# Chapter 1

**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

*probability space*

Characteristic
* don't know what result we will get in the future
* the best we can do is to describe/calculate the probability of these possible results.

— 樂透開獎号碼
— waiting time
— rain tomorrow?
- - -

Website of My Probability Course

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

---

Textbook page ⟶          LNp. (Lecture Note page)

**Definition** (sample space, TBp. 2)

A **sample space** $\Omega$ is the set of all possible outcomes in a random phenomenon.

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

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

h : head
t : tail

— $\Omega$ is a finite set

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

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

— $\Omega$ is an infinite, but countable, set

*discrete random variable*

cf.

*continuous random variable*

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

$$\Omega = \{t \mid t \geq 0\} = [0, \infty)$$

— $\Omega$ is an infinite, but uncountable, set

**Question**

What are the differences between the $\Omega$ in these examples?

**Definition** (event, TBp. 2)

A particular <u>subset</u> of $\Omega$ is called an **event.** ← collection of all "well-defined" events ⇒ $\sigma$-field.
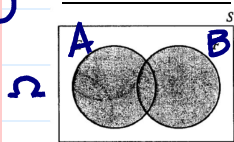
**Example 1.4** (cont. Ex. 1.1)

Let $A$ be the event that total <u>number of heads equals 2</u>, then $A = \{hht, hth, thh\}$.
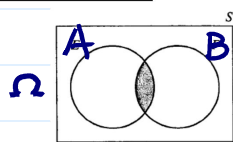
**Example 1.5** (cont. Ex. 1.2)

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

- **<u>union.</u>** $C = \underline{A \cup B} \Rightarrow C$: <u>at least one</u> of $A$ and $B$ <u>occur</u>.
- **<u>intersection.</u>** $C = \underline{A \cap B} \Rightarrow C$: <u>both</u> $A$ and $B$ <u>occur</u>.
- **<u>complement.</u>** $C = \underline{A^c} \Rightarrow C$: $A$ does <u>not occur</u>.
- mutually exclusive ⊙ **<u>disjoint.</u>** $\underline{A \cap B = \emptyset} \Rightarrow A$ and $B$ have <u>no outcomes in common</u>.
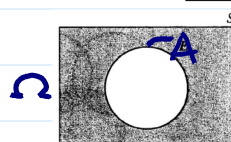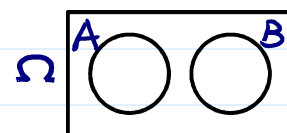
(a) Shaded region: $E \cup F$.    (b) Shaded region: $EF$.    (c) Shaded region: $E^c$.

**Definition** (probability measure, TBp. 4)

A **probability measure** on $\Omega$ is a function $P$ from <u>subsets of</u> $\underline{\Omega}$ to the <u>real numbers</u> that satisfies the following axioms:

1. $\underline{P(\Omega) = 1}$. ← total prob. = 1
2. If $A \subset \Omega$, then $\underline{P(A) \geq 0}$. ← non-negativity
3. If $A_1$ and $A_2$ are <u>disjoint</u>, then ← additivity

$$P(\underline{A_1 \cup A_2}) = \underline{P(A_1) + P(A_2)}.$$

More generally, if $A_1$, $A_2$, ... are <u>mutually disjoint</u>, then

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

Axioms of probability

$\mathcal{F}$

$P: \mathcal{F} \rightarrow [0, 1]$

**Example 1.6** (cont. Ex. 1.1)

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

$\Omega = \{hhh, hht, hth, thh, htt, tht, tth, ttt\}$
     1/8   1/8   1/8   1/8   1/8   1/8   1/8   1/8   $P: \Omega \rightarrow [0,1]$

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

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

**Property C.** If $A \subset B$, then $P(A) \leq P(B)$.
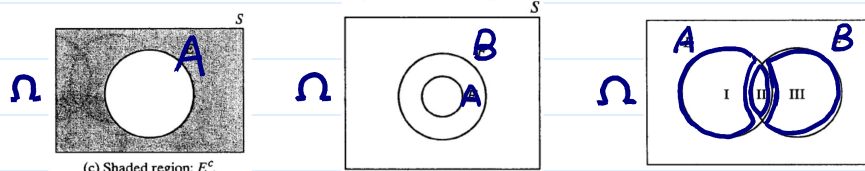
**Property D.** $P(A \cup B) = P(A) + P(B) - P(A \cap B)$.

*generalization:*
$P(A_1 \cup \cdots \cup A_n)$
$\quad = \sum P(A_i)$
$\quad - \sum P(A_i \cap A_j)$
$\quad + \sum P(A_i \cap A_j \cap A_k)$
$\quad - \cdots$
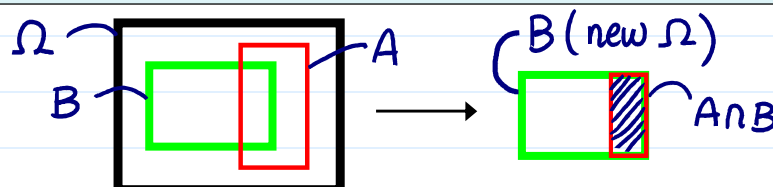


(c) Shaded region: $E^c$.

**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

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

Q: Why cond. prob. important in statistics?

Ans: update information.

**Example 1.7** (cont. Ex. 1.6)

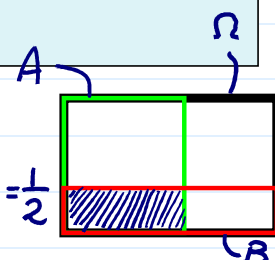Suppose that the first throw is $h$. What is the probability that we can get exact two $h$'s in the three trials?

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

$P(A|B) = \dfrac{P(A \cap B)}{P(B)} = \dfrac{2/8}{4/8} = \dfrac{1}{2}$



**Theorem** (Multiplication Law, TBp. 17)

Let $A$ and $B$ be events and assume $P(B) > 0$. Then
$$P(A \cap B) = P(A|B)P(B).$$ ← intuition

generalization
$P(A_1 \cap A_2 \cap \cdots \cap A_n)$
$= P(A_1) \cdot P(A_2|A_1) \cdot P(A_3|A_1 \cap A_2) \cdot \cdots$

→ Sometimes, this is easier to obtain ($\because \Omega \to B$)

**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 $P(B) = 0.2$. The probability that it is cloudy and raining is
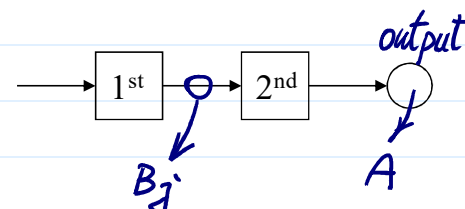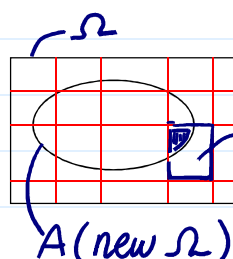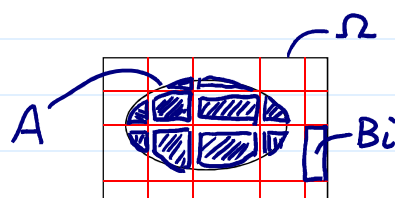$P(A \cap B) = P(A|B)P(B) = 0.3 \times 0.2 = 0.06$.

**Theorem** (Law of Total Probability, TBp. 18)

Let $B_1, B_2, \ldots, B_n$ 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 $A$,

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

平均 ← 權重 ← intuition

$P(A \cap B_i)$

a partition of $\Omega$



$A$ — $B_i$　　$\Omega$　　$B_j$　　$A$ (new $\Omega$)　　1st　2nd　output　$B_j$　$A$

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

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

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

update

**Definition** (independence, TBp. 24)

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

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

獨立

definition of independence

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

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

cf.　generalization of multiplication Law in LNp. 6

cf.

When $A$ and $B$ are independent,
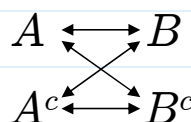
intuition of independence

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

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

$A \longleftrightarrow B$
$A^c \longleftrightarrow B^c$

independence & complement

required
optional

⊗ **Reading**: textbook, Sections 1.1, 1.2, 1.3, 1.5, 1.6, 1.7
⊗ **Further Reading**: Roussas, Chapters 1 and 2
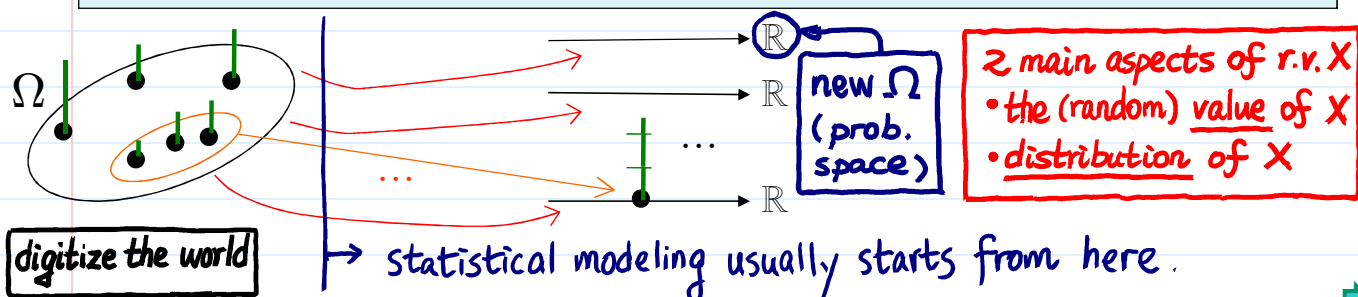
# Chapters 2 and 3

## Outline

- random variables (隨機變數)
- 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 $\Omega$ to the real numbers.



$\Omega$

digitize the world

R
R
...
R

new $\Omega$ (prob. space)

2 main aspects of r.v. X
- the (random) value of X
- distribution of X

$\rightarrow$ statistical modeling usually starts from here.

**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

update probability space

$$1/8 \quad 1/8 \quad 1/8 \quad 1/8 \quad 1/8 \quad 1/8 \quad 1/8 \quad 1/8$$
$$\Omega = \{hhh, hht, hth, thh, htt, tht, tth, ttt\}$$

$X_1 :$  3,  2,  2,  2,  1,  1,  1,  0.  $\leftarrow$ new $\Omega$
$\frac{1}{8}$  $\frac{3}{8}$  $\frac{3}{8}$  $\frac{1}{8}$  $\leftarrow$ new probability measure

$X_2 :$  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?

We need random variable because

Data $\rightarrow$ in $\mathbb{R}^n$ space

Uncertainty $\rightarrow$ need probability measure

$\rightarrow$ can do "+", "-", "×", "÷", exp, log, ---

- distribution ← *probability measure of r.v.* ⌐*don't know what value will appear*
  (分配,分布)        *• For r.v., it value : random, but its distribution : fixed*

> **Question 2.2**
>
> A random variable have a sample space on **real line**. Does it bring some special ways to characterize its probability measure?

*⌐ finite or countable infinity*

|  | discrete | continuous ← *uncountable* |
|---|---|---|
| uni-variate r.v. | • pmf ← → • cdf • mgf/chf | • pdf • cdf • mgf/chf |
| multi-variate r.v.'s | • joint pmf ←→ • joint cdf • joint mgf/chf | • joint pdf • joint cdf • joint mgf/chf |

*one r.v.* ←
*at least two r.v.*

*when any of them is known, the other 2 can be obtained*

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

> **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. ← *uncountable*

*e.g.*

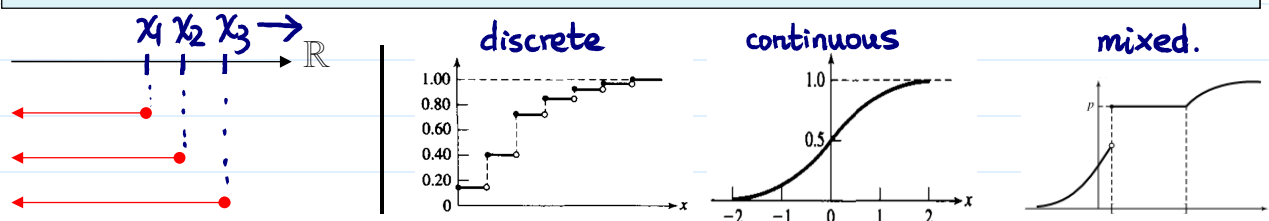| Discrete | Continuous |
|---|---|
| $X \in \{0,1,2,3\}$ | $X \in [0,1]$ |
| $X \in \mathbb{Z}_+$ | $X \in (-\infty, \infty)$ |

> **Definition 2.3** (cumulative distribution function, TBp. 36)
>
> A function $F$ is called the **cumulative distribution function (cdf)** of a random variable $X$ if
>
> $$F(x) = P(X \leq x), \ x \in \mathbb{R}.$$

$x_1 \ x_2 \ x_3 \rightarrow \mathbb{R}$

discrete    continuous    mixed.

**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) $\sum_{x \in \mathcal{X}} p(x) = 1$.

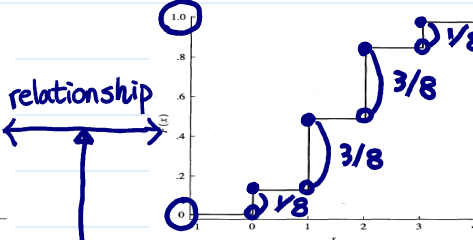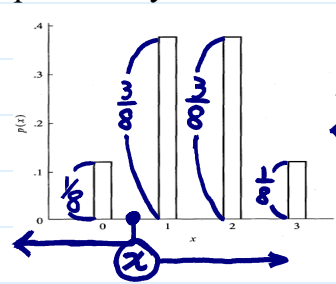For a discrete random variable $X$ with pmf $p(x)$,
$$P(X = x) = p(x),$$
and
$$P(X \in A) = \sum_{x \in A} p(x).$$

$\mathcal{X}$: a finite or countably infinite set.

$\rightarrow A \cap \mathcal{X}$

probability mass function          cumulative distribution function



relationship

$\frac{3}{8}$   $\frac{3}{8}$

$\frac{1}{8}$   $\frac{1}{8}$

$x$

$\frac{1}{8}$
$\frac{3}{8}$
$\frac{3}{8}$
$\frac{1}{8}$

$P(X \leq 1) = \frac{1}{8} + \frac{3}{8} = \frac{4}{8} = F(1)$

$P(X < 1) = \frac{1}{8} = F(1-)$

$P(X = 1) = P(X \leq 1) - P(X < 1)$
$\qquad = F(1) - F(1-)$

- $\boxed{\begin{aligned} F(x) &= \sum_{t \leq x} P(X = t) = \sum_{t \leq x} p(t) \\[2mm] p(x) &= P(X = x) = F(x) - F(x-) \end{aligned}}$    $= \lim_{t \uparrow x} F(t)$

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

Note. pdf plays a similar role as pmf, but
$$\Sigma \rightarrow \int$$

pdf of Uniform(0, 1)          cdf of Uniform(0.1)



area$^*$          relationship          $\} = $ area$^*$

$x$          $x$

area $= P(A)$

pdf

$A$

- $\boxed{\begin{aligned} F(x) &= \int_{-\infty}^{x} f(t)\, dt \\[2mm] f(x) &= \frac{d}{dx} F(x) \end{aligned}}$

The value of a pdf can be larger than one (c.f. pmf)

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

$f(x)$

$x - \frac{dx}{2}$   $x + \frac{dx}{2}$

**Question 2.3**

How to interpret $f(x)$?

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

proportional to prob.

**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. $\lim_{x \to \infty} F(x) = 1$ and $\lim_{x \to -\infty} F(x) = 0.$
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. $\rightarrow F(x-) = P(X < x)$
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.

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

$\hookrightarrow (X_1, X_2, \cdots, X_n) \in \mathbb{R}^n$

**Example 2.2** (cont. Ex. 2.1)

Why several marginal distributions not enough?

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

When $X_1 = 1$ occurs,
$$P(X_2 = 0 \mid X_1 = 1) = \frac{2/8}{3/8} = \frac{2}{3}$$
$$P(X_2 = 1 \mid X_1 = 1) = \frac{1/8}{3/8} = \frac{1}{3}$$

$(X_1, X_2) \in \mathbb{R}^2$

$X_2$: # of head on 1st toss    $X_1$: total # of heads

| | $0(1/8)$ | $1(3/8)$ | $2(3/8)$ | $3(1/8)$ |
|---|---|---|---|---|
| $(1/2)$ $0$ | $\frac{1}{8}\left(\frac{1}{16}\right)$ | $\frac{2}{8}\left(\frac{3}{16}\right)$ | $\frac{1}{8}\left(\frac{3}{16}\right)$ | $0\left(\frac{1}{16}\right)$ |
| $(1/2)$ $1$ | $0\left(\frac{1}{16}\right)$ | $\frac{1}{8}\left(\frac{3}{16}\right)$ | $\frac{2}{8}\left(\frac{3}{16}\right)$ | $\frac{1}{8}\left(\frac{1}{16}\right)$ |

marginal distribution
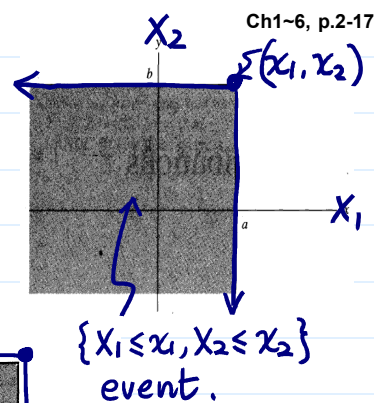
joint distribution

Note: two marginal distributions are not enough to describe their joint distribution.

**Question 2.5**

When we know the joint distribution, we can obtain every marginal distributions. Is the reverse statement true?
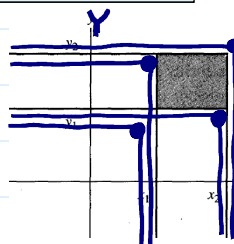
**Definition 2.6** (joint cumulative distribution function, TBp. 71)

The **joint cdf** of $X_1, X_2, \ldots, X_n$ is

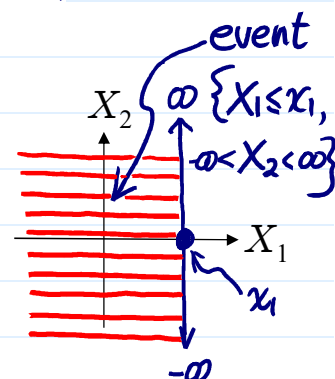$$F(x_1, x_2, \cdots, x_n) = P(X_1 \le x_1, X_2 \le x_2, \ldots, X_n \le x_n)$$

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

can be generalized to more than 2 r.v.'s

$$P(x_1 < X \le x_2, y_1 < Y \le y_2)$$
$$= F(x_2, y_2) - F(x_2, y_1)$$
$$\quad - F(x_1, y_2) + F(x_1, y_1)$$

$\{X_1 \le x_1, X_2 \le x_2\}$ event.

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

The **marginal cdf** of $X_1$ is

$$F_{X_1}(x_1) = P(X_1 \le x_1) = \lim_{x_2, x_3, \ldots, x_n \to \infty} F(x_1, x_2, \cdots, x_n)$$

event
$\{X_1 \le x_1, -\infty < X_2 < \infty\}$

- discrete case: marginal pmf  $p_{X_1}(x) = F_{X_1}(x) - F_{X_1}(x-)$.

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

○ discrete multivariate case

cf. the similarity between pmf & pdf

$$p(x_1, x_2, \cdots, x_n) = P(X_1 = x_1, X_2 = x_2, \ldots, X_n = x_n)$$
$$\Rightarrow \textbf{joint pmf} \text{ of } X_1, X_2, \ldots, X_n$$

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

$$F(x_1, x_2, \cdots, x_n) = \sum_{t_1 \le x_1, \, t_2 \le x_2, \, \ldots, \, t_n \le x_n} p(t_1, t_2, \ldots, t_n)$$

relationship b/w joint cdf & pmf

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

relationship b/w marginal & joint pmfs

○ continuous multivariate case

$$f(x_1, x_2, \cdots, x_n) = \frac{\partial^n}{\partial x_1 \cdots \partial x_n} F(x_1, x_2, \cdots, x_n)$$
$$\Rightarrow \textbf{joint pdf} \text{ of } X_1, X_2, \ldots, X_n$$

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

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

relationship b/w joint cdf & pdf

relationship b/w marginal & joint pdfs

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

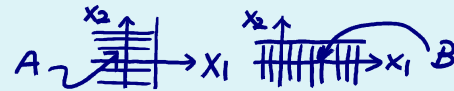• **independent random variables** ← Recall. independent events (LNp.8)

---

**Definition 2.8** (independent random variables, TBp. 84)

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

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

for all $x_1, x_2, \ldots, x_n$.

$A \cap B$ ⟋

$A$ →$X_1$　　$B$ →$X_1$

$$(\Rightarrow)\ f = \frac{\partial^n}{\partial x_1 \cdots \partial x_n} F = \frac{\partial^n}{\partial x_1 \cdots \partial x_n} F_{X_1} \cdots F_{X_n} = f_{X_1} \cdots f_{X_n}$$

$$(\Leftarrow)\ F = \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_n} f = \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_n} f_{X_1} \cdots f_{X_n} = F_{X_1} \cdots F_{X_n}$$

joint can be determined by marginals

---

**Theorem 2.2 (TBp. 85-86)**

1. For continuous case,

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

For discrete case,

Note: similarity between pdf & pmf.

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

---

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

For any A & B

for interpretation → $P(Y \in B | X \in A) = P(Y \in B)$

↳ No matter what data X occurs, it has no impact on the appearance probability of data Y.

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

indep. & transformation

intuition

e.g. $X$: 生日　—— indep. ——　$Y$: 身高

$g(X)$ 星座　—— indep. ——　$R(Y)$: 免費車票

generalization

$X_1, \ldots, X_n$ are independent
$1 < i_0 < i_1 < \cdots < i_k = n$
$Y_1 = g_1(X_1, \ldots, X_{i_1})$,
$Y_2 = g_2(X_{i_1+1}, \ldots, X_{i_2})$,
$\ldots$
$Y_k = g_k(X_{i_{k-1}+1}, \ldots, X_{i_k})$.
$Y_1, \ldots, Y_k$ are independent

**✸✸4.** marginal distributions of $\dot{X}_1, X_2, \ldots, X_n +$ independence $\Rightarrow$ joint distribution of $X_1, X_2, \ldots, X_n$

---

• **conditional distribution** ← conditional probability (LNp.5)

---

**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

$$p_{Y|X}(y|x) \equiv P(\underbrace{Y=y}_{\substack{\| \\ \text{event} \\ A}}|\underbrace{X=x}_{\substack{\| \\ \text{event} \\ B}}) = \frac{P(X=x, Y=y)}{P(X=x)} \quad \overset{A \cap B}{\underset{B}{}}$$

$$= \frac{p_{XY}(x, y)}{p_X(x)} = \frac{\text{joint}}{\text{marginal}}$$

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



$P_{x,y}$　$P_x(x)$: marginal

$\Sigma P_i = 1 (?)$

Ans: No.

LNp.16

---

**Example 2.3** (cont. Ex 2.2)

$p_{X_2|X_1}(0|1) = 2/3$, and $p_{X_2|X_1}(1|1) = 1/3$ ← update $\begin{cases} P_{X_2}(0) = 1/2 \\ P_{X_2}(1) = 1/2 \end{cases}$

---

**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

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

Notice the similarity between pmf & pdf.

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



$f_{XY}$　→ area = $f_X(x)$

$f_X(x)$

$= \int_{\mathbb{R}} f_{XY}(x, y) dy = 1 (?)$

---

**Theorem 2.3**

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

$$P(a \le Y \le b, x-\tfrac{\Delta x}{2} \le X \le x+\tfrac{\Delta x}{2}) \Big/ P(x-\tfrac{\Delta x}{2} \le X \le x+\tfrac{\Delta x}{2})$$

$$\|$$

$$P(a \le Y \le b | x-\Delta x/2 \le X \le x+\Delta x/2) = \frac{\int_a^b \int_{x-\Delta x/2}^{x+\Delta x/2} f_{XY}(u, v) du dv}{\int_{x-\Delta x/2}^{x+\Delta x/2} f_X(t) dt}$$

$$\approx \frac{\int_a^b f_{XY}(x, y) \Delta x \, dy}{f_X(x) \Delta x} = \int_a^b \frac{f_{XY}(x, y)}{f_X(x)} dy$$

---

2. For each <u>fixed</u> $x$, $p_{Y|X}(y|x)$ is a <u>pmf</u> for $y$ and $f_{Y|X}(y|x)$ is a <u>pdf</u> for $y$. ← *Notice the different roles of $x$ & $y$*

③ $\underline{p_{XY}(x,y) = p_{Y|X}(y|x)\,p_X(x)}$, and $\underline{f_{XY}(x,y) = f_{Y|X}(y|x)\,f_X(x)}$

  — <u>multiplication law</u> ← *cf. LNp.6*

④ $\underline{p_Y(y) = \sum_x p_{Y|X}(y|x)p_X(x)}$, and $\underline{f_Y(y) = \int_{-\infty}^{\infty} f_{Y|X}(y|x)f_X(x)dx}$

  — <u>law of total probability</u> ← *cf. LNp.7*

*update*
⑤ $\underline{p_{X|Y}(x|y) = \dfrac{p_{Y|X}(y|x)p_X(x)}{\sum_x p_{Y|X}(y|x)p_X(x)}}$, and *update* $\underline{f_{X|Y}(x|y) = \dfrac{f_{Y|X}(y|x)f_X(x)}{\int_{-\infty}^{\infty} f_{Y|X}(y|x)f_X(x)dx}}$,

*items 3,4,5 can be generalized to more than 2 r.v.'s*

*LNp.7 cf.* → — <u>Bayes' rule</u>    *intuition (graphs in LNp 21 & 22)*

6. $X, Y$ are <u>independent</u> ⇔ $\underline{p_{Y|X}(y|x) = p_Y(y)}$ or $\underline{f_{Y|X}(y|x) = f_Y(y)}$

---

• **functions of random variables**

*Raw Data*                   <u>Transformations</u>                    *r.v.*
                             $g_1(X_1, \ldots, X_n) = Y_1$ ←
$\left(\begin{array}{c} X_1, \\ \ldots, \\ X_n \end{array}\right)$    $\ldots,$
                             $g_k(X_1, \ldots, X_n) = Y_k$ ←
                             ⋮
                                                          Θ
           <u>Extract Information</u>    *unknown parameters in the statistical model*

**Question 2.6**

For given r.v.'s $X_1, \ldots, X_n$, how to derive the <u>distributions</u> of <u>their</u> <u>transformations</u>?

---

1. **<u>method of events</u>** → *discrete r.v.'s (pmf)*

**Theorem 2.7**

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

**Example 2.4** (<u>univariate</u> <u>discrete</u> random variable)

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

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

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

**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, $Z=X-Y$) ← Ans. $p_Z(z) = \sum_y p(z+y, y)$

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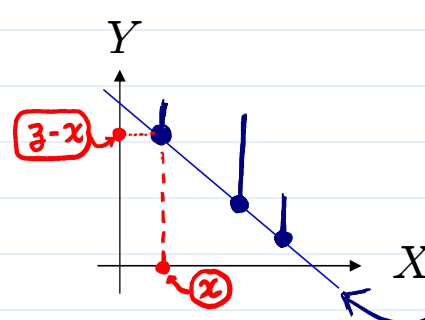
When $X$, $Y$ independent, $p(x,y) = p_X(x)p_Y(y)$,

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

cf. $\begin{cases} \text{value of r.v.} \\ \text{distribution of r.v.} \end{cases}$

$X+Y=z \Rightarrow Y=z-X$



---

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

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

1. Find the region $Y \leq y$ in the $(x_1, x_2, \ldots, x_n)$ space.
   — $A_y$ —

2. Find $F_Y(y) = P(Y \leq y)$ by summing the joint pmf or integrating the joint pdf of $X_1, X_2, \ldots, X_n$ 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 $\mathbf{Y} = (Y_1, Y_2, \ldots, Y_m)$.

$$P(A_y) = \begin{cases} \sum_{x \in A_y} P(x) \\ \int_{A_y} f(x)\, dx \end{cases}$$

$$F_Y(y_1, \cdots, y_m) = P(Y_1 \leq y_1, \cdots, Y_m \leq y_m)$$
$$= P(\underline{X} \in A_{y_1, \cdots, y_m})$$
$$f_Y(y_1, \cdots, y_m) = \frac{\partial^n}{\partial y_1 \cdots \partial y_m} F_Y(y_1, \cdots y_m)$$

**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 $Y = X^2$. $\quad\hookrightarrow X$ is a continuous r.v.

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

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

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

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

**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 $Z = X + Y$. $\quad\hookrightarrow X, Y$: continuous r.v.'s

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

Ans. $f_Z(z) = \int_{-\infty}^{\infty} f(z+y, y)\,dy$

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

$$\begin{aligned}
F_Z(z) &= P(Z \leq z) = P(X+Y \leq z) = \int\int_{R_z} f(x,y)\,dx\,dy \\
&= \int_{-\infty}^{\infty}\int_{-\infty}^{z-x} f(x,y)\,dy\,dx \\
&= \int_{-\infty}^{z}\int_{-\infty}^{\infty} f(x, v-x)\,dx\,dv \quad (\text{set } y = v - x) \\
& \hspace{8cm} \{x = x
\end{aligned}$$

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

$x+y=3$

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

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x)f_Y(z-x)\,dx \quad \Rightarrow \textbf{convolution} \text{ of } f_X \text{ and } f_Y$$
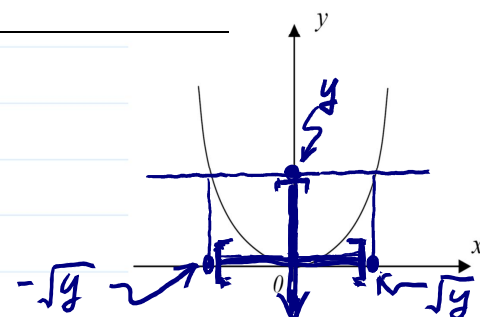
cf. ▸ the convolution for discrete r.v.'s ( LNp. 25)

**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, $Z=XY$)

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

P(Z ≤ z)
‖

$$F_Z(z) = \int\int_{Q_z} f(x,y)dxdy = \int_{-\infty}^0 \int_{xz}^\infty + \int_0^\infty \int_{-\infty}^{xz} f(x,y)dydx$$

$$= \int_{-\infty}^0 \int_z^{-\infty} + \int_0^\infty \int_{-\infty}^z xf(x,xv)dvdx \quad (\text{set } \begin{cases} y = xv \\ x = x \end{cases})$$

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

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

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

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

y = zx

**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 $[0,1]$.

→ ① no jump ② strictly increasing ⇒ X : a continuous r.v.

**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 $F$.

pdf of uniform distribution with α=0 and β=1

F(x)

• larger slope ⇒ more $X_i$'s
• smaller slope ⇒ fewer $X_i$'s

slope = pdf

an $X_i$ or $F^{-1}(U_i)$

an $F(x_i)$ or $U_i$

**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) : [check its proof in textbook]

> **Theorem 2.6** (univariate continuous case, TBp. 62)
>
> Let $X$ be a continuous random variable with pdf $f_X(x)$. Let $Y = g(X)$, where $g$ is differentiable, strictly monotone. Then,
>
> [can be relaxed to **piecewise strictly monotone**]
>
> $$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{dg^{-1}(y)}{dy} \right|$$
>
> [cf. Example 2.4 in LNp 24  Q: What's the role of the term?]
>
> for $y$ s.t. $y = g(x)$ for some $x$, and $f_Y(y) = 0$ otherwise.

> **Example 2.9**
>
> $X$ is a random variables with pdf $f_X(x)$. Find the distributon of $Y = 1/X$.

For $x > 0$ (or $x < 0$),
$$y = 1/x \equiv g(x) \quad \Rightarrow \quad x = g^{-1}(y) = 1/y$$
$$dg^{-1}/dy = -1/y^2 \quad \text{and} \quad \left| dg^{-1}/dy \right| = 1/y^2$$

hence
$$f_Y(y) = f_X(1/y)(1/y^2)$$

> **Theorem 2.7** (multivariate continuous case, TBp. 102-103)
>
> $\mathbf{X} = (X_1, X_2, \ldots, X_n)$ multivariate continuous, $\mathbf{Y} = (Y_1, Y_2, \ldots, Y_n) \equiv \mathbf{g}(\mathbf{X})$. $\mathbf{g}$ is one-to-one, so that its inverse exists and is denoted by ①
>
> $$\mathbf{x} = \mathbf{g}^{-1}(\mathbf{y}) = \mathbf{w}(\mathbf{y}) = (\underset{X_1}{w_1(\mathbf{y})}, \underset{X_2}{w_2(\mathbf{y})}, \ldots, \underset{X_n}{w_n(\mathbf{y})}).$$
>
> Assume $\mathbf{w}$ have continuous partial derivatives, and let ②
>
> $$J \underset{\text{Jacobin}}{=} \begin{vmatrix} \frac{\partial w_1(\mathbf{y})}{\partial y_1} & \frac{\partial w_1(\mathbf{y})}{\partial y_2} & \cdots & \frac{\partial w_1(\mathbf{y})}{\partial y_n} \\ \frac{\partial w_2(\mathbf{y})}{\partial y_1} & \frac{\partial w_2(\mathbf{y})}{\partial y_2} & \cdots & \frac{\partial w_2(\mathbf{y})}{\partial y_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial w_n(\mathbf{y})}{\partial y_1} & \frac{\partial w_n(\mathbf{y})}{\partial y_2} & \cdots & \frac{\partial w_n(\mathbf{y})}{\partial y_n} \end{vmatrix}$$
>
> [— determinant —]   [interpretation: similar to $\left| \frac{dg^{-1}}{dy} \right|$]
>
> Then   [absolute value]
>
> $$f_\mathbf{Y}(\mathbf{y}) = f_\mathbf{X}(\mathbf{g}^{-1}(\mathbf{y}))|J|.$$
>
> for $\mathbf{y}$ s.t. $\mathbf{y} = \mathbf{g}(\mathbf{x})$ for some $\mathbf{x}$, and $f_\mathbf{Y}(\mathbf{y}) = 0$, otherwise.

**Note.** When the dimensionality of $\mathbf{Y}$, denoted by $k$, is less than $n$, we can choose another $n - k$ transformations $\mathbf{Z}$ such that $(\mathbf{Y}, \mathbf{Z})$ satisfy the above assumptions. By integrating out the last $n-k$ arguments in the pdf of $(\mathbf{Y}, \mathbf{Z})$, the pdf of $\mathbf{Y}$ 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: $Y_1 = X_1 X_2$)

Let $Y_2 = X_1$. Then

$$x_1 = \quad y_2 \quad \equiv w_1(y_1, y_2)$$
$$x_2 = \quad y_1 y_2 \quad \equiv w_2(y_1, y_2).$$

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

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

Therefore,

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

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

cf. Ex2.8 in LNp.29

---

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

---

• extrema and order statistics → 順序統計量 → quantile (分位數)

**Definition 2.11** (order statistics, sec 3.7)

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

$$X_{(1)} = \min(X_1, X_2, \ldots, X_n) \quad \text{is the } \textbf{minimum}$$
$$X_{(n)} = \max(X_1, X_2, \ldots, X_n) \quad \text{is the } \textbf{maximum}$$
$$R \equiv X_{(n)} - X_{(1)} \quad \text{is called } \textbf{range}$$
$$S_j \equiv X_{(j)} - X_{(j-1)}, j = 2, \ldots, n \quad \text{are called } j\text{th } \textbf{spacings}$$

transformation $g$

$$X_4 \quad X_2 \quad X_3 \qquad X_6 \quad X_1 \qquad X_5 \longrightarrow \mathbb{R} \quad g^{-1} \text{ not exist}$$

$$X_{(1)} \, X_{(2)} \, X_{(3)} \qquad X_{(4)} \, X_{(5)} \qquad X_{(6)}$$

$R$

$S_2 \quad S_3 \qquad\qquad\qquad S_6$

**Note.** In the section, we only consider the <u>case</u> that $X_1, X_2, \ldots, X_n$ are <u>i.i.d</u> <u>continuous r.v.'s</u> with <u>cdf</u> $F$ and <u>pdf</u> $f$. Although $X_1, X_2, \ldots, X_n$ are <u>independent</u>, their <u>order statistics</u> are <u>not independent</u> in general. $\boxed{X_{(1)}, \cdots, X_{(n)}}$

---

**Definition 2.12** (i.i.d.)

$\underline{X_1, X_2, \ldots, X_n}$ are **i.i.d.** (independent, identically distributed) with cdf $F$/pmf $p$/pdf $f$ $\Rightarrow X_1, X_2, \ldots, X_n$ are <u>independent</u> and have a <u>common marginal</u> cdf $F$/pmf $p$/pdf $f$. $\quad \hookrightarrow$ joint = $\pi$ marginal
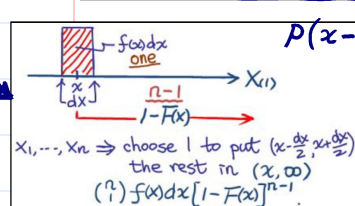
$\boxed{\text{but not } \underline{\text{common}} \text{ value}}$

---

**Theorem 2.8** (TBp. 104)

The <u>cdf</u> of $\underline{X_{(1)}}$ is $\underline{1-[1-F(x)]^n}$ and its <u>pdf</u> is $\underline{nf(x)[1-F(x)]^{n-1}}$.

The <u>cdf</u> of $\underline{X_{(n)}}$ is $\underline{[F(x)]^n}$ and its <u>pdf</u> is $\underline{nf(x)[F(x)]^{n-1}}$.

$P(x - \frac{dx}{2} < X_{(1)} < x + \frac{dx}{2}) \approx f_{X_{(1)}}(x) dx$



$x_1, \cdots, x_n \Rightarrow$ choose 1 to put $(x - \frac{dx}{2}, x + \frac{dx}{2})$ the rest in $(x, \infty)$
$\binom{n}{1} f(x) dx [1 - F(x)]^{n-1}$

$$\begin{aligned} F_{X_{(n)}}(x) &= P(X_{(n)} \le x) = P(X_1 \le x, \ldots, X_n \le x) \\ &= P(X_1 \le x) \cdots P(X_n \le x) \qquad \frac{d F_{X_{(n)}}(x)}{dx} \\ &= [F(x)]^n. \end{aligned}$$



$x_1, \cdots, x_n \Rightarrow$ choose 1 to put $(x - \frac{dx}{2}, x + \frac{dx}{2})$ the rest in $(-\infty, x]$
$\binom{n}{1} f(x) dx [F(x)]^{n-1}$

$$\begin{aligned} 1 - F_{X_{(1)}}(x) &= P(X_{(1)} > x) = P(X_1 > x, \ldots, X_n > x) \\ &= P(X_1 > x) \cdots P(X_n > x) \qquad \frac{d F_{X_{(1)}}(x)}{dx} \\ &= [1 - F(x)]^n. \end{aligned}$$
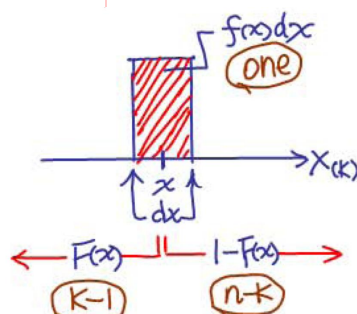
$P(x - \frac{dx}{2} < X_{(n)} < x + \frac{dx}{2}) \approx f_{X_{(n)}}(x) dx$

---

one

$\leftarrow F(x) \quad\vdash\!\!\vdash \quad 1-F(x) \rightarrow$
$(k-1) \qquad (n-k)$

$X_1, \cdots, X_n \Rightarrow$ choose 1 to place in $(x - \frac{dx}{2}, x + \frac{dx}{2})$
$\qquad = k-1 : \quad : \quad = (-\infty, x)$
$\qquad = n-k : \quad : \quad = (x, \infty)$
$\binom{n}{1 \ \ k-1 \ \ n-k} f(x) dx [F(x)]^{k-1} [1-F(x)]^{n-k}$
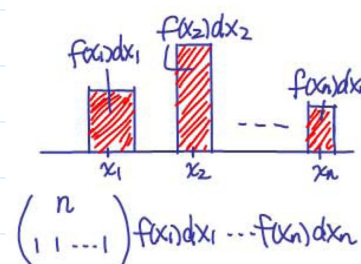
**Theorem 2.9** (TBp. 105)

The <u>pdf</u> of the $k$th order statistic $\underline{X_{(k)}}$ is

$P(x - \frac{dx}{2} < X_{(k)} < x + \frac{dx}{2}) \approx f_{X_{(k)}}(x) \cdot dx$

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



$\binom{n}{1 \ 1 \cdots 1} f(x_1) dx_1 \cdots f(x_n) dx_n$

---

**Theorem 2.10** (TBp. 114, Problem 73)

The <u>joint pdf</u> of $\underline{X_{(1)}, X_{(2)}, \cdots, X_{(n)}}$ is

$P(x_i - \frac{dx_i}{2} < X_{(i)} < x_i + \frac{dx_i}{2}, \ i=1, \cdots, n) \approx f_{X_{(1)} \cdots X_{(n)}}(x_1, \cdots, x_n) dx_1 \cdots dx_n$

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

for $\underline{x_1 \le x_2 \le \cdots \le x_n}$, and $f_{X_{(1)} X_{(2)} \ldots X_{(n)}} = 0$ otherwise.



$X_{(2)}$

$\leftarrow X_{(1)} = X_{(2)}$

$\rightarrow X_{(1)}$

not a product set.

**Question:** Are $X_{(1)}, X_{(2)}, \ldots, X_{(n)}$ <u>independent</u>, judged from the from of its joint pdf? $\leftarrow$ c.f. Thm 2.2, item 1 (LN$_p$.19)

**Example 2.11** (range, TBp. 105-106)

The joint pdf of $X_{(1)}$ and $X_{(n)}$ is $P(s-\frac{ds}{2} < X_{(1)} < s+\frac{ds}{2}, t-\frac{dt}{2} < X_{(n)} < t+\frac{dt}{2})$
$\approx f_{X_{(1)},X_{(n)}}(s,t)\,ds\,dt.$

$$f_{X_{(1)}X_{(n)}}(s,t) = n(n-1)f(s)f(t)[F(t)-F(s)]^{n-2}, \quad \text{for } s \leq t,$$
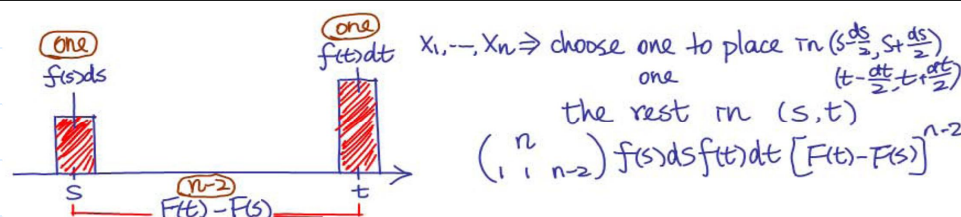
and 0 otherwise. Therefore, the pdf of $R = X_{(n)} - X_{(1)}$ is
— check exercise in Ex2.7(LNp.28)

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

one $f(s)ds$    one $f(t)dt$    $X_1, \cdots, X_n \Rightarrow$ choose one to place in $(s-\frac{ds}{2}, s+\frac{ds}{2})$
one    $(t-\frac{dt}{2}, t+\frac{dt}{2})$
the rest in $(s,t)$
$\binom{n}{1\ 1\ n-2} f(s)ds\,f(t)dt\,[F(t)-F(s)]^{n-2}$
$s$    $(n-2)$ $F(t)-F(s)$    $t$

**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

# Chapter 4

**Outline**

➢ expectation ← 期望值.
  • 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?

pmf    $X$    $3X$    $3X+3$

$\ell_1:\ell_2=1:3$    $\ell_1$    $\ell_2$    3

• expectation

---

**Definition 3.1** (expectation, TBp. 122, 123)

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

$\mathbb{R}^n \to \mathbb{R}^1$

$$E(Y) \equiv \sum_{-\infty < y < \infty} y p_Y(y) = E[g(X_1, \ldots, X_n)]$$

weighted average
加權平均
平均: $y$
權重: $P_Y / f_Y$

$$\equiv \sum_{-\infty < x_1 < \infty, \ldots, -\infty < x_n < \infty} g(x_1, \ldots, x_n) p(x_1, \ldots, x_n),$$

if $X_1, X_2, \ldots, X_n$ are discrete random variables, or

$Y$: *random*
$E(Y)$: *fixed value*

$$E(Y) \equiv \int_{-\infty}^{\infty} y f_Y(y) dy = E[g(X_1, \ldots, X_n)]$$

$$\equiv \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(x_1, \ldots, x_n) f(x_1, \ldots, x_n) dx_1 \cdots dx_n,$$

if $Y$ and $X_1, X_2, \ldots, X_n$ are continuous random variables.

---

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

1. (TBp.116&118)  $g(x) = x \Rightarrow E[g(X)] = E(X)$ is called **mean** of $X$, usually denoted by $E(X)$ or $\mu_X$.

constant

2. (TBp.131)  $g(x) = (x - \mu_X)^2 \Rightarrow E[g(X)] = E[(X - E(X))^2]$ is called **variance** of $X$, usually denoted by $Var(X)$ or $\sigma_X^2$. The square root of variance, i.e., $\sigma_X$, is called **standard deviation**.

constant, not random

3. (TBp.138)  $g(x,y) = (x - \mu_X)(y - \mu_Y) \Rightarrow E[g(X,Y)] = E[(X - E(X))(Y - E(Y))]$ is called **covariance** of $X$ and $Y$, usually denoted by $Cov(X,Y)$ or $\sigma_{XY}$.

4. (TBp.142)  The **correlation coefficient** of $X, Y$ is defined as $\sigma_{XY}/(\sigma_X \sigma_Y)$, usually denoted by $Cor(X,Y)$ or $\rho_{XY}$. $X$ and $Y$ are called **uncorrelated** if $\rho_{XY} = 0$. $\Longleftrightarrow \sigma_{XY} = 0$

**Notes.** (intuitive explanation of mean)

from its definition

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.
     ← center of gravity (重心)

3. Mean can be interpreted as a long-run average. (see Chapter 5.) → LLN

**Notes.** (intuitive explanation of variance and standard deviation)

from its definition

how the dist. is spread out

① variance is the average value of the squared deviation of $X$ from $\mu_X$.

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 $a, b_1, \ldots, b_n \in \mathbb{R}$,
$$E\left(a + \sum_{i=1}^{n} b_i X_i\right) = a + \sum_{i=1}^{n} b_i E(X_i).$$
$$\Rightarrow E(a+bX) = a + b \cdot E(X)$$

g: convex
$E[g(x)] \geq g[E(x)]$

g: concave
$E[g(x)] \leq g[E(x)]$

② (TBp.124)    If $X$, $Y$ are independent, then
$$E(g(X)h(Y)) = E(g(X))E(h(Y)).$$
   $\underbrace{\phantom{E(g(X))}}_{W} \underbrace{\phantom{E(h(Y))}}_{Z}$   W & Z are independent.

independent ⇒ uncorrelated

In particular, $E(XY) = E(X)E(Y)$.
(**Question 3.2**: $E(X/Y) = E(X)/E(Y)$? ← $E\left(\frac{x}{Y}\right) = E\left(x \cdot \frac{1}{Y}\right) = E(x) \cdot E\left(\frac{1}{Y}\right)$
**Note**. $E[g(X)] \neq g[E(X)]$ in general.   $\frac{1}{E(Y)}$?

**Theorem 3.2** (properties of variance and standard deviation)

① (TBp.132)    $\sigma_X^2 = Var(X) = E[(X - \mu_X)^2] = E(X^2) - \mu_X^2.$
→ for calculation purpose              $\bot [E(x)]^2$

② (TBp.131)    $Var(a + bX) = b^2 Var(X),$ $a, b \in \mathbb{R}$, and $\sigma_{a+bX} = |b|\sigma_X.$

• location shift ⇒ no impact on $\sigma^2$
• scale change ⇒ $\sigma^2 \to b^2\sigma^2$

$[b_1 \cdots b_n] \begin{bmatrix} a_{ij} = cov(x_i, x_j) \end{bmatrix} \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$

covariance matrix   var(x_i)

3. (TBp.140)
$$Var\left(a + \sum_{i=1}^{n} b_i X_i\right) = \sum_{i=1}^{n} b_i^2 Var(X_i) + 2 \sum_{1 \leq i < j \leq n} b_i b_j Cov(X_i, X_j).$$

gone

In particular, $Var(X + Y) = Var(X) + Var(Y) + 2Cov(X, Y).$

④ (TBp.140)    If $X_1, \ldots, X_n$ are independent,
         imply → $cov(x_i, x_j) = 0,$ i.e., uncorrelated, $\forall i \neq j$

cf.

mean of sum item 1, Thm 3.1 (LNp.41)

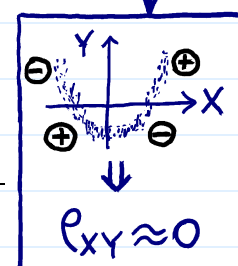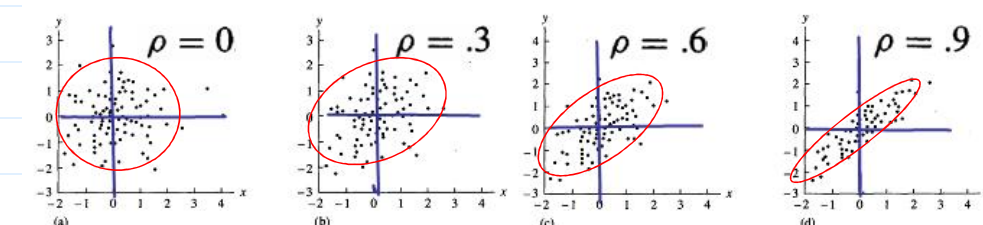$$Var\left(\sum_{i=1}^{n} X_i\right) = \sum_{i=1}^{n} Var(X_i).$$
$-\mu_x + \mu_x$

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

**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.

   *might not be causal relation* → *i.e., when $X$ (r.v.) is large (or small), will $Y$ tend to be larger or smaller?*

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. ← *from its definition.*

3. positive covariance and negative covariance → *drawback: cov depends on the scale/unit of $X$ & $Y$*

4. correlation coefficient is unit free

5. correlation coefficient measures the strength of the linear relationship between $X$ and $Y$.

$$COV(X,Y) = E[(X-\mu_X)(Y-\mu_Y)]$$

X: $X-\mu_X$ positive
X: $X-\mu_X$ negative
X: $Y-\mu_Y$ positive
X: $Y-\mu_Y$ negative

$(\mu_X, \mu_Y)$

positive covariance    zero covariance    negative covariance

$\rho_{XY} \approx 0$

$\rho = 0$    $\rho = .3$    $\rho = .6$    $\rho = .9$

(a)    (b)    (c)    (d)

**Theorem 3.4** (properties of covariance and correlation coefficient)

1. (TBp.138) $Cov(X,Y) = E[(X - \mu_X)(Y - \mu_Y)] = E(XY) - \mu_X\mu_Y$
   (Note. $Cov(X,X) = Var(X)$.)
   → for calculation purpose

2. (TBp.140)

$$[b_1 \cdots b_n]\begin{bmatrix} \sigma_{ij} = Cov(X_i, Y_j) \end{bmatrix}\begin{bmatrix} d_1 \\ \vdots \\ d_m \end{bmatrix}$$

$$Cov\left(a + \sum_{i=1}^{n} b_i X_i, c + \sum_{j=1}^{m} d_j Y_j\right) = \sum_{i=1}^{n}\sum_{j=1}^{m} b_i d_j Cov(X_i, Y_j)$$

gone

3. (TBp.140) If $X$, $Y$ are independent then $Cov(X,Y) = 0$, i.e., **independent ⇒ uncorrelated**. But, the converse statement is not necessarily true.

$\rho = +1 \Leftrightarrow a > 0$
$\rho = -1 \Leftrightarrow a < 0$

4. (TBp.143) $-1 \le \rho_{XY} \le 1$ and $\rho_{XY} = \pm 1$ if and only if $Y = aX + b$ with probability one for some $a$, $b \in \mathbb{R}$.

standardization (標準化)
After standardization, mean = 0, Var = 1

5. $\rho_{XY} = E\left[\left(\dfrac{X - \mu_X}{\sigma_X}\right)\left(\dfrac{Y - \mu_Y}{\sigma_Y}\right)\right]$

6. $|Cor(a + bX, c + dY)| = |Cor(X,Y)|$ → • location shift • scale change ⇒ no impact on |cor|

**• moment generating function & characteristics function**

**Definition 3.3** (moment generating function, TBp. 155)

The **moment generating function (mgf)** of a random variable $X$ is
$$M_X(t) = E(e^{tX}), \quad t \in \mathbb{R}$$
Laplace transformation of

$M_X(t) = \begin{cases} \int e^{tx} f_X(x)\, dx \\ \sum e^{tx} p_X(x) \end{cases}$

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 $t$.
   $t = 0 \Rightarrow E(e^{0 \cdot x}) = 1$ ← always exists
   i.e., $E(e^{tX}) < \infty$

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.

→ know mgf ⇒ know distribution.

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

*the reason why it's called moment generating function.*

*know all moments*
$\Rightarrow$ *know* $M_X(t) = \sum_{k=0}^{\infty} \frac{M_k^{(k)}(0)}{k!} t^k$
$\Rightarrow$ *know dist.*

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

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

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

↳ *useful for identifying the dist.* of $X_1 + \cdots + X_n$

6. <u>continuity theorem</u> (see <u>Chapter 5</u>)　↳ *generalization:* indep. $X_1, \cdots, X_n$

$$M_{X_1 + \cdots + X_n}(t) = \prod_{i=1}^{n} M_{X_i}(t)$$

---

**Definition 3.4** (moment, TBp. 155)

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

➤ Some Notes.

*$\mu_k'$: a linear combination of $\mu_1, \cdots, \mu_k$*

- $\underline{\mu_k'} = \underline{\sum_{i=0}^{k} \binom{k}{i} (-\mu_X)^{n-i} \mu_i}$.

*$\mu_k$: a linear combination of $\mu_1, \mu_2', \cdots, \mu_k'$*

- $\underline{\mu_k} = \underline{\sum_{i=0}^{k} \binom{k}{i} (\mu_X)^{n-i} \mu_i'}$.

- In particular, $\underline{E(X) = \mu_X = \mu_1}$, and,

$$\underline{Var(X) = \sigma_X^2 = \mu_2 - \mu_1^2 = \mu_2'}.$$

$$\begin{array}{c|ccc} \mu_1 & \mu_2 & \cdots & \mu_k & \cdots \\ \hline \mu_1 & \mu_2' & \cdots & \mu_k' & \cdots \end{array}$$

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

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

$M_{X_1, \cdots, X_n}(t, \cdots, t)$
$= M_{X_1 + \cdots + X_n}(t)$

$$M_{\underline{X_1 X_2 \cdots X_n}}(t_1, t_2, \ldots, t_n) \equiv \underline{E(e^{t_1 X_1 + t_2 X_2 + \cdots + t_n X_n})}$$

c.f. ↳ *mgf of* $X_1 + X_2 + \cdots + X_n = Y$
$= E(e^{tX_1 + tX_2 + \cdots + tX_n})$

if the <u>expection exists</u>.

---

**Theorem 3.6** (<u>properties</u> of joint mgf)

1. $M_{\underline{X_1}}(\underline{t_1}) = M_{\underline{X_1 X_2 \cdots X_n}}(\underline{t_1}, \underline{0}, \ldots, 0)$ ← *relationship between* <u>joint</u> *mgf &* <u>marginal</u> *mgf.*

2. <u>uniqueness theorem</u>

★3. $X_1, X_2, \ldots, X_n$ are <u>independent</u> if and only if

*LNp.19,*
<u>joint</u> ⟨ cdf / pmf / pdf
$= \prod_{i=1}^{n}$ marginal ⟨ cdf / pmf / pdf　*c.f.*

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

*c.f.*

*the mgf of the* <u>sum</u> *of* <u>indep.</u> $X_1, \cdots, X_n$
$= \prod_{i=1}^{n} M_{X_i}(t)$

★4. $\left. \dfrac{\partial^{r_1 + \cdots + r_n}}{\partial t_1^{r_1} \cdots \partial t_n^{r_n}} M_{\underline{X_1 X_2 \cdots X_n}}(t_1, t_2, \ldots, t_n) \right|_{t_1 = t_2 = \cdots = t_n = 0}$

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

• **conditional expectation** ← — Recall: conditional distribution (LNp.21~23)

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

The **conditional expectation** of $h(Y)$ given $X = x$ is

（平均: Y or $h(Y)$
權重: $p_{Y|X}(y|x)$
$f_{Y|X}(y|x)$）

— a function of $x$

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

In particular, $\underline{E(Y|X = x)} = \sum_y y\, p_{Y|X}(y|x)$ — a pmf for $y$

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

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

function of $x$ with unit of Y

$f(x, y)$: joint pdf

$E_{Y|X}(Y|\underline{x})$

$f(x^*, y)$ — a function of $y$ only

$E_{Y|X}(Y|\underline{x}^*)$

$f_{Y|X}(y|x^*) = \dfrac{f(x^*, y)}{f_X(x^*)}$

e.g.,
$h(Y) = Y$
$X$: height (cm)
$Y$: weight (kg)

$E(Y|X=\underline{170})$
= average weight of <u>people</u> whose <u>height</u> = 170



area = $f_X(x^*)$

$X = \underline{x}^*$

• a curve on (X, Y) plane
• a map from X to Y

made by S.-W. Cheng (NTHU, Taiwan)

**Theorem 3.8** (properties of conditional expectation)

1. $E_{Y|X}(h(Y)|x)$ is a function of $x$ and is free of $Y$.

fixed values → the Y part has been integrated or summed

② If $X$ and $Y$ are independent then $E_{Y|X}(h(Y)|x) = E_Y(h(Y))$.

By Thm2.3, item 6, LNp.23, $\begin{cases} P_{Y|X}(y|x)=P_Y(y) \\ f_{Y|X}(y|x)=f_Y(y) \end{cases}$

(intuition) $E_{Y|X}(h(Y)|x)$ is a constant function of $x$ ⇒ X offers no information of Y

3. $E(h(X)|X = x) = h(x)$

e.g. X:height Y:weight g(x):function that maps from height to average weight

④ Let $g(x) = E_{Y|X}(h(Y)|x)$, then $g(X)$ is a random variable (transformation of $X$) and usually denoted by $E_{Y|X}(h(Y)|X)$.
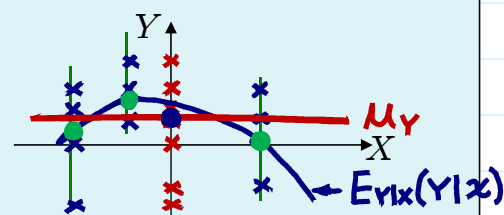
It's a function of X only. But, its random value reflects h(Y)

5. **law of total expection** (TBp.149)
$$E_X[E_{Y|X}(h(Y)|X)] = E_Y[h(Y)].$$
In particular,

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

$E_{X,Y} = E_X E_{Y|X} = E_Y E_{X|Y}$

$\sum_x \sum_y h(y) P_{XY}(x,y)$ $\parallel$ $P_{Y|X}(y|x)·P_X(x)$ $\uparrow$ $\sum_y \sum_x$

$\int\int h(y) f_{XY}(x,y)dy dx$ $\parallel$ $f_{Y|X}(y|x)f_X(x)$ $\downarrow$ $\int_y\int_x$

generalization $E_{XY}[h(X,Y)] = E_Y E_{X|Y}[h(X,Y)|Y] = E_X E_{Y|X}[h(X,Y)|X]$

4. **variance decomposition** (TBp.151)
$$Var_Y(Y) = Var_X[E_{Y|X}(Y|X)] + E_X[Var_{Y|X}(Y|X)]$$

**Note.**

1. $Var_Y(Y) \geq E_X[Var_{Y|X}(Y|X)]$ and the equality holds if and only if $E_{Y|X}(Y|X) = E_Y(Y)$ with probability one.
$Var_X[E_{Y|X}(Y|X)] = 0$

2. $Var_Y(Y) \geq Var_X[E_{Y|X}(Y|X)]$ and the equality holds if and only if $Var_{Y|X}(Y|X) = 0$ with probability one; i.e., $Y = E_{Y|X}(Y|X)$ with probability one.
$E_X[Var_{Y|X}(Y|X)] = 0$

**• 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?)

  | $X$ | 身高 | 雨量 |
  |---|---|---|
  | $Y$ | 體重 | 米產量 |

- **statistical modeling**: assign $(X, Y)$ a (known) joint distribution (cdf $F(x, y)$, pdf $f(x, y)$, or pmf $p(x, y)$)

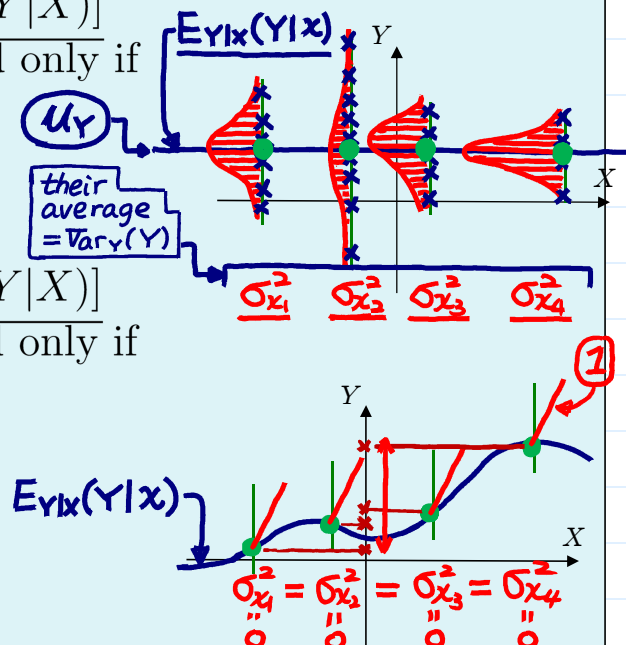- **objective**: Predict $Y$ by using a function of $X$, i.e., $g(X)$.

  We consider the following three groups of $g$'s:

  (i) $G_1 = \{g(x) : g(x) = c, \text{ where } c \in \mathbb{R}\}$ — *not use the information of $X$*

  (ii) $G_2 = \{g(x) : g(x) = a + bx, \text{ where } a, b \in \mathbb{R}\}$, and

  (iii) $G_3 = \{g(x) : g \text{ is arbitrary}\}$.

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

- **question**: Within each group, what is the "best" prediction?

  *i.e., how to choose $c$ for $G_1$*
  *: : : $a, b$ : $G_2$*
  *: : : $g$ for $G_3$*

- **criterion**: minimizing mean square error:

  *meaning?* $\to$ $\text{MSE} \equiv E_{X,Y}\{[Y - g(X)]^2\}$. *predicted value*

  *true value* — *error*

$G_1$ → **Example 3.2** ("best" constant prediction, TBp. 153)

$$E_{X,Y}(Y - c)^2 = E_Y(Y - c)^2 \geq E_Y[Y - E_Y(Y)]^2 = Var_Y(Y) \quad \boxed{min}$$

$G_3$

The equality holds if and only if $c = E_Y(Y)$. *— only need to know $\mu_Y$*

**Example 3.3** ("best" prediction of $Y$ using $X$, TBp. 153)

*mean: best predictor under MSE*

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

*cf. check LNp. 52*

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

**Notes for the best predictor in $G_3$.**

*best in $G_1: E_Y(Y)$*

- $E_{Y|X}(Y|X)$ is the best predictor of $Y$ based on $X$, in the mean squared prediction error sense. *intuition* ← *check the graph in LNp.50* | *median: best predictor under $E|Y - g(X)|$*

- need to know the joint distribution of $X$ and $Y$, or at least $E_{Y|X}(Y|x)$

- $E_{Y|X}(Y|x)$ is called the regression function of $Y$ on $X$. *迴歸*

$G_2$ → **Example 3.4** ("best" linear prediction of $Y$ using $X$, TBp. 153-154)

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

*$0 \leq |\rho| \leq 1$*    $\boxed{min}$

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

### 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  [best in $G_2$]  [more information better predictor]

  [best in $G_3$]  $\hookrightarrow$ linear regression analysis

- ⊙ needs to know only the means, variances and covariances

  [c.f.] $\to$ the best in $G_1$ & $G_3$ $\longrightarrow$ Which one require more information?

- ⊙ $\sigma_Y^2(1 - \rho^2)$ is small if $\rho$ is close to $+1$ or $-1$, and large if $\rho$ is close to $0$  [intuition] $\leftarrow$ check the plot in LNp.44

### Notes.

Collect data of X, Y to estimate their joint dist.

(1) $\min_{a,b} E[Y - (a + bX)]^2 \leq \min_c E(Y - c)^2$ and the equality holds if and only if $\rho = 0$.  $\because G_1 \subset G_2 \subset G_3$

(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?

❖ **Reading**: textbook, Chapter 4
❖ **Further Reading**: Roussas, 5.1, 5.3, 5.4, 5.5, 6.1, 6.2, 6.4, 6.5

# 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? ⎿ statistical modeling

• discrete distributions

**Definition 4.1** (Uniform distribution $U(a_1,...,a_m)$ )

Equal probability to obtain $a_1, a_2, \ldots, a_m$.

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

        a pmf ? (Ec)

                                               U(1,2,3,4)

• **mgf:** $\frac{\sum_{j=1}^{m} e^{a_j t}}{m}$ ← by definition (Ec)

• **mean:** $\frac{\sum_{j=1}^{m} a_j}{m} \equiv \bar{a}$          • **parameter:** $a_i \in \mathbb{R}, \ m = 1, 2, \ldots$

• **variance:** $\frac{\sum_{j=1}^{m}(a_j - \bar{a})^2}{m}$       • **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:** $p(x) = \begin{cases} p^x(1-p)^{(1-x)}, & \text{if } x = 0 \text{ or } x = 1 \\ 0, & \text{otherwise} \end{cases}$

  *a pmf ? (Ec)*

- **mgf:** $pe^t + 1 - p$ — *by definition* **(Ec)**

- **mean:** $p$ — *by definition, use mgf* **(Ec)** $\overset{P}{=}$

- **variance:** $p(1-p)$ — $Var(X) = E(x^2) - [E(x)]^2$ **(Ec)**
  $Var(X) = EX(X-1) + E(X) - [E(X)]^2$
  *use mgf* ⊙

- **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.

$\hookrightarrow p = P(A)$

$I_A: \Omega \to \mathbb{R}, \quad I_A(\omega) = \begin{cases} 1, & \text{if } \omega \in A \\ 0, & \text{if } \omega \notin A \end{cases}$

**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 $p$.

*Shape*   *explanation*

- **pmf:** $p(x) = \begin{cases} \dbinom{n}{x} p^x(1-p)^{(n-x)}, & x = 0, 1, \ldots, n \\ 0, & \text{otherwise} \end{cases}$

  *a pmf ?* **(Ec)** ← use **(✱)**

  $0 + 1 + 1 + \cdots + 0 = x$
  $\boxed{1}\ \boxed{2}\ \cdots\ \boxed{n}$

- **mgf:** $(pe^t + 1 - p)^n, \ t \in \mathbb{R}$ — *by definition, sum of i.i.d. B(p)* **(Ec)**

- **mean:** $np$ — *use definition, use mgf, sum of iid B(p)* **(Ec)**

  *intuition*

- **variance:** $np(1-p)$ ← *max at $p = \frac{1}{2}$, min at $p = 0$ or 1*

  $E(X) = \sum_{x=0}^{n} x \binom{n}{x} p^x (1-p)^{n-x}$
  $= \left[ \sum_{x=1}^{n} \binom{n-1}{x-1} p^{x-1} (1-p)^{(n-1)-(x-1)} \right] np$

  *pmf of $B(n-1, p)$*

  ↳ *STO (sum-to-one) method*

- **parameter:** $p \in [0, 1], \ n = 1, 2, \ldots$

- **example:** # of heads, toss a coin $n$ times

*Find $E(x^2)$ using mgf*
*Find $E[X(X-1)]$ using STO* **(Ec)**
*sum of i.i.d. B(p)*



$n = 10 \text{ and } p = .1$ (a)

$n = 10 \text{ and } p = .5$ (b)

**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　　*(intuition)*

2. Let $X_1, \ldots, X_n$ be i.i.d. $B(p)$, then $Y = X_1 + \cdots + X_n \sim B(n, p)$.— prove using ① mgf ( $M_Y(t) = \prod_{i=1}^{n} M_{X_i}(t)$ ) ② convolution & induction **(Ec)**

3. Let $X_i \sim B(n_i, p), i = 1, \ldots, k$, and $X_1, \ldots, X_k$ are independent. Then, $Y = X_1 + \cdots + X_k \sim B(n_1 + \cdots + n_k, p)$. 　*(intuition)*

prove using mgf
prove using convolution & induction **(Ec)**

$\underbrace{\substack{\frac{1}{6}\;\frac{1}{6} \\ 0\;1}}_{X_1}$ $\underbrace{\substack{\frac{1}{6}\;\frac{1}{6}\;\frac{1}{6}\;\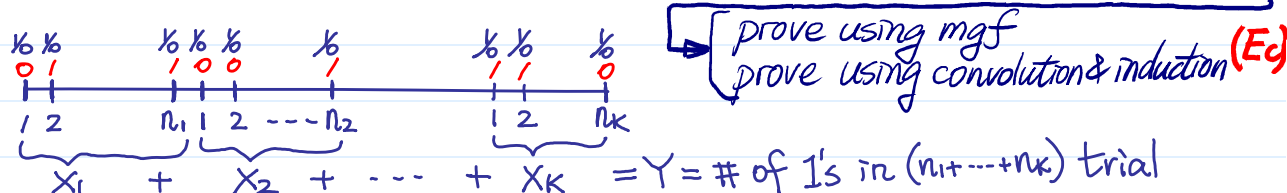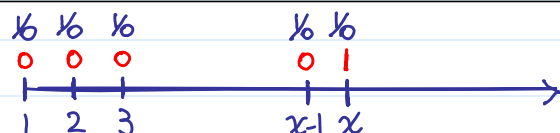;\;\frac{1}{6} \\ 1\;0\;0\;\;\;1}}_{X_2}$ $+ \cdots +$ $\underbrace{\substack{\frac{1}{6}\;\frac{1}{6}\;\;\;\frac{1}{6} \\ 1\;1\;\;\;0}}_{X_K}$ $= Y = \#$ of $1$'s in $(n_1 + \cdots + n_k)$ trial

---

**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.

$\substack{\frac{1}{6}\;\frac{1}{6}\;\frac{1}{6}\\ 0\;\;0\;\;0}\qquad\substack{\frac{1}{6}\;\frac{1}{6}\\ 0\;\;1}$
$\;\;1\;\;2\;\;3\qquad\qquad x{-}1\;\;x$

---

*(explanation)*

- **pmf:** $p(x) = \begin{cases} (1-p)^{(x-1)}p, & x = 1, 2, 3, \ldots \\ 0, & \text{otherwise} \end{cases}$

a pmf ? **(Ec)** ← use **(**✱✱**)**

- **cdf:** $F(x) = \begin{cases} 1 - (1-p)^{[x]}, & 1 \le [x] \le x < [x] + 1 \\ 0, & x < 1 \end{cases}$ — Find $P(X > x)$ using **(**✱✱**)** **(Ec)**

- **mgf:** $\frac{pe^t}{1-(1-p)e^t}$, $t < -\log(1-p)$. — use **(**✱✱**)** / use STO **(Ec)**

- **mean:** $\frac{1}{p}$ — use $E(X) = \sum_{k=1}^{\infty} P(X \ge k)$ or use **(**✱✱**)** / use mgf / use differentiation method (TBp.117, Example B) **(Ec)** *(intuition)*

- **variance:** $\frac{1-p}{p^2}$ — Find $E(X^2)$ using mgf / Find $E[X(X-1)]$ using differentiation method **(Ec)**

- **parameter:** $p \in [0, 1]$

- **example:** lottery, # of tickets a person must purchase up to and including the first winning ticket
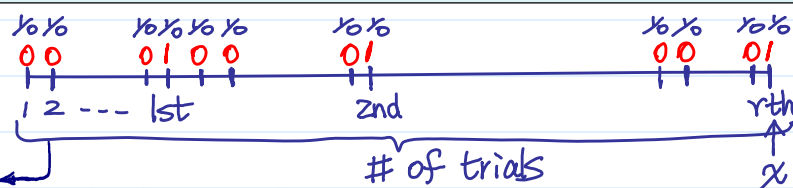
**Note:** a memoryless distribution ← *(intuition)*
└ check its definition (LNp.74) and prove **(Ec)**

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



*explanation*

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

  *a pmf? (Ec)*　*use (✗✗✗)*

- **mgf:** $\frac{p^r e^{rt}}{[1-(1-p)e^t]^r}, \quad t < -\log(1-p).$ ⎱ *use STO / use (✗✗✗) / sum of i.i.d G(p)* **(Ec)**

- **mean:** $\frac{r}{p}$ ⎱ *use mgf / use STO / sum of i.i.d. G(p)* **(Ec)**

  *intuition*

- **variance:** $\frac{r(1-p)}{p^2}$ ⎱ *Find E(X²) using mgf / Find E(X(X+1)) using STO / sum of i.i.d. G(p)* **(Ec)**

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

- **example:** lottery, <u># of tickets</u> a person must purchase up to and including the <u>$r$th winning</u> ticket

  **Note:** (✗✗✗)
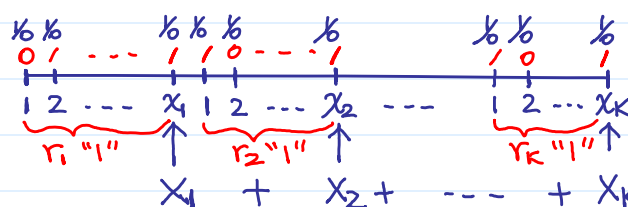  $\sum_{x=0}^{\infty} \binom{n+x-1}{x} t^x = \frac{1}{(1-t)^n},$
  for $-1 < t < 1$.

---

**Note.**

1. <u>negative binomial</u> distribution is a <u>generalization</u> of <u>geometric</u> distribution from <u>1st success to $r$th</u> success

   *intuition*

2. Let $X_1, X_2, \dots, X_r$ be i.i.d. $G(p)$, then $Y = X_1 + \dots + X_r \sim NB(r, p)$ — *prove using ① mgf ( $M_Y(t) = \prod_{i=1}^{n} M_{X_i}(t)$ ) ② convolution & induction* **(Ec)**

3. Let $X_i \sim NB(r_i, p), i = 1, \dots, k$, and $X_1, \dots, X_k$ are <u>independent</u>. Then, $Y = X_1 + \dots + X_k \sim NB(r_1 + \dots + r_k, p)$. ⎱ *prove using mgf / prove using convolution & induction* **(Ec)**

   *intuition*



$X_1 + X_2 + \dots + X_K = Y = $ # of trial until $(r_1 + \dots + r_K)$ one.

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

- **joint pmf:** *use (\*\*\*\*)*

  *a joint pmf? (Ec)*

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

  *explanation*

- **joint mgf:** $(p_1 e^{t_1} + \cdots + p_r e^{t_r})^n$, $t_1, \ldots, t_r \in \mathbb{R}$. ⎰ *use (\*\*\*\*)* ⎱ *use STO* **(Ec)**

- **marginal distribution:** $X_i \sim B(n, p_i)$, $i = 1, \ldots, r$ — *intuition* **(Ec)** — *prove using mgf*

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

- **variance:** $Var(X_i) = np_i(1 - p_i)$, $i = 1, \ldots, n$ — *Find $E(X_i X_j)$ using STO*

- **covariance:** $Cov(X_i, X_j) = -np_i p_j$, $i \neq j$ — *Find $E(X_i X_j)$ using mgf*

- **parameter:** $p_i \in [0, 1]$, and $\sum_{i=1}^{r} p_i = 1$.　　$n = 1, 2, \ldots$ **(Ec)**

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

  *why negative?*

> (\*\*\*\*)
> **Note:** $(a_1 + \cdots + a_k)^n = \sum\limits_{x_1 + \cdots + x_k = n} \begin{pmatrix} n \\ x_1, \cdots, x_k \end{pmatrix} a_1^{x_1} \cdots a_k^{x_k}$.

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

---

> **Definition 4.7** (Poisson distribution $P(\lambda)$, sec 2.1.5)
>
> Limit of binomial distributions $X_n \sim B(n, p_n)$, where $p_n \to 0$ as $n \to \infty$ in such a way that $\lambda_n \equiv np_n \to \lambda$.

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

$P_n = \dfrac{\lambda_n}{n}$

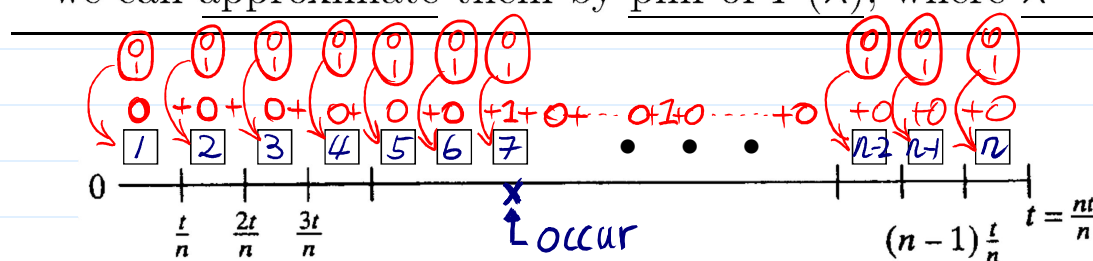**Note:** if $a_n \to a$, $\left(1 + \dfrac{a_n}{n}\right)^n \to e^a$.

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

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

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

$\frac{\lambda^x}{x!} \quad e^{-\lambda}$

**explanations.**

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



$t = \frac{nt}{n}$

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

*(handwritten)* $N_1 + N_2 + \cdots + N_n$
$\sim B(n, p)$ with large $n$ & small $p$

*(handwritten)* shape

*(handwritten)* a pmf? **(Ec)**    use **(★★★★)**

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

- **mgf:** $e^{\lambda(e^t - 1)}$, $t \in \mathbb{R}$.    *(handwritten)* use **(★★★★)** / use STO  **(Ec)**

- **mean:** $\lambda$    *(handwritten)* use STO / use mgf **(Ec)**    *(box)* meaning of parameter $\lambda$: average occurences

- **variance:** $\lambda$    *(handwritten)* Find $E[X(X-1)]$ using STO / Find $E(X^2)$ using mgf **(Ec)** / $np(1-p) \approx np$

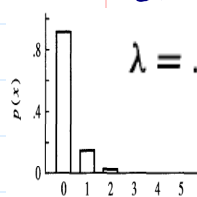- **parameter:** $\lambda > 0$

*(handwritten box)* **Note:** **(★★★★)** $e^\lambda = \sum_{x=0}^{\infty} \frac{\lambda^x}{x!}$

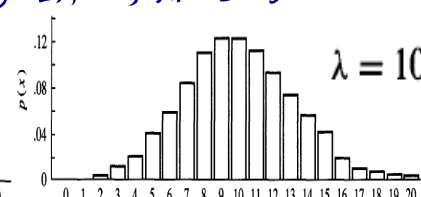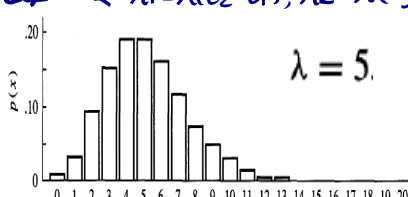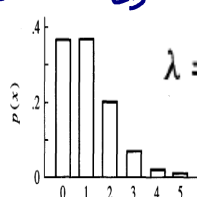- **example:** number of phone calls coming into an exchange during a unit of time

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

*(handwritten)* prove using mgf / prove using convolution & induction **(Ec)** / intuition

*(handwritten)* $X_1 + X_2 + X_3$
$t_1 \quad t_2 \quad t_3 \quad t_4$    $\leftarrow \lambda_1 = \lambda(t_2 - t_1), \lambda_2 = \lambda(t_3 - t_2), \cdots, \lambda_1 + \lambda_2 + \lambda_3 = \lambda(t_4 - t_1)$



$\lambda = .1$, $\lambda = 1$, $\lambda = 5$, $\lambda = 10$

**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.

*(handwritten)* c.f. with replacement $\Rightarrow X \sim B(r, \frac{n}{m+n})$

*(handwritten)* explanation

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

*(handwritten)* a pmf? **(Ec)**    use **(★★★★★)**

*(handwritten box)* **Note:** **(★★★★★)** $\binom{n+m}{r} = \sum_x \binom{n}{x}\binom{m}{r-x}$

- **mgf:** <u>exist</u>, but <u>no simple</u> expression

- **mean:** $\boxed{\frac{rn}{n+m}}$  ← use STO (Ec)

  (intuition) →

- **variance:** $\frac{rnm(n+m-r)}{(n+m)^2(n+m-1)}$  ← Find $E[X(X-1)]$ using STO (Ec)

- **parameter:** $r, n, m, = 1, 2, \ldots, r \leq n+m$

- **example:** <u>sampling industrial products</u> for <u>defect</u> inspection

**Notes.** a relationship between <u>hypergeometric</u> and <u>binomial</u> distributions: Let $\underline{m, n \to \infty}$ in such a way that

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

where $0 < p < 1$. Then,    (intuition): When $m, n$ are large, with replacement $\approx$ without replacement

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

---

- <u>continuous</u> distributions

> **Definition 4.9** (<u>Uniform</u> distribution $\underline{U(a, b)}$, sec 2.2)
>
> Choose a <u>number at random</u> between $\underline{a}$ and $\underline{b}$.

[shape]

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

  [a pdf? (Ec)]

  $f(x)$

  $\frac{1}{b-a}$ ————

  ____|_____
      $a$      $b$
        ↑
      $\frac{a+b}{2}$

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

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

- **mean:** $\frac{a+b}{2}$  ← { by definition / use mgf  (Ec)

  [intuition]

- **variance:** $\frac{(b-a)^2}{12}$  ← { Find $E(X^2)$ using definition / Find $E(X^2)$ using mgf  (Ec)

- **parameter:** $a, b \in \mathbb{R}$, $a < b$    [Thm 2.4, 2.5 (LNp.30)]

**Note:** $\underline{U(0,1)}$ is useful for <u>pseudo-random number</u> generation

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

*shape*

- **pdf:** $f(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$

  *a pdf ? (Ec)*

- **cdf:** $F(x) = \begin{cases} 1 - e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$ ← by definition (Ec)
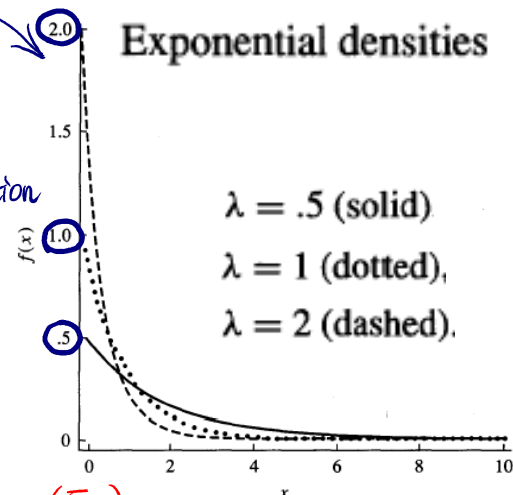
- **mgf:** $\frac{\lambda}{\lambda - t}$, $t < \lambda$. ← by definition, use STO (Ec)

- **mean:** $\frac{1}{\lambda}$ ← use STO, use mgf (Ec)

- **variance:** $\frac{1}{\lambda^2}$ ← Find $E(X^2)$ using STO, Find $E(X^2)$ using mgf (Ec)

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

- **example:** lifetime or waiting time

*meaning of parameter*
- $\frac{1}{\lambda}$: average waiting time $\left(\frac{時間}{次}\right)$
- $\lambda$: average occurence rate $\left(\frac{次}{時間}\right)$

**Exponential densities**

$\lambda = .5$ (solid)

$\lambda = 1$ (dotted),

$\lambda = 2$ (dashed).
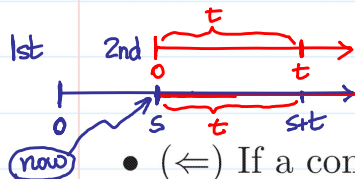
---

← **Notes:**

Ch1~6, p.2-72

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

   $T - s > t$

   $P(T > t + s | T > s) = \dfrac{P(T > t + s \text{ and } T > s)}{P(T > s)} = \dfrac{P(T > t + s)}{P(T > s)}$

   1st   2nd

   $= \dfrac{e^{-\lambda(t+s)}}{e^{-\lambda s}} = e^{-\lambda t} = P(T > t)$

   cdf of $T$: $F_T(t) = 1 - P(T > t)$

   *If discrete, then it is geometric*   *cf.*

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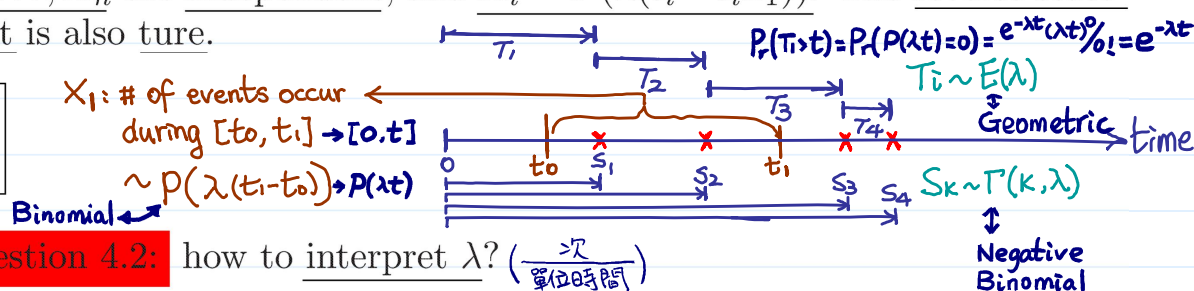   - It does not mean the two events $T > s$ and $T > t + s$ are independent.

2. relationship between exponential, gamma, and Poisson distributions

   Let $T_1, T_2, T_3, \ldots$ be i.i.d. $\sim E(\lambda)$ and $S_k = T_1 + \cdots + T_k$, $k = 1, 2, \ldots$.
   Let $X_i$ be the number of $S_k$'s that falls in $[t_{i-1}, t_i]$, $i = 1, \ldots, n$, then
   $X_1, \ldots, X_n$ are independent, and $X_i \sim P(\lambda(t_i - t_{i-1}))$. The reverse state-
   ment is also ture.

   $P_r(T_1 > t) = P_r(P(\lambda t) = 0) = e^{-\lambda t}(\lambda t)^0 / 0! = e^{-\lambda t}$

   $T_i \sim E(\lambda)$
   $\updownarrow$
   Geometric → time

   **Poisson Process**

   $X_1$: # of events occur during $[t_0, t_1] \to [0, t]$
   $\sim P(\lambda(t_1 - t_0)) \to P(\lambda t)$

   Binomial

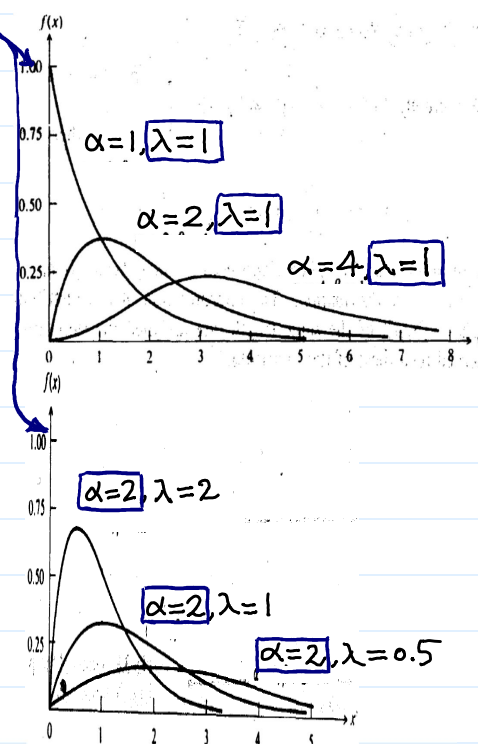   $S_k \sim \Gamma(k, \lambda)$
   $\updownarrow$
   Negative Binomial

   **Question 4.2:** how to interpret $\lambda$? $\left(\frac{次}{單位時間}\right)$

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

**Definition 4.11** (Gamma distribution $\underline{\Gamma(\alpha, \lambda)}$, sec 2.2.2)

[shape]

● **pdf:** $f(x) = \begin{cases} \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$

[a pdf? (Ec)] ← use gamma function (LNp.74)

● **mgf:** $(\frac{\lambda}{\lambda-t})^\alpha$, $t < \lambda$. ← use STO (Ec)
— sum of i.i.d. exponential

● **mean:** $\frac{\alpha}{\lambda}$ ← — use STO
— use mgf    (Ec)
[intuition]   — sum of iiid exponential

● **variance:** $\frac{\alpha}{\lambda^2}$

● **parameter:** $\alpha, \lambda > 0$

— Find $E(X^2)$ using STO
— Find $E(X^2)$ using mgf  (Ec)
— sum of iiid exponential

$f(x)$
1.00
0.75  $\alpha=1, \lambda=1$
0.50  $\alpha=2, \lambda=1$
0.25  $\alpha=4, \lambda=1$
0  1  2  3  4  5  6  7  8

$f(x)$
1.00
0.75  $\alpha=2, \lambda=2$
0.50  $\alpha=2, \lambda=1$
0.25  $\alpha=2, \lambda=0.5$
0  1  2  3  4  5

**Notes.**

1. $\underline{\alpha}$: $\underline{\text{shape}}$ parameter; $\underline{\lambda}$: $\underline{\text{scale}}$ parameter ( Question 4.3: how to $\underline{\text{inter-}}$ $\underline{\text{pret}}$ $\alpha$, $\lambda$ from the view point of $\underline{\text{Poisson process?}}$

(LNp.72) $\lambda$: occurence rate, $\alpha$: # of summed exponential r.v.'s

2. properties of $\underline{\text{gamma function}}$ $\Gamma(\alpha)$:

   • $\Gamma(\alpha) \equiv \underline{\int_0^\infty y^{\alpha-1} e^{-y} dy}$ (which is $\underline{\text{finite for } \alpha > 0}$)

   • $\underline{\Gamma(1) = 1}$ and $\underline{\Gamma(\frac{1}{2}) = \sqrt{\pi}}$

   • $\Gamma(\alpha) = \underline{(\alpha-1)\Gamma(\alpha-1)}$

   • $\Gamma(\alpha) = \underline{(\alpha-1)!}$ if $\underline{\alpha \text{ is an integer}}$

   • $\Gamma(\frac{\alpha}{2}) = \underline{\frac{\sqrt{\pi}(\alpha-1)!}{2^{\alpha-1}(\frac{\alpha-1}{2})!}}$ if $\underline{\alpha \text{ is an odd integer}}$

3. $\underline{\text{gamma distribution}}$ can be viewed as a $\underline{\text{generalization}}$ of $\underline{\text{exponential}}$ distribution, i.e., $\underline{\Gamma(1, \lambda) = E(\lambda)}$.   (Ec) [intuition]
prove using mgf

4. Let $\underline{X_1, \ldots, X_k}$ be $\underline{\text{i.i.d.}} \sim \underline{E(\lambda)}$, then $\underline{Y = X_1 + \cdots + X_k} \sim \underline{\Gamma(k, \lambda)}$.

5. Let $\underline{X_1, \ldots, X_k}$ be $\underline{\text{independent}}$, and $\underline{X_i \sim \Gamma(\alpha_i, \lambda)}$, then $\underline{Y = X_1 + \cdots + X_k} \sim \underline{\Gamma(\alpha_1 + \cdots + \alpha_k, \lambda)}$. [intuition]
prove using mgf (Ec)

6. Let $\underline{X \sim \Gamma(\alpha, \lambda)}$, then $\underline{cX} \sim \underline{\Gamma(\alpha, \lambda/c)}$, where $c > 0$. [intuition] (Ec)
prove using mgf

7. $\underline{X \sim \Gamma(\alpha, \lambda)} \Rightarrow \underline{E(X^k) = \frac{\Gamma(\alpha+k)}{\lambda^k \Gamma(\alpha)}}$, for $0 < k$ and $\underline{E(\frac{1}{X^k}) = \frac{\lambda^k \Gamma(\alpha-k)}{\Gamma(\alpha)}}$, for $\underline{0 < k < \alpha}$. — use ① STO ② mgf (Ec)   ∫ integration

**Definition 4.12** (Beta distribution $beta(\alpha, \beta)$, sec 15.3.2)

*shape*

- **pdf:** $f(x) = \begin{cases} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1}, & 0 \le x \le 1 \\ 0, & \text{otherwise} \end{cases}$
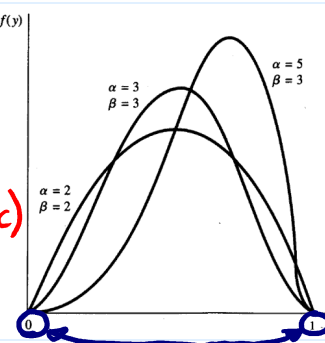
  *a pdf? (Ec)*

  *cf.* → pmf of $B(n,p) = \binom{n}{x} p^x (1-p)^{n-x}$

  $f(y)$

- **mgf:** $1 + \sum_{k=1}^{\infty} \left(\prod_{r=0}^{k-1} \frac{\alpha+r}{\alpha+\beta+r}\right) \frac{t^k}{k!}$ ← by definition

  (Note: $e^{tx} = \sum_{k=0}^{\infty} \frac{(tx)^k}{k!}$) (Ec)

- **mean:** $\frac{\alpha}{\alpha+\beta}$ ← [ use STO / use mgf (Ec) ]

  *intuition*

- **variance:** $\frac{\alpha\beta}{(\alpha+\beta+1)(\alpha+\beta)^2}$ ← [ Find $E(X^2)$ using STO / Find $E(X^2)$ using mgf ] (Ec)

- **parameter:** $\alpha, \beta > 0$

**Notes:**

① <u>Beta</u> 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. $\beta(1,1) = U(0,1)$     *meaning of $\alpha$ & $\beta$*

③ Let $X_1 \sim \Gamma(\alpha_1, \lambda), X_2 \sim \Gamma(\alpha_2, \lambda)$, and $X_1, X_2$ independent.
   Then, $\frac{X_1}{X_1+X_2} \sim beta(\alpha_1, \alpha_2)$. ← [ $Y_1 = X_1/X_1+X_2$ / $Y_2 = X_1+X_2$ ] find the joint pdf of $(Y_1, Y_2)$, then marginal pdf of $Y_1$ (Ec)

α = 5, β = 3   α = 3, β = 3   α = 2, β = 2