

Question

There are many random phenomena (example?) in our real life. What is the language/mathematical structure that we use to depict them?

Outline

- sample space
- event
- probability measure
 - conditional probability
 - independence
- three theorems
 - multiplication law
 - law of total probability
 - Bayes' rule

Website of My Probability Course

<http://www.stat.nthu.edu.tw/~swcheng/Teaching/math2810/index.php>

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Definition (sample space, TBp. 2)

A sample space Ω is the set of all possible outcomes in a random phenomenon.

Example 1.1 (throw a coin 3 times, TBp. 35)

$$\Omega = \{hhh, hht, hth, thh, htt, tht, tth, ttt\}$$

Ω is a finite set

Example 1.2 (number of jobs in a print queue, Ex. B, TBp. 2)

$$\Omega = \{0, 1, 2, \dots\}$$

Ω is an infinite, but countable, set

Example 1.3 (length of time between successive earthquakes, Ex. C, TBp. 2)

$$\Omega = \{t | t \geq 0\}$$

Ω is an infinite, but uncountable, set

Question

What are the differences between the Ω in these examples?

Definition (event, TBp. 2)

A particular subset of Ω is called an event.

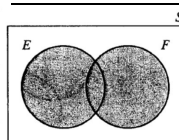
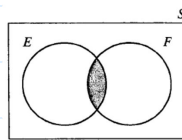
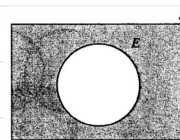
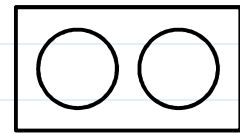
Example 1.4 (cont. Ex. 1.1)

Let A be the event that total number of heads equals 2, then $A = \{hht, hth, thh\}$.

Example 1.5 (cont. Ex. 1.2)

Let A be the event that fewer than 5 jobs in the print queue, then $A = \{0, 1, 2, 3, 4\}$.

- **union.** $C = A \cup B \Rightarrow C$: at least one of A and B occur.
- **intersection.** $C = A \cap B \Rightarrow C$: both A and B occur.
- **complement.** $C = A^c \Rightarrow C$: A does not occur.
- **disjoint.** $A \cap B = \emptyset \Rightarrow A$ and B have no outcomes in common.

(a) Shaded region: $E \cap F$.(b) Shaded region: EF .(c) Shaded region: E^c .

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Definition (probability measure, TBp. 4)

A probability measure on Ω is a function P from subsets of Ω to the real numbers that satisfies the following axioms:

1. $P(\Omega) = 1$.
2. If $A \subset \Omega$, then $P(A) \geq 0$.
3. If A_1 and A_2 are disjoint, then

$$P(A_1 \cup A_2) = P(A_1) + P(A_2).$$

More generally, if A_1, A_2, \dots are mutually disjoint, then

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i).$$

Example 1.6 (cont. Ex. 1.1)

Suppose the coin is fair. For every outcome $\omega \in \Omega$, $P(\omega) = \frac{1}{8}$.

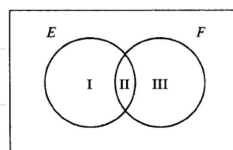
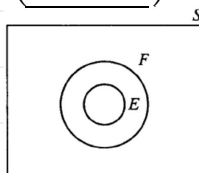
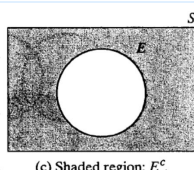
$$\Omega = \begin{matrix} hhh & hht & hth & thh & htt & tht & tth & ttt \\ 1/8 & 1/8 & 1/8 & 1/8 & 1/8 & 1/8 & 1/8 & 1/8 \end{matrix}$$

Property A. $P(A^C) = 1 - P(A)$.

Property B. $P(\emptyset) = 0$.

Property C. If $A \subset B$, then $P(A) \leq P(B)$.

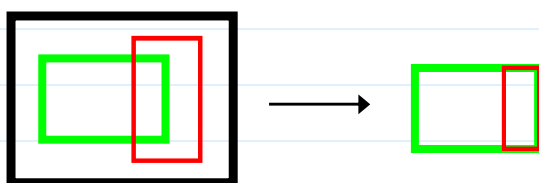
Property D. $P(A \cup B) = P(A) + P(B) - P(A \cap B)$.



Definition (conditional probability, TBp. 17)

Let A and B be two events with $P(B) > 0$. The conditional probability of A given B is defined to be

$$P(A|B) = \frac{P(A \cap B)}{P(B)}.$$



NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Example 1.7 (cont. Ex. 1.6)

Suppose that the first throw is h . What is the probability that we can get exact two h 's in the three trials?

$$\Omega = \{hhh, hht, hth, thh, htt, tht, tth, ttt\}$$

$$B = \{hhh, hht, hth, htt\}$$

$$A = \{hht, hth, thh\}$$



Theorem (Multiplication Law, TBp. 17)

Let A and B be events and assume $P(B) > 0$. Then

$$P(A \cap B) = \underline{P(A|B)}P(B).$$

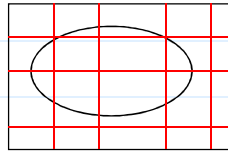
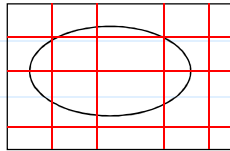
Example 1.7 (Ex. B, TBp. 18)

Suppose if it is cloudy (B), the probability that it is raining (A) is 0.3, and that the probability that it is cloudy is $P(B) = \underline{0.2}$. The probability that it is cloudy and raining is $P(A \cap B) = \underline{P(A|B)}P(B) = 0.3 \times 0.2 = 0.06$.

Theorem (Law of Total Probability, TBp. 18)

Let B_1, B_2, \dots, B_n be such that $\bigcup_{i=1}^n B_i = \Omega$ and $B_i \cap B_j = \emptyset$ for $i \neq j$, with $P(B_i) > 0$ for all i . Then, for any event A ,

$$P(A) = \sum_{i=1}^n P(A|B_i)P(B_i).$$

**Theorem (Bayes' Rule, TBp. 20)**

Let A and B_1, \dots, B_n be events where the B_i are disjoint, $\bigcup_{i=1}^n B_i = \Omega$ and $P(B_i) > 0$ for all i . Then

$$P(B_j|A) = \frac{P(A|B_j)P(B_j)}{\sum_{i=1}^n P(A|B_i)P(B_i)}.$$

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Definition (independence, TBp. 24)

Two events A and B are said to be **independent** if

$$P(A \cap B) = P(A)P(B).$$

A collection of events A_1, A_2, \dots, A_n are said to be **mutually independent** if for any subcollection, A_{i_1}, \dots, A_{i_m} ,

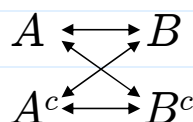
$$P(A_{i_1} \cap \dots \cap A_{i_m}) = P(A_{i_1}) \dots P(A_{i_m}).$$

When A and B are independent,

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)P(B)}{P(B)} = P(A),$$

and $P(A^c|B) = P(A^c)$.

Furthermore, $P(A|B^c) = P(A)$ and $P(A^c|B^c) = P(A^c)$.



❖ **Reading:** textbook, Sections 1.1, 1.2, 1.3, 1.5, 1.6, 1.7

❖ **Further Reading:** Roussas, Chapters 1 and 2

Chapters 2 and 3

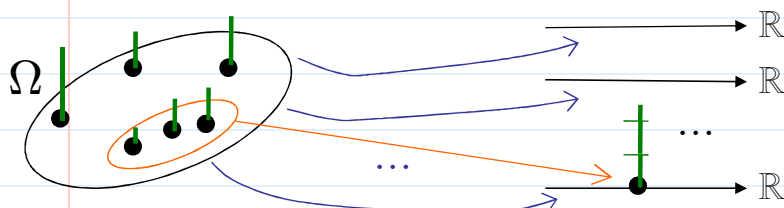
Outline

- random variables
- distribution
 - discrete and continuous
 - univariate and multivariate
 - cdf, pmf, pdf
- conditional distribution
- independent random variables
- function of random variables
 - distribution of transformed r.v.
 - extrema and order statistics

• random variable

Definition 2.1 (random variable, TBp. 33)

A random variable is a function from Ω to the real numbers.



NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Ch1-6, p.2-10

Example 2.1 (cont. Ex. 1.1)

- (1) X_1 = the total number of heads
- (2) X_2 = the number of heads on the first toss
- (3) X_3 = the number of heads minus the number of tails

	1/8	1/8	1/8	1/8	1/8	1/8	1/8	1/8
$\Omega =$	hhh	hht	hth	thh	htt	tht	tth	ttt
	↓	↓	↓	↓	↓	↓	↓	↓
$X_1 :$	3	2	2	2	1	1	1	0
$X_2 :$	1	1	1	0	1	0	0	0
$X_3 :$	3	1	1	1	-1	-1	-1	-3

Question 2.1

Why statisticians need random variables? Why they map to real line?

• distribution

Question 2.2

A random variable have a sample space on real line. Does it bring some special ways to characterize its probability measure?

	discrete	continuous
uni-variate r.v.	<ul style="list-style-type: none"> • pmf • cdf • mgf/chf 	<ul style="list-style-type: none"> • pdf • cdf • mgf/chf
multi-variate r.v.'s	<ul style="list-style-type: none"> • joint pmf • joint cdf • joint mgf/chf 	<ul style="list-style-type: none"> • joint pdf • joint cdf • joint mgf/chf

pmf: probability mass function, pdf: probability density function,
cdf: cumulative distribution function

mgf (moment generating function) and chf (characteristic function) will be defined in Chapter 4

NTHU MATH 2820, 2026, Lecture Notes
 made by S.-W. Cheng (NTHU, Taiwan)

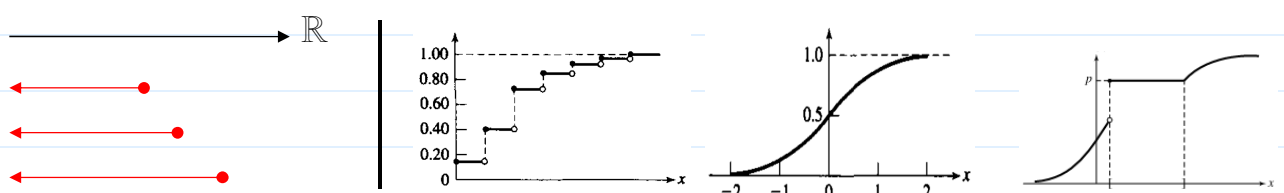
Definition 2.2 (discrete and continuous random variables, TBp. 35 and 47)

A discrete random variable can take on only a finite or at most a countably infinite number of values. A continuous random variable can take on a continuum of values.

Definition 2.3 (cumulative distribution function, TBp. 36)

A function \underline{F} is called the cumulative distribution function (cdf) of a random variable \underline{X} if

$$\underline{F}(x) = \underline{P}(\underline{X} \leq x), x \in \mathbb{R}.$$



Definition 2.4 (probability mass function/frequency function, TBp. 36)

A function $p(x)$ is called a **probability mass function (pmf)** or a **frequency function** if and only if (1) $p(x) \geq 0$ for all $x \in \mathcal{X}$, and (2) $\sum_{x \in \mathcal{X}} p(x) = 1$.

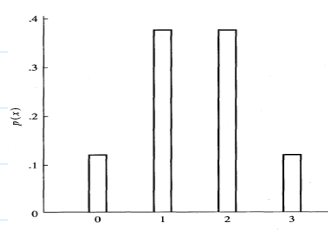
For a discrete random variable X with pmf $p(x)$,

$$P(X = x) = p(x),$$

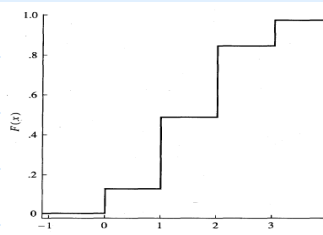
and

$$P(X \in A) = \sum_{x \in A} p(x).$$

probability mass function



cumulative distribution function



$$\bullet \quad F(x) = \sum_{t \leq x} P(X = t) = \sum_{t \leq x} p(t)$$

$$p(x) = P(X = x) = F(x) - F(x-)$$

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

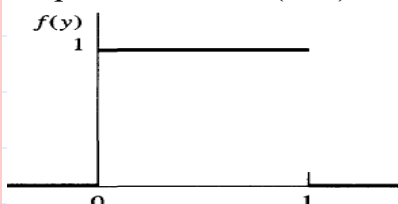
Definition 2.5 (probability density function, TBp. 46)

A function $f(x)$ is a **probability density function (pdf)** or **density function** if and only if (1) $f(x) \geq 0$ for all x , and (2) $\int_{-\infty}^{\infty} f(x) dx = 1$.

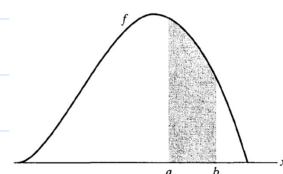
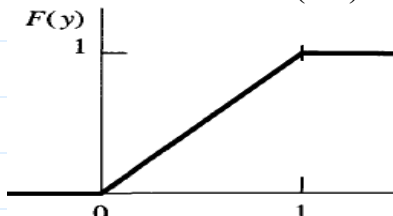
For a continuous random variable X with pdf f ,

$$P(X \in A) = \int_A f(x) dx.$$

pdf of Uniform(0, 1)



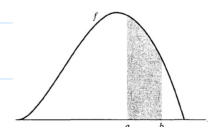
cdf of Uniform(0,1)



$$\bullet \quad F(x) = \int_{-\infty}^x f(t) dt$$

$$f(x) = \frac{d}{dx} F(x)$$

(**Note.** x st $f(x) > 0$, $P(X = x) = \int_x^x f(t) dt = 0$)

**Question 2.3**

How to interpret $f(x)$?

For small dx , $P\left(x - \frac{dx}{2} \leq X \leq x + \frac{dx}{2}\right) = \int_{x-\frac{dx}{2}}^{x+\frac{dx}{2}} f(t) dt \approx f(x) dx$

Theorem 2.1 (properties of cdf)

If $F(x)$ is a cumulative distribution function of some random variable X then the following properties hold.

1. $0 \leq F(x) \leq 1$
2. $F(x)$ is nondecreasing.
3. For any $x \in \mathbb{R}$, $F(x)$ is continuous from the right; i.e.

$$\lim_{t \downarrow x} F(t) = F(x).$$

4. $\lim_{x \rightarrow \infty} F(x) = 1$ and $\lim_{x \rightarrow -\infty} F(x) = 0$.
5. $P(X > x) = 1 - F(x)$ and $P(a < X \leq b) = F(b) - F(a)$.
6. For any $x \in \mathbb{R}$, $F(x)$ has left limit.
7. There are at most countably many discontinuity points of $F(x)$.

Conversely, if a function $F(x)$ satisfies properties 2, 3, 4 then $F(x)$ is a cdf.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Question 2.4 Why need joint distribution for the study of multivariate r.v.'s?

Example 2.2 (cont. Ex. 2.1)

$$\Omega = \{hhhh, hhht, hthh, thhh, hhtt, thht, tthh, tttt\}$$

X_2 : # of head on 1 st toss	X_1 : total # of heads			
	0(1/8)	1(3/8)	2(3/8)	3(1/8)
(1/2) 0	$\frac{1}{8}\left(\frac{1}{16}\right)$	$\frac{2}{8}\left(\frac{3}{16}\right)$	$\frac{1}{8}\left(\frac{3}{16}\right)$	$0\left(\frac{1}{16}\right)$
(1/2) 1	$0\left(\frac{1}{16}\right)$	$\frac{1}{8}\left(\frac{3}{16}\right)$	$\frac{2}{8}\left(\frac{3}{16}\right)$	$\frac{1}{8}\left(\frac{1}{16}\right)$

Question 2.5

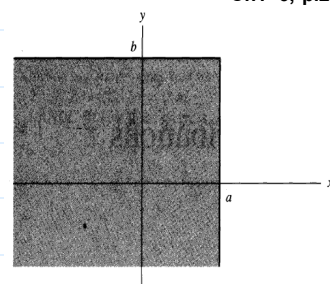
When we know the joint distribution, we can obtain every marginal distributions. Is the reverse statement true?

Definition 2.6 (joint cumulative distribution function, TBp. 71)

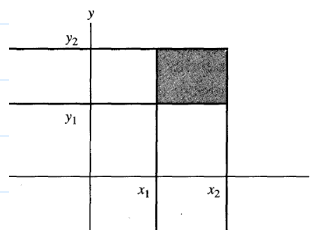
The joint cdf of X_1, X_2, \dots, X_n is

$$F(x_1, x_2, \dots, x_n) = P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n)$$

for $x_1, x_2, \dots, x_n \in \mathbb{R}$.

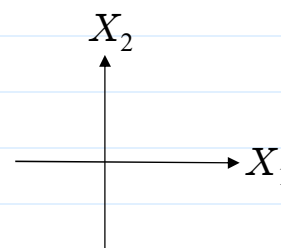


$$\begin{aligned} P(x_1 < X \leq x_2, y_1 < Y \leq y_2) \\ &= \frac{F(x_2, y_2) - F(x_2, y_1) \\ &\quad - F(x_1, y_2) + F(x_1, y_1)}{1} \end{aligned}$$

**Definition 2.7** (marginal cdf, TBp. 76)

The marginal cdf of X_1 is

$$F_{X_1}(x_1) = P(X_1 \leq x_1) = \lim_{x_2, x_3, \dots, x_n \rightarrow \infty} F(x_1, x_2, \dots, x_n)$$



- discrete case: marginal pmf $p_{X_1}(x) = F_{X_1}(x) - F_{X_1}(x-)$.
- continuous case: marginal pdf $f_{X_1}(x) = \frac{d}{dx} F_{X_1}(x)$.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

• discrete multivariate case

$$\begin{aligned} p(x_1, x_2, \dots, x_n) &= P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) \\ &\Rightarrow \text{joint pmf of } X_1, X_2, \dots, X_n \end{aligned}$$

$$P((X_1, \dots, X_n) \in A) = \sum_{(x_1, \dots, x_n) \in A} p(x_1, \dots, x_n)$$

$$F(x_1, x_2, \dots, x_n) = \sum_{t_1 \leq x_1, t_2 \leq x_2, \dots, t_n \leq x_n} p(t_1, t_2, \dots, t_n)$$

$$p_{X_1}(x_1) = P(X_1 = x_1) = \sum_{-\infty < t_2 < \infty, \dots, -\infty < t_n < \infty} p(x_1, t_2, \dots, t_n)$$

• continuous multivariate case

$$f(x_1, x_2, \dots, x_n) = \frac{\partial^n}{\partial x_1 \cdots \partial x_n} F(x_1, x_2, \dots, x_n)$$

$$\Rightarrow \text{joint pdf of } X_1, X_2, \dots, X_n$$

$$P((X_1, \dots, X_n) \in A) = \int \cdots \int_A f(x_1, \dots, x_n) dx_1 \cdots dx_n$$

$$F(x_1, x_2, \dots, x_n) = \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_n} f(t_1, t_2, \dots, t_n) dt_n \cdots dt_1$$

$$f_{X_1}(x_1) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(x_1, t_2, \dots, t_n) dt_2 \cdots dt_n$$

- independent random variables

Definition 2.8 (independent random variables, TBp. 84)

Random variables X_1, X_2, \dots, X_n are said to be independent if their joint cdf factors into the product of their marginal cdf's

$$F(x_1, x_2, \dots, x_n) = F_{X_1}(x_1)F_{X_2}(x_2) \cdots F_{X_n}(x_n)$$

for all x_1, x_2, \dots, x_n .

Theorem 2.2 (TBp. 85-86)

1. For continuous case,

$$F(x_1, \dots, x_n) = F_{X_1}(x_1) \cdots F_{X_n}(x_n) \Leftrightarrow f(x_1, \dots, x_n) = f_{X_1}(x_1) \cdots f_{X_n}(x_n)$$

For discrete case,

$$F(x_1, \dots, x_n) = F_{X_1}(x_1) \cdots F_{X_n}(x_n) \Leftrightarrow p(x_1, \dots, x_n) = p_{X_1}(x_1) \cdots p_{X_n}(x_n)$$

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

2. X, Y independent

$$\Leftrightarrow P(X \in A, Y \in B) = P(X \in A)P(Y \in B)$$

i.e. the events $\{X \in A\}$ and $\{Y \in B\}$ are independent

3. X, Y independent \Rightarrow $Z = g(X)$ and $W = h(Y)$ are independent

generalization

X_1, \dots, X_n are independent

$$1 < i_0 < i_1 < \cdots < i_k = n$$

$$Y_1 = g_1(X_1, \dots, X_{i_1}),$$

$$Y_2 = g_2(X_{i_1+1}, \dots, X_{i_2}),$$

\dots

$$Y_k = g_k(X_{i_{k-1}+1}, \dots, X_{i_k}).$$

Y_1, \dots, Y_k are independent

4. marginal distributions of X_1, X_2, \dots, X_n + independence \Rightarrow joint distribution of X_1, X_2, \dots, X_n

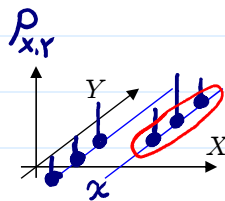
• conditional distribution

Definition 2.9 (conditional pmf for discrete case, TBp. 87)

X and Y are discrete random variables with joint pmf $p_{XY}(x, y)$, the conditional pmf of Y given X is

$$\begin{aligned} p_{Y|X}(y|x) &\equiv P(Y = y | X = x) = \frac{P(X = x, Y = y)}{P(X = x)} \\ &= \frac{p_{XY}(x, y)}{p_X(x)} \end{aligned}$$

if $p_X(x) > 0$. The probability is defined to be zero if $p_X(x) = 0$.



Example 2.3 (cont. Ex 2.2)

$$p_{X_2|X_1}(0|1) = 2/3, \text{ and } p_{X_2|X_1}(1|1) = 1/3$$

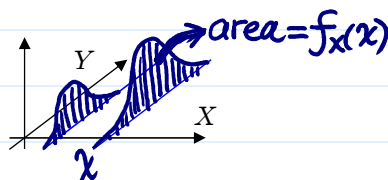
NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Definition 2.10 (conditional pdf for continuous case, TBp. 86)

X and Y are continuous random variables with joint pdf $f_{XY}(x, y)$, the conditional pdf of Y given X is defined by

$$f_{Y|X}(y|x) = \frac{f_{XY}(x, y)}{f_X(x)}, \quad y \in R,$$

if $0 < f_X(x) < \infty$ and 0 otherwise.



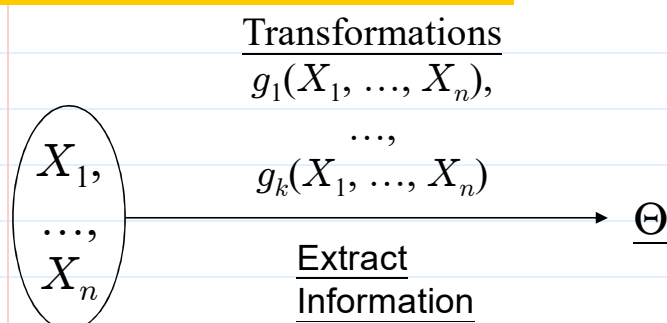
Theorem 2.3

1. The definition of $f_{Y|X}(y|x)$ comes from

$$\begin{aligned} P(a \leq Y \leq b | x - \Delta x/2 \leq X \leq x + \Delta x/2) &= \frac{\int_a^b \int_{x-\Delta x/2}^{x+\Delta x/2} f_{XY}(u, v) du dv}{\int_{x-\Delta x/2}^{x+\Delta x/2} f_X(t) dt} \\ &\approx \frac{\int_a^b f_{XY}(x, y) \Delta x dy}{f_X(x) \Delta x} = \int_a^b \frac{f_{XY}(x, y)}{f_X(x)} dy \end{aligned}$$

2. For each fixed x , $\underline{p_{Y|X}(y|x)}$ is a pmf for y and $\underline{f_{Y|X}(y|x)}$ is a pdf for y .
3. $\underline{p_{XY}(x, y)} = \underline{p_{Y|X}(y|x)} \underline{p_X(x)}$, and $\underline{f_{XY}(x, y)} = \underline{f_{Y|X}(y|x)} \underline{f_X(x)}$
— multiplication law
4. $\underline{p_Y(y)} = \sum_x \underline{p_{Y|X}(y|x)} \underline{p_X(x)}$, and $\underline{f_Y(y)} = \int_{-\infty}^{\infty} \underline{f_{Y|X}(y|x)} \underline{f_X(x)} dx$
— law of total probability
5. $\underline{p_{X|Y}(x|y)} = \frac{\underline{p_{Y|X}(y|x)} \underline{p_X(x)}}{\sum_x \underline{p_{Y|X}(y|x)} \underline{p_X(x)}}$, and $\underline{f_{X|Y}(x|y)} = \frac{\underline{f_{Y|X}(y|x)} \underline{f_X(x)}}{\int_{-\infty}^{\infty} \underline{f_{Y|X}(y|x)} \underline{f_X(x)} dx}$,
— Bayes' rule
6. X, Y are independent $\Leftrightarrow \underline{p_{Y|X}(y|x)} = \underline{p_Y(y)}$ or $\underline{f_{Y|X}(y|x)} = \underline{f_Y(y)}$

• functions of random variables



Question 2.6

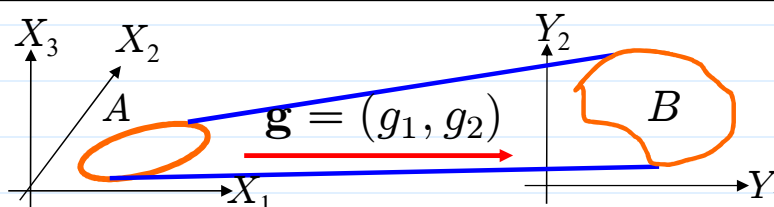
For given r.v.'s X_1, \dots, X_n , how to derive the distributions of their transformations?

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

1. method of events

Theorem 2.7

Let $\underline{\mathbf{X}} = (X_1, X_2, \dots, X_n)$ be random variables, and $\underline{\mathbf{Y}} = \underline{\mathbf{g}(\mathbf{X})}$. Then, the distribution of \mathbf{Y} is determined by the distribution of \mathbf{X} as follow: for any event B defined by $\underline{\mathbf{Y}}$, $P(\underline{\mathbf{Y}} \in B) = P(\underline{\mathbf{X}} \in \underline{A})$, where $\underline{A} = \underline{\mathbf{g}^{-1}(B)}$.



Example 2.4 (univariate discrete random variable)

Let \underline{X} be a discrete r.v. taking the values $\underline{x_i}, i = 1, 2, \dots$, and $\underline{Y} = \underline{g(X)}$. Then, \underline{Y} is also a discrete r.v. taking the values $\underline{y_j}, j = 1, 2, \dots$. To determine the pmf of Y , by taking $\underline{B} = \{y_j\}$, we have

$$\underline{A} = \{x_i : g(x_i) = y_j\} \text{ and hence}$$

$$\underline{p_Y(y_j)} = P(\{y_j\}) = P(A) = \sum_{x_i \in A} \underline{p_X(x_i)}.$$

Example 2.5 (sum of two discrete random variables, TBp. 96)

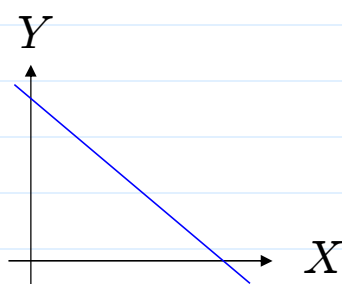
X and Y are random variables with joint pmf $p(x, y)$. Find the distribution of $Z = X + Y$.

(Exercise: difference of two random variables, $Z = X - Y$)

$$p_Z(z) = P(Z = z) = P(X + Y = z) = \sum_{x=-\infty}^{\infty} p(x, z - x)$$

When X, Y independent, $p(x, y) = p_X(x)p_Y(y)$,

$$p_Z(z) = \sum_{x=-\infty}^{\infty} \underline{p_X(x)p_Y(z-x)} \Rightarrow \underline{\text{convolution of } p_X \text{ and } p_Y}$$



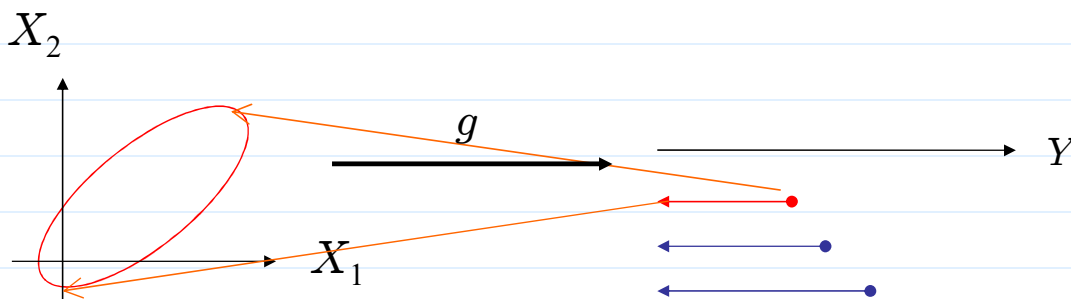
NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

2. method of cumulative distribution function (a special case of method 1)

Let Y be a function of the random variables X_1, X_2, \dots, X_n .

1. Find the region $Y \leq y$ in the (x_1, x_2, \dots, x_n) space.
2. Find $F_Y(y) = P(Y \leq y)$ by summing the joint pmf or integrating the joint pdf of X_1, X_2, \dots, X_n over the region $Y \leq y$.
3. (for continuous case) Find the pdf of Y by differentiating $F_Y(y)$, i.e., $f_Y(y) = \frac{d}{dy} F_Y(y)$.

Note. It can be generalized to multivariate $\mathbf{Y} = (Y_1, Y_2, \dots, Y_m)$.



Example 2.6 (square of a random variable, similar example see TBp. 61)

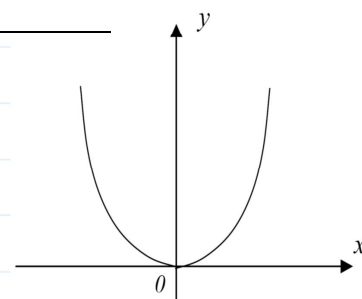
X is a random variables with pdf $f_X(x)$ and cdf $F_X(x)$. Find the distributon of $Y = X^2$.

For $y \geq 0$, $\{Y \leq y\} = \{-\sqrt{y} \leq X \leq \sqrt{y}\}$

$$F_Y(y) = P(Y \leq y) = P(-\sqrt{y} \leq X \leq \sqrt{y}) = F_X(\sqrt{y}) - F_X(-\sqrt{y})$$

$$\begin{aligned} f_Y(y) &= \frac{d}{dy} F_Y(y) = \frac{d}{dy} F_X(\sqrt{y}) - \frac{d}{dy} F_X(-\sqrt{y}) \\ &= f_X(\sqrt{y}) \frac{1}{2\sqrt{y}} - f_X(-\sqrt{y}) \left(-\frac{1}{2\sqrt{y}}\right) \\ &= \frac{1}{2\sqrt{y}} (f_X(\sqrt{y}) + f_X(-\sqrt{y})) \end{aligned}$$

and $f_Y(y) = 0$ for $y < 0$.



NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

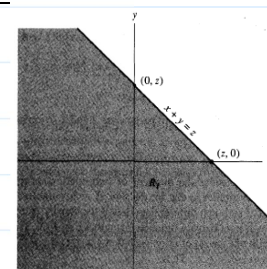
Example 2.7 (sum of two continuous random variables, TBp. 97)

X and Y are random variables with joint pdf $f(x, y)$. Find the distribution of $Z = X + Y$.

(Exercise: difference of two random variables, $Z = X - Y$)

Let R_z be $\{(x, y) : x + y \leq z\}$. Then,

$$\begin{aligned} F_Z(z) &= P(Z \leq z) = P(X + Y \leq z) = \iint_{R_z} f(x, y) dx dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{z-x} f(x, y) dy dx \\ &= \int_{-\infty}^z \int_{-\infty}^{\infty} f(x, v-x) dx dv \quad (\text{set } y = v-x) \\ f_Z(z) &= \frac{d}{dz} F_Z(z) = \int_{-\infty}^{\infty} f(x, z-x) dx \end{aligned}$$



When X, Y independent, $f(x, y) = f_X(x)f_Y(y)$,

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx \Rightarrow \text{convolution of } f_X \text{ and } f_Y$$

Example 2.8 (quotient of two continuous random variables, TBp. 98)

X and Y are r.v. with joint pdf $f(x, y)$. Find the distribution of $Z = Y/X$. (Exercise: product of two random variables, $Z=XY$)

$$Q_z = \{(x, y) : y/x \leq z\} = \{(x, y) : x < 0, y \geq zx\} \cup \{(x, y) : x > 0, y \leq zx\}$$

$$F_Z(z) = \iint_{Q_z} f(x, y) dx dy = \int_{-\infty}^0 \int_{xz}^{\infty} f(x, y) dy dx + \int_0^{\infty} \int_{-\infty}^{xz} f(x, y) dy dx$$

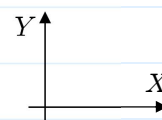
$$= \int_{-\infty}^0 \int_z^{-\infty} x f(x, xv) dv dx + \int_0^{\infty} \int_{-\infty}^z x f(x, xv) dv dx \quad (\text{set } y = xv)$$

$$= \int_{-\infty}^0 \int_{-\infty}^z (-x) f(x, xv) dv dx + \int_0^{\infty} \int_{-\infty}^z x f(x, xv) dv dx$$

$$= \int_{-\infty}^z \int_{-\infty}^{\infty} |x| f(x, xv) dx dv$$

$$f_Z(z) = \frac{d}{dz} F_Z(z) = \int_{-\infty}^{\infty} |x| f(x, xz) dx$$

$$\left(= \int_{-\infty}^{\infty} |x| f_X(x) f_Y(xz) dx \quad \text{when } X, Y \text{ independent} \right)$$



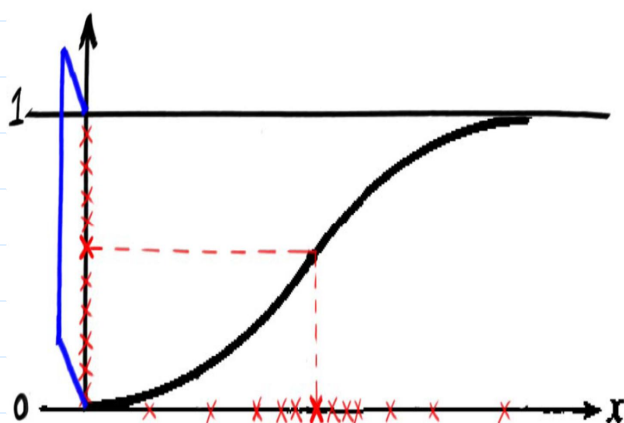
NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Theorem 2.4 (TBp. 63)

Let X be a random variable whose cdf F possesses a unique inverse F^{-1} . Let $Z = F(X)$, then Z has a uniform distribution on $[0, 1]$.

Theorem 2.5 (TBp. 63)

Let U be a uniform random variable on $[0, 1]$ and F is a cdf which possesses a unique inverse F^{-1} . Let $X = F^{-1}(U)$. Then the cdf of X is F .



Note. The 2 theorems are useful for generating pseudo-random numbers in computer simulation (the concepts can be generalized to any r.v.'s).

3. method of probability density function (for continuous r.v.'s and differentiable, one-to-one transformations, a special case of method 2) :

Theorem 2.6 (univariate continuous case, TBp. 62)

Let \underline{X} be a continuous random variable with pdf $f_X(x)$. Let $\underline{Y} = g(\underline{X})$, where g is differentiable, strictly monotone. Then,

$$\underline{f_Y(y)} = \underline{f_X(g^{-1}(y))} \left| \frac{dg^{-1}(y)}{dy} \right|$$

for y s.t. $y = g(x)$ for some x , and $f_Y(y) = 0$ otherwise.

Example 2.9

\underline{X} is a random variables with pdf $f_X(x)$. Find the distributon of $\underline{Y} = 1/\underline{X}$.

For $x > 0$ (or $x < 0$),

$$y = 1/x \equiv g(x) \Rightarrow x = g^{-1}(y) = 1/y$$

$$\underline{dg^{-1}/dy} = -1/y^2 \quad \text{and} \quad \left| dg^{-1}/dy \right| = 1/y^2$$

hence

$$\underline{f_Y(y)} = \underline{f_X(1/y)}(1/y^2)$$

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Theorem 2.7 (multivariate continuous case, TBp. 102-103)

$\underline{\mathbf{X}} = (X_1, X_2, \dots, X_n)$ multivariate continuous, $\underline{\mathbf{Y}} = (Y_1, Y_2, \dots, Y_n) \equiv \underline{\mathbf{g}}(\underline{\mathbf{X}})$. $\underline{\mathbf{g}}$ is one-to-one, so that its inverse exists and is denoted by

$$\underline{\mathbf{x}} = \underline{\mathbf{g}}^{-1}(\underline{\mathbf{y}}) = \underline{\mathbf{w}}(\underline{\mathbf{y}}) = (\underline{w_1}(\underline{\mathbf{y}}), \underline{w_2}(\underline{\mathbf{y}}), \dots, \underline{w_n}(\underline{\mathbf{y}})).$$

Assume $\underline{\mathbf{w}}$ have continuous partial derivatives, and let

$$\underline{J} = \begin{vmatrix} \frac{\partial w_1(\underline{\mathbf{y}})}{\partial y_1} & \frac{\partial w_1(\underline{\mathbf{y}})}{\partial y_2} & \dots & \frac{\partial w_1(\underline{\mathbf{y}})}{\partial y_n} \\ \frac{\partial w_2(\underline{\mathbf{y}})}{\partial y_1} & \frac{\partial w_2(\underline{\mathbf{y}})}{\partial y_2} & \dots & \frac{\partial w_2(\underline{\mathbf{y}})}{\partial y_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial w_n(\underline{\mathbf{y}})}{\partial y_1} & \frac{\partial w_n(\underline{\mathbf{y}})}{\partial y_2} & \dots & \frac{\partial w_n(\underline{\mathbf{y}})}{\partial y_n} \end{vmatrix}$$

Then

$$\underline{f_Y(y)} = \underline{f_X(g^{-1}(y))} |\underline{J}|.$$

for $\underline{\mathbf{y}}$ s.t. $\underline{\mathbf{y}} = \underline{\mathbf{g}}(\underline{\mathbf{x}})$ for some $\underline{\mathbf{x}}$, and $f_Y(\underline{\mathbf{y}}) = 0$, otherwise.

Note. When the dimensionality of $\underline{\mathbf{Y}}$, denoted by k , is less than n , we can choose another $n - k$ transformations $\underline{\mathbf{Z}}$ such that $(\underline{\mathbf{Y}}, \underline{\mathbf{Z}})$ satisfy the above assumptions. By integrating out the last $n - k$ arguments in the pdf of $(\underline{\mathbf{Y}}, \underline{\mathbf{Z}})$, the pdf of $\underline{\mathbf{Y}}$ can be obtained.

Example 2.10 (cont. Ex 2.8)

X_1 and X_2 are random variables with joint pdf $f_{X_1X_2}(x_1, x_2)$. Find the distribution of $Y_1 = X_2/X_1$. (Exercise: $Y_1 = X_1X_2$)

Let $Y_2 = X_1$. Then

$$x_1 = y_2 \equiv w_1(y_1, y_2)$$

$$x_2 = y_1 y_2 \equiv w_2(y_1, y_2).$$

$$\frac{\partial w_1}{\partial y_1} = 0, \quad \frac{\partial w_1}{\partial y_2} = 1, \quad \frac{\partial w_2}{\partial y_1} = y_2, \quad \frac{\partial w_2}{\partial y_2} = y_1.$$

$$J = \begin{vmatrix} 0 & 1 \\ y_2 & y_1 \end{vmatrix} = -y_2, \quad \text{and} \quad |J| = |y_2|$$

Therefore,

$$f_{Y_1Y_2}(y_1, y_2) = f_{X_1X_2}(y_2, y_1 y_2) |y_2|$$

$$f_{Y_1}(y_1) = \int_{-\infty}^{\infty} f_{Y_1Y_2}(y_1, y_2) dy_2 = \int_{-\infty}^{\infty} f_{X_1X_2}(y_2, y_1 y_2) |y_2| dy_2$$

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

4. method of moment generating function: based on the uniqueness theorem of moment generating function. To be explained later in Chapter 4.

Ch1-6, p.2-34

• extrema and order statistics

Definition 2.11 (order statistics, sec 3.7)

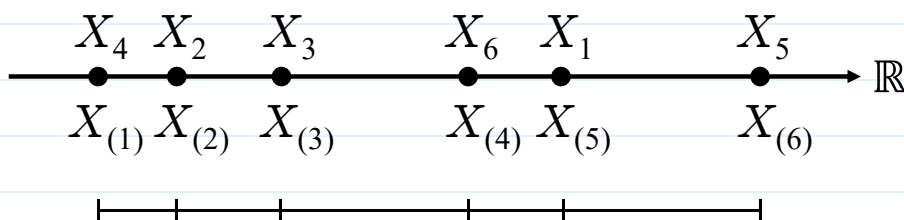
Let X_1, X_2, \dots, X_n be random variables. We sort the X_i 's and denote by $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ the **order statistics**. Using the notation,

$$X_{(1)} = \min(X_1, X_2, \dots, X_n) \quad \text{is the } \underline{\text{minimum}}$$

$$X_{(n)} = \max(X_1, X_2, \dots, X_n) \quad \text{is the } \underline{\text{maximum}}$$

$$R \equiv X_{(n)} - X_{(1)} \quad \text{is called } \underline{\text{range}}$$

$$S_j \equiv X_{(j)} - X_{(j-1)}, j = 2, \dots, n \quad \text{are called } \underline{j\text{th spacings}}$$



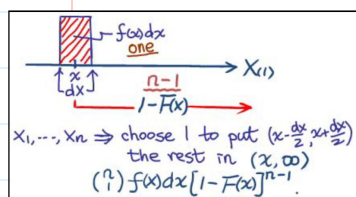
Note. In the section, we only consider the case that X_1, X_2, \dots, X_n are i.i.d continuous r.v.'s with cdf F and pdf f . Although X_1, X_2, \dots, X_n are independent, their order statistics are not independent in general.

Definition 2.12 (i.i.d.)

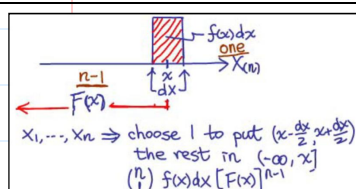
X_1, X_2, \dots, X_n are **i.i.d.** (i)ndependent, (i)dentically (d)istributed with cdf F /pmf p /pdf $f \Rightarrow X_1, X_2, \dots, X_n$ are independent and have a common marginal cdf F /pmf p /pdf f .

Theorem 2.8 (TBp. 104)

The cdf of $X_{(1)}$ is $1 - [1 - F(x)]^n$ and its pdf is $nf(x)[1 - F(x)]^{n-1}$.
The cdf of $X_{(n)}$ is $[F(x)]^n$ and its pdf is $nf(x)[F(x)]^{n-1}$.



$$\begin{aligned} F_{X_{(n)}}(x) &= P(X_{(n)} \leq x) = P(X_1 \leq x, \dots, X_n \leq x) \\ &= P(X_1 \leq x) \cdots P(X_n \leq x) \\ &= [F(x)]^n. \end{aligned}$$



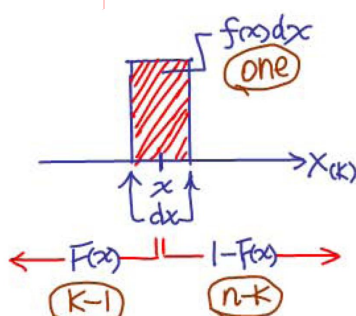
$$\begin{aligned} 1 - F_{X_{(1)}}(x) &= P(X_{(1)} > x) = P(X_1 > x, \dots, X_n > x) \\ &= P(X_1 > x) \cdots P(X_n > x) \\ &= [1 - F(x)]^n. \end{aligned}$$

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

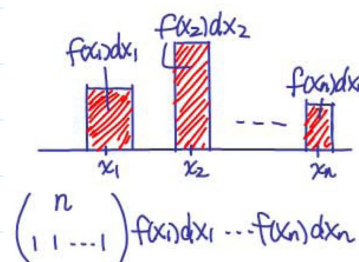
Theorem 2.9 (TBp. 105)

The pdf of the k th order statistic $X_{(k)}$ is

$$f_{X_{(k)}}(x) = \frac{n!}{(k-1)!(n-k)!} f(x) [F(x)]^{k-1} [1 - F(x)]^{n-k}.$$



$$\begin{aligned} X_1, \dots, X_n \Rightarrow \text{choose 1 to place in } (x - \frac{dx}{2}, x + \frac{dx}{2}) \\ = k-1 = (-\infty, x) \\ = n-k = (x, \infty) \\ \binom{n}{k-1, n-k} f(x)dx [F(x)]^{k-1} [1-F(x)]^{n-k} \end{aligned}$$

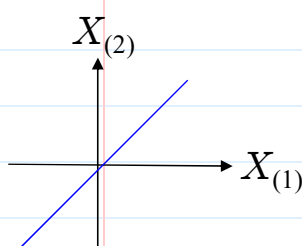


Theorem 2.10 (TBp. 114, Problem 73)

The joint pdf of $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ is

$$f_{X_{(1)}X_{(2)}\dots X_{(n)}}(x_1, x_2, \dots, x_n) = n! f(x_1) f(x_2) \cdots f(x_n),$$

for $x_1 \leq x_2 \leq \dots \leq x_n$, and $f_{X_{(1)}X_{(2)}\dots X_{(n)}} = 0$ otherwise.



Question: Are $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ independent, judged from the form of its joint pdf?

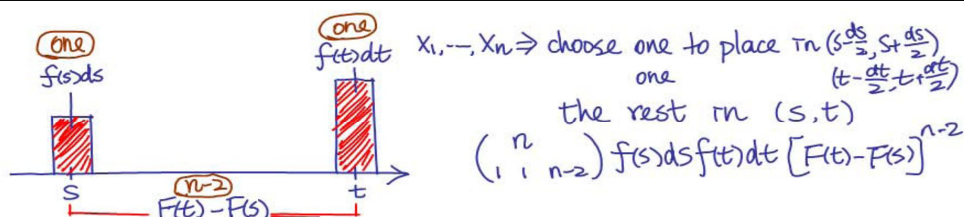
Example 2.11 (range, TBp. 105-106)

The joint pdf of $X_{(1)}$ and $X_{(n)}$ is

$$f_{X_{(1)}X_{(n)}}(s, t) = n(n-1)f(s)f(t)[F(t) - F(s)]^{n-2}, \quad \text{for } s \leq t,$$

and 0 otherwise. Therefore, the pdf of $R = X_{(n)} - X_{(1)}$ is

$$f_R(r) = \int_{-\infty}^{\infty} f_{X_{(1)}X_{(n)}}(s, s+r)ds \quad \text{for } r > 0, \text{ and } f_R(r) = 0, \text{ otherwise.}$$

**Exercise**

1. Find the joint pdf of $X_{(i)}$ and $X_{(j)}$, where $i < j$.
2. Find the joint pdf of $X_{(j)}$ and $X_{(j-1)}$, and derive the pdf of j th spacing $S_j = X_{(j)} - X_{(j-1)}$.

❖ **Reading:** textbook, 2.1 (not including 2.1.1~5), 2.2 (not including 2.2.1~4), 2.3, 2.4, Chapter 3

❖ **Further Reading:** Roussas, 3.1, 4.1, 4.2, 7.1, 7.2, 9.1, 9.2, 9.3, 9.4, 10.1

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

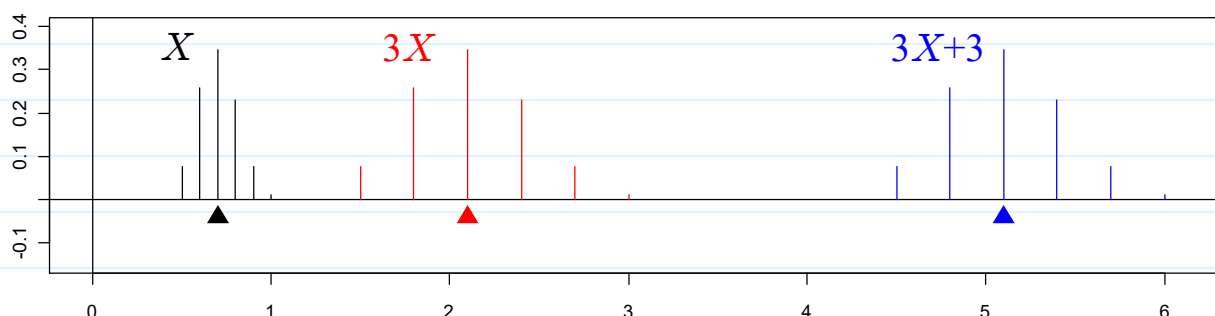
Chapter 4

Outline

- expectation
 - mean, variance, standard deviation, covariance, correlation coefficient
- moment generating function & characteristic function
- conditional expectation and prediction
- δ method

Question 3.1

Can we describe the characteristics of distributions by use of some intuitive and meaningful simple values?



- expectation

Definition 3.1 (expectation, TBp. 122, 123)

For random variables X_1, \dots, X_n , the **expectation** of a univariate random variable $\underline{Y} = g(X_1, \dots, X_n)$ is defined as

$$\begin{aligned}\underline{E(Y)} &\equiv \sum_{-\infty < y < \infty} yp_Y(y) = E[g(X_1, \dots, X_n)] \\ &\equiv \sum_{-\infty < x_1 < \infty, \dots, -\infty < x_n < \infty} \underline{g(x_1, \dots, x_n)} \underline{p(x_1, \dots, x_n)},\end{aligned}$$

if X_1, X_2, \dots, X_n are discrete random variables, or

$$\begin{aligned}\underline{E(Y)} &\equiv \int_{-\infty}^{\infty} yf_Y(y)dy = E[g(X_1, \dots, X_n)] \\ &\equiv \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \underline{g(x_1, \dots, x_n)} \underline{f(x_1, \dots, x_n)} dx_1 \cdots dx_n,\end{aligned}$$

if \underline{Y} and X_1, X_2, \dots, X_n are continuous random variables.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Definition 3.2 (mean, variance, standard deviation, covariance, correlation coefficient)

- (TBp.116&118) $\underline{g(x)} = \underline{x} \Rightarrow \underline{E[g(X)]} = \underline{E(X)}$ is called **mean** of \underline{X} , usually denoted by $\underline{E(X)}$ or $\underline{\mu_X}$.
- (TBp.131) $\underline{g(x)} = \underline{(x - \mu_X)^2} \Rightarrow \underline{E[g(X)]} = \underline{E[(X - E(X))^2]}$ is called **variance** of \underline{X} , usually denoted by $\underline{Var(X)}$ or $\underline{\sigma_X^2}$. The square root of variance, i.e., $\underline{\sigma_X}$, is called **standard deviation**.
- (TBp.138) $\underline{g(x, y)} = \underline{(x - \mu_X)(y - \mu_Y)} \Rightarrow \underline{E[g(X, Y)]} = \underline{E[(X - E(X))(Y - E(Y))]}$ is called **covariance** of \underline{X} and \underline{Y} , usually denoted by $\underline{Cov(X, Y)}$ or $\underline{\sigma_{XY}}$.
- (TBp.142) The **correlation coefficient** of $\underline{X}, \underline{Y}$ is defined as $\underline{\sigma_{XY}/(\sigma_X \sigma_Y)}$, usually denoted by $\underline{Cor(X, Y)}$ or $\underline{\rho_{XY}}$. \underline{X} and \underline{Y} are called **uncorrelated** if $\underline{\rho_{XY} = 0}$.

Notes. (intuitive explanation of mean)

1. Mean of a random variable parallels the notion of a weighted average.
2. It is helpful to think of the mean as the center of mass of the pmf/pdf.
3. Mean can be interpreted as a long-run average. (see Chapter 5.)

Notes. (intuitive explanation of variance and standard deviation)

1. variance is the average value of the squared deviation of X from μ_X .
2. If X has units, then mean and standard deviation have the same unit, and variance has unit squared.

Theorem 3.1 (properties of mean)

1. (TBp.125) For constants $\underline{a}, b_1, \dots, b_n \in \mathbb{R}$,

$$\underline{E}(\underline{a} + \sum_{i=1}^n b_i X_i) = \underline{a} + \sum_{i=1}^n b_i \underline{E}(X_i).$$
2. (TBp.124) If X, Y are independent, then

$$\underline{E}(g(X)h(Y)) = \underline{E}(g(X))\underline{E}(h(Y)).$$

In particular, $\underline{E}(XY) = \underline{E}(X)\underline{E}(Y)$.

(**Question 3.2:** $\underline{E}(X/Y) = \underline{E}(X)/\underline{E}(Y)$?)

Note. $\underline{E}[g(X)] \neq g[\underline{E}(X)]$ in general.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Theorem 3.2 (properties of variance and standard deviation)

1. (TBp.132) $\underline{\sigma}_X^2 = \underline{Var}(X) = \underline{E}[(X - \mu_X)^2] = \underline{E}(X^2) - \mu_X^2$.
2. (TBp.131) $\underline{Var}(\underline{a} + \underline{b}X) = \underline{b}^2 \underline{Var}(X)$, $\underline{a}, \underline{b} \in \mathbb{R}$, and $\underline{\sigma}_{\underline{a}+\underline{b}X} = |\underline{b}| \underline{\sigma}_X$.
3. (TBp.140)

$$\underline{Var}\left(\underline{a} + \sum_{i=1}^n b_i X_i\right) = \sum_{i=1}^n \underline{b}_i^2 \underline{Var}(X_i) + 2 \sum_{1 \leq i < j \leq n} \underline{b}_i \underline{b}_j \underline{Cov}(X_i, X_j).$$

In particular, $\underline{Var}(X + Y) = \underline{Var}(X) + \underline{Var}(Y) + 2\underline{Cov}(X, Y)$.

4. (TBp.140) If X_1, \dots, X_n are independent,

$$\underline{Var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \underline{Var}(X_i).$$

5. (TBp.136) $\underline{E}[(X - \theta)^2] = \underline{Var}(X) + (\mu_X - \theta)^2$ (Mean square error = variance + bias square)

Theorem 3.3 (Chebyshev's inequality, TBp. 133)

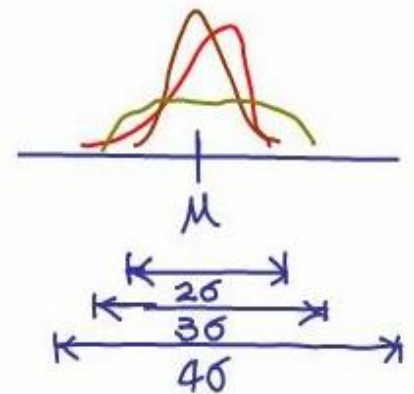
Let X be a random variable with mean μ and variance σ^2 . Then for any $t > 0$,

$$P(|X - \mu| > t) \leq \frac{\sigma^2}{t^2}.$$

Proof. Let $f(x)$ be the pdf of X . Let $R = \{x : |x - \mu| > t\}$.
Then

$$\begin{aligned} P(|X - \mu| > t) &= \int_R f(x) dx \leq \int_R \frac{(x - \mu)^2}{t^2} f(x) dx \\ &\leq \int_{-\infty}^{\infty} \frac{(x - \mu)^2}{t^2} f(x) dx = \frac{\sigma^2}{t^2}. \end{aligned}$$

No other restriction
on the functional
form of pdf/pmf



Note.

1. Setting $t = k\sigma$ we have

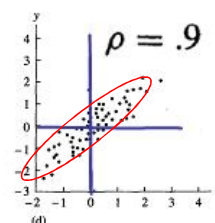
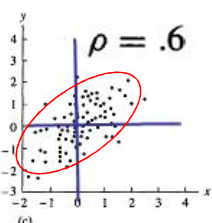
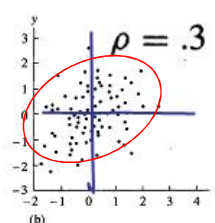
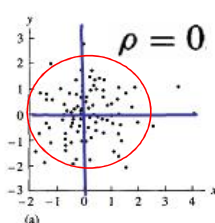
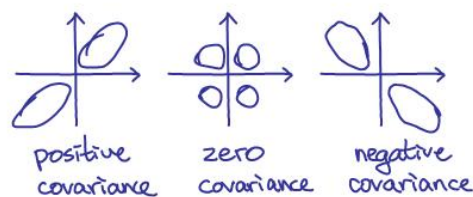
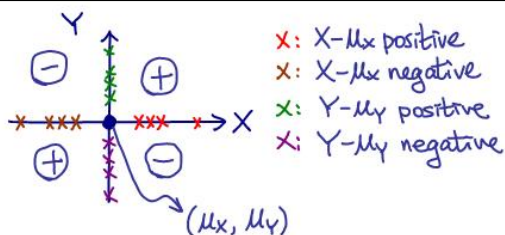
$$P(|X - \mu| > k\sigma) \leq \frac{1}{k^2}.$$

2. (TBp. 134) $\text{Var}(X) = 0 \Rightarrow P(X = \mu_X) = 1.$

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Notes. (intuitive explanation of covariance and correlation coefficient)

1. covariance is a measure of the joint variability of X and Y , or their degree of association.
2. covariance is the average value of the product of the deviation of X from its mean and the deviation of Y from its mean.
3. positive covariance and negative covariance
4. correlation coefficient is unit free
5. correlation coefficient measures the strength of the linear relationship between X and Y .



Theorem 3.4 (properties of covariance and correlation coefficient)

1. (TBp.138) $\underline{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)] = \underline{E}(XY) - \underline{\mu}_X \underline{\mu}_Y$
(Note. $\underline{Cov}(X, X) = \underline{Var}(X)$.)
2. (TBp.140)

$$\underline{Cov}\left(\underline{a} + \sum_{i=1}^n b_i X_i, \underline{c} + \sum_{j=1}^m d_j Y_j\right) = \sum_{i=1}^n \sum_{j=1}^m \underline{b}_i \underline{d}_j \underline{Cov}(X_i, Y_j)$$
3. (TBp.140) If X, Y are independent then $\underline{Cov}(X, Y) = 0$, i.e., **independent** \Rightarrow **uncorrelated**. But, the converse statement is not necessarily true.
4. (TBp.143) $\underline{-1} \leq \underline{\rho}_{XY} \leq \underline{1}$ and $\underline{\rho}_{XY} = \pm 1$ if and only if $\underline{Y} = \underline{a}X + \underline{b}$ with probability one for some $\underline{a}, \underline{b} \in \mathbb{R}$.
5. $\underline{\rho}_{XY} = \underline{E}\left[\left(\frac{X - \mu_X}{\sigma_X}\right)\left(\frac{Y - \mu_Y}{\sigma_Y}\right)\right]$
6. $|\underline{Cor}(\underline{a} + \underline{b}X, \underline{c} + \underline{d}Y)| = |\underline{Cor}(X, Y)|$

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

• moment generating function & characteristics function

Definition 3.3 (moment generating function, TBp. 155)

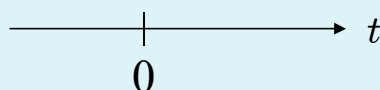
The moment generating function (mgf) of a random variable \underline{X} is

$$\underline{M}_X(t) = \underline{E}(e^{tX}), \quad t \in \mathbb{R}$$

if the expectation exists.

Theorem 3.5 (properties of moment generating function)

1. The moment generating function may or may not exist for any particular value of t .
2. **uniqueness theorem** (TBp.143). If the moment generating function exists for t in an open interval containing zero, it uniquely determines the probability distribution.



3. (TBp.156) If the moment generating function exists in an open interval containing zero, then

$$\underline{M_X^{(k)}(0)} = \underline{E(X^k)}.$$

4. (TBp.158) For any constants a, b , $\underline{M_{a+bX}(t)} = \underline{e^{at} M_X(bt)}$.

5. (TBp.159) X, Y independent $\Rightarrow \underline{M_{X+Y}(t)} = \underline{M_X(t) M_Y(t)}$.

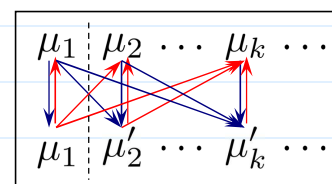
6. continuity theorem (see Chapter 5)

Definition 3.4 (moment, TBp. 155)

The k th **moment** of a random variable is $\underline{E(X^k)} \equiv \underline{\mu_k}$, and the k th **central moment** is $\underline{E[(X - \mu_X)^k]} \equiv \underline{\mu'_k}$.

➤ Some Notes.

- $\underline{\mu'_k} = \sum_{i=0}^k \binom{k}{i} (-\mu_X)^{n-i} \underline{\mu_i}$.
- $\underline{\mu_k} = \sum_{i=0}^k \binom{k}{i} (\mu_X)^{n-i} \underline{\mu'_i}$.
- In particular, $\underline{E(X)} = \underline{\mu_X} = \underline{\mu_1}$, and,
 $\underline{Var(X)} = \underline{\sigma_X^2} = \underline{\mu_2 - \mu_1^2} = \underline{\mu'_2}$.



NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Definition 3.5 (joint moment generating function, TBp. 161)

For random variables X_1, X_2, \dots, X_n , their **joint mgf** is defined as:

$$\underline{M_{X_1 X_2 \dots X_n}(t_1, t_2, \dots, t_n)} = \underline{E(e^{t_1 X_1 + t_2 X_2 + \dots + t_n X_n})}$$

if the expection exists.

Theorem 3.6 (properties of joint mgf)

1. $\underline{M_{X_1}(t_1)} = \underline{M_{X_1 X_2 \dots X_n}(t_1, 0, \dots, 0)}$
2. uniqueness theorem
3. X_1, X_2, \dots, X_n are independent if and only if

$$\underline{M_{X_1 X_2 \dots X_n}(t_1, t_2, \dots, t_n)} = \prod_{i=1}^n \underline{M_{X_i}(t_i)}.$$

4.
$$\left. \frac{\partial^{r_1 + \dots + r_n}}{\partial t_1^{r_1} \dots \partial t_n^{r_n}} M_{X_1 X_2 \dots X_n}(t_1, t_2, \dots, t_n) \right|_{t_1=t_2=\dots=t_n=0}$$

$$= \underline{E(X_1^{r_1} X_2^{r_2} \dots X_n^{r_n})}$$

Definition 3.6 (characteristic function, TBp. 161)

The characteristic function (chf) of a random variable \underline{X} is

$$\underline{\phi_X}(t) = \underline{E}(e^{itX}) = \underline{E}[\cos(tX)] + i \cdot \underline{E}[\sin(tX)],$$

where $i = \sqrt{-1}$, and the joint characteristic function of $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$ is

$$\underline{\phi_{X_1 X_2 \dots X_n}}(t_1, t_2, \dots, t_n) = \underline{E}(e^{it_1 X_1 + it_2 X_2 + \dots + it_n X_n}).$$

Theorem 3.7 (properties of characteristic function)

1. The characteristic function always exists.
2. If $\underline{M_X}(t)$ exists, then $\underline{\phi_X}(t) = \underline{M_X}(it)$.
3. uniqueness theorem
4. (FYI) inversion theorem:
 - discrete case: $\underline{p_X}(x) = \lim_{T \rightarrow \infty} \int_{-T}^T e^{-itx} \underline{\phi_X}(t) dt$
 - continuous case: $\underline{f_X}(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-itx} \underline{\phi_X}(t) dt$
5. The properties of characteristic function are similar to those of moment generating function.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

• conditional expectation

Definition 3.7 (conditional expectation, TBp. 135-136)

The conditional expectation of $\underline{h(Y)}$ given $\underline{X = x}$ is

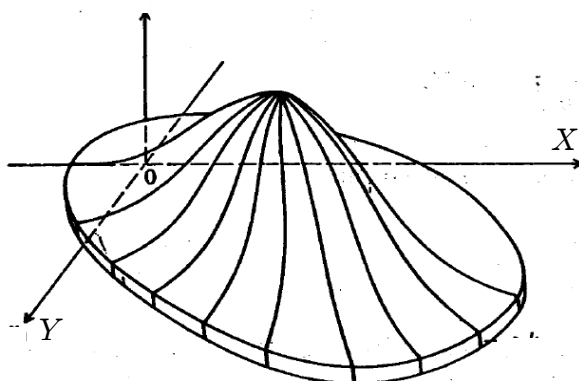
$$[\underline{\text{Discrete case}}] : \underline{E(h(Y)|X = x)} = \sum_y \underline{h(y)} \underline{p_{Y|X}(y|x)}$$

$$\text{In particular, } \underline{E(Y|X = x)} = \sum_y \underline{y} \underline{p_{Y|X}(y|x)}$$

$$[\underline{\text{Continuous case}}]: \underline{E(h(Y)|X = x)} = \int \underline{h(y)} \underline{f_{Y|X}(y|x)} dy$$

$$\text{In particular, } \underline{E(Y|X = x)} = \int \underline{y} \underline{f_{Y|X}(y|x)} dy$$

$f(x, y)$: joint pdf



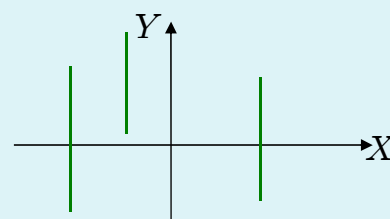
Theorem 3.8 (properties of conditional expectation)

1. $E_{Y|X}(h(Y)|x)$ is a function of x and is free of Y .
2. If X and Y are independent then $E_{Y|X}(h(Y)|x) = E_Y(h(Y))$.
3. $E(h(X)|X = x) = h(x)$
4. Let $g(x) = E_{Y|X}(h(Y)|x)$, then $g(X)$ is a random variable (transformation of X) and usually denoted by $E_{Y|X}(h(Y)|X)$.
5. **law of total expectation** (TBp.149)

$$E_X[E_{Y|X}(h(Y)|X)] = E_Y[h(Y)].$$

In particular,

$$E_Y(Y) = E_X[E_{Y|X}(Y|X)].$$



NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

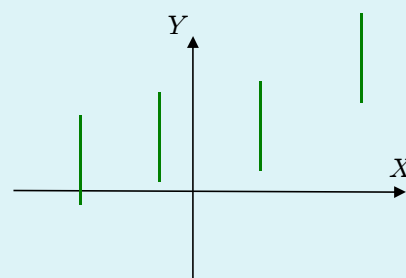
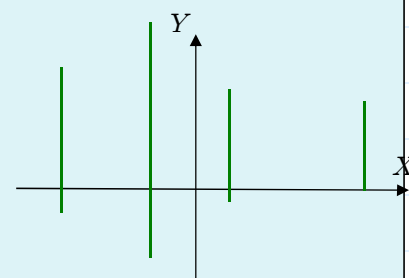
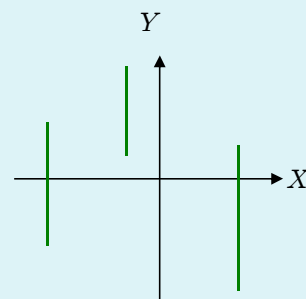
4. variance decomposition (TBp.151)

$$\begin{aligned} \text{Var}_Y(Y) = & \\ & \text{Var}_X[E_{Y|X}(Y|X)] + \\ & E_X[\text{Var}_{Y|X}(Y|X)] \end{aligned}$$

Note.

1. $\text{Var}_Y(Y) \geq E_X[\text{Var}_{Y|X}(Y|X)]$
and the equality holds if and only if $E_{Y|X}(Y|X) = E_Y(Y)$
with probability one.

2. $\text{Var}_Y(Y) \geq \text{Var}_X[E_{Y|X}(Y|X)]$
and the equality holds if and only if $\text{Var}_{Y|X}(Y|X) = 0$
with probability one; i.e.,
 $Y = E_{Y|X}(Y|X)$
with probability one.



• prediction

Example 3.1 (predicting the value of a r.v. Y from another r.v. X , TBp. 152-154)

- **data:** X and Y (example?)
- **statistical modeling:** assign (X, Y) a (known) joint distribution (cdf $F(x, y)$, pdf $f(x, y)$, or pmf $p(x, y)$)
- **objective:** Predict Y by using a function of X , i.e., $g(X)$.

We consider the following three groups of g 's:

- (i) $G_1 = \{g(x) : \underline{g(x) = c}, \text{ where } c \in \mathbb{R}\}$
- (ii) $G_2 = \{g(x) : \underline{g(x) = a + bx}, \text{ where } a, b \in \mathbb{R}\}$, and
- (iii) $G_3 = \{g(x) : \underline{g \text{ is arbitrary}}\}$.

Note. $G_1 \subset G_2 \subset G_3$.

- **question:** Within each group, what is the “best” prediction?
- **criterion:** minimizing mean square error:

$$\underline{\text{MSE}} \equiv \underline{E_{X,Y}\{[Y - \underline{g(X)}]^2\}}.$$

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Ch1-6, p.2-54

Example 3.2 (“best” constant prediction, TBp. 153)

$$\underline{E_{X,Y}(Y - \underline{c})^2} = \underline{E_Y(Y - \underline{c})^2} \geq \underline{E_Y[Y - \underline{E_Y(Y)}]^2} = \underline{\text{Var}_Y(Y)}$$

The equality holds if and only if $\underline{c = E_Y(Y)}$.

Example 3.3 (“best” prediction of Y using X , TBp. 153)

$$\underline{E_{X,Y}[Y - \underline{g(X)}]^2} \geq \underline{E_{X,Y}[Y - \underline{E_{Y|X}(Y|X)}]^2} = \underline{E_X[\text{Var}_{Y|X}(Y|X)]}$$

The equality holds if and only if $\underline{g(x) = E_{Y|X}(Y|x)}$.

Notes for the best predictor in G_3 .

- $\underline{E_{Y|X}(Y|X)}$ is the best predictor of Y based on X , in the mean squared prediction error sense.
- need to know the joint distribution of X and Y , or at least $E_{Y|X}(Y|x)$
- $\underline{E_{Y|X}(Y|x)}$ is called the regression function of Y on X .

Example 3.4 (“best” linear prediction of Y using X , TBp. 153-154)

$$\underline{E_{X,Y}[Y - \underline{(a + bX)}]^2} \geq \underline{E_{X,Y}\left\{Y - \left[\mu_Y + \rho \frac{\sigma_Y}{\sigma_X}(X - \mu_X)\right]\right\}^2} = \underline{\sigma_Y^2(1 - \rho^2)}$$

The equality holds if and only if $\underline{a = \mu_Y - b\mu_X}$ and $\underline{b = \rho \frac{\sigma_Y}{\sigma_X}}$.

Notes for the best predictor in G_2 .

- $E_{Y|X}(Y|x) = \mu_Y + \rho \frac{\sigma_Y}{\sigma_X}(x - \mu_X)$ if (X, Y) is distributed as bivariate normal
- needs to know only the means, variances and covariances
- $\sigma_Y^2(1 - \rho^2)$ is small if ρ is close to $+1$ or -1 , and large if ρ is close to 0

Notes.

1. $\min_{a,b} E[Y - (a + bX)]^2 \leq \min_c E(Y - c)^2$ and the equality holds if and only if $\rho = 0$.
2. $\min_g E(Y - g(X))^2 \leq \min_{a,b} E[Y - (a + bX)]^2$ and the equality holds if and only if $E_{Y|X}(Y|x) = \mu_Y + \rho(\sigma_Y/\sigma_X)(x - \mu_X)$.

Question 3.3

What if the joint distribution of X and Y is unknown?

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

• δ method

Question 3.3

Let $Y = g(X)$. Suppose we only know the mean μ_X and variance σ_X^2 of X , but **not** the entire distribution (i.e., do not know cdf, pdf/pmf of X). Can we derive the distribution of Y ? If not, can we “roughly” describe the mean and variance of Y ? (**Note.** $E[g(X)] \neq g[E(X)]$.)

Theorem 3.9 (δ method for univariate case, TBp. 162)

$$Y = g(X) \approx g(\mu_X) + (X - \mu_X)g'(\mu_X) \quad (\text{by Taylor expansion})$$

$$\Rightarrow E[g(X)] \approx g(\mu_X)$$

$$\text{Var}[g(X)] \approx \text{Var}(X)[g'(\mu_X)]^2$$

or
$$Y = g(X) \approx g(\mu_X) + (X - \mu_X)g'(\mu_X) + \frac{1}{2}(X - \mu_X)^2 g''(\mu_X)$$

$$\Rightarrow E[g(X)] \approx g(\mu_X) + \frac{1}{2}\sigma_X^2 g''(\mu_X)$$

Note. How good these approximations are depends on whether g can be reasonably well approximated by the 1st or 2nd order polynomials in a neighborhood of μ_X and on the size of σ_X .

Theorem 3.10 (δ method for multivariate case, TBp. 165)

Function of two univariate random variables $Z = g(X, Y)$:

Let $\underline{\mu} = (\mu_X, \mu_Y)$.

$$Z = g(X, Y) \approx g(\underline{\mu}) + (X - \mu_X) \frac{\partial g(\underline{\mu})}{\partial x} + (Y - \mu_Y) \frac{\partial g(\underline{\mu})}{\partial y}$$

$$\Rightarrow E(Z) \approx g(\underline{\mu})$$

$$\text{Var}(Z) \approx \underline{\sigma}_X^2 \left[\frac{\partial g(\underline{\mu})}{\partial x} \right]^2 + \underline{\sigma}_Y^2 \left[\frac{\partial g(\underline{\mu})}{\partial y} \right]^2 + 2\underline{\sigma}_{XY} \left[\frac{\partial g(\underline{\mu})}{\partial x} \right] \left[\frac{\partial g(\underline{\mu})}{\partial y} \right]$$

or

$$g(X, Y) \approx g(\underline{\mu}) + (X - \mu_X) \frac{\partial g(\underline{\mu})}{\partial x} + (Y - \mu_Y) \frac{\partial g(\underline{\mu})}{\partial y}$$

$$+ \frac{1}{2} (X - \mu_X)^2 \frac{\partial^2 g(\underline{\mu})}{\partial x^2} + (X - \mu_X)(Y - \mu_Y) \frac{\partial^2 g(\underline{\mu})}{\partial x \partial y}$$

$$+ \frac{1}{2} (Y - \mu_Y)^2 \frac{\partial^2 g(\underline{\mu})}{\partial y^2}$$

$$\Rightarrow E[g(X, Y)] \approx g(\underline{\mu}) + \frac{1}{2} \underline{\sigma}_X^2 \frac{\partial^2 g(\underline{\mu})}{\partial x^2} + \underline{\sigma}_{XY} \frac{\partial^2 g(\underline{\mu})}{\partial x \partial y} + \frac{1}{2} \underline{\sigma}_Y^2 \frac{\partial^2 g(\underline{\mu})}{\partial y^2}$$

Note. The general case of a function of n random variables can be worked out similarly.

❖ **Reading:** textbook, Chapter 4

❖ **Further Reading:** Roussas, 5.1, 5.3, 5.4, 5.5, 6.1, 6.2, 6.4, 6.5

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Some Commonly Used Distributions (from Chapters 2, 3, 6)

Question 4.1

For a given random phenomenon or data, what distribution (or statistical model) is more appropriate to depict it?

• discrete distributions

Definition 4.1 (Uniform distribution $U(a_1, \dots, a_m)$)

Equal probability to obtain a_1, a_2, \dots, a_m .

• **pmf:** $p(x) = \begin{cases} \frac{1}{m}, & x = a_1, \dots, a_m \\ 0, & \text{otherwise} \end{cases}$

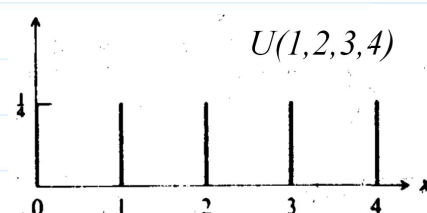
• **mgf:** $\frac{\sum_{j=1}^m e^{a_j t}}{m}$

• **mean:** $\frac{\sum_{j=1}^m a_j}{m} \equiv \bar{a}$

• **variance:** $\frac{\sum_{j=1}^m (a_j - \bar{a})^2}{m}$

• **parameter:** $a_i \in \mathbb{R}, m = 1, 2, \dots$

• **example:** throw a fair die once



Definition 4.2 (Bernoulli distribution $B(p)$, sec 2.1.1)

A Bernoulli distribution takes on only two values: 0 and 1, with probabilities $1 - p$ and p , respectively.

- **pmf:** $p(x) = \begin{cases} p^x(1-p)^{(1-x)}, & \text{if } x = 0 \text{ or } x = 1 \\ 0, & \text{otherwise} \end{cases}$
- **mgf:** $pe^t + 1 - p$
- **mean:** p
- **variance:** $p(1 - p)$
- **parameter:** $p \in [0, 1]$
- **example:** toss a coin once, p =probability that head occurs

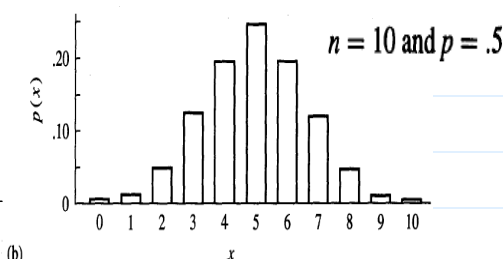
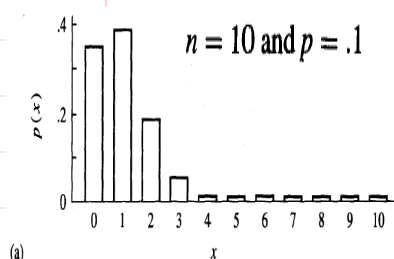
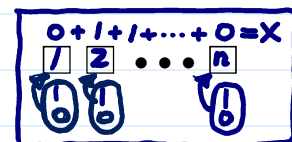
Note: If A is an event, then the indicator random variable I_A follows the Bernoulli distribution.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Definition 4.3 (Binomial distribution $B(n, p)$, sec 2.1.2)

Suppose that n independent Bernoulli trials are performed, where n is a fixed number. The total number of 1 appearing in the n trials follows a binomial distribution with parameters n and p .

- **pmf:** $p(x) = \begin{cases} \binom{n}{x} p^x (1-p)^{(n-x)}, & x = 0, 1, \dots, n \\ 0, & \text{otherwise} \end{cases}$
- **mgf:** $(pe^t + 1 - p)^n, t \in \mathbb{R}.$
- **mean:** np
- **variance:** $np(1 - p)$
- **parameter:** $p \in [0, 1], n = 1, 2, \dots$
- **example:** # of heads, toss a coin n times

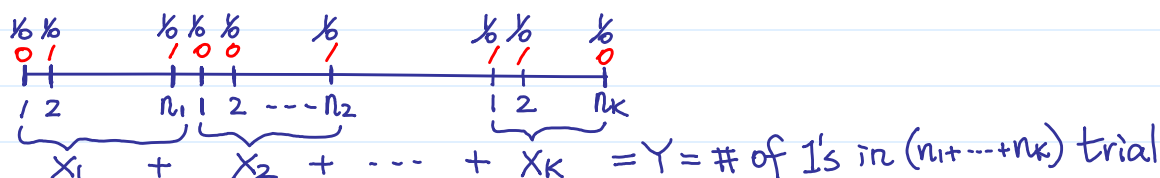


Note:

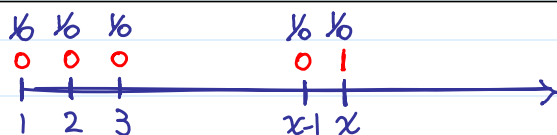
$$(a + b)^n = \sum_{x=0}^n \binom{n}{x} a^x b^{n-x}$$

Note.

1. binomial distribution is a generalization of bernoulli distribution from 1 trial to n trials
2. Let X_1, \dots, X_n be i.i.d. $B(p)$, then $Y = X_1 + \dots + X_n \sim B(n, p)$.
3. Let $X_i \sim B(n_i, p), i = 1, \dots, k$, and X_1, \dots, X_k are independent. Then, $Y = X_1 + \dots + X_k \sim B(n_1 + \dots + n_k, p)$.

**Definition 4.4** (Geometric distribution $G(p)$, sec 2.1.3)

The geometric distribution is constructed from an infinite sequence of independent Bernoulli trials. Let X be the total number of trials up to and including the first appearance of 1. Then, X follows the geometric distribution.



NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

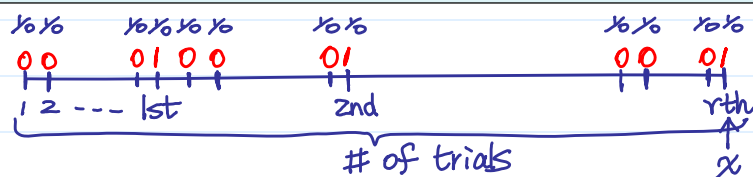
- **pmf:** $p(x) = \begin{cases} (1-p)^{(x-1)}p, & x = 1, 2, 3, \dots \\ 0, & \text{otherwise} \end{cases}$
- **cdf:** $F(x) = \begin{cases} 1 - (1-p)^{[x]}, & 1 \leq [x] \leq x < [x] + 1 \\ 0, & x < 1 \end{cases}$
- **mgf:** $\frac{pe^t}{1-(1-p)e^t}, t < -\log(1-p)$.
- **mean:** $\frac{1}{p}$
- **variance:** $\frac{1-p}{p^2}$
- **parameter:** $p \in [0, 1]$
- **example:** lottery, # of tickets a person must purchase up to and including the first winning ticket

Note: a memoryless distribution

Note:
 $\sum_{x=n}^{\infty} t^x = \frac{t^n}{1-t},$
for $-1 < t < 1$.

Definition 4.5 (Negative Binomial distribution $NB(r, p)$, sec 2.1.3)

An infinite sequence of independent Bernoulli trials is performed until the appearance of the r th 1. Let X denote the total number of trials. Then, X follows negative binomial distribution.



• pmf:
$$p(x) = \begin{cases} \binom{x-1}{r-1} p^r (1-p)^{(x-r)}, & x = r, r+1, \dots \\ 0, & \text{otherwise} \end{cases}$$

• mgf: $\frac{p^r e^{rt}}{[1-(1-p)e^t]^r}, \quad t < -\log(1-p).$

• mean: $\frac{r}{p}$

• variance: $\frac{r(1-p)}{p^2}$

• parameter: $p \in [0, 1], \quad r = 1, 2, \dots$

• example: lottery, # of tickets a person must purchase up to and including the r th winning ticket

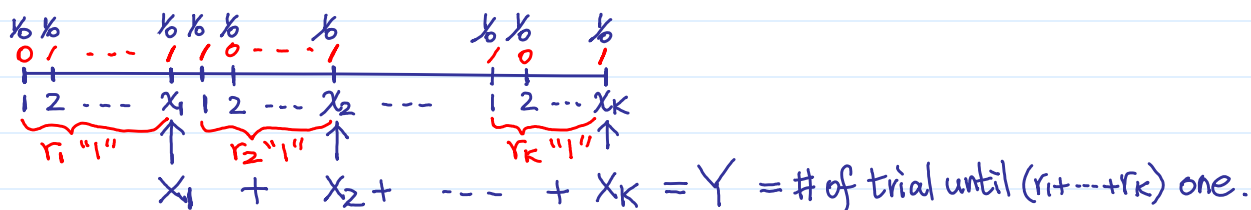
Note:

$$\sum_{x=0}^{\infty} \binom{n+x-1}{x} t^x = \frac{1}{(1-t)^n}, \quad \text{for } -1 < t < 1.$$

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Note.

1. negative binomial distribution is a generalization of geometric distribution from 1st success to r th success
2. Let X_1, X_2, \dots, X_r be i.i.d. $G(p)$, then $Y = X_1 + \dots + X_r \sim NB(r, p)$.
3. Let $X_i \sim NB(r_i, p), i = 1, \dots, k$, and X_1, \dots, X_k are independent. Then, $Y = X_1 + \dots + X_k \sim NB(r_1 + \dots + r_k, p)$.

**Definition 4.6** (Multinomial distribution $Multinomial(n, p_1, p_2, \dots, p_r)$, TBp.73-74)

Suppose that each of n independent trials can result in one of r types of outcomes, and that on each trial the probabilities of the r outcomes are p_1, p_2, \dots, p_r . Let X_i be the total number of outcomes of type i in the n trials, $i = 1, \dots, r$. Then, (X_1, \dots, X_r) follows a multinomial distribution.

- joint pmf:

$$p(x_1, \dots, x_r) = \begin{cases} \binom{n}{x_1 \dots x_r} p_1^{x_1} \dots p_r^{x_r}, & x_i = 0, 1, \dots, n, \text{ and } \sum_{i=1}^r x_i = n \\ 0, & \text{otherwise} \end{cases}$$

- joint mgf: $(p_1 e^{t_1} + \dots + p_r e^{t_r})^n$, $t_1, \dots, t_r \in \mathbb{R}$.

- marginal distribution: $X_i \sim B(n, p_i)$, $i = 1, \dots, r$

- mean: $E(X_i) = np_i$, $i = 1, \dots, n$

- variance: $Var(X_i) = np_i(1 - p_i)$, $i = 1, \dots, n$

- covariance: $Cov(X_i, X_j) = -np_i p_j$, $i \neq j$

- parameter: $p_i \in [0, 1]$, and $\sum_{i=1}^r p_i = 1$. $n = 1, 2, \dots$

- example: randomly choose n people, record the numbers of people with different religions

Note: $(a_1 + \dots + a_k)^n = \sum_{x_1 + \dots + x_k = n} \binom{n}{x_1, \dots, x_k} a_1^{x_1} \dots a_k^{x_k}$.

Notes: multinomial distribution is a generalization of the binomial distribution from 2 outcomes to r outcomes.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Definition 4.7 (Poisson distribution $P(\lambda)$, sec 2.1.5)

Limit of binomial distributions $X_n \sim B(n, p_n)$, where $p_n \rightarrow 0$ as $n \rightarrow \infty$ in such a way that $\lambda_n \equiv np_n \rightarrow \lambda$.

$$\binom{n}{x} p_n^x (1 - p_n)^{(n-x)}$$

Note: if $a_n \rightarrow a$, $(1 + \frac{a_n}{n})^n \rightarrow e^a$.

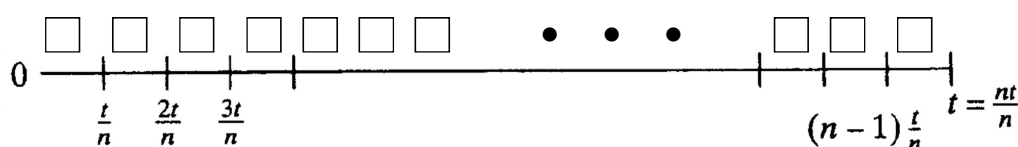
$$= \frac{n(n-1) \dots (n-x+1)}{x!} \left(\frac{\lambda_n}{n}\right)^x \left(1 - \frac{\lambda_n}{n}\right)^{n-x}$$

$$= \frac{n(n-1) \dots (n-x+1)}{n^x} \frac{1}{x!} \lambda_n^x \left(1 - \frac{\lambda_n}{n}\right)^{n-x}$$

$$= 1 \left(1 - \frac{1}{n}\right) \dots \left(1 - \frac{x-1}{n}\right) \frac{\lambda_n^x}{x!} \left(1 - \frac{\lambda_n}{n}\right)^n \left(1 - \frac{\lambda_n}{n}\right)^{-x} \rightarrow 1^x \cdot \frac{\lambda^x}{x!} \cdot e^{-\lambda} \cdot 1 = \frac{\lambda^x e^{-\lambda}}{x!}$$

explanations.

- if n large, the pmf of $B(n, p)$ is not easily calculated. Then, we can approximate them by pmf of $P(\lambda)$, where $\lambda = np$.



2. Let X be the number of times some event occurs in a given time interval I . Divide the interval into many small subintervals I_k , $k = 1, \dots, n$, of equal length. Let N_k be the number of events occurring in I_k . When we can assume N_1, \dots, N_n are independent and approximately $\sim B(p)$, X has a distribution near $P(\lambda)$, where $\lambda = np$.

- **pmf:** $p(x) = \begin{cases} \frac{\lambda^x}{x!} e^{-\lambda}, & x = 0, 1, 2, \dots \\ 0, & \text{otherwise} \end{cases}$

- **mgf:** $e^{\lambda(e^t - 1)}$, $t \in \mathbb{R}$.

- **mean:** λ

- **variance:** λ

- **parameter:** $\lambda > 0$

Note:

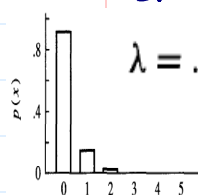
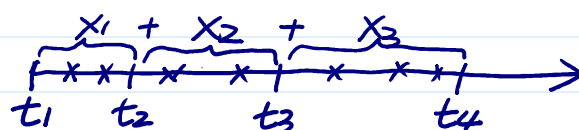
$$e^\lambda = \sum_{x=0}^{\infty} \frac{\lambda^x}{x!}.$$

- **example:** number of phone calls coming into an exchange during a unit of time

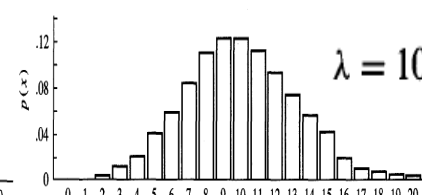
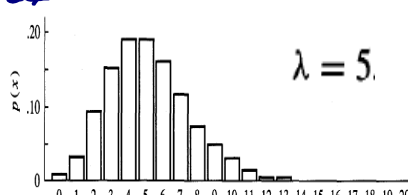
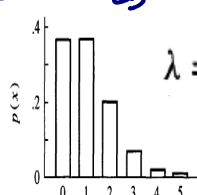
NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Note: Let $X_i \sim P(\lambda_i)$, $i = 1, \dots, k$, and X_1, \dots, X_k are independent.
Then, $Y = X_1 + \dots + X_k \sim P(\lambda_1 + \dots + \lambda_k)$.

Ch1-6, p.2-68



(a)



Definition 4.8 (Hypergeometric distribution $HG(r, n, m)$, sec 2.1.4)

Suppose that an urn contains n black balls and m white balls. Let X denote the number of black balls drawn when taking r balls without replacement. Then, X follows hypergeometric distribution.

- **pmf:** $p(x) = \begin{cases} \frac{\binom{n}{x} \binom{m}{r-x}}{\binom{n+m}{r}}, & x = 0, 1, \dots, \min(r, n), \\ 0, & \text{otherwise} \end{cases}$

Note:

$$\binom{n+m}{r} = \sum_x \binom{n}{x} \binom{m}{r-x}.$$

- **mgf:** exist, but no simple expression
- **mean:** $\frac{rn}{n+m}$
- **variance:** $\frac{rnm(n+m-r)}{(n+m)^2(n+m-1)}$
- **parameter:** $r, n, m, = 1, 2, \dots, r \leq n + m$
- **example:** sampling industrial products for defect inspection

Notes. a relationship between hypergeometric and binomial distributions: Let $m, n \rightarrow \infty$ in such a way that

$$\underline{p_{m,n}} \equiv \frac{n}{m+n} \rightarrow p,$$

where $0 < p < 1$. Then,

$$\frac{\binom{n}{x} \binom{m}{r-x}}{\binom{n+m}{r}} \rightarrow \underline{\binom{r}{x} p^x (1-p)^{r-x}}.$$

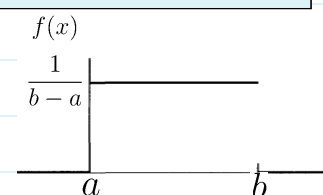
NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

• continuous distributions

Definition 4.9 (Uniform distribution $U(a, b)$, sec 2.2)

Choose a number at random between a and b .

- **pdf:** $f(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{otherwise} \end{cases}$



- **cdf:** $F(x) = \begin{cases} 0, & x < a \\ \frac{x-a}{b-a}, & a \leq x \leq b \\ 1, & x > b \end{cases}$

- **mgf:** $\frac{e^{bt}-e^{at}}{t(b-a)}, t \in \mathbb{R}.$

- **mean:** $\frac{a+b}{2}$

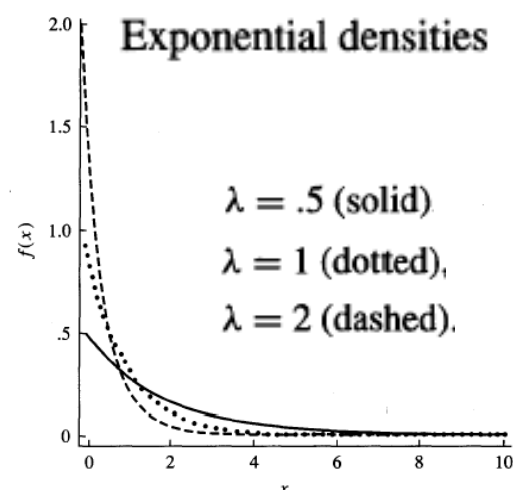
- **variance:** $\frac{(b-a)^2}{12}$

- **parameter:** $a, b \in \mathbb{R}, a < b$

Note: $U(0, 1)$ is useful for pseudo-random number generation

Definition 4.10 (Exponential distribution $E(\lambda)$, sec 2.2.1)

- pdf: $f(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$
- cdf: $F(x) = \begin{cases} 1 - e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$
- mgf: $\frac{\lambda}{\lambda - t}, t < \lambda$.
- mean: $\frac{1}{\lambda}$
- variance: $\frac{1}{\lambda^2}$
- parameter: $\lambda > 0$
- example: lifetime or waiting time



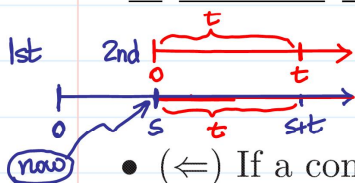
NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Notes:

Ch1-6, p.2-72

1. memoryless (future independent of past): Let $T \sim E(\lambda)$, then

$$\begin{aligned} P(T > t+s | T > s) &= \frac{P(T > t+s \text{ and } T > s)}{P(T > s)} = \frac{P(T > t+s)}{P(T > s)} \\ &= \frac{e^{-\lambda(t+s)}}{e^{-\lambda s}} = e^{-\lambda t} = P(T > t) \end{aligned}$$

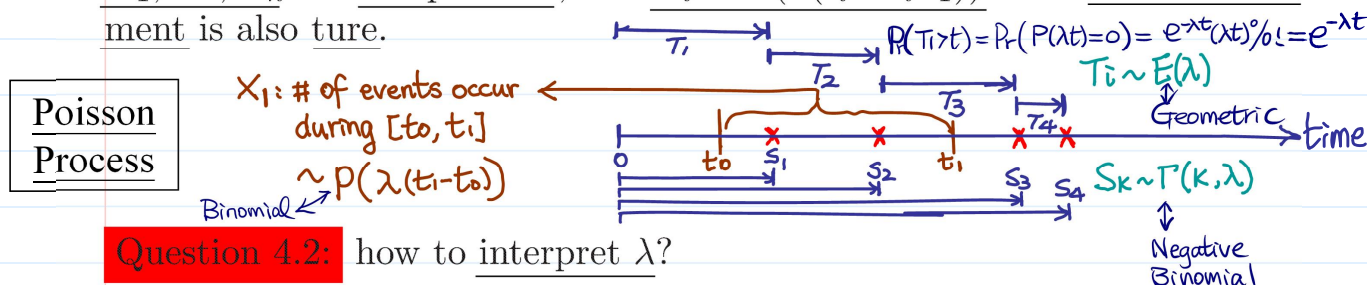


- (\Leftarrow) If a continuous distribution is memoryless, it is exponential.
- It does not mean the two events $T > s$ and $T > t+s$ are independent.

2. relationship between exponential, gamma, and Poisson distributions

Let T_1, T_2, T_3, \dots be i.i.d. $\sim E(\lambda)$ and $S_k = T_1 + \dots + T_k, k = 1, 2, \dots$

Let X_i be the number of S_k 's that falls in $[t_{i-1}, t_i], i = 1, \dots, n$, then X_1, \dots, X_n are independent, and $X_i \sim P(\lambda(t_i - t_{i-1}))$. The reverse statement is also true.

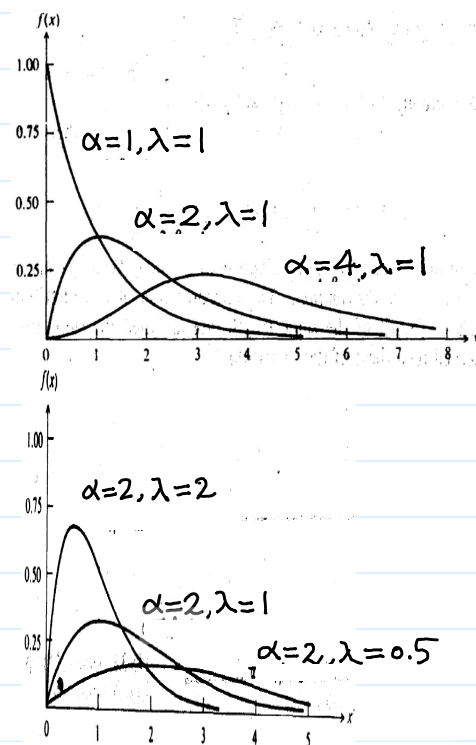


Question 4.2: how to interpret λ ?

3. Sometimes, the pdf is written as $\frac{1}{\lambda} e^{-\frac{x}{\lambda}}$. In the case, how to interpret λ ?

Definition 4.11 (Gamma distribution $\Gamma(\alpha, \lambda)$, sec 2.2.2)

- pdf: $f(x) = \begin{cases} \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$
- mgf: $(\frac{\lambda}{\lambda-t})^\alpha, t < \lambda$.
- mean: $\frac{\alpha}{\lambda}$
- variance: $\frac{\alpha}{\lambda^2}$
- parameter: $\alpha, \lambda > 0$



Notes.

1. α : shape parameter; λ : scale parameter (Question 4.3: how to interpret α, λ from the view point of Poisson process?)

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

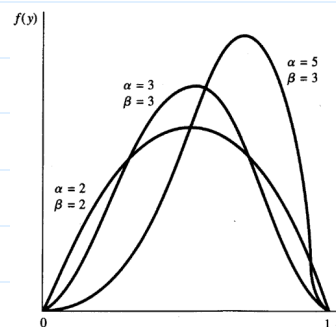
2. properties of gamma function $\Gamma(\alpha)$:

- $\Gamma(\alpha) \equiv \int_0^\infty y^{\alpha-1} e^{-y} dy$ (which is finite for $\alpha > 0$)
- $\Gamma(1) = 1$ and $\Gamma(\frac{1}{2}) = \sqrt{\pi}$
- $\Gamma(\alpha) = (\alpha-1)\Gamma(\alpha-1)$
- $\Gamma(\alpha) = (\alpha-1)!$ if α is an integer
- $\Gamma(\frac{\alpha}{2}) = \frac{\sqrt{\pi}(\alpha-1)!}{2^{\alpha-1}(\frac{\alpha-1}{2})!}$ if α is an odd integer

3. gamma distribution can be viewed as a generalization of exponential distribution, i.e., $\Gamma(1, \lambda) = E(\lambda)$.
4. Let X_1, \dots, X_k be i.i.d. $\sim E(\lambda)$, then $Y = X_1 + \dots + X_k \sim \Gamma(k, \lambda)$.
5. Let X_1, \dots, X_k be independent, and $X_i \sim \Gamma(\alpha_i, \lambda)$, then $Y = X_1 + \dots + X_k \sim \Gamma(\alpha_1 + \dots + \alpha_k, \lambda)$.
6. Let $X \sim \Gamma(\alpha, \lambda)$, then $cX \sim \Gamma(\alpha, \lambda/c)$, where $c > 0$.
7. $X \sim \Gamma(\alpha, \lambda) \Rightarrow E(X^k) = \frac{\Gamma(\alpha+k)}{\lambda^k \Gamma(\alpha)}$, for $0 < k$ and $E(\frac{1}{X^k}) = \frac{\lambda^k \Gamma(\alpha-k)}{\Gamma(\alpha)}$, for $0 < k < \alpha$.

Definition 4.12 (Beta distribution $\text{beta}(\alpha, \beta)$, sec 15.3.2)

- pdf: $f(x) = \begin{cases} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}, & 0 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases}$
- mgf: $1 + \sum_{k=1}^{\infty} \left(\prod_{r=0}^{k-1} \frac{\alpha+r}{\alpha+\beta+r} \right) \frac{t^k}{k!}$
- mean: $\frac{\alpha}{\alpha+\beta}$
- variance: $\frac{\alpha\beta}{(\alpha+\beta+1)(\alpha+\beta)^2}$
- parameter: $\alpha, \beta > 0$

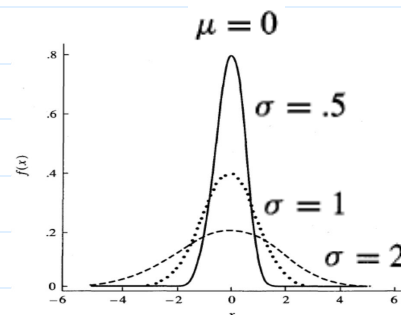

Notes:

1. Beta function: $B(\alpha, \beta) \equiv \int_0^1 t^{\alpha-1} (1-t)^{\beta-1} dt = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$
2. $\beta(1, 1) = U(0, 1)$
3. Let $X_1 \sim \Gamma(\alpha_1, \lambda)$, $X_2 \sim \Gamma(\alpha_2, \lambda)$, and X_1, X_2 independent. Then, $\frac{X_1}{X_1+X_2} \sim \text{beta}(\alpha_1, \alpha_2)$.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Definition 4.13 (Normal distribution $N(\mu, \sigma^2)$, sec. 2.2.3)

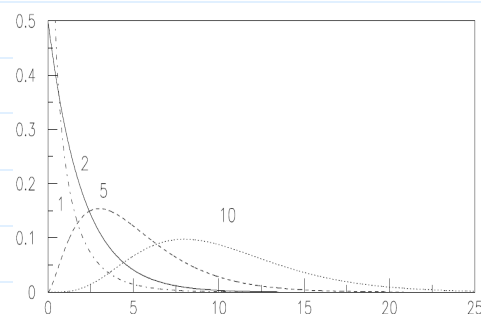
- pdf: $f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, x \in \mathbb{R}$.
- mgf: $e^{\mu t + \frac{\sigma^2 t^2}{2}}, t \in \mathbb{R}$
- mean: μ
- variance: σ^2
- parameter: $\mu \in \mathbb{R}, \sigma > 0$


Notes:

1. bell-shaped pdf (symmetric about μ , where it has maximum, and falls off in the rate determined by σ)
2. play a central role in probability and statistics (e.g., CLT, Chapter 5)
3. $X \sim N(\mu, \sigma^2) \Rightarrow$ for $a, b \in \mathbb{R}$, $aX + b \sim N(a\mu + b, a^2\sigma^2)$. In particular, $\frac{X-\mu}{\sigma} \sim N(0, 1)$.
4. Let X_1, \dots, X_k be independent, and $X_i \sim N(\mu_i, \sigma_i^2)$. Then $Y = X_1 + \dots + X_k \sim N(\sum_{i=1}^k \mu_i, \sum_{i=1}^k \sigma_i^2)$
5. by 3 and 4, let X_1, \dots, X_k be i.i.d $\sim N(\mu, \sigma^2)$, then $\bar{X}_k \sim N(\mu, \frac{\sigma^2}{k})$.

Definition 4.14 (Chi-square distribution $\chi^2_{\underline{n}}$, TBp.177)

- pdf: $f(x) = \begin{cases} \frac{1}{\Gamma(\frac{n}{2})2^{\frac{n}{2}}} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, & x \geq 0 \\ 0, & x \leq 0 \end{cases}$
- mgf: $(\frac{1}{1-2t})^{\frac{n}{2}}$
- mean: n
- variance: $2n$
- parameter: $n = 1, 2, 3, \dots$

**Notes:**

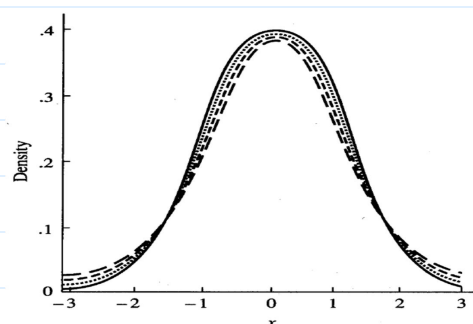
1. \underline{n} is called degree of freedom
2. $\chi^2_n = \Gamma(\frac{n}{2}, \frac{1}{2})$
3. Let $\underline{X}_1, \dots, \underline{X}_k$ be independent and $\underline{X}_i \sim \chi^2_{n_i}$, then $\underline{Y} = \underline{X}_1 + \dots + \underline{X}_k \sim \chi^2_{n_1 + \dots + n_k}$
4. Let $\underline{Z} \sim N(0, 1)$, then $\underline{X} = \underline{Z}^2 \sim \chi^2_1$.
5. By 3 and 4, let $\underline{Z}_1, \dots, \underline{Z}_n$ be i.i.d. $\sim N(0, 1)$, then $\underline{Y} = \underline{Z}_1^2 + \dots + \underline{Z}_n^2 \sim \chi^2_n$.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Definition 4.15 (t distribution \underline{t}_n , TBp.178)

- pdf: $f(x) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi}\Gamma(\frac{n}{2})} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}, x \in \mathbb{R}.$
- mgf: not exist, except at $\underline{t} = 0$
- mean: $0, (n > 1)$
- variance: $\frac{n}{n-2}, (n > 2)$
- moments:
$$E(X^k) = \begin{cases} \frac{\Gamma(\frac{k+1}{2})\Gamma(\frac{n-k}{2})}{\sqrt{\pi}\Gamma(\frac{n}{2})} n^{\frac{k}{2}}, & k < n \text{ and } \underline{\text{even}} \\ 0, & k < n \text{ and } \underline{\text{odd}} \end{cases}$$
- parameter: $n = 1, 2, 3, \dots$

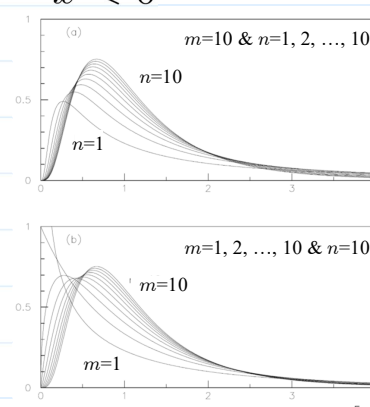
t_5 (long dash), t_{10} (short dash),
 t_{30} (dot), $N(0,1)$ (solid)

**Notes:**

1. Let $\underline{Z} \sim N(0, 1)$ and $\underline{U} \sim \chi^2_n$ be independent, then $\frac{\underline{Z}}{\sqrt{\underline{U}/n}} \sim \underline{t}_n$.
2. $f(x) = f(-x)$, i.e., \underline{t}_n distribution is symmetric about zero
3. as $\underline{n} \rightarrow \infty$, \underline{t}_n tends to $N(0, 1)$. (by LLN, Chapter 5)
4. \underline{t}_n has heavier tail than $\underline{N}(0, 1)$

Definition 4.16 (F distribution $F_{m,n}$, TBp.179)

- pdf: $f(x) = \begin{cases} \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} \left(\frac{m}{n}\right)^{-\frac{m}{2}} x^{\frac{m}{2}-1} \left(1 + \frac{m}{n}x\right)^{-\frac{m+n}{2}}, & x \geq 0 \\ 0, & x < 0 \end{cases}$
- mgf: not exist, except at $t = 0$
- mean: $\frac{n}{n-2}, (n > 2)$
- variance: $\frac{2n^2(m+n-2)}{m(n-2)^2(n-4)}, (n > 4)$
- moments: $E(X^k) = \frac{\Gamma(\frac{m+2k}{2})\Gamma(\frac{n-2k}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} \left(\frac{n}{m}\right)^k, k < \frac{n}{2}$
- parameter: $m, n = 1, 2, 3, \dots$


Notes:

1. Let $U \sim \chi_m^2$ and $V \sim \chi_n^2$ be independent, then $\frac{U/m}{V/n} \sim F_{m,n}$.
2. Let $X \sim t_n$, then $Y = X^2 \sim F_{1,n}$.
3. $X \sim F_{m,n} \Rightarrow X^{-1} \sim F_{n,m}$

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Theorem 4.1 (distributions of sample mean and sample variance of i.i.d. normal, sec. 6.3)

Let X_1, X_2, \dots, X_n be i.i.d. $\sim N(\mu, \sigma^2)$. Define

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad (\text{called } \underline{\text{sample mean}})$$

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \quad (\text{called } \underline{\text{sample variance}})$$

Then,

1. $\bar{X}_n \sim N(\mu, \frac{\sigma^2}{n})$ and $\sqrt{n}(\bar{X}_n - \mu)/\sigma \sim N(0, 1)$.
2. (TBp.195) The random variable \bar{X}_n and the random vector $(X_1 - \bar{X}_n, X_2 - \bar{X}_n, \dots, X_n - \bar{X}_n)$ are independent.
3. (TBp.196) \bar{X}_n and S_n^2 are independently.
4. (TBp.197) The distribution of $(n-1)S_n^2/\sigma^2$ is the chi-square distribution with $n-1$ degrees of freedom.
5. (TBp.198)
$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n} \sim t_{n-1}.$$

Proof of 2. The joint mgf of \bar{X}_n and $(X_1 - \bar{X}_n, X_2 - \bar{X}_n, \dots, X_n - \bar{X}_n)$ is

$$M(s, t_1, t_2, \dots, t_n) = E \left\{ e^{[s\bar{X}_n + \sum_{i=1}^n t_i(X_i - \bar{X}_n)]} \right\} = E \left\{ e^{[\sum_{i=1}^n (\frac{s}{n} + t_i - \bar{t})X_i]} \right\}.$$

Let $a_i = \frac{s}{n} + t_i - \bar{t}$, $i = 1, 2, \dots, n$. Then

$$\sum_{i=1}^n a_i = s, \quad \sum_{i=1}^n a_i^2 = \frac{s^2}{n} + \sum_{i=1}^n (t_i - \bar{t})^2.$$

Now we have

$$\begin{aligned} M(s, t_1, t_2, \dots, t_n) &= \prod_{i=1}^n M_{X_i}(a_i) = \prod_{i=1}^n \exp \left(\mu a_i + \frac{\sigma^2}{2} a_i^2 \right) \\ &= \exp \left(\mu \sum_{i=1}^n a_i + \frac{\sigma^2}{2} \sum_{i=1}^n a_i^2 \right) = \exp \left[\mu s + \frac{\sigma^2}{2} \frac{s^2}{n} + \frac{\sigma^2}{2} \sum_{i=1}^n (t_i - \bar{t})^2 \right] \\ &= \exp \left(\mu s + \frac{\sigma^2}{2n} s^2 \right) \exp \left[\frac{\sigma^2}{2} \sum_{i=1}^n (t_i - \bar{t})^2 \right] \end{aligned}$$

Thus, the joint mgf factorizes into product of the mgf of \bar{X}_n and the mgf of $(X_1 - \bar{X}_n, X_2 - \bar{X}_n, \dots, X_n - \bar{X}_n)$.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Question 4.4: Are $(X_1 - \bar{X}_n, \dots, X_n - \bar{X}_n)$ independent? [Hint: $\sum_{i=1}^n (X_i - \bar{X}_n) = 0$]

Proof of 4. First note that

$$\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 = \sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2 \sim \chi_n^2.$$

Also,

$$\underbrace{\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu)^2}_W = \underbrace{\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X}_n)^2}_U + \underbrace{\left(\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \right)^2}_V.$$

Since V and U are independent,

$$M_W(t) = M_U(t) M_V(t)$$

$$+ \frac{2}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X}_n)(\bar{X}_n - \mu) = 0$$

and then

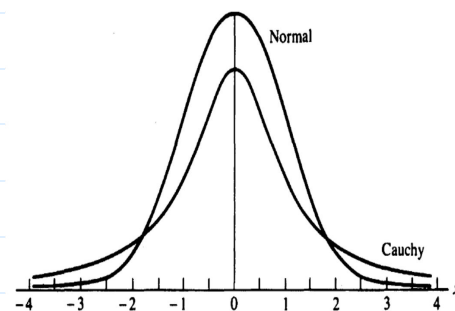
$$M_U(t) = \frac{M_W(t)}{M_V(t)} = \frac{(1 - 2t)^{-\frac{n}{2}}}{(1 - 2t)^{-\frac{1}{2}}} = (1 - 2t)^{-\frac{n-1}{2}},$$

which is the mgf of a χ_{n-1}^2 distribution. Thus $\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$

Question 4.5: Why degree of freedom = $n - 1$, rather than n ?

Definition 4.17 (Cauchy distribution $C(\mu, \sigma)$, TBp.95)

- pdf: $f(x) = \frac{\sigma}{\pi} \frac{1}{\sigma^2 + (x - \mu)^2}$, $x \in \mathbb{R}$.
- cdf: $\frac{1}{2} + \frac{1}{\pi} \tan^{-1}(\mu + \sigma x)$, $x \in \mathbb{R}$.
- mgf: not exist, except at $t = 0$
- chf: $e^{i\mu t - \sigma|t|}$
- mean: not exist
- variance: not exist
- parameter: $\mu \in \mathbb{R}, \sigma > 0$

**Notes:**

1. a heavy tail distribution
2. $C(0, 1) = t_1$
3. Let X, Y be i.i.d. $\sim N(0, 1)$. Then, $X/Y \sim C(0, 1)$.
4. $X \sim C(\mu, \sigma) \Rightarrow$ for $a, b \in \mathbb{R}$, $aX + b \sim C(a\mu + b, |a|\sigma)$
5. Let X_1, \dots, X_k be independent, and $X_i \sim C(\mu_i, \sigma_i)$. Then $Y = X_1 + \dots + X_k \sim C(\sum_{i=1}^k \mu_i, \sum_{i=1}^k \sigma_i)$
6. by 4 and 5, let X_1, \dots, X_k be i.i.d. $\sim C(\mu, \sigma)$, then $\bar{X}_k \sim C(\mu, \sigma)$.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Some other distributions

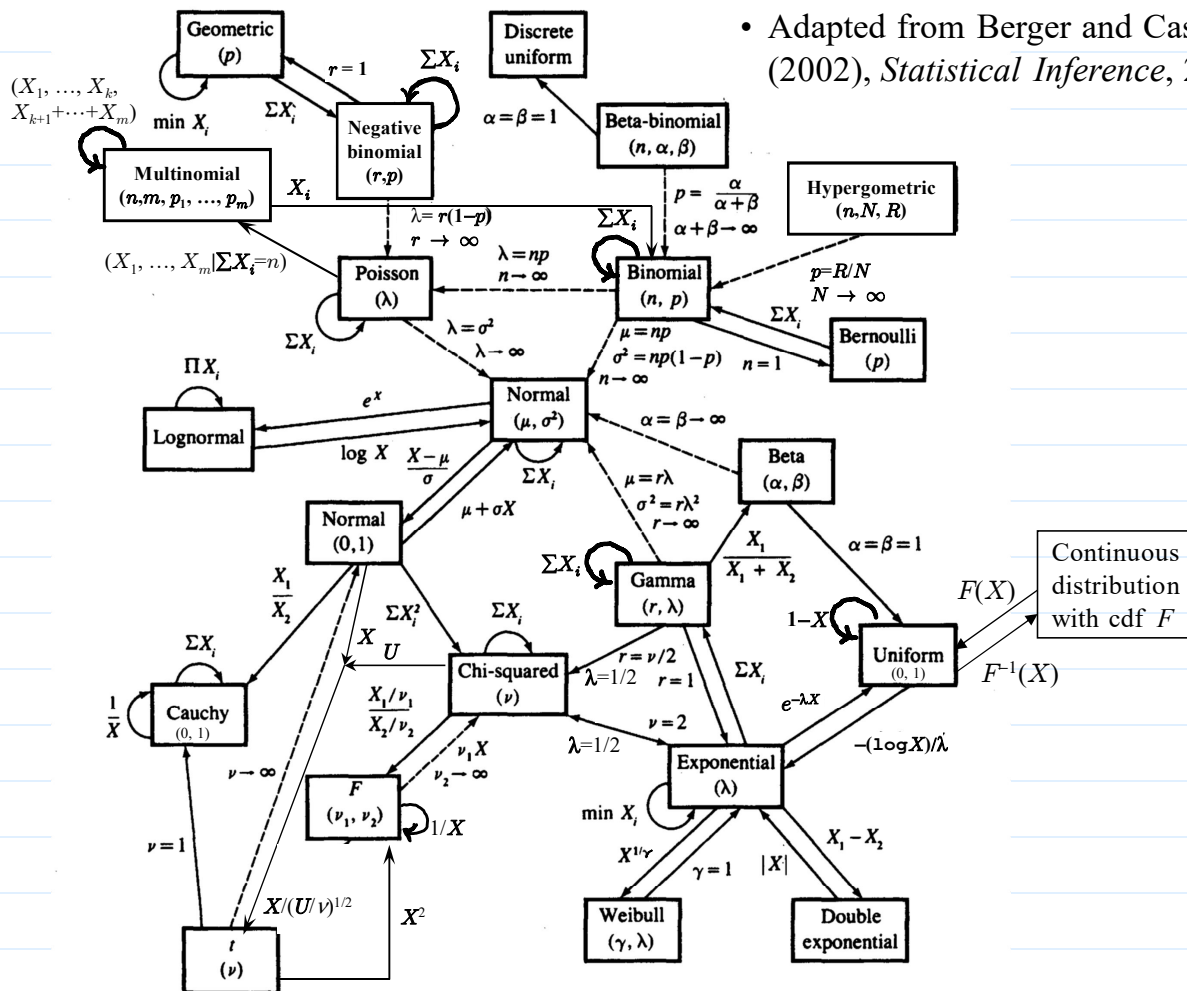
Log-normal (TBp. 69)
Weibull (TBp. 69)
Double exponential (TBp. 111)
Logistic
Pareto (TBp. 323)
Maxwell (TBp. 121)

In this chapter, you should learn

1. random phenomenon behind each distribution
2. statistical modeling (assigning a distribution) of data
3. relationship among distributions
4. meaning of parameters in each distribution
5. HOW to derive cdf/mgf/chf/mean/variance from pdf/pmf (optional)

but you are not necessary to

1. memorize their pdf/pmf/cdf/mgf/chf/mean/variance/...



Relationships among common distributions. Solid lines represent transformations and special cases, dashed lines represent limits. Adapted from Leemis (1986).

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Chapter 5

Outline

- 3 types of convergence
 - a.s., in prob., in dist.
- law of large number
- central limit theorem

Question 5.1

1. repeat flipping a fair coin **2 or 3 times**. Can you accurately predict the average appearance of heads?
2. repeat flipping a fair coin **many many times**. What will you predict the average appearance of heads?

Note.

1. Some deterministic patterns emerge from random phenomena when more and more data are collected, i.e., more and more information is gathered.
2. In the following, $n \rightarrow \infty$ can be interpreted as sample size of data is large enough.

Definition 5.1 (converge almost surely, TBp. 178)

A sequence of random variables $\{Z_n : \Omega \rightarrow \mathbb{R}\}$ is said to **converge almost surely** to a random variable $Z : \Omega \rightarrow \mathbb{R}$, and denoted as $Z_n \xrightarrow{\text{a.s.}} Z$, if for any $\epsilon > 0$,

$$P\left(\left\{\underline{\omega} \in \Omega : \lim_{n \rightarrow \infty} |Z_n(\omega) - Z(\omega)| < \epsilon\right\}\right) = \underline{1}.$$

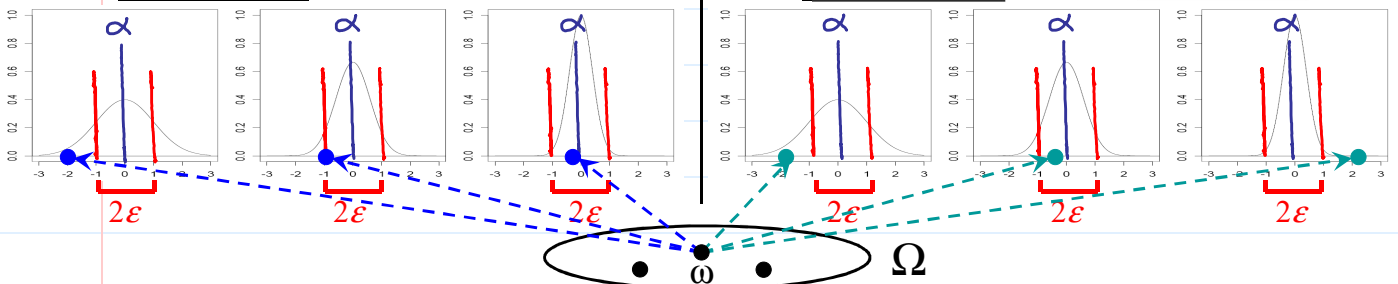
Definition 5.2 (converge in probability, TBp. 178)

A sequence of random variables $\{Z_n : \Omega \rightarrow \mathbb{R}\}$ is said to **converge in probability** to a random variable $Z : \Omega \rightarrow \mathbb{R}$, and denoted as $Z_n \xrightarrow{P} Z$, if for any $\epsilon > 0$,

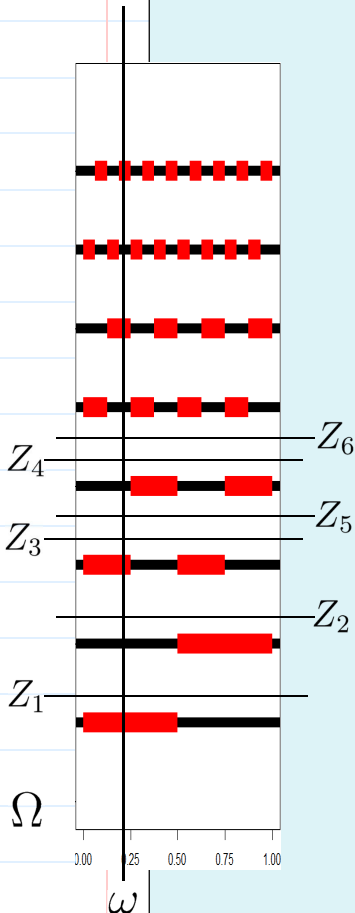
$$\lim_{n \rightarrow \infty} P(\{\omega \in \Omega : |Z_n(\omega) - Z(\omega)| < \epsilon\}) = \underline{1}.$$

$Z_n \xrightarrow{\text{a.s.}} \alpha$, α : a constant.

$Z_n \xrightarrow{P} \alpha$, α : a constant.



NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Example 5.1 (convergence in probability need not imply convergence a.s.)

- $\Omega = (0, 1]$
- P : uniform probability measure on Ω
- For $k = 1, 2, \dots$, divide $(0, 1]$ into 2^k subintervals of equal length. These intervals are given by

$$I_{k,j} = \left(\frac{j-1}{2^k}, \frac{j}{2^k}\right]$$

for $j = 1, 2, \dots, 2^k$.

- Let $Z_1, Z_2, \dots : \Omega \rightarrow \mathbb{R}$ be a sequence of r.v.'s defined as follows:

$$Z_n(\omega) = \begin{cases} 1, & \text{if } \omega \in I_{k,j}, \\ 0, & \text{if } \omega \notin I_{k,j}, \end{cases}$$

where $n = 2^k + j - 2$.

- $Z : \Omega \rightarrow \mathbb{R}$ such that $Z(\omega) = 0$ for $\omega \in \Omega$

- $Z_n \xrightarrow{P} Z$, because for $0 < \epsilon < 1$,

$$P(\{\omega \in \Omega : |Z_n(\omega) - Z(\omega)| < \epsilon\}) = 1 - \frac{1}{2^k} \rightarrow \underline{1}.$$

- Z_n not converge to Z a.s. because

$$P(\{\omega \in \Omega : Z_n(\omega) \rightarrow Z(\omega)\}) = P(\emptyset) = \underline{0}.$$

Definition 5.3 (converge in distribution, TBp. 181)

Let Z_1, Z_2, \dots be a sequence of random variables with cdf's F_1, F_2, \dots and let Z be a random variable with cdf F . Then Z_n converges in distribution to Z , denoted as $Z_n \xrightarrow{d} Z$, if

$$\lim_{n \rightarrow \infty} F_n(z) = F(z)$$

at every point z where F is continuous.

Theorem 5.1 (some properties about the 3 types of convergence)

1. $Z_n \xrightarrow{\text{a.s.}} Z \Rightarrow Z_n \xrightarrow{P} Z$
2. $Z_n \xrightarrow{P} Z \Rightarrow Z_n \xrightarrow{d} Z$
3. $Z_n \xrightarrow{d} c, c: \text{a constant} \Rightarrow Z_n \xrightarrow{P} c$
4. (**convergence of transformation**) Let $g: \mathbb{R}^k \mapsto \mathbb{R}$ be a continuous function.
 - (a) $Z_n^{(j)} \xrightarrow{\text{a.s.}} Z^{(j)}, j = 1, \dots, k \Rightarrow g(Z_n^{(1)}, \dots, Z_n^{(k)}) \xrightarrow{\text{a.s.}} g(Z^{(1)}, \dots, Z^{(k)})$.
 - (b) $Z_n^{(j)} \xrightarrow{P} Z^{(j)}, j = 1, \dots, k \Rightarrow g(Z_n^{(1)}, \dots, Z_n^{(k)}) \xrightarrow{P} g(Z^{(1)}, \dots, Z^{(k)})$.
 - (c) $(Z_n^{(1)}, \dots, Z_n^{(k)}) \xrightarrow{d} (Z^{(1)}, \dots, Z^{(k)}) \Rightarrow g(Z_n^{(1)}, \dots, Z_n^{(k)}) \xrightarrow{d} g(Z^{(1)}, \dots, Z^{(k)})$.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

5. (**Slutsky's theorem**) If $X_n \xrightarrow{d} X, Y_n \xrightarrow{P} a$, where a is a constant, then

$$(a) \ Y_n X_n \xrightarrow{d} aX$$

$$(b) \ X_n + Y_n \xrightarrow{d} X + a$$

$$(c) \ \frac{X_n}{Y_n} \xrightarrow{d} \frac{X}{a}, \text{ provided that } P(Y_n \neq 0) = 1 \text{ for all } n \text{ and } a \neq 0.$$

6. (**limit theorem for δ method**) Suppose

$$\frac{\sqrt{n}(X_n - \theta)}{\sigma} \xrightarrow{d} N(0, 1).$$

For a given function g , suppose that $g'(\theta) \neq 0$ exists. Then

$$\frac{\sqrt{n}[g(X_n) - g(\theta)]}{\sigma|g'(\theta)|} \xrightarrow{d} N(0, 1).$$

Theorem 5.2 (Continuity Theorem, TBp. 181)

Let $F_n(x)$ be a sequence of cdfs with the corresponding mgfs $M_n(t)$. Let $F(x)$ be a cdf with the mgf $M(t)$. If $M_n(t) \rightarrow M(t)$ as $n \rightarrow \infty$ for all t in an open interval containing zero, then $F_n(x) \rightarrow F(x)$ at all continuity point of F .

Notes.

1. The reverse of the continuity theorem also holds.
2. The continuity theorem still holds when the moment generating function is replaced by characteristics function (chf always exists).

(for your information)

F_n, F : cdf; f_n, f : pdf; p_n, p : pmf;

Q: $F_n \xrightarrow{d} F$ implies $\lim_{n \rightarrow \infty} f_n(x) = f(x)$?

or $\lim_{n \rightarrow \infty} p_n(x) = p(x)$?

Ans: In general, **NO**.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Example 5.2 (Convergence of Poisson to Normal, TBp. 181-182)

Let $X_n \sim P(\lambda_n)$, $n = 1, 2, \dots$ with $\lambda_n \rightarrow \infty$. We know that $E(X_n) = Var(X_n) = \lambda_n$ and $M_{X_n}(t) = e^{\lambda_n(e^t - 1)}$. Let

$$Z_n = (X_n - \lambda_n) / \sqrt{\lambda_n},$$

Then $M_{Z_n}(t) = e^{-t\sqrt{\lambda_n}} M_{X_n}\left(\frac{t}{\sqrt{\lambda_n}}\right) = e^{-t\sqrt{\lambda_n}} e^{\lambda_n(e^{t/\sqrt{\lambda_n}} - 1)}$. Because

$$\lim_{n \rightarrow \infty} \log M_{Z_n}(t) = \lim_{n \rightarrow \infty} -t\sqrt{\lambda_n} + \lambda_n(e^{t/\sqrt{\lambda_n}} - 1) = \frac{t^2}{2},$$

$M_{Z_n}(t) \rightarrow e^{t^2/2}$, which is the mgf of $N(0, 1)$. By continuity theorem, $Z_n \xrightarrow{d} N(0, 1)$, i.e., when λ is large, we can approximate the distribution of $P(\lambda)$ by $N(\lambda, \lambda)$.

Theorem 5.3 (Weak Law of Large Numbers (WLLN), TBp. 178)

Let $X_1, X_2, \dots, X_n, \dots$ be a sequence of independent random variables with $E(X_i) = \mu$ and $Var(X_i) = \sigma^2$.

Let $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. Then $\bar{X}_n \xrightarrow{P} \mu$.

Proof: $E(\bar{X}_n) = \mu, \quad Var(\bar{X}_n) = \sigma^2/n$

By Chebyshev's inequality,

$$P(|\bar{X}_n - \mu| > \epsilon) \leq \frac{Var(\bar{X}_n)}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2} \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

Notes. Under the same assumptions, a **strong law of large numbers (SLLN)**, which asserts that $\bar{X}_n \xrightarrow{\text{a.s.}} \mu$, can be proved.

Example 5.3 (Monte Carlo integration, TBp. 179)

To calculate $I(f) = \int_0^1 f(x)dx$, we can generate X_1, X_2, \dots, X_n i.i.d. $\sim U(0, 1)$ and compute $\hat{I}(f) = \frac{1}{n} \sum_{i=1}^n f(X_i)$. By the LLN, $\hat{I}(f)$ will be close to $E[f(X_i)] = \int_0^1 f(x) \times 1 dx = I(f)$ as n is large.

NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Example 5.4 (Repeated Measurements, TBp. 179-180)

Let X_1, \dots, X_n be i.i.d. with mean μ and variance σ^2 , then

$$\bar{X}_n \xrightarrow{P} \mu.$$

Let

$$S_n^2 \equiv \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \left(\frac{1}{n} \sum_{i=1}^n X_i^2 \right) - \bar{X}_n^2.$$

Because $g(x) = x^2$ is continuous, $\bar{X}_n^2 \xrightarrow{P} \mu^2$. Next, the r.v.'s X_1^2, \dots, X_n^2 are i.i.d. with mean $\sigma^2 + \mu^2$. By WLLN

$$\frac{1}{n} \sum_{i=1}^n X_i^2 \xrightarrow{P} \sigma^2 + \mu^2.$$

Therefore,

$$S_n^2 \xrightarrow{P} (\sigma^2 + \mu^2) - \mu^2 = \sigma^2.$$

(Note. $\frac{1}{n}$ in S_n^2 can be replaced by $\frac{1}{n-1}$.)

Example 5.5

If $X_n \sim t_n$, then $X_n \xrightarrow{d} N(0, 1)$.

(Ec) If $X_n \sim F_{m,n}$, then $mX_n \xrightarrow{d} \chi_m^2$ as $n \rightarrow \infty$.

- Let X_1, \dots, X_n be i.i.d. \sim Exponential(1),

then $\mu_X = E(X_i) = 1$, $\sigma_X^2 = \text{Var}(X_i) = 1$, for $i = 1, \dots, n$, and

$$n\bar{X}_n = \sum_{i=1}^n X_i \sim \text{Gamma}(n, 1) \Rightarrow \bar{X}_n \sim \frac{1}{n} \text{Gamma}(n, 1).$$

$$\Rightarrow E(\bar{X}_n) = \mu_X = 1, \text{Var}(\bar{X}_n) = \sigma_X^2/n = 1/n$$

LLN

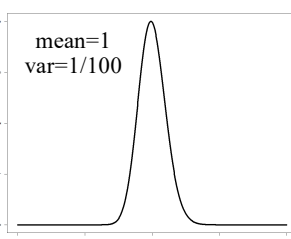
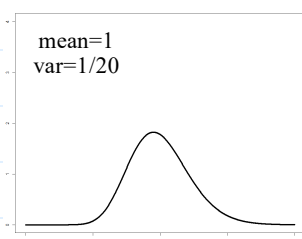
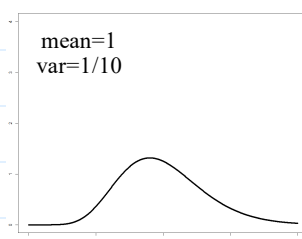
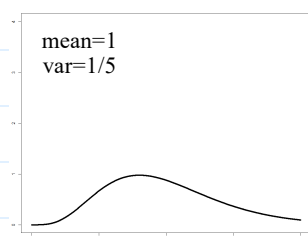
$n=5$

$n=10$

$n=20$

$n=100$

pdf
of
 \bar{X}_n

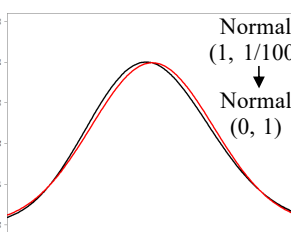
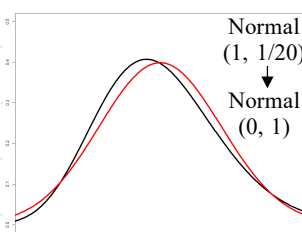
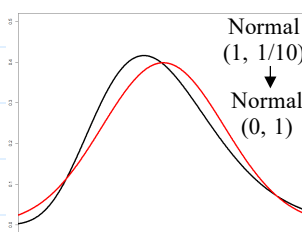
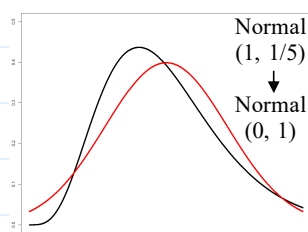


$$\frac{\bar{X}_n - \mu_X}{\sigma_X / \sqrt{n}}$$

CLT

pdf
of

$$\frac{\bar{X}_n - \mu_X}{\sigma_X / \sqrt{n}}$$



NTHU MATH 2820, 2026, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Theorem 5.4 (Central Limit Theorem, TBp. 169)

Let X_1, X_2, \dots be i.i.d. with mean μ and variance σ^2 . Let

$$\bar{X}_n (= \frac{1}{n} \sum_{i=1}^n X_i) \quad \text{and} \quad T_n (= n\bar{X}_n = \sum_{i=1}^n X_i)$$

be the average and the sum of data, respectively. Then,

$$\lim_{n \rightarrow \infty} P\left(\frac{T_n - n\mu}{\sigma\sqrt{n}} \leq x\right) = \lim_{n \rightarrow \infty} P\left(\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \leq x\right) = \Phi(x),$$

for $-\infty < x < \infty$, where $\Phi(x)$ is the cdf of $N(0, 1)$.

Proof. Let $W_i = \frac{X_i - \mu}{\sigma}$, then $E(W_i) = 0$ and $\text{Var}(W_i) = 1$. Let

$$Z_n \equiv \frac{T_n - n\mu}{\sigma\sqrt{n}} = \frac{1}{\sqrt{n}} \sum_{i=1}^n W_i.$$

Let $M(t)$ be the mgf of W_i 's and $M_{Z_n}(t)$ be the mgf of Z_n , then

$$\begin{aligned} M_{Z_n}(t) &= \left[M\left(\frac{t}{\sqrt{n}}\right) \right]^n = \left[M(0) + M'(0)\frac{t}{\sqrt{n}} + \frac{1}{2}M''(0)\left(\frac{t}{\sqrt{n}}\right)^2 + o\left(\frac{1}{n}\right) \right]^n \\ &= \left[1 + \frac{t^2}{2n} + o\left(\frac{1}{n}\right) \right]^n \rightarrow e^{t^2/2} \quad [\text{because if } a_n \rightarrow a, (1 + \frac{a_n}{n})^n \rightarrow e^a] \end{aligned}$$

Notes.

- When mgf's do not exist, we can use chf's to prove it instead.
- This is one of the simplest versions of CLT.

Example 5.6 (Normal approximation to Binomial distribution, TBp.187)

Let X_1, X_2, \dots, X_n be i.i.d. $\sim B(1, p)$, then $T_n \sim B(n, p)$. Note that $E(X_i) = p$, $Var(X_i) = p(1-p)$ and $E(T_n) = np$, $Var(T_n) = np(1-p)$. By CLT,

$$\frac{T_n - np}{\sqrt{np(1-p)}} \xrightarrow{d} N(0, 1),$$

i.e., when n is large enough, we can approximate the distribution of $B(n, p)$ by $N(np, np(1-p))$.

- Note:** 1. how about those distributions that can be generated from a sum of some i.i.d. random variables? (example?)
 2. (cf.) Poisson in Def. 4.7 (LNp.66) and Example 5.2 (LNp.92)

Example 5.7 (measurement error (or called sampling error), TBp. 186)

- Suppose that you want to know the average income of families living in Taipei.
- If you can ask every families their incomes, you will get the exact value of the average, denoted by μ .

NTHU MATH 2820, 2026, Lecture Notes
 made by S.-W. Cheng (NTHU, Taiwan)

- However, what if you only take a random sample of, say, 1000 families?
- The average income of the 1000 families, denoted by \bar{X}_{1000} , is a random variable. It has an error $\bar{X}_{1000} - \mu$, which is called measurement error or sampling error.
- By CLT, the error will be distributed normally, and we can approximate $P(|\bar{X}_{1000} - \mu| < c)$ using normal distribution no matter what the distribution of incomes is.

Example 5.8 (experimental error)

- It is usually true that an experimental error ϵ is a function of a number of component errors $\epsilon_1, \dots, \epsilon_n$.
- for example, errors in the settings of experimental conditions, errors due to variation in raw materials, and so on.
- If each individual component error is fairly small, it is possible to approximate the overall error ϵ as a linear function of independently distributed component errors

$$\epsilon \approx a_1\epsilon_1 + \dots + a_n\epsilon_n.$$

- By CLT, the distribution of ϵ will tend to normal as the number of component errors becomes large.
- This argument also offers a good justification for why in many statistical methods, such as in ANOVA or linear regression, the error part is assumed to be distributed normally.

Example 5.9 (cont. Ex.5.4 in LNp.94, Repeated Measurements)

Let X_1, X_2, \dots, X_n be i.i.d. with mean μ and variance σ^2 . Then, by LLN, CLT, and Slutsky's theorem,

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n} \xrightarrow{d} N(0, 1),$$

$$\text{because } \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} N(0, 1) \text{ and } S_n^2 \xrightarrow{P} \sigma^2 \left(\Rightarrow \frac{\sigma^2}{S_n^2} \xrightarrow{P} 1 \right).$$

❖ **Reading:** textbook, chapter 5

❖ **Further reading:** Roussas, chapter 8