\sigma^2 \in \mathbb{R}_{>0}$ is a distribution with the following PDF:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right). \quad (90)$$



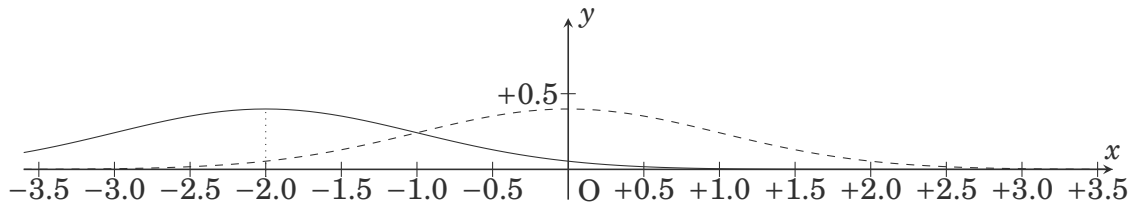
Normal distributions' PDF (Solid: $\mu = 1, \sigma = 1$, Dashed: $\mu = 0, \sigma = 1$).

The mean, the variance, and the standard deviation are μ , σ^2 , and $\sigma := \sqrt{\sigma^2}$, respectively.

PDF of the normal distribution

The **normal distribution**, also known as the **Gaussian distribution** with a mean parameter $\mu \in \mathbb{R}$ and a variance parameter $\sigma^2 \in \mathbb{R}_{>0}$ is a distribution with the following PDF:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right). \quad (90)$$



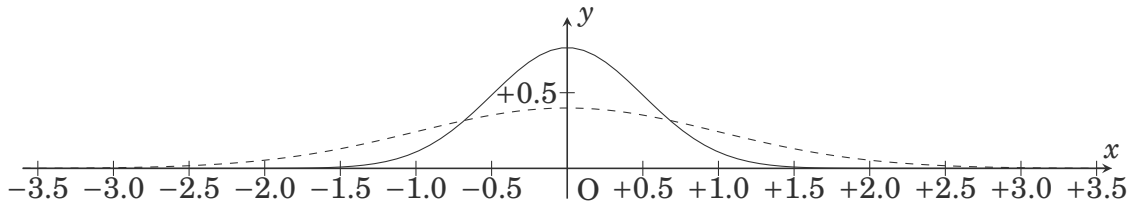
Normal distributions' PDF (Solid: $\mu = -2, \sigma = 1$, Dashed: $\mu = 0, \sigma = 1$).

The mean, the variance, and the standard deviation are μ , σ^2 , and $\sigma := \sqrt{\sigma^2}$, respectively.

PDF of the normal distribution

The **normal distribution**, also known as the **Gaussian distribution** with a mean parameter $\mu \in \mathbb{R}$ and a variance parameter $\sigma^2 \in \mathbb{R}_{>0}$ is a distribution with the following PDF:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right). \quad (90)$$



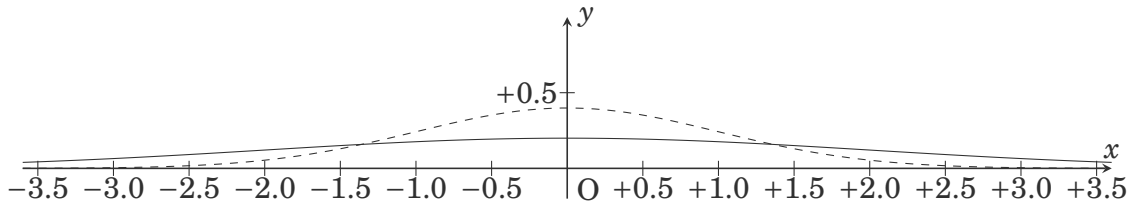
Normal distributions' PDF (Solid: $\mu = 0, \sigma = 0.5$, Dashed: $\mu = 0, \sigma = 1$).

The mean, the variance, and the standard deviation are μ , σ^2 , and $\sigma := \sqrt{\sigma^2}$, respectively.

PDF of the normal distribution

The **normal distribution**, also known as the **Gaussian distribution** with a mean parameter $\mu \in \mathbb{R}$ and a variance parameter $\sigma^2 \in \mathbb{R}_{>0}$ is a distribution with the following PDF:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right). \quad (90)$$



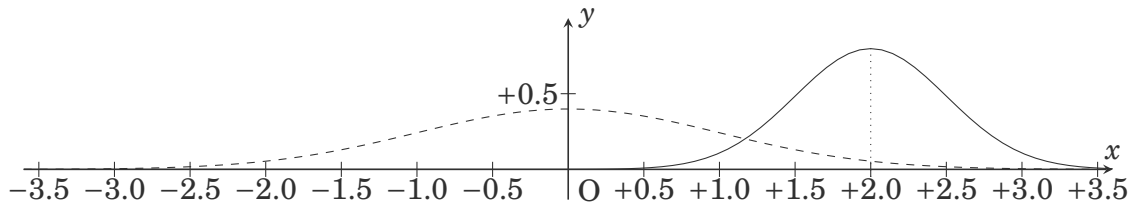
Normal distributions' PDF (Solid: $\mu = 0, \sigma = 2.0$, Dashed: $\mu = 0, \sigma = 1$).

The mean, the variance, and the standard deviation are μ , σ^2 , and $\sigma := \sqrt{\sigma^2}$, respectively.

PDF of the normal distribution

The **normal distribution**, also known as the **Gaussian distribution** with a mean parameter $\mu \in \mathbb{R}$ and a variance parameter $\sigma^2 \in \mathbb{R}_{>0}$ is a distribution with the following PDF:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right). \quad (90)$$



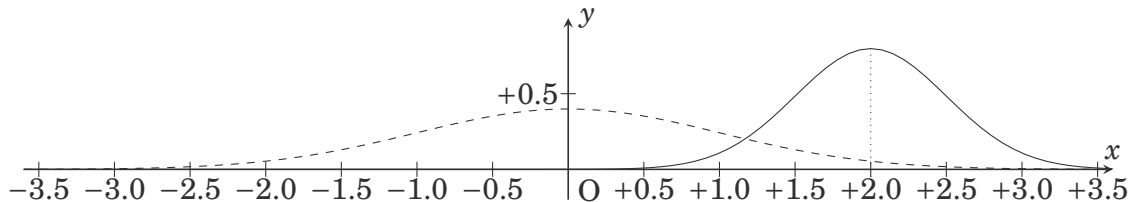
Normal distributions' PDF (Solid: $\mu = 2, \sigma = 0.5$, Dashed: $\mu = 0, \sigma = 1$).

The mean, the variance, and the standard deviation are μ , σ^2 , and $\sigma := \sqrt{\sigma^2}$, respectively.

PDF of the normal distribution

The **normal distribution**, also known as the **Gaussian distribution** with a mean parameter $\mu \in \mathbb{R}$ and a variance parameter $\sigma^2 \in \mathbb{R}_{>0}$ is a distribution with the following PDF:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right). \quad (90)$$



Normal distributions' PDF (Solid: $\mu = 2, \sigma = 0.5$, Dashed: $\mu = 0, \sigma = 1$).

The PDF is symmetric about $x = \mu$ and it is dense around $x = \mu$.

Outline

4. Sample Statistics

4.1 Introduction: why do we learn sample statistics?

4.2 Terminology

4.3 Sample mean, law of large numbers, and central limit theorem

4.4 Estimation of distribution and parametric model

4.5 Likelihood

4.6 Maximum likelihood estimator

4.7 Exercises

Likelihood

To determine a parameter of a parametric model from data points, we quantify how “likely” the distribution indicated by a parameter is correct.

When we have a PMF or PDF of a distribution, we simply define the value of the PMF or PDF of the data points as the **likelihood** of the distribution.

Definition (Likelihood of a discrete parametric model)

Let $P(\cdot; \cdot)$ be a discrete parametric model with a parameter set Θ and $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$ be values of data points.

Then the **likelihood** of $P(\cdot; \boldsymbol{\theta})$ (or often called the likelihood of the parameter $\boldsymbol{\theta}$) is defined as the following product.

$$P(\mathbf{x}_1; \boldsymbol{\theta}) \cdot P(\mathbf{x}_2; \boldsymbol{\theta}) \cdot \dots \cdot P(\mathbf{x}_m; \boldsymbol{\theta}). \quad (91)$$

Likelihood

To determine a parameter of a parametric model from data points, we quantify how “likely” the distribution indicated by a parameter is correct.

When we have a PMF or PDF of a distribution, we simply define the value of the PMF or PDF of the data points as the **likelihood** of the distribution.

Definition (Likelihood of a continuous parametric model)

Let $p(\cdot; \cdot)$ be a continuous parametric model with a parameter set Θ and $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$ be values of data points.

Then the **likelihood** of $p(\cdot; \boldsymbol{\theta})$ (or often called the likelihood of the parameter $\boldsymbol{\theta}$) is defined as the following product.

$$p(\mathbf{x}_1; \boldsymbol{\theta}) \cdot p(\mathbf{x}_2; \boldsymbol{\theta}) \cdots p(\mathbf{x}_m; \boldsymbol{\theta}). \quad (91)$$

Examples of likelihood calculation

Suppose that we have data points $(x_1, x_2, x_3, x_4) = (1, 1, 0, 1)$, and consider the Bernoulli distribution $P(0; \theta) = 1 - \theta, P(1; \theta) = \theta$.

The likelihood of the Bernoulli distribution with θ on the data is given by

$$P(x_1; \theta)P(x_2; \theta)P(x_3; \theta)P(x_4; \theta) = P(1; \theta)P(1; \theta)P(0; \theta)P(1; \theta) = \theta \cdot \theta \cdot (1 - \theta) \cdot \theta. \quad (92)$$

- The likelihood of $\theta = 0$ is $0 \cdot 0 \cdot (1 - 0) \cdot 0 = 0$.

Examples of likelihood calculation

Suppose that we have data points $(x_1, x_2, x_3, x_4) = (1, 1, 0, 1)$, and consider the Bernoulli distribution $P(0; \theta) = 1 - \theta, P(1; \theta) = \theta$.

The likelihood of the Bernoulli distribution with θ on the data is given by

$$P(x_1; \theta)P(x_2; \theta)P(x_3; \theta)P(x_4; \theta) = P(1; \theta)P(1; \theta)P(0; \theta)P(1; \theta) = \theta \cdot \theta \cdot (1 - \theta) \cdot \theta. \quad (92)$$

- The likelihood of $\theta = 0$ is $0 \cdot 0 \cdot (1 - 0) \cdot 0 = 0$.
- The likelihood of $\theta = \frac{1}{4}$ is $\frac{1}{4} \cdot \frac{1}{4} \cdot \left(1 - \frac{1}{4}\right) \cdot \frac{1}{4} = \frac{3}{256}$.

Examples of likelihood calculation

Suppose that we have data points $(x_1, x_2, x_3, x_4) = (1, 1, 0, 1)$, and consider the Bernoulli distribution $P(0; \theta) = 1 - \theta, P(1; \theta) = \theta$.

The likelihood of the Bernoulli distribution with θ on the data is given by

$$P(x_1; \theta)P(x_2; \theta)P(x_3; \theta)P(x_4; \theta) = P(1; \theta)P(1; \theta)P(0; \theta)P(1; \theta) = \theta \cdot \theta \cdot (1 - \theta) \cdot \theta. \quad (92)$$

- The likelihood of $\theta = 0$ is $0 \cdot 0 \cdot (1 - 0) \cdot 0 = 0$.
- The likelihood of $\theta = \frac{1}{4}$ is $\frac{1}{4} \cdot \frac{1}{4} \cdot \left(1 - \frac{1}{4}\right) \cdot \frac{1}{4} = \frac{3}{256}$.
- The likelihood of $\theta = \frac{1}{2}$ is $\frac{1}{2} \cdot \frac{1}{2} \cdot \left(1 - \frac{1}{2}\right) \cdot \frac{1}{2} = \frac{1}{16} = \frac{16}{256}$.

Examples of likelihood calculation

Suppose that we have data points $(x_1, x_2, x_3, x_4) = (1, 1, 0, 1)$, and consider the Bernoulli distribution $P(0; \theta) = 1 - \theta, P(1; \theta) = \theta$.

The likelihood of the Bernoulli distribution with θ on the data is given by

$$P(x_1; \theta)P(x_2; \theta)P(x_3; \theta)P(x_4; \theta) = P(1; \theta)P(1; \theta)P(0; \theta)P(1; \theta) = \theta \cdot \theta \cdot (1 - \theta) \cdot \theta. \quad (92)$$

- The likelihood of $\theta = 0$ is $0 \cdot 0 \cdot (1 - 0) \cdot 0 = 0$.
- The likelihood of $\theta = \frac{1}{4}$ is $\frac{1}{4} \cdot \frac{1}{4} \cdot \left(1 - \frac{1}{4}\right) \cdot \frac{1}{4} = \frac{3}{256}$.
- The likelihood of $\theta = \frac{1}{2}$ is $\frac{1}{2} \cdot \frac{1}{2} \cdot \left(1 - \frac{1}{2}\right) \cdot \frac{1}{2} = \frac{1}{16} = \frac{16}{256}$.
- The likelihood of $\theta = \frac{3}{4}$ is $\frac{3}{4} \cdot \frac{3}{4} \cdot \left(1 - \frac{3}{4}\right) \cdot \frac{3}{4} = \frac{27}{256}$.

Examples of likelihood calculation

Suppose that we have data points $(x_1, x_2, x_3, x_4) = (1, 1, 0, 1)$, and consider the Bernoulli distribution $P(0; \theta) = 1 - \theta, P(1; \theta) = \theta$.

The likelihood of the Bernoulli distribution with θ on the data is given by

$$P(x_1; \theta)P(x_2; \theta)P(x_3; \theta)P(x_4; \theta) = P(1; \theta)P(1; \theta)P(0; \theta)P(1; \theta) = \theta \cdot \theta \cdot (1 - \theta) \cdot \theta. \quad (92)$$

- The likelihood of $\theta = 0$ is $0 \cdot 0 \cdot (1 - 0) \cdot 0 = 0$.
- The likelihood of $\theta = \frac{1}{4}$ is $\frac{1}{4} \cdot \frac{1}{4} \cdot \left(1 - \frac{1}{4}\right) \cdot \frac{1}{4} = \frac{3}{256}$.
- The likelihood of $\theta = \frac{1}{2}$ is $\frac{1}{2} \cdot \frac{1}{2} \cdot \left(1 - \frac{1}{2}\right) \cdot \frac{1}{2} = \frac{1}{16} = \frac{16}{256}$.
- The likelihood of $\theta = \frac{3}{4}$ is $\frac{3}{4} \cdot \frac{3}{4} \cdot \left(1 - \frac{3}{4}\right) \cdot \frac{3}{4} = \frac{27}{256}$.
- The likelihood of $\theta = 1$ is $1 \cdot 1 \cdot (1 - 1) \cdot 1 = 0$.

Examples of likelihood calculation

Suppose that we have data points $(x_1, x_2, x_3, x_4) = (1, 1, 0, 1)$, and consider the Bernoulli distribution $P(0; \theta) = 1 - \theta, P(1; \theta) = \theta$.

The likelihood of the Bernoulli distribution with θ on the data is given by

$$P(x_1; \theta)P(x_2; \theta)P(x_3; \theta)P(x_4; \theta) = P(1; \theta)P(1; \theta)P(0; \theta)P(1; \theta) = \theta \cdot \theta \cdot (1 - \theta) \cdot \theta. \quad (92)$$

- The likelihood of $\theta = 0$ is $0 \cdot 0 \cdot (1 - 0) \cdot 0 = 0$.
- The likelihood of $\theta = \frac{1}{4}$ is $\frac{1}{4} \cdot \frac{1}{4} \cdot \left(1 - \frac{1}{4}\right) \cdot \frac{1}{4} = \frac{3}{256}$.
- The likelihood of $\theta = \frac{1}{2}$ is $\frac{1}{2} \cdot \frac{1}{2} \cdot \left(1 - \frac{1}{2}\right) \cdot \frac{1}{2} = \frac{1}{16} = \frac{16}{256}$.
- The likelihood of $\theta = \frac{3}{4}$ is $\frac{3}{4} \cdot \frac{3}{4} \cdot \left(1 - \frac{3}{4}\right) \cdot \frac{3}{4} = \frac{27}{256}$.
- The likelihood of $\theta = 1$ is $1 \cdot 1 \cdot (1 - 1) \cdot 1 = 0$.

Hence, among the above three, the distribution given by $\theta = \frac{3}{4}$ most likely generates the data sequence $(x_1, x_2, x_3, x_4) = (1, 1, 0, 1)$.

Probability and likelihood

The value of the product

$$P(\mathbf{x}_1; \boldsymbol{\theta}) \cdot P(\mathbf{x}_2; \boldsymbol{\theta}) \cdots P(\mathbf{x}_m; \boldsymbol{\theta}) \quad (93)$$

can be interpreted as either

- the probability of the random variable sequence taking the value sequence $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$, i.e., a function of a value sequence, or
- the likelihood of the distribution determined by the parameter $\boldsymbol{\theta}$, i.e., a function of a distribution (or parameter).

In other words, the above product is the probability (or the probability density for continuous distribution case) if we interpret it as a function of a value sequence, and the likelihood if we interpret it as a function of a distribution (or a parameter).

Outline

4. Sample Statistics

4.1 Introduction: why do we learn sample statistics?

4.2 Terminology

4.3 Sample mean, law of large numbers, and central limit theorem

4.4 Estimation of distribution and parametric model

4.5 Likelihood

4.6 Maximum likelihood estimator

4.7 Exercises

Maximum likelihood estimator

Once we define the likelihood of a distribution, all we need to do is find a parameter that maximizes the likelihood.

The parameter vector that maximizes the likelihood is called the **maximum likelihood estimator (MLE)**.

Definition (Maximum likelihood estimator)

Let $P(\cdot; \cdot)$ be a discrete parametric model with a parameter set Θ and $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$ be values of data points.

The parameter vector $\boldsymbol{\theta}$ is called a maximum likelihood estimator (MLE) if it maximizes the likelihood

$$P(\mathbf{x}_1; \boldsymbol{\theta}) \cdot P(\mathbf{x}_2; \boldsymbol{\theta}) \cdots P(\mathbf{x}_m; \boldsymbol{\theta}). \quad (94)$$

If there is a unique MLE, we often denote it by $\hat{\boldsymbol{\theta}}$.

MLE maximizes the score and minimizes the negative log likelihood

For a parameter vector θ , the following is equivalent¹³.

- The parameter vector θ maximizes the likelihood function

$$P(\mathbf{x}_1; \theta) \cdot P(\mathbf{x}_2; \theta) \cdots P(\mathbf{x}_m; \theta). \quad (95)$$

- The parameter vector θ maximizes the **log-likelihood** function

$$\log P(\mathbf{x}_1; \theta) + \log P(\mathbf{x}_2; \theta) + \cdots + \log P(\mathbf{x}_m; \theta). \quad (96)$$

- The parameter vector θ minimizes the **negative log likelihood** function

$$-\log P(\mathbf{x}_1; \theta) - \log P(\mathbf{x}_2; \theta) - \cdots - \log P(\mathbf{x}_m; \theta). \quad (97)$$

¹³It follows since log is an increasing function. It holds regardless of the base of the logarithm.

Why do we consider the logarithm of the likelihood?

- The likelihood is a product and its logarithm is a sum. When we maximize it in a computer, we rely on its derivative (gradient descent methods). Differentiation of a sum is much easier than that of a product, so the (negative) log-likelihood has an advantage over the original likelihood from the optimization viewpoint.
- If the data size m is large, the absolute value of the likelihood, the product of many small values, tends to be too small to represent in a computer (underflow). Since the logarithm sees the power index, it can handle extremely small likelihood.
- The negative log-likelihood can be interpreted as the sum of the errors. For example, we can interpret the negative log-likelihood of the normal distribution as the squared error.

The MLE of the normal distribution minimizes the square error.

Let $p(x; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$. Then, the negative (natural) log-likelihood of the data sequence is given by

$$\begin{aligned} & \log(2\pi\sigma^2) + \frac{(x_1 - \mu)^2}{2\sigma^2} + \log(2\pi\sigma^2) + \frac{(x_2 - \mu)^2}{2\sigma^2} + \cdots + \log(2\pi\sigma^2) + \frac{(x_m - \mu)^2}{2\sigma^2} \\ &= m \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \left[(x_1 - \mu)^2 + (x_2 - \mu)^2 + \cdots + (x_m - \mu)^2 \right]. \end{aligned} \quad (98)$$

The MLE of the normal distribution minimizes the square error.

Let $p(x; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$. Then, the negative (natural) log-likelihood of the data sequence is given by

$$\begin{aligned} & \log(2\pi\sigma^2) + \frac{(x_1 - \mu)^2}{2\sigma^2} + \log(2\pi\sigma^2) + \frac{(x_2 - \mu)^2}{2\sigma^2} + \cdots + \log(2\pi\sigma^2) + \frac{(x_m - \mu)^2}{2\sigma^2} \\ &= m \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \left[(x_1 - \mu)^2 + (x_2 - \mu)^2 + \cdots + (x_m - \mu)^2 \right]. \end{aligned} \quad (98)$$

When we minimize the above with respect to μ , we can ignore the gray parts.

In this sense, the MLE of the mean parameter of the normal distribution model is equivalent to minimizing the squared error.

MLE example: Bernoulli case

Example

Suppose that we have data points x_1, x_2, \dots, x_m , and consider the Bernoulli distribution $P(0; \theta) = 1 - \theta, P(1; \theta) = \theta$.

The negative log-likelihood of the Bernoulli distribution with θ on the data is given by

$$-\log P(x_1; \theta) P(x_2; \theta) \dots P(x_m; \theta) = m_0 \log(1 - \theta) + m_1 \log \theta, \quad (99)$$

where m_0 and m_1 are the numbers of zeros and ones in the data sequence. Obviously, $m_0 + m_1 = m$, and the sample mean $\bar{x} = \frac{m_1}{m}$. Let l denote the above negative log-likelihood.

Suppose that $m_0 \neq 0$ and $m_1 \neq 0$, then l takes the minimum¹⁴ if and only if $\theta = \frac{m_1}{m} = \bar{x}$.

Hence, the MLE $\hat{\theta} = \frac{m_1}{m} = \bar{x}$.

For example, if $(x_1, x_2, x_3, x_4) = (1, 1, 0, 1)$, then $\hat{\theta} = \frac{m_1}{m} = \bar{x} = \frac{3}{4}$.

¹⁴To prove it, differentiate the loss by θ and apply the first derivative test.

Why can we justify the maximum likelihood estimator (MLE)?

Similar to the sample mean, if data points are generated by a distribution indicated by a parameter vector in the parameter set of a parametric vector, the MLE has the following properties:

- **Consistency**: The MLE converges to the true parameter as $m \rightarrow \infty$.
- **Asymptotic normality**: An appropriately scaled MLE's distribution converges to a normal distribution, and its error is proportional to $\frac{1}{\sqrt{m}}$ for sufficiently large m .

Outline

4. Sample Statistics

4.1 Introduction: why do we learn sample statistics?

4.2 Terminology

4.3 Sample mean, law of large numbers, and central limit theorem

4.4 Estimation of distribution and parametric model

4.5 Likelihood

4.6 Maximum likelihood estimator

4.7 Exercises

Exercise (Standard normal distribution)

Write down the standard normal distribution's (probability density function) PDF.

Exercise (The central limit theorem (CLT))

Let X_1, X_2, \dots be an infinite sequence of independently identically distributed RVs. Assume the distribution of each random variable X_i is given by one of the following. For each case, apply the CLT and find what random variable converges to which normal distribution.

- The probability mass function P_{X_i} defined by $P_{X_i}(-1) = \frac{1}{4}$ and $P_{X_i}(+1) = \frac{3}{4}$.
- The probability density function p_{X_i} defined by $p_{X_i}(x) = \frac{1}{2\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-3}{2}\right)^2\right)$.

Exercise (Exercise: likelihood calculation)

Suppose that we have data points $(x_1, x_2, x_3, x_4) = (1, 0, 0, 1)$, and consider the Bernoulli distribution $P(0; \theta) = 1 - \theta, P(1; \theta) = \theta$.

Find the likelihoods of the Bernoulli distribution given by $\theta = 0, \frac{1}{4}, \frac{2}{4}, \frac{3}{4}, 1$.

Also, answer which distribution most likely generates the data points.

Outline

5. Statistical Test

5.1 Introduction: why do we learn statistical tests?

5.2 The logic of statistical tests

5.3 Example test statistics

5.4 p-value

5.5 Failure of statistical test

5.6 Exercises

Outline

5. Statistical Test

5.1 Introduction: why do we learn statistical tests?

5.2 The logic of statistical tests

5.3 Example test statistics

5.4 p-value

5.5 Failure of statistical test

5.6 Exercises

Statistical tests support our judgements

In real applications (e.g., physical, engineering, medical, etc.), we need to judge from data whether a phenomenon happens or not.

Specifically, for some summary statistics or parameter θ and a set \mathcal{H}_1 , we often want to judge from data points whether $\theta \in \mathcal{H}_1$ or not.

Statistical tests support our judgements

In real applications (e.g., physical, engineering, medical, etc.), we need to judge from data whether a phenomenon happens or not.

Specifically, for some summary statistics or parameter θ and a set \mathcal{H}_1 , we often want to judge from data points whether $\theta \in \mathcal{H}_1$ or not.

For example, if we investigate the purity of a factory's chemical product, we might want to know whether the true expectation μ of the purity is the same as the purity μ_0 of the natural material or not.

In this case, $\mathcal{H}_1 = [0, 1] \setminus \{\mu_0\}$, and we want to discuss whether $\theta \in \mathcal{H}_1$ or not.

Statistical tests give us a framework to make such a judgement.

Learning outcomes

By the end of this section, you should be able to:

- Explain the logic of statistical tests
- Explain the definitions of p-value, significance level, type-I error, and type-II error.
- Make a judgment from data using statistical tests

Outline

5. Statistical Test

5.1 Introduction: why do we learn statistical tests?

5.2 The logic of statistical tests

5.3 Example test statistics

5.4 p-value

5.5 Failure of statistical test

5.6 Exercises

We cannot directly prove that “the hypothesis is correct.”

What we want to “prove” is the following statement: “if the data points’ values are x_1, x_2, \dots, x_m , then $\theta \in \mathcal{H}_1$,” in some probability theory sense.

A naïve idea is to evaluate the “probability” of $\theta \in \mathcal{H}_1$ when the data points’ values are x_1, x_2, \dots, x_m .

We cannot directly prove that “the hypothesis is correct.”

What we want to “prove” is the following statement: “if the data points’ values are x_1, x_2, \dots, x_m , then $\theta \in \mathcal{H}_1$,” in some probability theory sense.

A naïve idea is to evaluate the “probability” of $\theta \in \mathcal{H}_1$ when the data points’ values are x_1, x_2, \dots, x_m .

However, in (frequentism) statistics, we cannot discuss the probability of a parameter θ being in a set since a parameter θ is not a random variable, while it regards data points x_1, x_2, \dots, x_m as values of random variables.

We cannot directly prove that “the hypothesis is correct.”

What we want to “prove” is the following statement: “if the data points’ values are x_1, x_2, \dots, x_m , then $\theta \in \mathcal{H}_1$,” in some probability theory sense.

A naïve idea is to evaluate the “probability” of $\theta \in \mathcal{H}_1$ when the data points’ values are x_1, x_2, \dots, x_m .

However, in (frequentism) statistics, we cannot discuss the probability of a parameter θ being in a set since a parameter θ is not a random variable, while it regards data points x_1, x_2, \dots, x_m as values of random variables.

In contrast, we can discuss the other direction, that is, given a parameter θ , we can discuss the probability of the random variables taking the given values x_1, x_2, \dots, x_m .

So, we take the **contraposition** of the statement that we originally wanted to prove.

The fundamental logic of statistical test

The contraposition of “if the data points’ values are x_1, x_2, \dots, x_m , then $\theta \in \mathcal{H}_1$,” is:

“If $\theta \notin \mathcal{H}_1$, then the data points’ values are NOT x_1, x_2, \dots, x_m .”

Hence, discussing the event $\theta \notin \mathcal{H}_1$ is essential.

The fundamental logic of statistical test

The contraposition of “if the data points’ values are x_1, x_2, \dots, x_m , then $\theta \in \mathcal{H}_1$,” is:

“If $\theta \notin \mathcal{H}_1$, then the data points’ values are NOT x_1, x_2, \dots, x_m .”

Hence, discussing the event $\theta \notin \mathcal{H}_1$ is essential.

Let \mathcal{H} be the set of all the possible values that θ can take and define $\mathcal{H}_0 := \mathcal{H} \setminus \mathcal{H}_1$.

The event $\theta \notin \mathcal{H}_1$, which we focus on, is equivalent to $\theta \in \mathcal{H}_0$.

Hence, \mathcal{H}_0 plays an essential role in statistical tests. \mathcal{H}_0 is called the **null hypothesis** and \mathcal{H}_1 is called the **alternative hypothesis**.

In statistical tests, a **hypothesis** is a set of values that the variable θ , which we are interested in, may take.

Test statistics

Our starting point is to assume $\theta \notin \mathcal{H}_1$, or equivalently, $\theta \in \mathcal{H}_0$. Our objective is that the data points x_1, x_2, \dots, x_m “contradict in a probability theory sense” the assumption.

Test statistics

Our starting point is to assume $\theta \notin \mathcal{H}_1$, or equivalently, $\theta \in \mathcal{H}_0$. Our objective is that the data points x_1, x_2, \dots, x_m “contradict in a probability theory sense” the assumption.

To judge whether a “contradiction” happens, we evaluate a summary statistic of the empirical distribution. Such a summary statistic is called a **test statistic**. The test statistic is a RV since it is a function of the data points, which are the values of RVs. Hence, the distribution of a test statistic is determined if we fix a distribution of the data points.

Test statistics

Our starting point is to assume $\theta \notin \mathcal{H}_1$, or equivalently, $\theta \in \mathcal{H}_0$. Our objective is that the data points x_1, x_2, \dots, x_m “contradict in a probability theory sense” the assumption.

To judge whether a “contradiction” happens, we evaluate a summary statistic of the empirical distribution. Such a summary statistic is called a **test statistic**. The test statistic is a RV since it is a function of the data points, which are the values of RVs. Hence, the distribution of a test statistic is determined if we fix a distribution of the data points.

For a distribution corresponding to \mathcal{H}_0 , if the value of the test statistic is unlikely taken on the distribution (i.e. if a “probabilistic contradiction” happens), then we can conclude that the data points are not generated by the distribution. That is, we can conclude $\theta \notin \mathcal{H}_0$, i.e., $\theta \in \mathcal{H}_1$. This is the basic idea of the statistical test.

Terminology: rejecting and accepting a hypothesis

- We say that we **reject** a hypothesis when we conclude that the true distribution is **not in** the distributions corresponding to the hypothesis.
- We say that we **accept** a hypothesis when we conclude that the true distribution is **in** the distributions corresponding to the hypothesis.

Outline

5. Statistical Test

5.1 Introduction: why do we learn statistical tests?

5.2 The logic of statistical tests

5.3 Example test statistics

5.4 p-value

5.5 Failure of statistical test

5.6 Exercises

Example: are our products better?

We are going to compose a component purer than a natural one. Suppose that the purity of a natural one is 92% on average.

Our factory composed a component 8 times and the purity was the following:

Trial	1	2	3	4	5	6	7	8
Purity	95	93	94	94	92	93	91	96

8 trial results of our factory

Are our factory's products better than natural ones on average?

The sample mean of the factory's products is 93.5, which is better than 92, the natural components average. Could we conclude that our factory's products are better than natural components?

What's our concern?

The sample mean of the factory's products is 93.5, which is better than 92, the natural components average.

A possible bad story is that the true mean μ_0 is not larger than 92, but the sample mean was "luckily" 93.5, better than 92, owing to its stochastic behavior. This is our concern.

Hence, we consider how likely this bad story can happen by "luck."

***t*-test about the true expectation**

Suppose that X_1, X_2, \dots, X_m are random variables independently and identically following the normal distribution with an unknown true expectation μ and variance σ^2 .

We want to see whether or not the true mean equals a value μ_0 . That is, the null hypothesis is $\mathcal{H}_0 = \{\mu_0\}$.

Following the idea of the statistical test, we evaluate whether or not those random variables' values are extreme under the null hypothesis $\mu = \mu_0$. For this purpose, we consider the following value, called ***t*-statistic**.

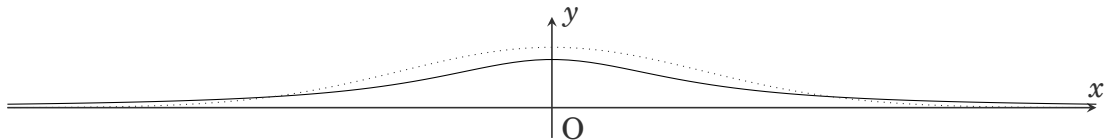
$$t := \frac{\bar{X} - \mu_0}{\frac{s}{\sqrt{m}}}, \quad (100)$$

where \bar{X} and s are the sample mean and sample standard deviation defined by

$$\bar{X} := \frac{1}{m} \sum_{i=1}^m X_i, \quad s := \sqrt{\frac{1}{m} \sum_{i=1}^m (X_i - \bar{X})^2}. \quad (101)$$

t -distribution

Suppose that X_1, X_2, \dots, X_m are independently and identically following a normal distribution. Then, t follows the t -distribution with $m - 1$ degree of freedom, whose PDF p_{m-1} is illustrated as follows.



Black solid curve: the PDF of the t -distribution with 1 degree of freedom.

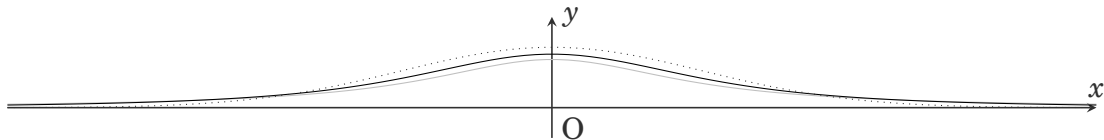
Black dotted curve: the PDF of the standard normal distribution.

The t -distribution's PDF is symmetric and similar to the standard normal distribution's PDF but has a larger probability of taking extremely large or small values.

As m increases, the PDF converges to the standard normal distribution's PDF.

t -distribution

Suppose that X_1, X_2, \dots, X_m are independently and identically following a normal distribution. Then, t follows the t -distribution with $m - 1$ degree of freedom, whose PDF p_{m-1} is illustrated as follows.



Black solid curve: the PDF of the t -distribution with 2 degree of freedom.

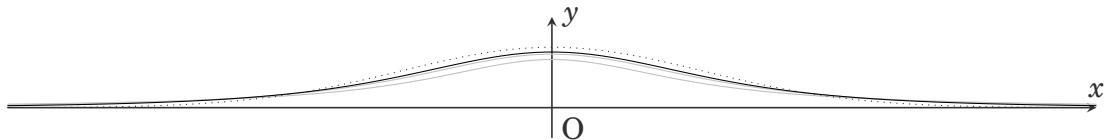
Black dotted curve: the PDF of the standard normal distribution.

The t -distribution's PDF is symmetric and similar to the standard normal distribution's PDF but has a larger probability of taking extremely large or small values.

As m increases, the PDF converges to the standard normal distribution's PDF.

t -distribution

Suppose that X_1, X_2, \dots, X_m are independently and identically following a normal distribution. Then, t follows the t -distribution with $m - 1$ degree of freedom, whose PDF p_{m-1} is illustrated as follows.



Black solid curve: the PDF of the t -distribution with 3 degree of freedom.

Black dotted curve: the PDF of the standard normal distribution.

The t -distribution's PDF is symmetric and similar to the standard normal distribution's PDF but has a larger probability of taking extremely large or small values.

As m increases, the PDF converges to the standard normal distribution's PDF.

Note: The specific form of the t -distribution.

The PDF $p_{m-1}(x)$ of the t -distribution with $m - 1$ degree of freedom is given by

$$p_{m-1}(x) = \frac{\Gamma(\frac{m}{2})}{\sqrt{(m-1)\pi}\Gamma(\frac{m-1}{2})} \left(1 + \frac{x^2}{m-1}\right)^{-\frac{m}{2}} \quad (102)$$

where $\Gamma(z) := \int_0^\infty s^{z-1} \exp(-s) ds$.

t -test is not limited to the one about the true expectation.

We have focused on a statistical test about the true expectation.

In general, a statistic is called a t statistic if it follows the t distribution. Also, a statistical test using a t statistic is called a t test. Hence, if you find a t test in another context, it might not be about the true expectation. It is always essential to confirm what the null hypothesis is and what the alternative hypothesis is in the context you are interested in.

Outline

5. Statistical Test

5.1 Introduction: why do we learn statistical tests?

5.2 The logic of statistical tests

5.3 Example test statistics

5.4 p-value

5.5 Failure of statistical test

5.6 Exercises

p-value

How do we determine the unlikeliness of the value of the test statistic?

As a criterion of the unlikeliness of the statistic's value, we consider the probability of the statistic taking a more extreme value ¹⁵. The probability is called the **p-value**. A small p-value indicates that the value of the statistic takes an extreme value.

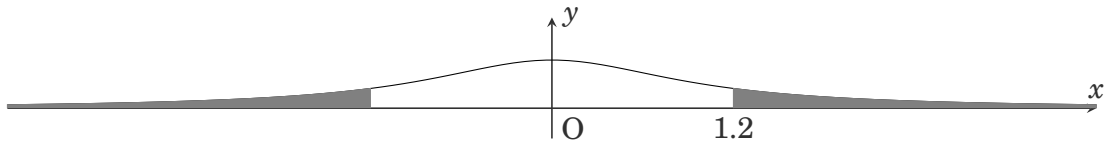
¹⁵Hence, we need to define in which case the value of the statistic is extreme. Although it is intuitive for well-known cases, there does not seem to be a way to mathematically decide it.

p-value in t-test

The t -statistic takes zero if $\bar{X} = \mu$. In non-extreme cases, where the sample mean \bar{X} is around the mean μ , t is around zero. In extreme cases, where the sample mean \bar{X} is distant from the mean μ , $|t|$ takes a large value. The larger $|t|$, the more extreme.

Here, when t -statistic takes a value t_0 , we define its p-value by

$$p = \Pr(|t| > |t_0|). \quad (103)$$



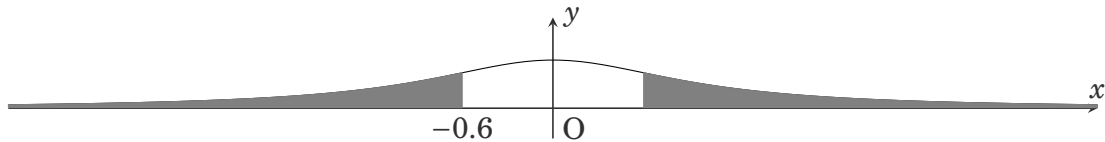
The p-value (the gray area) when t takes $t_0 = 1.2$.

p-value in t-test

The t -statistic takes zero if $\bar{X} = \mu$. In non-extreme cases, where the sample mean \bar{X} is around the mean μ , t is around zero. In extreme cases, where the sample mean \bar{X} is distant from the mean μ , $|t|$ takes a large value. The larger $|t|$, the more extreme.

Here, when t -statistic takes a value t_0 , we define its p-value by

$$p = \Pr(|t| > |t_0|). \quad (103)$$



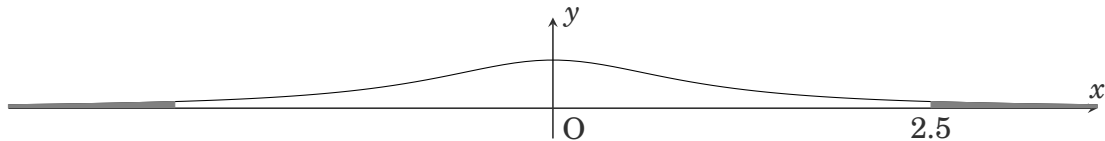
The p-value (the gray area) when t takes $t_0 = -0.6$.

p-value in t-test

The t -statistic takes zero if $\bar{X} = \mu$. In non-extreme cases, where the sample mean \bar{X} is around the mean μ , t is around zero. In extreme cases, where the sample mean \bar{X} is distant from the mean μ , $|t|$ takes a large value. The larger $|t|$, the more extreme.

Here, when t -statistic takes a value t_0 , we define its p-value by

$$p = \Pr(|t| > |t_0|). \quad (103)$$



The p-value (the gray area) when t takes $t_0 = 2.5$.

Significance level

We reject a hypothesis consisting of a single distribution if the p-value of the distribution on the data points is small¹⁶.

Now, how small should the threshold, called the ***significance level*** be?

There is no mathematical reason to determine it.

There is a convention to set the threshold at 0.05.

That is,

- If p-value is larger than 0.05, then we do not reject the null hypothesis \mathcal{H}_0 .
- If p-value is smaller than 0.05, then we reject the null hypothesis and accept the alternative hypothesis \mathcal{H}_1 .

¹⁶We reject a hypothesis consisting of multiple distributions if we can reject the hypothesis consisting of any distribution in the original hypothesis

Statistical test procedure

The standard procedure of the statistical test is the following.

- **Step 1:**

- **Step 2:**

- **Step 3:**

- **Step 4:**

Statistical test procedure

The standard procedure of the statistical test is the following.

- **Step 1:** Set the null hypothesis and alternative hypothesis. Also, fix the significance level α (usually 0.05 or 0.005) and determine which statistic to use.
- **Step 2:**
- **Step 3:**
- **Step 4:**

Statistical test procedure

The standard procedure of the statistical test is the following.

- **Step 1:** Set the null hypothesis and alternative hypothesis. Also, fix the significance level α (usually 0.05 or 0.005) and determine which statistic to use. For example, if we are interested in the true expectation, the null hypothesis is $\mu = \mu_0$, where μ is the unknown true expectation and μ_0 is a value, which we decide. The alternative hypothesis is $\mu \neq \mu_0$. We can use t statistic.
- **Step 2:**
- **Step 3:**
- **Step 4:**

Statistical test procedure

The standard procedure of the statistical test is the following.

- **Step 1:** Set the null hypothesis and alternative hypothesis. Also, fix the significance level α (usually 0.05 or 0.005) and determine which statistic to use. For example, if we are interested in the true expectation, the null hypothesis is $\mu = \mu_0$, where μ is the unknown true expectation and μ_0 is a value, which we decide. The alternative hypothesis is $\mu \neq \mu_0$. We can use t statistic.
- **Step 2:** Calculate the statistic.
- **Step 3:**
- **Step 4:**

Statistical test procedure

The standard procedure of the statistical test is the following.

- **Step 1:** Set the null hypothesis and alternative hypothesis. Also, fix the significance level α (usually 0.05 or 0.005) and determine which statistic to use. For example, if we are interested in the true expectation, the null hypothesis is $\mu = \mu_0$, where μ is the unknown true expectation and μ_0 is a value, which we decide. The alternative hypothesis is $\mu \neq \mu_0$. We can use t statistic.
- **Step 2:** Calculate the statistic. For example, in the t -test, calculate $t = \frac{\sqrt{m}(\bar{X} - \mu_0)}{s}$.
- **Step 3:**
- **Step 4:**

Statistical test procedure

The standard procedure of the statistical test is the following.

- **Step 1:** Set the null hypothesis and alternative hypothesis. Also, fix the significance level α (usually 0.05 or 0.005) and determine which statistic to use. For example, if we are interested in the true expectation, the null hypothesis is $\mu = \mu_0$, where μ is the unknown true expectation and μ_0 is a value, which we decide. The alternative hypothesis is $\mu \neq \mu_0$. We can use t statistic.
- **Step 2:** Calculate the statistic. For example, in the t -test, calculate $t = \frac{\sqrt{m}(\bar{X} - \mu_0)}{s}$.
- **Step 3:** Evaluate the p -value from the value of the statistic.
- **Step 4:**

Statistical test procedure

The standard procedure of the statistical test is the following.

- **Step 1:** Set the null hypothesis and alternative hypothesis. Also, fix the significance level α (usually 0.05 or 0.005) and determine which statistic to use. For example, if we are interested in the true expectation, the null hypothesis is $\mu = \mu_0$, where μ is the unknown true expectation and μ_0 is a value, which we decide. The alternative hypothesis is $\mu \neq \mu_0$. We can use t statistic.
- **Step 2:** Calculate the statistic. For example, in the t -test, calculate $t = \frac{\sqrt{m}(\bar{X} - \mu_0)}{s}$.
- **Step 3:** Evaluate the p -value from the value of the statistic. For example, in the t -test, we can evaluate p -value by referring to t -tables.
- **Step 4:**

Statistical test procedure

The standard procedure of the statistical test is the following.

- **Step 1:** Set the null hypothesis and alternative hypothesis. Also, fix the significance level α (usually 0.05 or 0.005) and determine which statistic to use. For example, if we are interested in the true expectation, the null hypothesis is $\mu = \mu_0$, where μ is the unknown true expectation and μ_0 is a value, which we decide. The alternative hypothesis is $\mu \neq \mu_0$. We can use t statistic.
- **Step 2:** Calculate the statistic. For example, in the t -test, calculate $t = \frac{\sqrt{m}(\bar{X} - \mu_0)}{s}$.
- **Step 3:** Evaluate the p -value from the value of the statistic. For example, in the t -test, we can evaluate p -value by referring to t -tables.
- **Step 4:** If $p < \alpha$, then we reject the null hypothesis and accept the alternative hypothesis. If $p \leq \alpha$, we can **neither reject nor accept a hypothesis**.

t-test example

Example

Our factory composed a component 8 times and the purity was (95, 93, 94, 94, 92, 93, 91, 96).

Suppose that the purity of a natural one is 92% on average.

Are our factory's products better than natural ones on average?

t-test example

Example

Our factory composed a component 8 times and the purity was (95, 93, 94, 94, 92, 93, 91, 96).

Suppose that the purity of a natural one is 92% on average.

Are our factory's products better than natural ones on average?

Step 1: Set the null hypothesis and alternative hypothesis. Also, fix the significance level α (usually 0.05 or 0.005).

The null hypothesis is $\mu = \mu_0 = 92$. The alternative hypothesis is $\mu \neq \mu_0 = 92$. Let's use the significance level $\alpha = 0.05$.

t-test example

Example

Our factory composed a component 8 times and the purity was (95, 93, 94, 94, 92, 93, 91, 96).

Suppose that the purity of a natural one is 92% on average.

Are our factory's products better than natural ones on average?

Step 2: Calculate the t -statistic.

The sample mean and standard deviation are $\bar{X} = 93.5$ and $s \approx 1.60$.

The t -statistic is $t = \frac{\sqrt{m}(\bar{X} - \mu_0)}{s} \approx \frac{93.5 - 92}{\frac{1.6}{2\sqrt{2}}} = 2.65$.

t-test example

Example

Our factory composed a component 8 times and the purity was (95, 93, 94, 94, 92, 93, 91, 96).

Suppose that the purity of a natural one is 92% on average.

Are our factory's products better than natural ones on average?

Step 3: Evaluate the p -value from the value of the t -statistic.

Here, under the null hypothesis, t follows the t -distribution with 7 degrees of freedom.

Then, if $t \approx 2.65$, the p -value is $p \approx 0.032$, according to an online calculator.

t-test example

Example

Our factory composed a component 8 times and the purity was (95, 93, 94, 94, 92, 93, 91, 96).

Suppose that the purity of a natural one is 92% on average.

Are our factory's products better than natural ones on average?

Step 4: Conclude from the p -value.

Since $p \approx 0.032 < \alpha = 0.05$, we reject the null hypothesis and accept the alternative hypothesis.

Hence, we can **statistically conclude that our factory produces better components than natural ones.**

Outline

5. Statistical Test

5.1 Introduction: why do we learn statistical tests?

5.2 The logic of statistical tests

5.3 Example test statistics

5.4 p-value

5.5 Failure of statistical test

5.6 Exercises

False positive (Type I error) and false negative (Type II error)

Statistical tests behave stochastically, so they may make a mistake. We may make two types of mistakes:

- **False positive (Type I error):** Accepts the alternative hypothesis \mathcal{H}_1 when the null hypothesis \mathcal{H}_0 is actually correct.
- **False negative (Type II error):** Fails to reject the null hypothesis \mathcal{H}_0 when the alternative hypothesis \mathcal{H}_1 is actually correct.

In the simple t -test case, the type I error probability equals to the significance level α .

Significance level, false-positive, false-negative

The false-positive rate, the possibility of accepting the alternative hypothesis when the data points are generated by a distribution in the null hypothesis, is determined by the significance level.

So, is it better to use a smaller significance level?

The answer is NO. It is because it increases the false-negative rate, the possibility of failing to accept the alternative hypothesis when the data points are generated by a distribution in the alternative hypothesis.

Outline

5. Statistical Test

5.1 Introduction: why do we learn statistical tests?

5.2 The logic of statistical tests

5.3 Example test statistics

5.4 p-value

5.5 Failure of statistical test

5.6 Exercises

Exercise (Statistical test 1)

Suppose that we apply a statistical test.

If the p-value of the test statistic is **lower** than the significance level, what of the following are the correct actions? Select all that apply.

- Accept the null hypothesis.
- Reject the null hypothesis.
- Accept the alternative hypothesis.
- Reject the alternative hypothesis.

Exercise (Statistical test 2)

Suppose that we apply a statistical test.

If the p-value of the test statistic is **higher** than the significance level, what of the following are the correct actions? Select all that apply.

- Accept the null hypothesis.
- Reject the null hypothesis.
- Accept the alternative hypothesis.
- Reject the alternative hypothesis.

Exercise (t-test 1)

Our factory composed a component 24 times and the purity was (95, 93, 94, 94, 92, 93, 91, 96, 93, 95, 96, 91, 92, 93, 94, 94, 91, 95, 96, 93, 94, 93, 94, 92).

Suppose that the purity of a natural one is 92% on average.

Are our factory's products better than natural ones on average?

Set the significance level $\alpha = 0.05$.

You can use the fact that the t -distribution can be approximated by the standard normal distribution if the degree of freedom is larger than 20.

Note that $\Pr(Z \geq 2) \approx 0.025$, where Z is a random variable following the standard normal distribution.

Exercise (t-test 2)

Our factory composed a component 24 times and the purity was (95, 93, 94, 94, 92, 93, 91, 96, 93, 95, 96, 91, 92, 93, 94, 94, 91, 95, 96, 93, 94, 93, 94, 92).

Suppose that the purity of a natural one is 93.25% on average.

Are our factory's products better than natural ones on average?

Set the significance level $\alpha = 0.05$.

You can use the fact that the t -distribution can be approximated by the standard normal distribution if the degree of freedom is larger than 20.

Note that $\Pr(Z \geq 2) \approx 0.025$, where Z is a random variable following the standard normal distribution.