

MAS

---

---

---

---

---





# Chapter 4:

## Continuous random variables and Probability distribution

### 1.) Continuous random variable.

- Definition: A continuous random variables is a random variable whose possible values includes in an interval of real numbers.  
một khoảng

E.g: ) Height: 180 cm - 190 cm

) Weight: 45.6 → 70.2 Kg

### 2.) Probability density function (pdf)

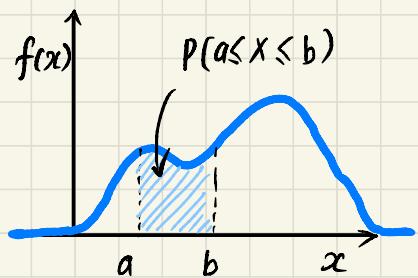
(Hàm mật độ xác suất) : Denote as:  $f(x)$

- Definition: The probability density function (pdf) of continuous random variables  $X$  is a function such that:

$$1. f(x) \geq 0, \forall x$$

$$2. \int_{-\infty}^{+\infty} f(x) dx = 1$$

$$3. P(a \leq X \leq b) = \int_a^b f(x) dx$$



So P is an area from a to b ← Total S = 1  
(area)

**Remark:** Property :  $X$  is a continuous random variable  
 (tính chất)

$$\begin{aligned} P(x_1 \leq X \leq x_2) &= P(x_1 \leq X < x_2) \\ &= P(x_1 < X \leq x_2) \\ &= P(x_1 < X < x_2) \end{aligned}$$

**Example 1:** Suppose that the probability density function of a continuous random variable  $X$  is:

$$f(x) = e^{-(x-3)}, x \geq 3$$

Determine  $P(1 \leq X < 5); P(X < 8); P(X \geq 0)$ .

(img 1)

**Example 2:** The probability density function of the length of a metal rod is

$$f(x) = cx^2 \text{ for } 2 \leq x < 3$$

- a. What is the value of  $c$ ?
- b. Find  $P(X < 2.5 \text{ or } X \geq 2.8)$

Eg1: pdf:  $f(x) = \begin{cases} e^{-(x-3)} & , \text{if } x > 3 \\ 0 & , \text{elsewhere} \end{cases}$

$$\text{a.) } P(1 \leq X < 5) = \int_a^b f(x) dx = \int_1^5 f(x) dx$$

$$= \int_3^5 e^{-(x-3)} dx = 0.865$$

$$\text{b.) } P(X < 8) = \int_3^8 e^{-(x-3)} dx = 0.993262$$

$$\text{c.) } P(X > 0) = \int_3^{+\infty} e^{-(x-3)} dx = \dots$$

$$C_2: P(X \geq 0) = 1 - P(X < 0) = 1 - 0 = 1$$

### 3.) Cumulative distribution function (cdf) (Hàm phân bố tích luỹ)

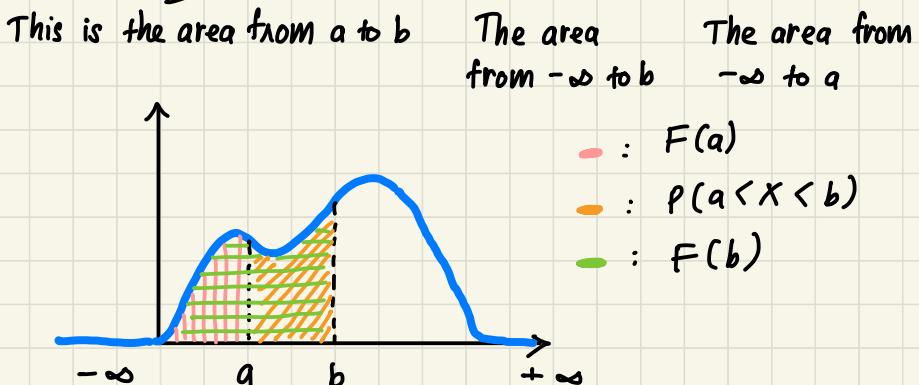
- Definition: The cumulative distribution function (cdf) of a continuous random variable  $X$  is

$$F(x) := \int_{-\infty}^x f(t) dt (= P(X \leq x) = P(X < x))$$

(for  $-\infty < x < +\infty$ )

- Remark: If  $X$  is continuous random variable with cumulative distribution function  $F(x)$  then we can use:

$$P(a < X < b) = F(b) - F(a)$$



- Remark(2):

$$\Rightarrow F'(x) = f(x)$$

$$\Rightarrow 0 \leq F(x) \leq 1, \forall x$$

( $F$  is non-decreasing :  $x < y \Rightarrow F(x) < F(y)$ )

⊕  $P(X = a) = 0$  ( $\# a$ )

⊕  $P(a < X < b) = \int_a^b f(x) dx$

⊕  $P(X > a) = 1 - P(X < a)$

### Cumulative distribution function (cdf)

Example 1: Suppose the cumulative distribution function of the random variable  $X$  is

$$F(x) = \begin{cases} 0, & \text{if } x < 1 \\ 0.5x - 0.5, & \text{if } 1 \leq x < 3 \\ 1, & \text{if } x \geq 3 \end{cases}$$

Find  $P(X < 2.8); P(0 < X < 1.5)$ ?

(img 2)

Example 2: Suppose the probability density function of the random variable  $X$  is

$$f(x) = e^{-(x-3)}, x \geq 3$$

Find the cumulative distribution function of  $X$ .

Eg 1:

$$\text{cdf: } F(x) = \begin{cases} 0, & \text{if } x < 1 \\ 0.5x - 0.5, & \text{if } 1 \leq x < 3 \\ 1, & \text{if } x \geq 3 \end{cases}$$

a.)  $P(X < 2.8) = F(2.8)$   
 $= 0.5 \times (2.8) - 0.5$   
 $= 0.9$

b.)  $P(0 < X < 1.5) = F(1.5) - F(0)$   
 $= 0.25 - 0 = 0.25$

$$\text{Eq2: pdf: } f(x) = Cx^2 \quad (2 \leq x \leq 3)$$

(img 1)

a.)  $\int_{-\infty}^{+\infty} f(x) dx = 1$

$$\Rightarrow \int_2^3 f(x) dx = 1 \Leftrightarrow C \int_2^3 x^2 dx = 1$$

$$\Leftrightarrow C \cdot \frac{19}{3} = 1 \Rightarrow C = \frac{3}{19}$$

b.)  $P(X < 2.5 \text{ or } X \geq 2.8)$

$$= P(X < 2.5) + P(X \geq 2.8)$$

$$= \int_2^{2.5} \frac{3}{19} x^2 dx + \int_{2.8}^3 \frac{3}{19} x^2 dx$$

Eq2 (img 2)

$$C_1: \text{pdf: } f(x) = e^{-(x-3)}, x \geq 3$$

$$\text{We have } F'(x) = f(x) = e^{-(x-3)}$$

$$\Rightarrow F(x) = C + [-e^{-(x-3)}]$$

$\Rightarrow$  check the answer

$$C_2: F(x) = \int_{-\infty}^x f(t) dt = \int_{-\infty}^x e^{-(t-3)} dt$$

...

#### 4.) Mean and Variance of a continuous random variable

- **Definition:** Suppose  $X$  is a continuous random variable with probability density function  $f(x)$ .

The  $\Rightarrow$  Mean or expected value of  $X$  is defined by

$$\mu = E(X) := \int_{-\infty}^{+\infty} x f(x) dx$$

$\Rightarrow$  The Variance of  $X$  is defined by

$$\begin{aligned}\sigma^2 &= V(X) := \int_{-\infty}^{+\infty} (x - \mu)^2 f(x) dx \\ &= \int_{-\infty}^{+\infty} x^2 f(x) dx - \mu^2\end{aligned}$$

$\Rightarrow$  The standard deviation of  $X$  is:

$$\sigma = \sqrt{V(X)}$$

**Example 1:** Assume that  $X$  is a continuous random variable with the following probability density function

$$f(x) = \frac{x^2}{18}, \text{ for } -3 < x < 3$$

Determine the mean and variance of  $X$ .

**Example 2:** The cumulative distribution function of the random variable  $X$  is

$$F(x) = \begin{cases} 0, & \text{if } x < 1 \\ 0.5x - 0.5, & \text{if } 1 \leq x < 3 \\ 1, & \text{if } x \geq 3 \end{cases}$$

Find the standard deviation of  $X$

Eg1: pdf:  $f(x) = \frac{x^2}{18} \quad (-3 < x < 3)$

$$\Rightarrow M = E(X) = \int_{-3}^3 x \cdot \frac{x^2}{18} dx = 0$$

$$\Rightarrow \sigma^2 = V(X) = \int_{-3}^3 x^2 \cdot \frac{x^2}{18} dx - M^2 = \frac{27}{5}$$

Eg2: cdf:  $F(x) := \begin{cases} 0, & \text{if } x < 1 \\ 0.5x - 0.5, & \text{if } 1 \leq x < 3 \\ 1, & \text{if } x \geq 3 \end{cases}$

$$f(x) = F'(x) = 0.5 \quad (1 \leq x < 3)$$

$$\Rightarrow M = \int_1^3 x \cdot 0.5 dx = 2$$

$$\Rightarrow \sigma^2 = V(X) = \int_1^3 x^2 \cdot 0.5 dx - M^2 = \frac{1}{3}$$

$$\Rightarrow \sigma = \sqrt{V(X)} = \sqrt{\frac{1}{3}}$$

## 5.) Continuous uniform distribution

- Definition: Suppose  $X$  has a continuous uniform distribution over the interval  $[a, b]$  if

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

- Mean & Variance:

$$\begin{aligned} + M = E(X) &= \int_{-\infty}^{+\infty} x \cdot f(x) dx \\ &= \int_a^b \frac{1}{b-a} \cdot x dx = \frac{a+b}{2} \\ + \sigma^2 = V(X) &= E(X^2) - (E(X))^2 \\ &= \frac{(b-a)^2}{12} \end{aligned}$$

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

**Example 1:** Suppose X has a continuous uniform distribution over the interval [1;10].

- Determine the mean, variance and standard deviation of X.
- Find  $P(X < 6.5)$ .
- Determine the cumulative distribution function of X.

**Example 2:** Suppose X has a continuous uniform distribution over [5;15]. What is the mean and variance of  $Y=8X$ ?

**Eg 1:** X has continuous uniform distribution on [1;10]

$$\Rightarrow [a, b] = [1, 10] \Rightarrow \begin{cases} a = 1 \\ b = 10 \end{cases}$$

a.)

$$\hat{)} \text{ Mean} = \mu = \frac{a+b}{2} = 5.5$$

$$\hat{)} \text{ Variance} = \sigma^2 = \frac{(b-a)^2}{12} = 6.75$$

$$\hat{)} \text{ Standard deviation} = \delta = \sqrt{\delta^2} = 2.6$$

b.)

pdf:  $f(x) = \begin{cases} \frac{1}{b-a} = \frac{1}{9}, & \text{if } x \in [1, 10] \\ 0, & \text{if } x \notin [1, 10] \end{cases}$

$$\Rightarrow P(X < 6.5) = \int_1^{6.5} f(x) dx = 0.6(1)$$

c.) cdf:  $F(x) := \begin{cases} 0, & \text{if } x < 1 \\ \frac{x-1}{9}, & \text{if } 1 \leq x < 10 \\ 1, & \text{if } x \geq 10 \end{cases}$

Eg2:  $Y = 8X$

$$\hat{)} E(Y) = E(8X) = 8E(X) = 8 \cdot \frac{a+b}{2}$$

$$\hat{)} V(Y) = V(8X) = 8^2 V(X)$$

$$\hat{)} \delta(Y) = \sqrt{V(Y)}$$

## 6.) Normal distribution (Phân phối chuẩn)

- Definition:  $X$  has a normal random variable with parameters  $\mu$  and  $\sigma^2$  if

( Notation:  $X = N(\mu, \sigma^2)$  )

- pdf:  $f(x) = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ ,  $-\infty < x < +\infty$

When  $-\infty < \mu < +\infty$  and  $\sigma > 0$

- Mean:  $\mu = E(X)$

- Variance:  $V(X) = \sigma^2$

- Standard deviation:  $\sigma = \sqrt{V(X)}$

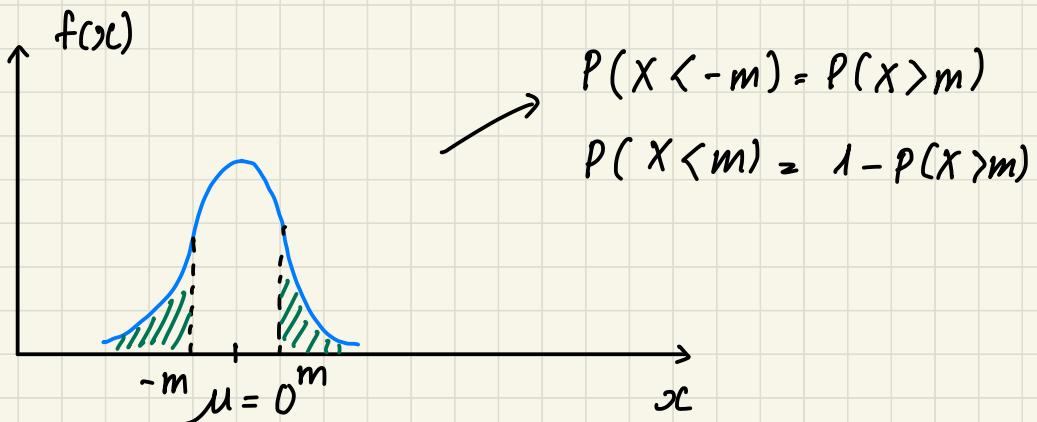
## 7.) Standard normal distribution (Phân phối chuẩn hóa)

- Notation:  $Z = N(0, 1)$

$\Rightarrow \mu = 0$  and  $\sigma^2 = 1 \Rightarrow \sigma = 1$

- pdf:  $f(x) = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{x^2}{2}}$ ,  $-\infty < x < +\infty$

- cdf:  $\Phi(z) = \int_{-\infty}^z \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}} dx$



### 8.) Normal Standardizing (Chuẩn hóa)

- Definition: If  $X \sim N(\mu, \sigma^2)$  then

$$Z = \frac{X - \mu}{\sigma}$$

is a standard normal random variable  $N(0, 1)$

$$(X \sim N(\mu, \sigma^2) \Rightarrow Z = \frac{X - \mu}{\sigma} \sim N(0, 1))$$

$$\begin{aligned} \text{D) } P(X < x) &= P\left(\frac{X - \mu}{\sigma} < \frac{x - \mu}{\sigma}\right) \\ &= P(Z < \frac{x - \mu}{\sigma}) \end{aligned}$$

$$\text{D) } P(a < X < b) = P(X < b) - P(X < a)$$

Eg: The weekly salaries are normally distributed with a mean of \$490 and standard deviation of \$45.

$$\text{Find } P(\text{a teacher earns} > \$525)$$

Solution:  $\mu = 190$ ,  $\delta = 45$

$$\begin{aligned} P(X > 525) &= 1 - P(X < 525) \\ &= 1 - P\left(Z < \frac{525 - 190}{45}\right) \\ &= 1 - P\left(Z < \frac{7}{9}\right) \\ &= 0.2183 \end{aligned}$$

**Example 1:** A machine pours beer into 16 oz. bottles. Experience has shown that the number of ounces poured is normally distributed with a standard deviation of 1.3 ounces. Find the probabilities that the amount of beer the machine will pour into the next bottle will be more than 16.25 ounces.

**Example 2:** The tread life of a particular brand of tire is a random variable best described by a normal distribution with a mean of 65,000 miles and a standard deviation of 1,500 miles. What warranty should the company use if they want 95% of the tires to outlast the warranty?

Eg1:  $\mu = 16$ ;  $\delta = 1.3 \Rightarrow \delta^2 = 1.69$   
 $\Rightarrow X \sim N(16, 1.69)$

$$\begin{aligned} P(X > 16.25) &= 1 - P(X < 16.25) \\ &= 1 - P\left(Z < \frac{5}{26}\right) = 0.42375 \end{aligned}$$

Eg2:  $\mu = 65,000$ ;  $\delta = 1,500$

$X$ : the tread life of a brand of tire

$x$ : the warranty

$$P(X > x) = 0.95$$

$$\Rightarrow P(X < x) = 0.05$$

$$\Leftrightarrow P\left(\frac{Z < \frac{x - 65,000}{1500}}{1500}\right) = 0.05$$

$$\Rightarrow \frac{x - 65,000}{1500} = -1.645$$

$$\Rightarrow x = 62,532.5 \text{ (miles)}$$

### g.) Approximation using Normal Distribution

a. Normal approximation to the Binomial distribution  
(Xấp xỉ phân phối nhị thức bởi phân phối)

**Definition:** If  $X$  has a Binomial distribution  $B(n, p)$  then random variable:

$$Z = \frac{X - n \cdot p}{\sqrt{n \cdot p \cdot (1-p)}} \quad \begin{cases} E(Z) = n \cdot p \\ \sigma(Z) = \sqrt{n \cdot p \cdot (1-p)} \end{cases}$$

$\approx N(0, 1)$

is approximately a standard random variable  $(0, 1)$

We have:

$$P(X \leq x) \approx P(Z \leq \frac{x + 0.5 - n \cdot p}{\sqrt{n \cdot p \cdot (1-p)}})$$

$$P(X \geq x) \approx P(Z \geq \frac{x - 0.5 - n \cdot p}{\sqrt{n \cdot p \cdot (1-p)}})$$

**Remark:** The approximation is good for  $n \cdot p > 5$  and  $n(1-p) > 5$

**Example:** The manufacturing of semiconductor chips produces 3% defective chips. Assume the chips are independent and that a lot contains 800 chips. Approximate the probability that more than 30 chips are defective.

**Hint:**  $X \sim B(n = 800; p = 0.03)$ . Use Normal approximation.

**Remark:** Can use BINOM.DIST in Excel to find actual value.

$$\text{Eg: } X \sim B(n = 800; p = 0.03)$$

$$\begin{aligned} P(X > 30) &= 1 - P(X \leq 30) \\ &\approx 1 - P(Z \leq \frac{30 + 0.5 - 800 \times 0.03}{\sqrt{800 \times 0.03 \times 0.97}}) \\ &= 1 - P(Z \leq 1.347) \\ &= 0.9111 \end{aligned}$$

### b.) Normal Approximation to the Poisson Distribution

- \*  $\lambda$  : mean number of events in a unit interval
- $X$  : number of event in a unit interval

**Definition:** If  $X$  has Poisson distribution  $P(\lambda)$  then random variable :

$$Z = \frac{X - \lambda}{\sqrt{\lambda}}$$

$$\left\{ \begin{array}{l} E(X) = \mu = \lambda \\ \sigma^2 = V(X) = \lambda \end{array} \right.$$

is approximately a standard random variable  $N(0,1)$

We have :

$$\Rightarrow P(X \leq x_c) \approx P(Z \leq \frac{x_c + 0.5 - \lambda}{\sqrt{\lambda}})$$

$$\Rightarrow P(X \geq x_c) \approx P(Z \geq \frac{x_c - 0.5 - \lambda}{\sqrt{\lambda}})$$

**Remark:** The approximation is good for  $\lambda > 5$

Example: The number of customers that arrive at a fast-food business during a one-hour period is known to be Poisson distributed with a mean equal to 9.6. What is the probability that more than 10 customers will arrive in a one-hour period?

**Hint:**  $X \sim P(9.6)$ . Use Normal approximation

**Remark:** Can use POISSON.DIST in Excel to find actual value.

Eg:  $\mu = 9.6 = \lambda$

$$P(X > 10) = 1 - P(X \leq 10)$$

$$\approx 1 - P(Z \leq \frac{10 + 0.5 - 9.6}{\sqrt{9.6}})$$

$$= 0.386$$

## 10.) Exponential distribution (phân phối mũ)

- [
  - $\lambda$ : mean number of event in a unit interval
  - $X$ : the distance between two consecutive event

**Definition:** The random variable  $X$  that equals the distance between two consecutive events with mean number of event  $\lambda > 0$  per unit interval is an exponential random variable with parameter  $\lambda$ . The probability density function of  $X$  is:

$$\text{pdf: } f(x) = \lambda e^{-\lambda x}, x \geq 0$$

$$\text{cdf: } F(x) = 1 - e^{-\lambda x}, x \geq 0$$

**Remark:** If random variable  $X$  has exponential distribution with parameter  $\lambda$  then  $M = E(X) = \frac{1}{\lambda}$ ;  $S^2 = V(X) = \frac{1}{\lambda^2}$

**Example 1:** The time between customer arrivals at a furniture store has an approximate exponential distribution with mean of 9 minutes. If a customer just arrived, find the probability that the next customer will not arrive for at least 15 minutes.

**Example 2:** The time between patients arriving at an outpatient clinic follows an exponential distribution at a rate of 15 patients per hour. What is the probability that a randomly chosen arriving interval will not exceed 6 minutes?

**Eg1:**  $X = \text{the time between customer arrivals}$   
 $\Rightarrow X \sim \exp(\lambda)$

We have mean of  $X = 9$  minute  $\Rightarrow \mu = \frac{1}{\lambda} = 9$

$$\Rightarrow \lambda = \frac{1}{9} \text{ (arrivals / minute)}$$

$$\begin{aligned} P(X > 15) &= 1 - P(X \leq 15) \\ &= 1 - \int_0^{15} \frac{1}{9} \cdot e^{-\frac{1}{9}x} dx \\ &= 0.189 \end{aligned}$$

Eg2:  $X$ : the time between arrivals (in hours)

$$\lambda = 15 \text{ (patient / hour)}$$

$$\begin{aligned} P(X \leq 0.1) &= \int_0^{0.1} 15 \cdot e^{-15x} dx \\ &= 0.777 \end{aligned}$$

# Chapter 6:

## Data description

Introduction to statistic

Statistic



Descriptive Statistic  
(Thống kê mô tả)

Use numerical summaries  
and visual displays to  
describe data  
(Chapter 6)



Inferential Statistic  
(Thống kê suy luận)

Use information from samples  
(data) to estimate for  
Population  
(Chapter 8, 9, 10, 11)

1.) Numerical summarise of data

a.) Sample mean (Trung bình mẫu)

Definition: If the  $n$  observations in a sample are denoted by  $x_1, x_2, \dots, x_n$ , the sample mean is

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

( $n$  is the size of sample)

**Example:** Let's consider the weight of the eight observations collected from the prototype engine connectors: 12.6, 12.9, 13.4, 12.3, 13.6, 13.5, 12.6 and 13.1

Find the sample mean.

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{12.6 + 12.9 + \dots + 13.1}{8} = 13.0$$

### b.) Sample median (Trung vị mẫu)

**Definition:** - The value that lies in the middle of the data when the data set is ordered.

- Measures the center of an ordered data set by dividing it into two equal part.
- If the data set has an:
  - (a) even number of entries: median is the average of the two middle data entries.
  - (b) odd number of entries: median is the middle data entries.

**Example:** The prices (in dollars) for a sample of round-trip flights from Chicago, Illinois to Cancun, Mexico are listed. Find the median of the flight prices.

872 432 397 427 388 782 397

**Eg:** Order the data:

388, 397, 397, 427, 432, 782, 872

The sample median is 427.

### c.) Sample mode

- The data entry that occurs with the greatest frequency
- If no entry is repeated, the data set has no mode
- If two entries occur with the same greatest frequency, each entry is a mode (bimodal)

### d.) Sample variance and sample standard deviation

- If  $x_1, x_2, \dots, x_n$  is a sample of  $n$  observations, the sample variance is

$$S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$
$$= \frac{\sum_{i=1}^n x_i^2 - \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2}{n-1}$$

(  $S$  = sample standard deviation )

Example: Let's consider the weight of the eight observations collected from the prototype engine connectors: 12, 13, 9, 12, 10 and 12.

Find the sample standard deviation.

## e.) Sample range (Khoảng biến thiên)

$$\text{- Sample range} = R = \max\{x_i\} - \min\{x_i\}$$

## 2.) Stem and leaf diagram (Sơ đồ thân và lá)

### Stem and leaf diagram

A **stem-and-leaf diagram** is a good way to obtain an informative visual display of a data set where each number  $x_i$  consists of at least two digits. To construct a stem-and-leaf diagram, use the following steps:

- Divide each number  $x_i$  into two parts: a stem, consisting of one or more of the leading digits, and a leaf, consisting of the remaining digit.
- List the stem values in a vertical column.
- Record the leaf for each observation beside its stem.
- Write the units for stems and leaves on the display.

\* Remark: - Có bao nhiêu số thì có bấy nhiêu lá'

⇒ the number of stems  $\leq n$

⇒ the number of leaves  $= n$

\* Note: If leaf unit = 0.1

$\Rightarrow 3|6$  represent 3.6

unit = 100 ( $6 \times 100 = 600$ ) ↙ Put 3 in here

$\Rightarrow 3|6$  represent 3600

## Stem and leaf diagram

Example: The listening scores of 12 students in a TOEIC test are listed below

55 115 225 240 330 335 385 400 405 405 495 495

Stem	Leaves
5	5
11	5
22	5
24	0
33	0 5
38	5
40	0 5 5
49	5 5

The stem and leaf diagram:

## 3.) Box-plots (Biểu đồ hộp)

phân tử

### Three quartiles

An ordered set of data is divided into four equal parts, the division points are called **quartiles**:

- The **first quartile**,  $q_1$  or  $Q_1$ : is a value that has approximately 25% of the observations below.
- The **sample median** or **second quartile**,  $q_2$  or  $Q_2$ , has approximately 50% of the observations below its value.
- The **third quartile**,  $q_3$  or  $Q_3$ , has approximately 75% of the observations below its value.
- The **interquartile range**,  $IQR = Q_3 - Q_1$  (chênh lệch từ phân vị)

Example: Use the given sample data to find the sample quartiles, the sample mode and the IQR.

55, 52, 52, 52, 49, 74, 67, 55.

\* Way to determine  $Q_1$ ,  $Q_2$ ,  $Q_3$ :

- 1.) Find  $Q_2$  = sample median ( $Q_2$  là trung vi mẫu)
- 2.) Find  $Q_1$  = sample median of the  $Q_2$ 's left part
- 3.) Find  $Q_3$  = \_\_\_\_\_ right part
- 4.) Compute IQR (interquartile range = chênh lệch tứ phân vị)

$$IQR = Q_3 - Q_1$$

- 5.) Compute  $\begin{cases} Q_1 - 1.5 * IQR & (1) \\ Q_3 + 1.5 * IQR & (2) \end{cases}$

6.) Outlier (phân tử ngoại lai) : an entry is called outlier if it smaller than (1) or greater than (2)

7.) Minimum : is the smallest entry (not an outlier)

Maximum : is the biggest entry (not an outlier)

Eg: Data set: 49, 52, 52, 52 | 55, 55, 67, 71

$$Q_2 = \text{sample median} = \frac{52 + 55}{2} = 53.5$$

$$Q_1 = \frac{52 + 52}{2} = 52$$

$$Q_3 = \frac{55 + 67}{2} = 61$$

$$\text{IQR} = Q_3 - Q_1 = 61 - 52 = 9$$

$$Q_1 - 1.5 * \text{IQR} = 38.5$$

$$Q_3 + 1.5 * \text{IQR} = 74.5$$

$\Rightarrow$  doesn't have outlier

$$\Rightarrow \begin{cases} \text{Minimum} = 49 \\ \text{Maximum} = 71 \end{cases}$$

## 4.) Frequency distribution (phân bố tần suất)

### Frequency Distribution

Construction of frequency distribution: divide the range of the data into intervals (called class intervals, cells, or bins). The bins should be of equal width.

#### Example:

Data = Grades = {2.4, 4.4, 4.6, 5.0, 5.0, 5.8, 6.0, 7.4, 8.2, 9.0}

- Divide grade ranges into 5 bins:  
0 - 2, 2 - 4, 4 - 6, 6 - 8, 8 - 10.
- Count the number of data values in each bin:

Bin	Frequency
0 - 2	0
2 - 4	1
4 - 6	6
6 - 8	1
8 - 10	2

### \* Remark:

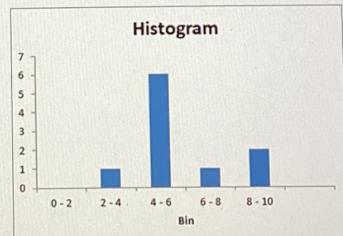
f) frequency: tần số ( $n_i$ )

f) relative frequency: tần suất ( $f_i = \frac{n_i}{n}$ )

# Histogram

The **histogram** is a visual display of the frequency distribution.

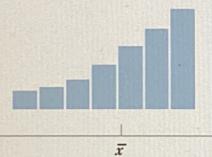
- Label the bin (class interval) boundaries on a horizontal scale.
- Mark and label the vertical scale with the frequencies or the relative frequencies.
- Above each bin, draw a rectangle where height is equal to the Frequency (or relative frequency) corresponding to that bin.



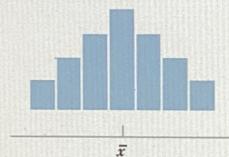
# Histogram

Remark:

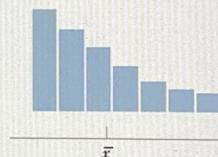
1. Histograms are very useful to explore the **distribution** of data.



Negative or left skew

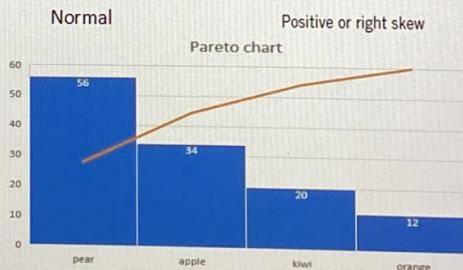


Normal



Positive or right skew

2. Pareto chart:  
(frequencies are ordered decreasingly)



## Times sequence plots

A **time series plot** is a graph in which the vertical axis denotes the observed value of the variable (say,  $x$ ) and the horizontal axis denotes the time (which could be minutes, days, years, etc.)

Years    Annual salary

Years	Annual salary
2005	20614
2006	21175
2007	24766
2008	28100
2009	30189
2010	24618
2011	22006
2012	30912
2013	32523
2014	35285



# Chapter 7:

## Sampling distribution and Point estimates of parameter

### 1. Point estimates of parameters.

( Các ước lượng điểm của các tham số )

Parameters of population  
( Các t特sô' cua tông the' )

$\mu$ : population mean

$\sigma^2$ : variance

$\delta$ : standard deviation

$p$ : population proportion  
( tí lê cua tông the' )

Statistics of a sample  
( Các thông kê cua mâu )

$\bar{x}$ : sample mean

$s^2$ : sample variance

$s$ : " standard deviation

$\hat{p}$ : sample proportion  
( tí lê mâu )

Definition: A point estimate for a parameter is a single value computed from a sample

⇒ The best point estimate for  $\mu$  is the  $\bar{x}$

$\sigma^2$  is the  $s^2$

$\delta$  is the  $s$

**Example 1:** In a sample of 73 products, there are 7 defective products. So a point estimate for proportion  $p$  of all defective products is:

$$\hat{p} = \frac{x}{n} = \frac{7}{73} = 9.6\%$$

**Example 2:** Market researchers use the number of sentences per advertisement as a measure of readability for magazine advertisements. The following represents a random sample of the number of sentences found in 15 advertisements.

9 20 18 16 9 9 11 13 22 16 5 18 6 6 5

- Find a point estimate of the population mean,  $\mu$ .
- Find a point estimate of the population standard deviation,  $\sigma$ .

**Hint:**

$$\bar{x} = \frac{\sum x_i}{n} = 12.2; s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = 5.77$$

**Eg1:** A point estimate for proportion  $p$  of all defective product is  $\hat{p}$

$$\hat{p} = \frac{x}{n} = \frac{7}{73} = 0.096 = 9.6\%$$

**Eg2:** a.) A p/e of the population mean  $\mu$  is  $\bar{x}$

$$\bar{x} = \frac{\sum x_i}{n} = 12.2$$

b.) A p/e of the population standard deviation  $\sigma$  is  $s$

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1} = 5.77$$

$$\Rightarrow s = \sqrt{s^2}$$

## 2.) Central limit theorem (Định lý giới hạn trung tâm)

Theorem 1. (CLT for 1 population)

- Consider a population with mean  $= \mu$ ; variance  $= \sigma^2$
- A random sample of size  $n$ :  $X_1, X_2, \dots, X_n$

$$\text{Sample mean: } \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

Then: 1) The sampling distribution (phân bố mẫu) of sample mean is

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right), \text{ for } n \geq 30$$

$\downarrow$   
normal distribution

$\Rightarrow$  It means that mean  $= \mu$ , variance of  $\bar{X}$  is  $\frac{\sigma^2}{n}$

Example: One year, professional players salaries averaged 1.5 million with a standard deviation of 0.9 million. Suppose a sample of 100 players was taken. Find the approximate probability that the average salary of these 100 players does not exceed 1.4 million.

Eg: Population mean  $= \mu = 1.5$ ; standard deviation  $= \sigma = 0.9$

Size of a sample:  $n = 100$ .

$$P(\bar{X} < 1.4) = P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < \frac{1.4 - \mu}{\sigma/\sqrt{n}}\right)$$

$$= P(Z < \frac{1.4 - 1.5}{0.9/\sqrt{100}}) = P(Z < -1.11) = 0.1333$$

$$\text{Remark: } \textcircled{1} Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

↑  
Phân bố chuẩn hóa  
(n > 30)

⇒ If the population already has normal distribution then for any sample size:

$$\bar{X} \sim N(\mu; \frac{\sigma^2}{n})$$

**Example:** An electrical firm manufactures light bulbs that have a length of life that is approximately normally distribution, with mean equal to 800 hours and a standard deviation of 40 hours. Find the probability that a random sample of 16 bulbs will have an average life of less than 775 hours.

**Eg 2:** Suppose that population is normal with mean  $\mu = 800$ ,  
Standard deviation  $= \sigma = 40$ .

Sample size :  $n = 16$ .

We have :

$$\begin{aligned} P(\bar{X} < 775) &= P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < \frac{775 - \mu}{\sigma/\sqrt{n}}\right) \\ &= P(Z < \frac{775 - 800}{40/\sqrt{16}}) = P(Z < -2.5) \\ &= 0.006 \quad (\text{Don't need "n > 30"}) \end{aligned}$$

\* **Chú ý:**

**Chương 4:**  $P(X < b) = P\left(\frac{X - \mu}{\sigma} < \frac{b - \mu}{\sigma}\right)$

X : is normal random variable,  $X \sim N(\mu, \sigma^2)$

$$\text{Choosing } F: P(\bar{X} < b) = P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < \frac{b - \mu}{\sigma/\sqrt{n}}\right)$$

$$= P\left(Z < \frac{b - \mu}{\sigma/\sqrt{n}}\right)$$

$\bar{X}$  is sample mean of a random sample taken from a population with mean =  $\mu$ , Standard deviation =  $\sigma$   
 $n$ : sample size

Theorem 2: (CTL for two population)

Population 1	Population 2
mean: $\mu_1$	mean: $\mu_2$
Standard deviation: $\sigma_1$	— — — ; $\sigma_2$
Sample size: $n_1$	Sample size: $n_2$
Sample mean: $\bar{X}_1$	Sample mean: $\bar{X}_2$

We have:

$$\bar{X}_1 - \bar{X}_2 \sim N\left(\underbrace{\mu_1 - \mu_2}_{\text{mean}}, \underbrace{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}_{\text{variance}}\right)$$

(for  $n_1 \geq 30, n_2 \geq 30$ )

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1)$$

$$\text{Note: } P\left(\overline{X}_1 - \overline{X}_2 < b\right) = P\left(\frac{\overline{X}_1 - \overline{X}_2 - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} < \frac{b - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}\right)$$

⇒ Các xs  $P(\overline{X}_1 - \overline{X}_2 > a)$  or  $P(a < \overline{X}_1 - \overline{X}_2 < b)$   
 điều kiện trong đó

Example: The television picture tubes of manufacturer A have a mean lifetime of 6.5 years and a standard deviation of 0.9 year, while those of manufacturer B have mean lifetime of 6.5 years and a standard deviation of 0.8 year. What is the probability that a random sample of 36 tubes from manufacture A will have mean lifetime that is at least 1 year more than the mean lifetime of a sample of 49 tubes from manufacturer B?

Population 1	Population 2
$\mu_1 = 6.5$	$\mu_2 = 6.5$
$\sigma_1 = 0.9$	$\sigma_2 = 0.8$
$n_1 = 36$	$n_2 = 49$

Eg: Population 1 (Manufacture A)  $\mu_1 = 6.5$ ,  $\sigma_1 = 0.9$ ,  $n_1 = 36$

Population 2 (Manufacture B)  $\mu_2 = 6.5$ ,  $\sigma_2 = 0.8$ ,  $n_2 = 49$

We have

$$P(\overline{X}_1 - \overline{X}_2 > 1) = P\left(\frac{\overline{X}_1 - \overline{X}_2 - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} >\right)$$

$$\geq \frac{1 - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \Big)$$

$$= P(Z \geq \frac{1 - (6.5 - 6.5)}{\sqrt{\frac{0.9^2}{36} + \frac{0.8^2}{49}}}) = P(Z \geq 5.302)$$

$$= 1 - P(Z < 5.302) \approx 1 - 1 = 0$$

## MAS PT2: (Ryu Yamada)

- 1.)  $P(X = 2018) = 0$ . Because  $P(X = a) = 0$  ( $\#a$ ) (D)
- 2.)  $P(\bar{X} > \frac{1458}{86}) = P(\bar{X} > 40.5) = 1 - P(\bar{X} \leq 40.5)$   
 $= 1 - P(Z \leq \frac{40.5 - 40}{\frac{2}{\sqrt{50}}})$   
 $= 1 - P(Z \leq 1.5) = 0.0668$   
 $\approx 0.067$  (A)
- 3.) Standard error =  $\frac{\sigma}{\sqrt{n}} = \frac{30}{\sqrt{86}} = 10$  (E)

4.) Cumulative relative frequency: fan so' furing atsi' ticks huy

$$\Rightarrow \text{the number of...} = [P(21-30) - P(11-20)] * 80$$

$$= (72.5\% - 27.5\%) * 80$$

$$= 35 \text{ (members)} \quad \text{span style="color: blue;">(C)}$$

5.) 32 numbers

$$Q_2 = \frac{x_{16} + x_{17}}{2} = 114 \quad \left. \begin{array}{l} \\ \end{array} \right\} \Rightarrow \text{span style="color: blue;">(B)}$$

$$Q_1 = \frac{x_8 + x_9}{2} = 108$$

$$Q_3 = \frac{x_{24} + x_{25}}{2} = 118$$

6.) Mode number (the number has the biggest frequency)  
= 16 (2 times) B

7.)  $\mu = 2.6$ ,  $\sigma = 0.4$

$$\begin{aligned} P(2.2 < X < 3) &= P(X < 3) - P(X < 2.2) \\ &= P(Z < \frac{3-2.6}{0.4}) - P(Z < \frac{2.2-2.6}{0.4}) \\ &= P(Z < 1) - P(Z < -1) \\ &= 0.68 = 68\% \quad \text{A} \end{aligned}$$

8.) continuous uniform distribution  $[1, 5]$

$$\Rightarrow \mu = \frac{a+b}{2} = 3$$

$$\Rightarrow \sigma^2 = \frac{(b-a)^2}{12} \Rightarrow \sigma = 1.1547.$$

C

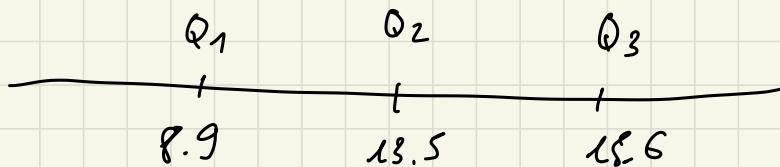
$$\begin{aligned} 9.) \quad P(X < 980) &= 1 - P(X \geq 950) \\ &\approx 1 - P(Z \geq \frac{950-0.5-\lambda}{\sqrt{\lambda}}) \\ &= 1 - P(Z \geq 1.5969) \\ &= 0.8548 \quad \text{D} \end{aligned}$$

10.)  $\lambda = 5$  uniform distribution  $[4.5, 7.5]$   
 $a \quad b$

$$\begin{aligned}
 P(5.0 < X < 6.0) &= P(X < 6.0) - P(X < 5.0) \\
 &= \int_0^6 f(x) dx - \int_0^5 f(x) dx \\
 &= \int_0^6 \frac{1}{b-a} dx - \int_0^5 \frac{1}{b-a} dx = \frac{1}{3}
 \end{aligned}$$

(C)

13.)  $\bar{x} = \frac{\text{Sum}}{1500} \Rightarrow \text{Sum} = 21300 \text{ in } /$



$$IQR = Q_3 - Q_1 = 6.7$$

= 0

$$Q_1 - 1.5 IQR = -1.15 \quad ) \text{ 2 kisay}$$

$$Q_3 + 1.5 \times IQR = 25.65$$

=) (E)

$$11.) \quad \begin{array}{c|cc} & 8 & 8 \\ g & 1 & 8 \\ 10 & 5 & 5 \end{array} \Rightarrow 85, 88, 91, 98, 105, 105$$

(C)

$$12.) \quad \text{expected value} = \text{mean} = \mu = \int_{-\infty}^{+\infty} xf(x) dx$$

$$= 0.8 \quad \text{(A)}$$

$= 1$   
 $\int_{-\infty}^{+\infty} f(x) dx = 0$

$$14.) \quad \text{exponential distribution. } \mu = 3 \Rightarrow \lambda = \frac{1}{3}$$

$$P(X > a) \stackrel{? \text{ time}}{=} 0.1$$

$$\Rightarrow P(X \leq a) = 0.9 = \int_{-\infty}^{+\infty} f(x) dx$$

$$\Leftrightarrow \int_0^a \frac{1}{3} \cdot e^{-\frac{1}{3}x} dx = 0.9 \quad (0.899)$$

$$\Rightarrow a \approx 6.9 \quad (\text{thúk tápán}) \quad \text{(D)}$$

MAS PT<sub>2</sub>: Nguyên tử

1) \* (Đã làm) (D)

$$\begin{aligned}2.) P(-1 < X < 2) &= F(2) - F(-1) \\&= 1 - \left[ \frac{1}{2} (-1)^3 + \frac{1}{2} \right] \\&= 1\end{aligned}$$

$$\begin{aligned}3.) P(X \leq 31) &\approx P(Z \leq \frac{31 + 0.5 - 4.1}{\sqrt{n \cdot p \cdot (1-p)}}) \\&= P(Z \leq 0.327) \\&= 0.6293\end{aligned}$$

4.) \* (C)

$$5.) \text{ex dis } \mu = 3.4 \Rightarrow \lambda = \frac{1}{3.4}$$

$$\begin{aligned}P(X > 5) &= 1 - P(X \leq 5) \\&\approx 1 - \int_0^5 \lambda \cdot e^{-\lambda x} dx \\&= 1 - 0.77 = 0.229790\end{aligned}$$

$$6.) \text{mean} = \frac{4 + 14 + 3 + 16 + 9 + 8 + 16}{7} = 10$$

(D)

$$7.) \quad F(x) = 1 - e^{-x/5} \Rightarrow f(x) = F'(x)$$

$$= \frac{1}{5} e^{-x/5}$$

$$\Rightarrow \mu = E(X) = \int_0^{+\infty} x \cdot \frac{1}{5} e^{-x/5} dx \approx 5$$

(C)

$$8.) \quad P(X < 48) = P(Z < \frac{48-45}{6})$$

$$= P(Z < \frac{1}{2})$$

$$= 0.6915 \quad (\text{C})$$

9.) \* (A)

$$10.) \quad \bar{x} ? \quad n = 15, \mu = 6, \delta = 2.5$$

$\bar{x} = \text{mean}$

$$\text{Variance} = \frac{\delta^2}{n} \Rightarrow \text{(B)}$$

11.) (B)

$$12.) \quad \mu = 1; \quad \delta = 0.1$$

$$P(X > a) = 0.025$$

$$\Rightarrow P(X \leq a) = 0.975$$

$$\Rightarrow P(Z \leq \frac{a-1}{0.1}) = 0.975$$

$$\Rightarrow a = 1.05649$$

(C)

$$13) P(X < 1.5) = \int_0^{1.5} 0.5 dx = 0.75 \quad \textcircled{B}$$

$$14.) \quad \textcircled{A}$$

$$15.) P(X < 0 \text{ or } X > 0.5) = 1 - P(0 \leq X \leq 0.5)$$
$$= 1 - \int_0^{0.5} 1.5x^2 = 0.9375 \quad \textcircled{A}$$

MAS PT2: Kết phán hùng

$$\int_{-\infty}^{\lambda} \lambda \cdot e^{-\lambda x} dx$$

1.)  $\mu = 3.4 \Rightarrow P(X < 1) = \int_0^1 \frac{1}{3.4} \cdot e^{-\frac{1}{3.4}x} dx$   
 $\Rightarrow \lambda = \frac{1}{3.4}$   
 $= 0.2548 \quad \textcircled{B}$

3.)  $s^2 = V(x) = \int_0^2 (x - \mu)^2 f(x) dx = \frac{2}{9} \quad \textcircled{D}$

2.)  $\textcircled{C.} \quad \mu = \int_0^2 x f(x) dx = \frac{4}{3}$

4) Mean =  $\mu$   
Variance =  $\frac{s^2}{n} (x^2)$

Standard deviation =  $\sqrt{\frac{s^2}{n}} \rightarrow \sigma \text{ đor.}$

$\textcircled{D}$

$n \uparrow \Rightarrow \sigma \downarrow$

5.)  $s = \sqrt{\frac{s^2}{n}} = \sqrt{\frac{20}{100}} = 0.447 \approx 0.45 \quad \textcircled{B}$

6.)  $\bar{x} = \frac{64.9 + 65 + 65.5}{3} = \frac{97.7}{15}$

$\approx s^2 = \frac{(64.9 - \bar{x})^2 + (65 - \bar{x})^2 + (65.5 - \bar{x})^2}{n-1}$

$$= 0.103337$$

$$\Rightarrow \varsigma = \sqrt{\varsigma^2} = 0.321 \quad (\text{A})$$

7.) \* A

8.) A

9.)  $f(x) = F'(x) = \frac{1}{\alpha} e^{-\frac{x}{\alpha}}$

$$P(X > 1) = 1 - P(X < 1)$$

$$= 1 - \int_0^1 \frac{1}{2} e^{-\frac{x}{2}} dx \\ = 0.6065 \quad (\text{C})$$

10.)  $P(13 < X < 65) = P(X < 65) - P(X < 13)$

$$= \int_0^{65} \underbrace{\frac{1}{b-a} dx}_{70-10} - \int_0^{13} \frac{1}{b-a} dx$$

$$= 0.86 \quad (\text{C})$$

11.) (?)

12.)  $4 + 5 + 8 + 6 + 2 = 25 \quad (\text{A})$

13.) (?)

14.) B

$$15) P(X > a) = 75.8\%$$

$$\Rightarrow P(X \leq a) = 0.242$$

$$\Leftrightarrow P(Z \leq \frac{a - \mu}{\sigma}) = 0.242$$

$$\Rightarrow \frac{a - 3.5}{1} = -0.7$$

$$\Rightarrow a = 2.8 \quad \text{B}$$

# Chapter 8:

## Statistical intervals for a single sample

⇒ Confidence interval (CI) : Khoảng tin cậy

⇒ Confidence level :  $100 \cdot (1-\alpha) \% = 1-\alpha$  (Độ tin cậy)

Ex: Confidence level = 95% = 0.95 =  $1-\alpha \Rightarrow \alpha = 0.05$

⇒ Critical value (a percentage point) : giá trị giới hạn :  $Z_\alpha$

$$P(Z > z_\alpha) = \alpha$$

$$\Rightarrow P(Z < z_\alpha) = 1 - \alpha$$

Ex: Find  $z_{0.025} \Rightarrow P(Z > z_{0.025}) = \alpha$

$$\Rightarrow P(Z < z_{0.025}) = 0.979$$

$$\Rightarrow z_{0.025} = 1.96$$

### Confidence interval for $\mu$ ( $\sigma$ is known)

Theorem (C.I. for  $\mu$  if  $\sigma$  is known)

If  $\bar{x}$  is the sample mean of a random sample of size  $n$  from a normal population with known variance  $\sigma^2$ , a  $100(1-\alpha)\%$  CI on  $\mu$  is given by:

$$\bar{x} - z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

Example: In a sample of 36 randomly selected women, it was found that their mean height was 65.3 inches. From previous studies, it is assumed that the standard deviation of all women heights  $\sigma = 2.5$  (in). Construct a 90% confidence interval for the mean height of all women.

1. Continuous interval for  $\mu$  of a normal distribution  $\delta$  is known

$\mu$ : population mean (unknown) } parameters  
 $\delta$ : population standard deviation }

a.) A  $100(1-\alpha)\%$  two-sided CI on  $\mu$  is given by:

$$\bar{x} - z_{\frac{\alpha}{2}} \cdot \frac{\delta}{\sqrt{n}} < \mu < \bar{x} + z_{\frac{\alpha}{2}} \cdot \frac{\delta}{\sqrt{n}}$$

b.) A  $100(1-\alpha)\%$  upper confidence bound on  $\mu$  is:  
(chỗn trên của  $\mu$ )

$$\mu \leq \bar{x} + z_{\alpha} \cdot \frac{\delta}{\sqrt{n}}$$

c.) A  $100(1-\alpha)\%$  lower confidence bound on  $\mu$  is:  
(chỗn dưới của  $\mu$ )

$$\mu \geq \bar{x} - z_{\alpha} \cdot \frac{\delta}{\sqrt{n}}$$

\* Here:  $\bar{x}$  is sample mean;  $n$  is sample size

Ex1: Suppose population is normal with standard deviation

$$\delta = 0.010 \text{ (}\delta\text{ is known)}$$

A 99% two-sided CI on  $\mu$  is:

$$\bar{x} - z_{\frac{\alpha}{2}} \cdot \frac{\delta}{\sqrt{n}} < \mu < \bar{x} + z_{\frac{\alpha}{2}} \cdot \frac{\delta}{\sqrt{n}}$$

$$n = 10 ; \bar{x} = 1.545$$

Confidence level:  $99\% = 1 - \alpha \Rightarrow \alpha = 0.01$

$$\Rightarrow Z_{\frac{\alpha}{2}} = Z_{0.005} = 2.58$$

$$1.545 - 2.58 \times \frac{0.010}{\sqrt{10}} < \mu < 1.545 + 2.58 \times \frac{0.010}{\sqrt{10}}$$

$$\Leftrightarrow 1.537 < \mu < 1.553$$

Answer:  $\mu \in (1.537, 1.553)$

b.) A 90% CI on  $\mu$  is: ( $Z_{\frac{\alpha}{2}} = Z_{0.05} = 1.645$ )

$$1.545 - 1.645 \times \frac{0.010}{\sqrt{10}} < \mu < 1.545 + 1.645 \times \frac{0.010}{\sqrt{10}}$$

$$\Rightarrow \mu \in (1.540; 1.550)$$

\* Remark: Nếu  $n$ ,  $\delta$  cố định thì độ tin cậy càng cao

→ Khoảng tin cậy (CI) càng dài (longer)

c.) Construct a 95% upper confidence bound on  $\mu$ :

$$\bar{x} + Z_{\alpha} \frac{\sigma}{\sqrt{n}}$$

$$= 1.545 + Z_{0.05} \times \frac{0.010}{\sqrt{10}} , Z_{0.05} = 1.645$$

$$= 1.55$$

d.) Construct a 90% lower confidence bound on  $\mu$ :

$$\bar{x} - z_{\alpha} \frac{\sigma}{\sqrt{n}} ; z_{\alpha} = z_{0.1}$$

\* Remark:

1.) A length of a CI is:  $2 \cdot z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$

2.) The error:  $|\bar{x} - \mu|$   
(Sai  $\delta\delta'$ )

It is clear that: the error  $= |\bar{x} - \mu| \leq z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$

Ex:  $\sigma$  is known,  $\delta = 6$ .  
(SL8/18)

She desires to be 98% confidence that  $|\bar{x} - \mu| \leq 4$

Solve:  $z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} = 4$

$$\Leftrightarrow z_{0.01} \times \frac{6}{\sqrt{n}} = 4$$

$$\Rightarrow n = 12.173$$

The required sample size is:  $n = 13$ .  
(Kích thước mẫu cần thiết)

Solution 2:  $n = \left\lceil \left( \frac{z_{\alpha}}{2} \cdot \frac{\sigma}{E} \right)^2 \right\rceil = \left\lceil \left( \frac{2.326 \times 6}{4} \right)^2 \right\rceil$   
 $= \lceil 12.173 \rceil = 13$

\* Remark: (Interpreting a 95% CI)

If an infinite number of random samples are collected and a 95% CI for  $\mu$  is computed from each sample, 95% of these intervals will contain  $\mu$ .

2.) Confidence on  $\mu$  of a normal distribution,  $\sigma$  is unknown

$\sigma$ is known	$\sigma$ is unknown
use $\mu$	use $s$ (sample standard dev.)
use $z_{\frac{\alpha}{2}}$	use $t_{\frac{\alpha}{2}, n-1}$
use $z_{\alpha}$	use $t_{\alpha, n-1}$

$t_{n-1}$ : t-distribution with  $n-1$  degrees of freedom

$$t_{0.05, 18} = t_{\alpha, n-1} \Rightarrow \begin{cases} \alpha = 0.05 \\ n-1 = 18 \end{cases}$$

a.) A  $100(1-\alpha)\%$  two-sided CI on  $\mu$  is:

$$\bar{x} - t_{\frac{\alpha}{2}, n-1} \cdot \frac{s}{\sqrt{n}} < \mu < \bar{x} + t_{\frac{\alpha}{2}, n-1} \cdot \frac{s}{\sqrt{n}}$$

b.) A  $100(1 - \alpha)\%$  upper confidence bound is :

$$\mu \leq \bar{x} + t_{\alpha, n-1} \frac{s}{\sqrt{n}}$$

c.) A  $100(1 - \alpha)\%$  lower confidence bound is :

$$\mu \geq \bar{x} - t_{\alpha, n-1} \frac{s}{\sqrt{n}}$$

Ex: The population has normal distribution,  $\sigma$  is unknown

A 95% two-sided CI on  $\mu$  is:

$$\bar{x} - t_{\frac{\alpha}{2}, n-1} \cdot \frac{s}{\sqrt{n}} < \mu < \bar{x} + t_{\frac{\alpha}{2}, n-1} \cdot \frac{s}{\sqrt{n}}$$

$$\bar{x} = 3450, \quad n = 28, \quad s = 600$$

Confidence level = 95% =  $1 - \alpha \Rightarrow \alpha = 0.05$

$$\Rightarrow t_{\alpha/2, n-1} = t_{0.025, 27} = 2.052$$

$$\Rightarrow 3218 < \mu < 3682, \text{ or } \mu \in (3218, 3682)$$

$$\mu \in (3218, 3282)$$

3) Confidence interval for  $p$  - Khoảng tin cho tỷ lệ  $p$ .

$p$ : population proportion (tỷ lệ của tổng thể - chưa biết)

$$P = \hat{P} = \frac{x}{n} \quad \text{Sample proportion (tỷ lệ mẫu)}$$

If  $\begin{cases} np \geq 5 \\ n(1-p) \geq 5 \end{cases}$  then we have

$$Z = \frac{\hat{P} - P}{\sqrt{\frac{P(1-P)}{n}}} \sim N(0, 1)$$

$$\text{Since, } P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) = 1 - \alpha.$$

$$\Rightarrow P(-z_{\alpha/2} \leq \frac{\hat{P} - P}{\sqrt{\frac{P(1-P)}{n}}} \leq z_{\alpha/2}) = 1 - \alpha.$$

$$\Rightarrow P(-z_{\alpha/2} \leq \frac{P - \hat{P}}{\sqrt{\frac{P(1-P)}{n}}} \leq z_{\alpha/2}) = 1 - \alpha.$$

$$\Rightarrow P(\hat{P} - z_{\alpha/2} \sqrt{\frac{P(1-P)}{n}} \leq P \leq \hat{P} + z_{\alpha/2} \sqrt{\frac{P(1-P)}{n}}) = 1 - \alpha$$



a) A  $100(1-\alpha)\%$  CI for  $p$  is given by:

$$\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} < p < \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$\hat{p}$ : is the proportion of a sample of size  $n$  ( $\hat{p} = \frac{x}{n}$ )

ex)

b) A  $100(1-\alpha)\%$  upper confidence bound for  $p$  is:

$$P \leq \hat{p} + z_{\alpha} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

c) A  $100(1-\alpha)\%$  lower confidence bound for  $p$  is:

$$P \geq \hat{p} - z_{\alpha} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

VĐ: Population homeless people

$\hat{p}$ : the proportion of homeless people who are veterans

A 95% upper confidence bound for  $p$  is:

$$P \leq \hat{p} + z_{\alpha} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$\Rightarrow 250 = n$$

$$\hat{p} = \frac{x}{n} = \frac{47}{250} = 0.188$$

$$\text{Confidence level} = 1 - \alpha = 95\% \rightarrow \alpha = 0.05.$$

$$\rightarrow \text{Critical solve: } z_2 - z_{0.05} = 1.645$$

$$\rightarrow p \leq 0.188 + 1.645 \times \sqrt{\frac{0.188(1-0.188)}{250}}$$

$$\Rightarrow p \leq 0.029.$$

VD: 1000 randomly selected cases of lung cancer, 750 resulted in death within 10 years  
 Calculate a 95% two-sided confidence interval of the death rate from lung cancer.

$$z_2 = z_{0.05} = 1.645$$

$$z_{2/2} = z_{0.025} = 1.96.$$

$$n = 1000$$

$$\hat{p} = \frac{x}{n} = \frac{750}{1000} = 0.75.$$

A 95% two-sided CI interval of the death rate from lung cancer :

$$0.75 - 1.96 \times \sqrt{\frac{0.75(1-0.75)}{1000}} \leq p \leq 0.75 + 1.96 \times \sqrt{\frac{0.75(1-0.75)}{1000}}$$

$$0.723 \leq p \leq 0.776.$$

Problem: Find the required sample size such that the error:  $= |\hat{p} - p| \leq E$

Case 1:  $p$  is known

$$\text{Solve: } z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} = E \Rightarrow n = \left( \frac{z_{\alpha/2}}{E} \right)^2 p(1-p)$$

Then choose  $n = \left[ \left( \frac{z_{\alpha/2}}{E} \right)^2 p(1-p) \right]$

Case 2:  $p$  is unknown but  $\hat{p}$  is known.

$$\text{Solve: } z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = E \Rightarrow \text{choose } n = \left[ \left( \frac{z_{\alpha/2}}{E} \right)^2 \hat{p}(1-\hat{p}) \right]$$

Case 3: Both  $p$  and  $\hat{p}$  is unknown

$$n = \left[ \left( \frac{z_{\alpha/2}}{E} \right)^2 \times 0.25 \right]$$

VP:  $p = 0.08$  ( $p$  is known)

Find min of  $n$  you should use to be 99% confident that the point estimate of  $p$  (is  $\hat{p}$ ) will be within 0.04 around  $p$ .

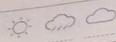
$$\text{Solve: } |\hat{p} - p| \leq 0.04$$

$$\Rightarrow n = \left[ \left( \frac{z_{\alpha/2}}{E} \right)^2 p(1-p) \right] = \left[ \left( \frac{2.58}{0.04} \right)^2 0.08(1-0.08) \right]$$

$$= [306.1944] = 307$$

$$E = 0.04$$

$$z_{\alpha/2} = z_{0.005} = 2.58$$



VD: (Both  $p$  and  $\hat{p}$  is unknown)  $Z_{0.01} = 2.33$

Find  $n$  such that  $|\hat{p} - p| \leq 6\%$  with  $1 - \alpha = 98\%$

$$n = \left\lceil \left( \frac{Z_{0.01}}{E} \right)^2 \times 0.25 \right\rceil = \left\lceil \left( \frac{2.33}{0.06} \right)^2 \times 0.25 \right\rceil = \left\lceil 377.006 \right\rceil$$

$$= 378$$

VD:

$$n = 4000 \rightarrow \hat{p} = \frac{x}{n} = 0.625$$

$$x = 2500$$

$$E = 5\% = 0.05$$

$$\times 1 - \alpha = 94\%$$

$$\rightarrow \alpha = 0.06.$$

## Chapter 9: Test of hypotheses for a single sample kiểm định các giả thuyết cho một mẫu đơn

Statistical hypothesis (Giả thuyết thống kê)

Hypothesis: is a assertion about the parameters  $\theta$   
 $\theta: \mu$  (mean);  $p$  (proportion)

Null hypothesis  $H_0$  (Giả thuyết ko) - Giả thuyết gốc  $\theta = \theta_0$   
 Alternative hypothesis  $H_1$ : (Giả thuyết đối)



Thứ ngày

Case 1: (two-tails test / Two-sided test)

$$\begin{cases} H_0 : \mu = \mu_0 \\ H_1 : \mu \neq \mu_0 \end{cases}$$

Case 2: (Left-tailed test / Lower-sided test)

$$\begin{cases} H_0 : \mu = \mu_0 \\ H_1 : \mu < \mu_0 \end{cases}$$

Case 3 (Right-tailed test / Upper-sided test)

$$\begin{cases} H_0 : \mu = \mu_0 \\ H_1 : \mu > \mu_0 \end{cases}$$

## § Type of error (các loại sai lầm)

\*  $\alpha = P(\text{type I error}) = P(\text{reject } H_0 \text{ when } H_0 \text{ is true})$

$$= P(\bar{x} \text{ is in critical region} \mid \mu = \mu_0)$$

$\alpha$ : the significance level (mức ý nghĩa)

$\beta = P(\text{type II error}) = P(\text{Fail to reject } H_0 \text{ when } H_0 \text{ is false})$

## § Hypothesis testing procedure for $\mu$ ( $\sigma$ is known)

### A. Traditional method

#### 1. Case 1 (two-tails test)

$$\begin{cases} H_0 : \mu = \mu_0 \\ H_1 : \mu \neq \mu_0 \end{cases}$$

⇒ Determine a test statistic (thông kê kiểm định) and its value

$$Z_0 = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$$

⇒ Determine the critical value (percentage point):  $Z_{\alpha/2}$

⇒ Make conclusion:

If  $Z_0 < -Z_{\alpha/2}$  or  $Z_0 > Z_{\alpha/2}$  ( $|Z_0| > Z_{\alpha/2}$ )

→ Reject  $H_0$ .

If  $-z_{\alpha/2} \leq z_0 \leq z_{\alpha/2}$  ( $|z_0| \leq z_{\alpha/2}$ )

$\rightarrow$  Fail to reject  $H_0$   
(Do not reject  $H_0$ )

2.) Case 2: (Left-tailed test, lower-tailed test)

$$\begin{cases} H_0: \mu = \mu_0 \\ H_1: \mu < \mu_0 \end{cases}$$

$\Rightarrow$  The critical value:  $z_\alpha$

$\Rightarrow$  Conclusion: If  $z_0 < -z_\alpha$ , then reject  $H_0$

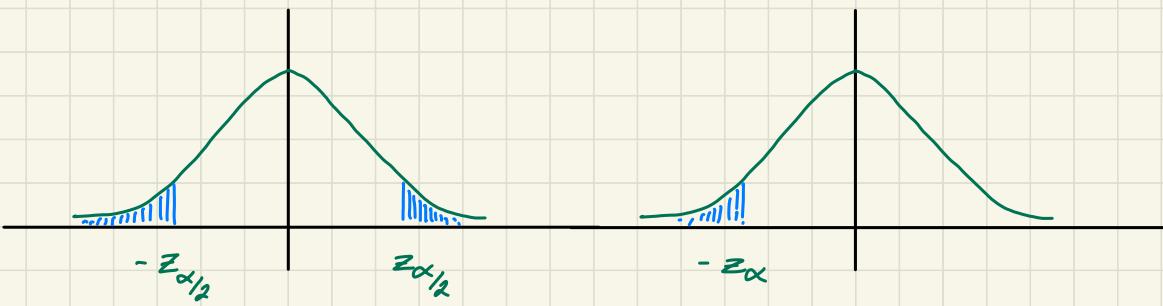
If  $z_0 > -z_\alpha \rightarrow$  Fail to reject  $H_0$

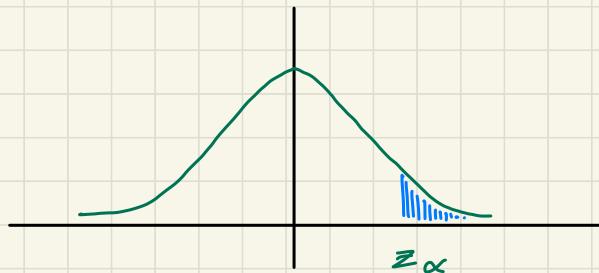
3.) Case 3:

$$\begin{cases} H_0: \mu = \mu_0 \\ H_1: \mu > \mu_0 \end{cases}$$

If  $z_0 > z_\alpha \rightarrow$  Reject to  $H_0$

If  $z_0 \leq z_\alpha \rightarrow$  Fail to reject  $H_0$





**Ex:** Population is normally distributed  
(SL 8/18)

$$\sigma = 0.03 \quad (\sigma \text{ is known})$$

$n = 43$ ,  $\bar{x} = 1.64$ , the significance level is  $\alpha = 0.05$

⇒ The hypotheses :  $\begin{cases} H_0: \mu = 1.7 \\ H_1: \mu \neq 1.7 \end{cases}$  ( $H_0 = 1.7$ )

⇒ The test statistic is

$$Z_0 = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}} = \frac{1.64 - 1.7}{0.03 / \sqrt{43}} = -13.115$$

(Note: In this problem, we will reject  $H_0$  if  $Z_0 < -z_{\alpha/2}$  or  $Z_0 > z_{\alpha/2}$ )

⇒ The critical value:  $z_{\alpha/2} = z_{0.025} = 1.96$

⇒ Since  $Z_0 = -13.115 < -z_{\alpha/2} = -1.96$   
⇒ Reject  $H_0$  and accept  $H_1$

\*  $\begin{cases} H_0: \mu = \dots \\ H_1: \mu \neq \dots \quad (\text{not } \geq \text{ or } \leq) \\ > \\ < \end{cases}$

Ex 2: Population is normally distributed  
(SL 12/18)

$$\sigma = 0.03 \quad (\sigma \text{ is known})$$

$n = 43$ ,  $\bar{x} = 1.64$ , the significance level is  $\alpha = 0.05$

⇒ The hypotheses:  $\begin{cases} H_0: \mu = 1.6 \\ H_1: \mu > 1.6 \end{cases} \quad (\mu_0 = 1.6)$

⇒ The test statistic:

$$z_0 = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}} = \frac{1.64 - 1.6}{0.03 / \sqrt{43}} = 8.743$$

⇒ The critical value:  $z_\alpha = z_{0.05} = 1.645$

⇒  $z_0 = 8.743 > z_\alpha = 1.645 \rightarrow \text{Reject } H_0 \text{ (Accept } H_1\text{)}$

⇒ We conclude that the mean height is greater than 1.6m!

### B.) P-value method

1.) Two-tailed test:  $\begin{cases} H_0: \mu = \mu_0 \\ H_1: \mu \neq \mu_0 \end{cases}$

$$\text{P-value} = P(Z < -|z_0| \text{ or } Z > |z_0|)$$

$$= 2 \cdot P(Z < -|z_0|) = 2 \cdot P(Z > |z_0|)$$

Conclusion: ⇒ If P-value  $< \alpha$ , then reject  $H_0$ .

⇒ If P-value  $> \alpha$ , then fail to reject  $H_0$ .

2.) Left-tailed test:  $\begin{cases} H_0: \mu = \mu_0 \\ H_1: \mu < \mu_0 \end{cases}$

$$P\text{-value} = P(Z < z_0)$$

3.) Right-tailed test:  $\begin{cases} H_0: \mu = \mu_0 \\ H_1: \mu > \mu_0 \end{cases} \Rightarrow P\text{-value} = P(Z > z_0)$

Ex (SL 12/18)

$$z_0 = 8.743$$

$$\begin{cases} H_0: \mu = 1.6 \\ H_1: \mu > 1.6 \end{cases}$$

$$\begin{aligned} P\text{-value} &= P(Z > z_0) = P(Z > 8.743) \\ &= 1 - P(Z < 8.743) \approx 1 - 1 = 0 \end{aligned}$$

Thus  $P\text{-value} = 0 < \alpha = 0.05 \rightarrow \text{Reject } H_0$

§ Test of hypotheses for population mean  $\mu$  ( $\sigma$  is unknown)

$\sigma$  is known

$\sigma$  is unknown

The test statistic:

$$z_0 = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$$

use  $\sigma$

use  $z_{\alpha/2}$

use  $z_\alpha$

The test statistic:

$$t_0 = \frac{\bar{x} - \mu_0}{s / \sqrt{n}}$$

use  $s$

use  $t_{\alpha/2, n-1}$

use  $t_{\alpha, n-1}$

Note:  $n-1 = \text{degrees of freedom}$

Ex: ( $\sigma$  is unknown)

Population is normally distributed

$$n = 49, \bar{x} = 5.1, s = 1.2$$

The significance level is  $\alpha = 0.05$

⇒ The value of the test statistic is:

$$t_0 = \frac{\bar{x} - \mu_0}{S / \sqrt{n}} = \frac{5.1 - 4.5}{1.2 / \sqrt{49}} = 3.5$$

⇒  $\begin{cases} H_0 : \mu = 4.5 \\ H_1 : \mu > 4.5 \end{cases}$

⇒ The critical value is:  $t_{\alpha, n-1} = t_{0.05, 48} = 1.677$

Note: Since  $t_0 = 3.5 > t_{\alpha, n-1} = 1.677 \rightarrow$  Reject  $H_0$ , accept  $H_1$

### f) Test of hypotheses for population proportion $p$

\* Tính toán bài toán kiểm định cho  $\mu$  khi 3 điều kiện:  
lưu ý là:

The test statistic:  $Z_0 = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$

Ex (Slide 17)

Test the hypotheses:  $\begin{cases} H_0: p = 3\% = 0.03 \\ H_1: p < 0.03 \end{cases} \quad (p_0 = 0.03)$

⇒ The test statistic is:

$$Z_0 = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{\frac{6}{135} - 0.03}{\sqrt{\frac{0.03(1-0.03)}{135}}} = 0.984$$

$$z_{\alpha} = z_{0.05} = 1.645$$

Because  $z_0 > z_{\alpha} \rightarrow$  Fail to reject  $H_0$

Solution (Use P-value)

$$P\text{-value} = P(Z < z_0) = P(Z < 0.984) = 2.144$$

Since  $P\text{-value} > \alpha \rightarrow$  Fail to reject  $H_0$ .

# Chapter 10:

## Statistical inference for Two Samples

Population 1: normally distribution  
mean:  $\mu_1$ , standard deviation:  $\sigma_1$   
Sample:  $\bar{x}_1, n_1$

Population 2: normally distribution  
mean:  $\mu_2$ , standard deviation:  $\sigma_2$   
Sample:  $\bar{x}_2, n_2$

Remark:  $\bar{x}_1 \sim N(\mu_1, \frac{\sigma_1^2}{n_1})$ ;  $\bar{x}_2 \sim N(\mu_2, \frac{\sigma_2^2}{n_2})$

$$\Rightarrow \bar{x}_1 - \bar{x}_2 \sim N(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2})$$

$$Z = \frac{\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1)$$

I. Inference on the difference in means of two normal distribution variances known.  $(\mu_1 - \mu_2) \uparrow$

1. CI on  $\mu_1 - \mu_2$

(a) A  $100(1-\alpha)\%$  two-sided CI on  $\mu_1 - \mu_2$  is given by

$$(\bar{x}_1 - \bar{x}_2) - z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} < \mu_1 - \mu_2$$

$$< (\bar{x}_1 - \bar{x}_2) + z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

(b) A  $100(1-\alpha)\%$  upper confidence bound on  $\mu_1 - \mu_2$  is:

$$\mu_1 - \mu_2 \leq (\bar{x}_1 - \bar{x}_2) + z_\alpha \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

(c) A  $100(1-\alpha)\%$  lower confidence bound on  $\mu_1 - \mu_2$  is

$$\mu_1 - \mu_2 \geq (\bar{x}_1 - \bar{x}_2) - z_{\alpha} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Ex 1: Population 1:  $\sigma_1^2 = 1.5$  Population 2:  $\sigma_2^2 = 1.2$

$$n_1 = 15$$

$$n_2 = 20$$

$$\bar{x}_1 = 89.6$$

$$\bar{x}_2 = 92.5$$

Confidence level:  $1 - \alpha = 95\% \Rightarrow \alpha = 5\% = 0.05$

A 95% Confidence interval on the difference in mean  $(\mu_1 - \mu_2)$  is:

$$(\bar{x}_1 - \bar{x}_2) - z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} < \mu_1 - \mu_2$$

$$\leq (\bar{x}_1 - \bar{x}_2) + z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

$$\Rightarrow -3.684 < \mu_1 - \mu_2 < -2.116$$

$$\text{Or } \mu_1 - \mu_2 \in (-3.684, -2.116)$$

b) A 99% lower confidence bound on  $\mu_1 - \mu_2$  is

$$\mu_1 - \mu_2 \geq (\bar{x}_1 - \bar{x}_2) - z_{\alpha} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

$$z_{\alpha} = z_{0.01} = 2.326$$

2. Test of hypotheses for difference in means  $\mu_1 - \mu_2$ .

$$\text{The test statistic is } Z_0 = \frac{(\bar{x}_1 - \bar{x}_2) - \Delta_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Case 1:  $\begin{cases} H_0: \mu_1 - \mu_2 = \Delta_0 \\ H_1: \mu_1 - \mu_2 \neq \Delta_0 \end{cases}$  (two-tailed test)

If  $Z_0 < -Z_{\alpha/2}$  or  $Z_0 > Z_{\alpha/2} \rightarrow \text{Reject } H_0$

$$\Rightarrow P\text{-value} = 2P(Z < -|Z_0|) = 2\underline{P}(Z > |Z_0|)$$

Case 2:  $\begin{cases} H_0: \mu_1 - \mu_2 = \Delta_0 \\ H_1: \mu_1 - \mu_2 < \Delta_0 \end{cases}$  (Left-tailed test)  
 lower —

$\Rightarrow$  Reject  $H_0$  if  $Z_0 < -Z_\alpha$

$$\Rightarrow P\text{-value} = \underline{P}(Z < Z_0)$$

Case 3:  $\begin{cases} H_0: \mu_1 - \mu_2 = \Delta_0 \\ H_1: \mu_1 - \mu_2 > \Delta_0 \end{cases}$  (Right-tailed test)

$\Rightarrow$  Reject  $H_0$  if  $Z_0 > Z_\alpha$

$$\Rightarrow P\text{-value} = \underline{P}(Z > Z_0)$$

Ex: Formulation 1:  $s_1 = 8$ ,  $n_1 = 10$ ,  $\bar{x}_1 = 121$ .  
 $(\mu_1 \text{ is unknown})$

Formulation 2:  $\{$  has a new drying ingredient  $\} (\mu_2 \text{ is unknown})$

$$s_2 = 8, n_2 = 10, \bar{x}_2 = 112$$

Significant level:  $\alpha = 0.05$

The hypotheses :  $\begin{cases} H_0 : \mu_1 - \mu_2 = 0 \\ H_1 : \mu_1 - \mu_2 > 0 \end{cases}$  ( $\Delta_0 = 0$ )

The value of the test statistic is :

$$z_0 = \frac{(\bar{x}_1 - \bar{x}_2) - \Delta_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = 2.516$$

†  $z_\alpha = z_{0.05} = 1.645$

† Since  $z_0 > z_\alpha \rightarrow$  Reject  $H_0$ , accept  $H_1$ .

→ We conclude that adding the new gradient reduce the drying time.

Remark: P-value =  $P(Z > z_0) = P(Z > 2.516)$   
 $= 1 - P(Z < 2.516) = 0.006$

Since P-value = 0.006 <  $\alpha = 0.05$

→ Reject  $H_0$

II. Inference on  $\mu_1 - \mu_2$  of two normal distributions

(Assume  $\sigma_1, \sigma_2$  are unknown and  $\sigma_1 = \sigma_2$ )

Sample 1:  $n_1, \bar{x}_1, s_1$  ( $s_1$  is sample standard deviation)

Sample 2:  $n_2, \bar{x}_2, s_2$

Dkt: Pooled variance (phuong sai gop)

$$S_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

$S_p^2 = \frac{s_1^2 + s_2^2}{n_1 + n_2 - 2}$ $\sigma_1, \sigma_2$ are known Use $\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$ Use $t_{\alpha/2}$ Use $\bar{z}_2$ The test statistic	$\sigma_1, \sigma_2$ are unknown Use $\sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}} = S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$ Use $t_{\alpha/2, n_1 + n_2 - 2}$ $n_1 + n_2 - 2$ degrees of freedom $t_{\alpha/2, n_1 + n_2 - 2}$ $t_{\alpha/2, n_1 + n_2 - 2}$ $t_0 = \frac{(\bar{x}_1 - \bar{x}_2) - \Delta_0}{\sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}}}$
---	---

Ex 1: Two population are normally distributed and have the same variance

Catalyst 1:  $n_1 = 8$ ,  $\bar{x}_1 = 92.255$ ,  $s_1 = 2.39$

Catalyst 2:  $n_2 = 8$ ,  $\bar{x}_2 = 92.753$ ,  $s_2 = 2.98$

⇒ Pooled variance:

$$S_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} = 7.3$$

A 95% Confidence interval for difference in means  $\mu_1 - \mu_2$

$$(\bar{x}_1 - \bar{x}_2) - t_{\alpha/2, n_1 + n_2 - 2} \sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}} < \mu_1 - \mu_2 <$$

$$(\bar{x}_1 - \bar{x}_2) + t_{\alpha/2, n_1+n_2-2} \sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}}$$

We have  $t_{\alpha/2, n_1+n_2-2} = t_{0.025, 14} = 2.145$

$$\Rightarrow -3.376 < \mu_1 - \mu_2 < 2.42$$

**Ex2:** Sample 1:  $n_1 = 14$ ,  $\bar{x}_1 = 3$ ,  $s_1 = 2.5$

" " 2:  $n_2 = 12$ ,  $\bar{x}_2 = 4$ ,  $s_2 = 2.8$

The Pooled variance is:

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} = 6.98$$

We want to test:  $\mu_1 - \mu_2 > 0 \rightarrow \Delta_0 = 0$

The test statistic is:

$$\text{to: } \frac{(\bar{x}_1 - \bar{x}_2) - \Delta_0}{\sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}}} = \frac{(3 - 4) - 0}{\sqrt{\frac{6.98}{14} + \frac{6.98}{12}}} = -0.96.$$