

## CHAPTER 1:

\* Obj : Statistics :

- collecting
- organizing
- analyzing
- interpreting }  $\Rightarrow$  make conclusion

\* purpose : change "Data" into conclusion

\* Classification:

Statistics

Descriptive : organizing, summarizing and display

Inferential : using sample data to draw conclusion about to population

\* Voc

- population : ( $N$ ) : everything, everyone
- parameter : characteristics of population
- sample : ( $n$ ) : a portion of population
- statistics : description of sample
- variable : ~~number~~, ~~des~~ characteristics of elements
- data : actual value of variables

\* ExP :

A		\$ 1000
B		\$ 1500
C		\$ 2000

- population : all of companies

- sample : 3

- statistics : profits of

- variable : profits of each company

Thứ

Ngày

No.

## \* collect data

1, Retrospective (Historical) : lây data ở quá khứ

2, Observational Study

3, Experiment

4, Simulation

5, Survey [ sample : 50 ]

population : CENSUS

## \* Type of data

- Qualitative :

- Quantitative :

+ Counting : # students, # cars, ...

+ Measure : temperature, weight, height, salary

## A Sampling

- Representative

- Replacement / Without replacement

( hoán lại )

- Non random : not representative

- Random :

- simple random ( same chance )

- Stratified, phân tầng

- cluster : khảo sát toàn bộ 1 nhóm

## CHAPTER 2: Probability

Ex: Random experiment

toss a coin : - outcomes [ head  
tail ]

- sample space: all of outcomes  $\Omega = \{H, T\}$

- event: + simple:  $B_1 = \{H\}$

+ compound:  $B_2 = \{HH\}$

(+) event is a subset of the sample space  
of a random / is a collection of outcomes / is  
a collection of elementary events

Ex: Roll a dice

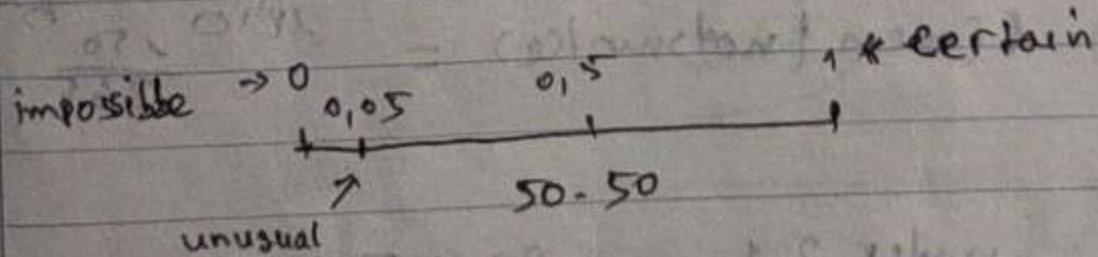
$$\Omega = \{1, 2, 3, 4, 5, 6\}$$

$$A = \{1, 3, 5\}$$

$$* P(A) = \frac{|A|}{|\Omega|}$$

\* Properties:

$$0 \leq P(A) \leq 1$$



\* Some probability formulas

① Addit. (OR)

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

Thứ

Ngày

Nº.

- Special :  $\begin{array}{c} A \\ \cap \\ \bigcirc \bigcirc \end{array}$   $P(A \cup B) = P(A) + P(B)$

## ② Multiplication (AND)

- If  $A$  and  $B$  are independent  $\Rightarrow P(A \cap B) = P(A) \cdot P(B)$

- General case :

$$\begin{aligned} P(A \cap B) &= P(A) P(B|A) \\ &= P(B) P(A|B) \end{aligned}$$

③  $| P(A|B) = P(A) \text{ given that } B$

## ④ Complement :

$$P(\bar{A}) = 1 - P(A)$$

Ex :

	smoke	not smoke
men	200	150
women	85	115

a)  $P(\text{men}) = \frac{200}{385} = \frac{5}{11}$

b)  $P(\text{smoke}) = \frac{200}{385} = \frac{40}{77}$

c)  $P(\text{men} | \text{smoke}) = \frac{200}{285} = \frac{40}{57}$

d)  $P(\text{men and smoke}) = \frac{200}{385} = \frac{40}{77} = \frac{9}{11}$

e)  $P(\text{men} | \text{not smoke}) = \frac{120}{265} = \frac{24}{53}$

Ex :

	under 21	21-25	over 25
ticket	82	30	18
no ticket	17	27	61
(%)	(90)	(60)	(29)

$$(82+17) - (17+27) = 114 - 44 = 70$$

$$114 - (17+27) = 114 - 44 = 70$$

(a)

$$b) P(\text{no ticket}) = \frac{105}{299}$$

$$c) P(21-25 \text{ and no ticket}) = \frac{22}{299}$$

$$d) P(21-25 \text{ or no ticket}) = \frac{199}{299}$$

$$e) P(21-25 | \text{no ticket}) = \frac{27}{105}$$

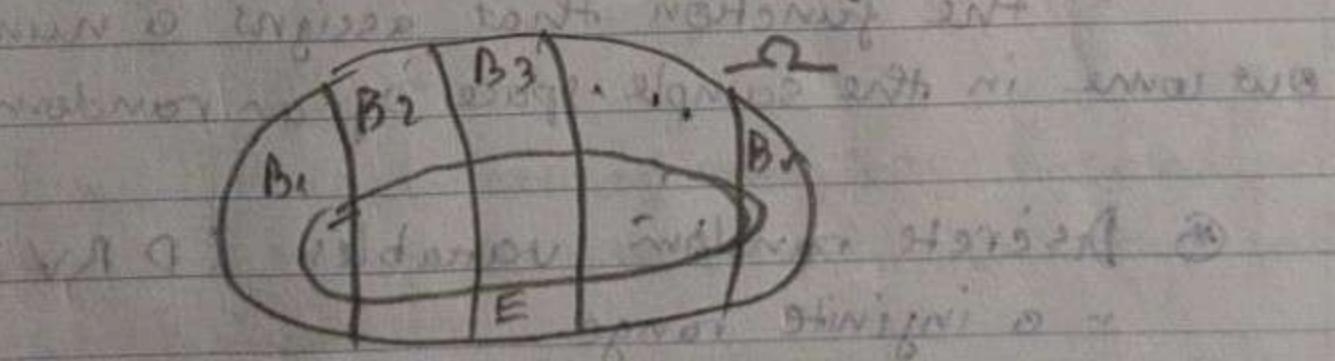
$$f) P(\text{no ticket} | 21-25) = \frac{27}{66}$$

Note:

A, B independent

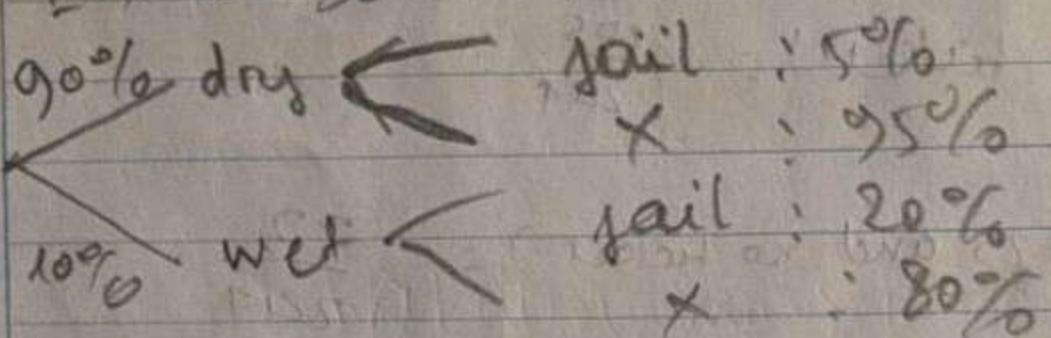
$$\begin{aligned} P(A \cap B) &= P(A) P(B) \\ \text{and } \Rightarrow & \left\{ \begin{array}{l} P(A|B) = P(A) \\ P(B|A) = P(B) \end{array} \right. \end{aligned}$$

### \* Total probability



$$P(E) = P(E|B_1)P(B_1) + P(E|B_2)P(B_2) + \dots + P(E|B_n)P(B_n)$$

Ex: Build tree



$$P(\text{fail}) = 90\% \cdot 5\% + 10\% \cdot 20\% \\ = 6.5\%$$

④ Bayes's

$$P(A|B) = P(B|A) \cdot P(A) / P(B)$$

## CHAPTER 3:

### Discrete random variables

#### \* Random variables

- describle the outcome  $x$  of an experiment
- the function that assigns a number to each outcome in the sample space of an random experiment

#### ④ Discrete random variables (DRV)

- a infinite range

#### ⑤ Probability distribution function (PDF)

$X$	0	1	2	3	4	5
$P(x)$	0.2	0.1	X	0.3	0.1	0.1

④ Mean and variance (ki' rong và phu'ng sai)

$$\mu = E(x) = \sum x_i P(X=x_i) = \sum \frac{x_i}{N}$$

$$s^2 = V(x) = E(x-\mu)^2 = \sum_{i=1}^N (x_i - \mu)^2 P(X=x_i)$$

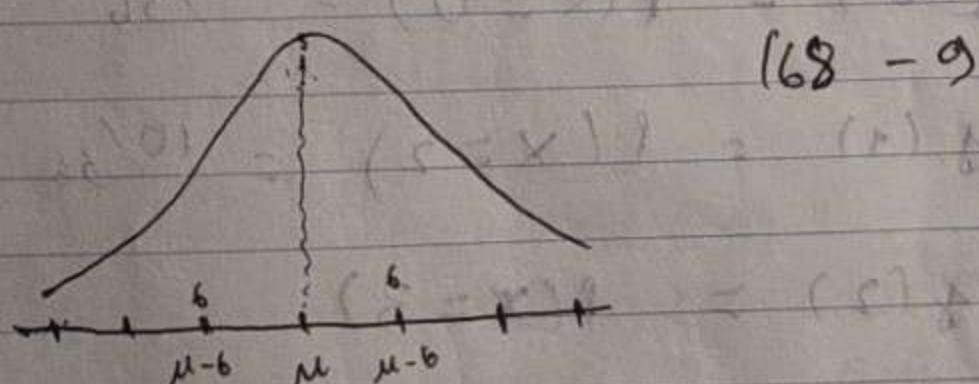
$$= \sum_{i=1}^N x_i^2 P(X=x_i) - \mu^2$$

$$s = \sqrt{s^2} = \sqrt{V(x)}$$

Ex: 5, 6, 7, 8

$$s = \sqrt{5^2 \cdot \frac{1}{4} + 6^2 \cdot \frac{1}{4} + 7^2 \cdot \frac{1}{4} + 8^2 \cdot \frac{1}{4}} = 6,5$$

⑤ Empirical Rule



$$(68 - 95 - 99, 78)$$

Ex: The number of messages sent per hour over a computer network has the following distribution

X	9	10	11	12	13	14
P(x)	0.08	0.15	0.3	0.2	0.2	0.07

a) Find avg / mean / expected value

b) variance

c) standard deviation

$$a) \mu = 11,5$$

$$b) \sigma = 1,85$$

$$c) \delta = 1,36$$

\* Probability mass function

$$f(x_i) = P(X=x_i), 0 \leq P(X=x_i) \leq 1$$

↓

$$\left\{ \begin{array}{l} 0 \leq f(x_i) \leq 1 \\ \sum f(x_i) \end{array} \right.$$

$E_x :$

$x$	0	1	2
$P(x)$	$25/36$	$10/36$	$1/36$

$$f(0) = P(X=0) = 25/36$$

$$f(1) = P(X=1) = 10/36$$

$$f(2) = P(X=2)$$

- Cumulative distribution

$$F(a) = P(X \leq a) = \sum_{x_i \leq a} P(X=x_i)$$

$$= \sum_{x_i \leq a} f(x_i)$$

↓

$$\left\{ \begin{array}{l} 0 \leq f(x) \leq 1 \\ F : \text{not continuous} \end{array} \right.$$

Thứ

Ngày

Nº.

$$\text{Ex 1: } F(0) = P(x \leq 0) = P(x=0) = \frac{25}{36}$$

$$F(1) = P(x \leq 1) = P(x=1) - P(x=0)$$

$$- P(a \leq x \leq b) = P(x \leq b) - P(x < a)$$

Ex 2: Determine the probability mass function of  $x$  from the following cumulative distribution function ( $F(x) \rightarrow g(x)$ )

$$F(x) = \begin{cases} 0, & x < -2 \\ 0.2, & -2 \leq x < 0 \\ 0.7, & 0 \leq x < 2 \\ 1, & x \geq 2 \end{cases}$$

$$\text{a) find } P(0 \leq x \leq 2)$$

$$\text{b) find } P(x < 1.2)$$

$$g(-2) = 0.2 - 0 = 0.2$$

$$g(0) = 0.2 - 0.2 = 0.5$$

$$g(2) = 1 - 0.7 = 0.3$$

Thứ

Ngày

No.

\* Review: - Mass function

$$f(x_1) = P(x = x_1)$$

$$\rightarrow \left\{ \begin{array}{l} 0 \leq f(x_1) \leq 1 \\ \sum_i f(x_i) = 1 \end{array} \right.$$

- Cumulative distribution function

$$F(a) = P(X \leq a) = \sum_{x \leq a} P(x = x_i) = \sum_{x \leq a} f(x_i)$$

$$- E(ax + by) = aE(X) + b(EY)$$

$$- V(ax + by) = a^2 V(X) + b^2 V(Y)$$

$$(!) V(ax + by) = a^2 V(X)$$

\* Special distribution

① Discrete Uniform Distribution

$$X = \{x_1, x_2, \dots, x_n\}$$

$$f(x_i) = 1/n = P(x = x_i)$$



Thứ

Ngày

Nº.

Ex: Rolling a dice  $\rightarrow x \in \{1, 2, 3, 4, 5, 6\}$

$$E(x) = \sum x_i P(x = x_i) = 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + \dots + 6 \cdot \frac{1}{6} = \mu$$

$$V(x) = \sum x_i^2 P(x = x_i) - \mu^2 \\ = 1^2 \cdot \frac{1}{6} + 2^2 \cdot \frac{1}{6} + \dots + 6^2 \cdot \frac{1}{6} - \mu^2$$

1.2 ~~\*)~~  $E(X) = \frac{a+b}{2}$  (chỉ dùng trong tập  
 $V(X) = \frac{(b-a+1)^2-1}{12}$  rời rạc (khi tiếp)

$X$ : discrete uniform random variable  
on the consecutive integer  $[a, b]$

Ex 1:  $X$ : uniform discrete random over  $[10, 60]$

2 (standard deviation) of  $X$ ?

$$\sigma = \sqrt{V(X)} = \sqrt{E(X-\mu)^2} = \sqrt{\frac{(60-10+1)^2-1}{12}}$$
$$= \sqrt{5 \cdot \frac{1}{3} \cdot 28}$$

Ex 2:

$x \sim U[2, 5]$ .  $V(ax)$ ?

$$V(ax) = 16 V(x) = 16 \cdot \frac{(5-2)^2-1}{12}$$

Ex 3:  $x \sim \{1, 2, 3, 5, 6, 7, 8, 9\}$  equally  
Mean = 10  
find  $\mu$ ?

$$\frac{1+2+3+6+u}{5} = 10$$

$$\rightarrow u = 38$$

### ② Binomial Distribution ( $X \sim B(n, p)$ )

(What is the probability of  $k$  success out of  $n$  trials)

Ex: 9 student,  $p$  (success) = 0.3

$P(3 \text{ students: success}) = ?$

$\Rightarrow X$ : number of  $k$  success ~~student~~ trial,  
 $p$ : prob of success,  $q$ : prob of fail

$$P(X=k) = C_n^k \cdot p^k \cdot q^{n-k}$$

\* Mean and variance ( $X \sim B(n, p)$ )

$$E(X) = np$$

$$V(X) = npq$$

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

$$P(X \geq k) = 1 - P(X \leq k-1)$$

Ex: dice rolled 20 times,  $X$  be the numbers of  $\boxed{\square}$  come up. Find the  $G$

$$G = \sqrt{V(X)} = \sqrt{npq} = \sqrt{20 \cdot \frac{1}{6} \cdot \frac{5}{6}}$$

### ③ Poisson

(what is the prob of  $n$  of success in a fixed interval: time, area, ...)

-  $X$  = number of events that can appear in a fixed interval with parameter  $\lambda > 0$

-  $\lambda$  = The average of events / interval  
 (mean) (per)

EXCEL:  $P(X=k) = \text{POISSON.DIST}(k, \lambda, 1)$

$$\Rightarrow P(X = k) = \frac{e^{-\lambda T}}{k!} \cdot (\lambda T)^k \quad k=0, 1, \dots$$

### 3.1 Mean and Variance

$$E(X) = \lambda T$$

$$V(X) = \lambda T$$

### ④ Hypergeometric distribution

(what is the probability of the number of success in sample)

$N$  total population  
 Success:  $k$   
 Fail:  $N-k$

sample  $n$   
 without replacement

$x = \text{num of success in } n \text{ sample}$   
 $x = \{1, 2, 3, \dots, n\}$

$$(l \leq n) \quad P(X=l) = \frac{C^R_l C^{N-R}_{N-l}}{C^n_N}$$

$$E(X) = np \quad (p = \frac{R}{N})$$

$$V(X) = npq \cdot \frac{n-r}{N-1}$$

Ex: 30 red, 20 green, without replacement  
 choose 5 balls,  $E(X)$  and  $\sigma$  at the number of red balls

$$E(X) = np = 5 \cdot \frac{30}{50} = 3$$

$$\sigma = \sqrt{V(X)} = \sqrt{5 \cdot \frac{3}{5} \cdot \frac{2}{5} \cdot \frac{50-5}{50-1}} \approx 1.05$$

Ex:  $N = 36$ ,  $n = 7$ ,  $P(3) ?$

$$P(3) = \frac{C_3^3 \cdot C_{12}^{7-3}}{C_{29}^7}$$

⑤ Geometric Distribution

What is the probability of the number of trials until the 1st success before

$X$  = number of trials until the first success

$$X \sim \text{geometric} \rightarrow P(X=k) = (1-p)^{k-1} \cdot p$$

$$E(X) = \frac{1}{p}$$

$$V(X) = \frac{1-p}{p \cdot p}$$

$$E(X) = P(\text{error}) = 0.03$$

$$\text{a)} P(X=5) = (1-0.03)^5 \cdot 0.03 \\ = 0.0257$$

$$\text{b)} P(X \leq 5) = P(X=1) + P(X=2) + \\ P(X=3) + P(X=4) + P(X=5) \\ = (1+p) + (1-p)p + (1-p)^2 p \\ + (1-p)^3 p + (1-p)^4 p$$

⑥ Negative Binomial distribution  
 (What is the probability of the number of trials until  $r^{\text{th}}$  success)

$X = \text{number of trials until } r^{\text{th}} \text{ success}$

$$X = \underbrace{(\text{first trial}, \dots, \text{th } r^{\text{th}} \text{ trial})}_{\text{successes}}, \alpha \gamma$$

$$P(X = n) = C_{n-1}^{r-1} P^r q^{n-r}$$

+ Mean and variance

$$E(X) = \frac{r}{p}$$

$$V(X) = \frac{r(1-p)}{p^2}$$

## CHAPTER 4: Continuous random

!!) continuos ≈ measure (weight; height; temperature)

⑦ Definition: A continuous random variable is a random variable  $S$  whose possible values includes in an interval of real numbers

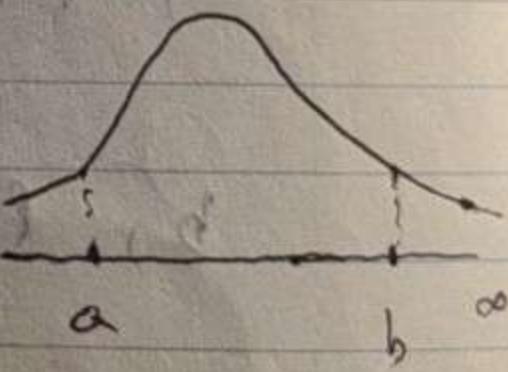
⑧ Probability density function

- $f(x) \geq 0$

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

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

- $P(a < X < b) = \int_a^b f(x) dx$



Shift

End

2

3

PgDn

Enter

Thứ

Ngày

No.

$$\text{Ex 1: } f(x) = e^{-(x-7)} \quad x > 7$$

$$\text{a) } P(0 \leq x \leq 8)$$

$$\text{b) } P(x < 9)$$

$$\text{c) } P(x \geq 9)$$

$$\text{d) } P(x \geq 5)$$

Ex 2: Find  $c$  such that  $f(x)$  is a probability density function  $f(x) = c(x^2 + 2x)$ ,  $0 \leq x < 2$

Ex 1:

$$\text{a) } P(0 \leq x \leq 8) = \int_0^8 f(x) dx = \int_7^8 e^{-(x-7)} dx$$

$$= 0.6321$$

$$\text{b) } P(x < 9) = 1 - P(x \geq 9)$$

$$\text{c) } 1$$

Ex 2:

$$\int_0^2 f(x) dx = 1$$

$$\cdot c \cdot \left( \frac{1}{3}x^3 + x^2 \right) \Big|_0^2 = 1 \Rightarrow c = 1$$

$$c \left( \frac{8}{3} + 4 \right) = 1 \Rightarrow c = 0.15$$

$$\text{f(x) \geq 0} \rightarrow c \geq 0$$

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

$$\rightarrow c \int_0^2 f(x) dx = 1$$

$$\rightarrow c \cdot \frac{2}{3} = 1$$

$$\rightarrow c = 0.15 (\geq 0)$$

Thứ  
Ngày

No.

## \* Probability - function

\* Probability function

$$F(a) = P(X \leq a) = \int_{-\infty}^a f(x) d(x)$$

or  $F(x) = P(\cancel{X} \leq x)$

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

$$\delta(x) = F'(x)$$

$$\text{Ex 1: } f(x) = e^{-(x-a)}$$

Find  $F(x)$

$$F(x) = \begin{cases} e^{-(x-2)} & = -e^{(2-x)} & x > 2 \\ 0 & & x \leq 2 \end{cases}$$

$$X > 2, F(x) = P(X < x) = \int_{-\infty}^x e^{-(t-x)}$$

$$\int_7^x e^{-(t-\gamma)} dt = -e^{-(x-\gamma)} + 1$$

Ex 2: Suppose the cumulative distribution function of the random variable  $x$  is  $F(x) = \begin{cases} 0, & x < -1 \\ -x^2 + 1, & -1 \leq x < 1 \\ 1, & x \geq 1 \end{cases}$

variable  $x$  is  $F(x) = \begin{cases} 0, & x < -1 \\ -x^2 + 1, & -1 \leq x < 1 \\ 1, & x \geq 1 \end{cases}$

$$a) P(x \leq 0)$$

b)  $P(X \leq 2)$

C) D( $x > 0, 5$ )

$$d_1 D (-1 < x < 2)$$

~~John~~

$$\begin{aligned}
 a) P(X \leq 0) &= F(0) = 0^3 + 1 = 1 \\
 b) P(X < -2) &= F(-2) = 0 \\
 c) P(X \geq 0.5) &= 1 - P(X \leq 0.5) = 1 - F(0.5) \\
 &= 1 - 0.5^2 \cdot 1 = 0.5^2 \\
 d) P(-1 < X < 2) &= P(X < 2) - P(X < -1) \\
 &= F(2) - F(-1) \\
 &\quad \text{---} \quad \text{---} \\
 &= 1 + (1)^2 - 1 = 1
 \end{aligned}$$

### ④ Mean and Variance

[Mean]

expected value

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

$$V(X) = \int_{-\infty}^{\infty} x^2 f(x) dx - \mu^2 = 6^2$$

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

$$\text{Ex 1: } f(x) = e(x^2 + 2x), \quad 0 \leq x < 2$$

$$\text{Find } E(X), V(X) = ?$$

$$c = 0.15$$

$$E(X) = \int_0^2 x \cdot 0.15(x^2 + 2x) dx = 1.4 = \mu$$

$$V(X) = 0.15 \int_0^2 x^2(x^2 + 2x) dx - 1.4^2 = \mu^2$$

$$= 0.15 \int_0^2 x^4 + 2x^3 dx - 1.4^2$$

$$= 0.15 \left[ \frac{x^5}{5} + \frac{2x^4}{4} \right]_0^2 - 1.4^2$$

## CHAPTER 9: CONTINUOUS DISTRIBUTIONS

### ① Continuous uniform distribution

$$x \sim U[a, b]$$

$f(x) = \text{constant over } [a, b]$

$\Rightarrow f(x) = 0, \text{ elsewhere}$

$$\therefore f(x) = \frac{1}{b-a}, a \leq x \leq b$$

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

### ② Mean and variance

$$E(x) = \frac{a+b}{2}, V(x) = \frac{(b-a)^2}{12}$$

Ex 1: A plumber estimates that service call are uniformly distributed between 0.5 and 8 hour

a) What are mean and variance

b) Find prob call < 3 hours

c) II II between 2 - 4 hours

d) II II more than 5 hours

e) II II more than 5 given

that less than 7 hours

$$E(x) = \frac{a+b}{2} = \frac{0.5 + 8}{2}$$

$$V(x) = \frac{(b-a)^2}{12} = \frac{7.5^2}{12}$$

$$b) P(x < 3) = \int_{-0.5}^3 \frac{1}{8-0.5} dx$$

$$c) P(2 < x < 9) = \int_2^9 \frac{1}{8-0.5} dx$$

$$d) P(x \geq 5) = \int_5^{\infty} \frac{1}{8-0.5} dx = \int_5^8 \frac{1}{8-0.5} dx$$

$$e) P(5 < x | x < 7) = \int_5^7 \frac{1}{8-0.5} dx / \int_{-0.5}^7 \frac{1}{8-0.5} dx$$

$$E(ax+by) = aE(x) + bE(y)$$

$$E(ax+b) = aE(x) + b$$

$$V(ax+by) = a^2 V(x) + b^2 V(y)$$

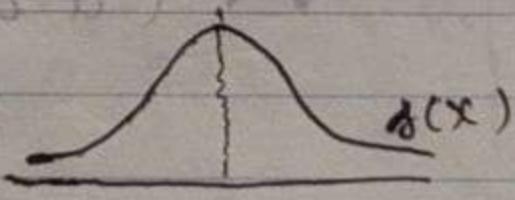
$$V(ax+b) = a^2 V(x)$$

$$\Sigma(h(x)) = \sum_{i=1}^n h(x_i) P(x=x_i) \text{ (dis)}$$

$$E(h(x)) = \int_{-\infty}^{\infty} h(x) g(x) dx \text{ (con)}$$

## ② Continuous normal distribution

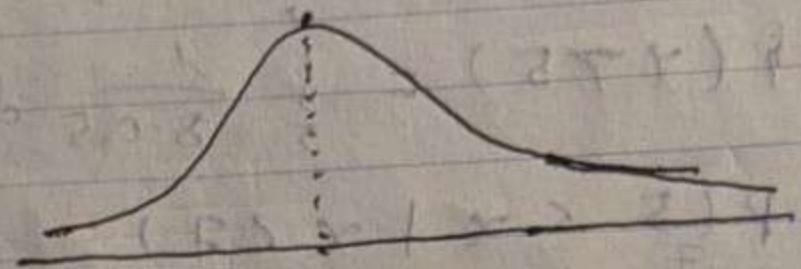
$\left\{ \begin{array}{l} f(x) \geq 0 \\ \text{bell-shaped} \\ \text{symmetric} \end{array} \right.$



$$f(x) = \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

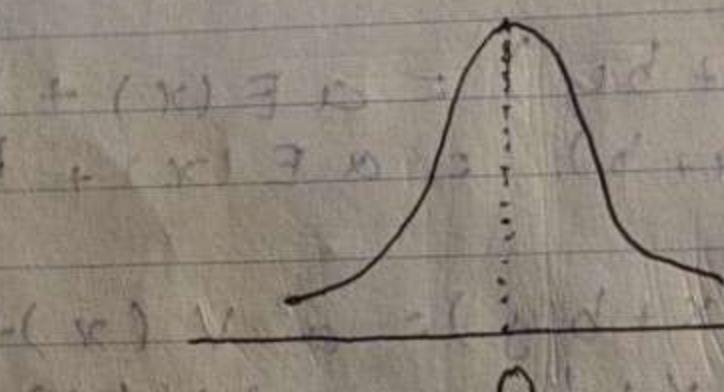
$$* X \sim N(\mu, \sigma^2)$$

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



$$\downarrow Z = \frac{x - \mu}{\sigma}$$

z



$$P(X < a) = P(Z < \frac{a - \mu}{\sigma}) = \Phi\left(\frac{a - \mu}{\sigma}\right)$$

$$P(Z < a) = \Phi(a) \quad (\text{nâng cao} \rightarrow \text{tách riêng})$$

$$\Phi(-a) = 1 - \Phi(a)$$

$$\therefore P(Z > a) = 1 - P(Z < a) = 1 - \Phi(a)$$

$$\therefore P(a < Z < b) = P(Z < b) - P(Z < a)$$

$$= \Phi(b) - \Phi(a)$$

Ex1:  $X \sim N(2000, 100^2)$  Find  $\sigma$

$$\text{a), } P(X > 2000) = 1 - P(X < 2000)$$

$$= 1 - P\left(Z < \frac{2000 - 2000}{100}\right) =$$

$$= 1 - P(Z < 0) = \underline{\underline{\Phi(0)}} = \underline{\underline{\Phi(1)}}$$

$$\therefore 0.6108$$

$$b, P(1700 < x < 2000)$$

$$= P\left(z < \frac{2000 - 1700}{200}\right) - P\left(z < \frac{1700 - 1100}{200}\right)$$

$$= P(z < -0.5) - P(z < -2)$$

$$= 0.285$$

c) Find  $a$  such that  $P(x > a) = 0.03$

$$P(x > a) = 1 - P(x < a) = 0.97$$

$$\Rightarrow P\left(z < \frac{a - 2100}{200}\right) = 0.97$$

$$\Rightarrow \Phi\left(\frac{a - 2100}{200}\right) = \Phi(0.97)$$

$$4. P(0.22 < z < a) \quad z \sim N(0, 1^2)$$

$$\Rightarrow P(z < a) - P(z < 0.22) = 0.24$$

$$\Rightarrow \Phi(a) = 0.29 + 0.59 = 0.83$$

$$5. X \sim N(16, 13^2)$$

$$P(X > 16.25) = 1 - P(X < 16.25)$$

$$= 1 - P\left(z < \frac{16.25 - 16}{1.3}\right)$$

$$(0 = 1 - P(z < 0.192))$$

$$= 0.423$$

Thứ  
Hàng

No.

3, Normal approximation to the Binomial and Poisson distribution

① Normal approximation to Binomial

$n \ggggg$   $p \lllll$   $\Rightarrow$  Normal  $\xrightarrow{\text{approximation}}$  Binomial

$$Z = \frac{X - np}{\sqrt{npq}}$$

$$\begin{aligned} P(X_{\text{Binomial}} \leq a) &= P(X_{\text{NORM}} \leq a + 0.5) \\ &= P(Z \leq \frac{a + 0.5 - np}{\sqrt{np(1-p)}}) \end{aligned}$$

$$\begin{aligned} P(X_{\text{Binomial}} \geq a) &= P(X_{\text{NORM}} \geq a - 0.5) \\ &= 1 - P(X_{\text{NORM}} < a - 0.5) \\ &= 1 - P(Z < \frac{a - 0.5 - np}{\sqrt{npq}}) \end{aligned}$$

② Normal approximation to Poisson

$X \sim \text{Poisson } (\lambda)$ ; Normal  $\xrightarrow{\text{approximation}}$  Poisson

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

$$\begin{aligned} P(X_{\text{Poisson}} \leq a) &= P(X_{\text{NORM}} \leq a + 0.5) \\ &\approx P(Z \leq \frac{a + 0.5 - \lambda}{\sqrt{\lambda}}) \end{aligned}$$

$$P(X_{\text{Poisson}} \geq a) = P(X_{\text{NORM}} \geq a - 0.5)$$

$$\approx P(Z \geq \frac{a - 0.5 - \lambda}{\sqrt{\lambda}})$$

5 10 15 20



Ex 11:

$$P = 0.7, n = 800$$

$$P(X_{\text{BINORM}} > 579) = P(X_{\text{BINORM}} \geq 580)$$

$$= P(X_{\text{NORM}} \geq 579.5) = 1 - P(X_{\text{NORM}} \leq 579.5)$$

$$= 1 - P(Z < \frac{579.5 - 800 \cdot 0.7}{\sqrt{800 \cdot 0.7 \cdot 0.3}})$$

$$= 1 - P(Z < -1.504)$$

Ex 12:

$$P = 0.03, n = 800$$

$$\text{a)} P(X > 30) = P(X_{\text{BINORM}} \geq 31)$$

$$= P(X_{\text{NORM}} \geq 30.5)$$

$$= 1 - P(X_{\text{NORM}} \leq 30.5)$$

$$= 1 - P(Z < \frac{30.5 - 800 \cdot 0.03}{\sqrt{800 \cdot 0.03 \cdot 0.97}})$$

$$\text{b)} P(20 < X_{\text{BINORM}} < 30) = P(21 < X_{\text{BIN}} \leq 29)$$

$$= P(20.5 \leq X_{\text{NORM}} \leq 29.5)$$

$$= P(X_{\text{NORM}} \leq 29.5) - P(X_{\text{NORM}} \leq 20.5)$$

Ex 13:

$$E(X) = \lambda = 1000$$

$$P(X_{\text{POISSON}} = 950) = P(950 \leq X_{\text{POI}} \leq 950)$$

$$= P(949.5 \leq X_{\text{NORM}} \leq 950.5)$$

$$= P(X_{\text{NORM}} \leq 950.5) - P(X_{\text{NORM}} \leq 949.5)$$

Ex 14:

$E(Y) = \lambda = 9.6$  ( $\lambda > 5 \Rightarrow$  using Normal approximation)

$$P(X_{\text{poisson}} > 10) = P(X_{\text{poisson}} > 11)$$

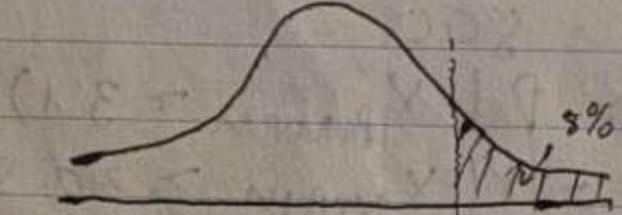
$$= 1 - P\left(Z < \frac{10.5 - 9.6}{\sqrt{9.6}}\right)$$

Ex 7:

$$E(X) = \mu = 8.21, \sigma = 2.14$$

$$\text{Find } a, P(X > a) = 8\%$$

$$P(X < a) = 0.92$$



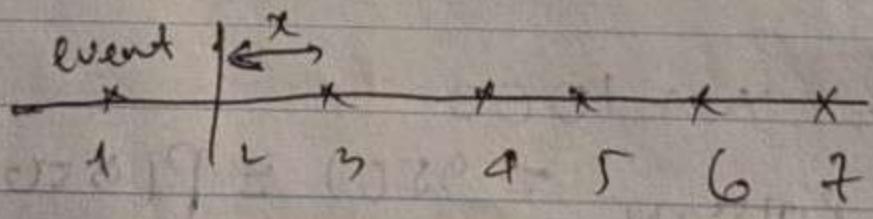
$$P\left(Z < \frac{a - 8.21}{2.14}\right) = 0.92$$

$$\phi\left(\frac{a - 8.21}{2.14}\right) = \phi(\text{NORM.INV}(0.92))$$

$$\rightarrow \frac{a - 8.21}{2.14} = \text{NORM.S.INV}(0.92)$$

a) Exponential distribution

x starting point



$\lambda$  = mean number of events per units

$x$  = distance between | 2 and 3 | successive event



Thứ

Ngày

No.

$$P(X > x) = e^{-\lambda x} \quad (x > 0)$$

$$F(x) = P(X \leq x) = 1 - e^{-\lambda x}$$

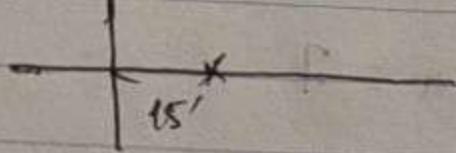
\* Mean and Variance

$$E(X) = \frac{1}{\lambda}$$

$$\text{Var}(X) = \frac{1}{\lambda^2}$$

$$\text{Ex 15, } E(X) = 9 = \frac{1}{\lambda} \Rightarrow \lambda = \frac{1}{9}$$

Standard



$$P(X > 15) = e^{-1/9 \cdot 15}$$

$$\text{Ex 16, } E(X) = 35 \rightarrow \lambda = 1/35$$

$$P(X > 6) = e^{-\lambda 6} = e^{-1/35 \times 6}$$

$$\text{Ex 17, } \lambda = 15/n$$

$$P(X \leq 6) = P(X \leq 1/10(n)) = 1 - e^{-15/10}$$

## CHAPTER 6: Descriptive Statistics

Take a sample of population (data) »

$$x_1, x_2, \dots, x_n$$

From a population  $N$

\* Definition:

$$\textcircled{1} \text{ Sample mean } \bar{x} = \frac{\sum x_i}{n}$$

$$\textcircled{2} \text{ Sample median (sorted data)}$$

KOKUYO

$$L = \frac{n+1}{2} \rightarrow \text{median} = \frac{x_{\lceil \frac{n+1}{2} \rceil} + x_{\lfloor \frac{n+1}{2} \rfloor}}{2}$$

Ex 1: 8 2 8 8 7 7  
 $\rightarrow 2 7 7 8 8 8$  (sorted),  $n=6$

$$L = \frac{6+1}{2} = 3.5$$

$$\rightarrow \text{median} = \frac{x_3 + x_4}{2} = \frac{7+8}{2}$$

Ex 2: 1 2 3 7 8 9 10

$$L = \frac{7+1}{2} = 4$$

$$\rightarrow \text{median} = x_4 = 7$$

### ③ Sample mode

Ex 1:

1, 8 2 8 8 7 7

mode: 2, 8

2, 5 6 8 & 10 10

mode: 8, 10

3, 1 2 3 7 8 9 10

mode: 3

### ④ Variance

1, Sample range

$$\Rightarrow r = \max(x_i) - \min(x_i)$$

2, Variance with direct method

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

$$S^2 = \frac{\sum x_i^2 - n\bar{x}^2}{n-1}$$

Standard deviation:  $s = \sqrt{s^2}$

Ex: 32 41 47 50 56 58 61  
62 68 75 82 85 90 92 98

Find sample mean and standard deviation

$$\bar{x} = 66.4665.09$$

$$s \approx 18.7832$$

$$s = 19.4$$

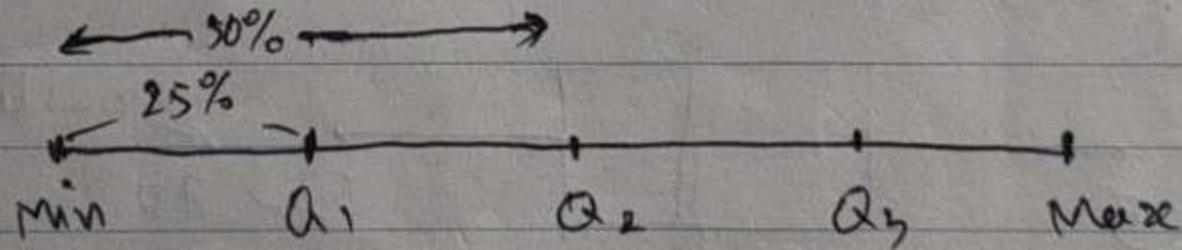
### ⑤ Stem and leaf diagram

Ex1: 55 115 225 290 330  
335 385 400 405 405

stem		leaf
5		3 5
11		5
22		5
24		0
33		0 5
38		5
40		0 5 5

(!) X. Y : stem X, leaf Y

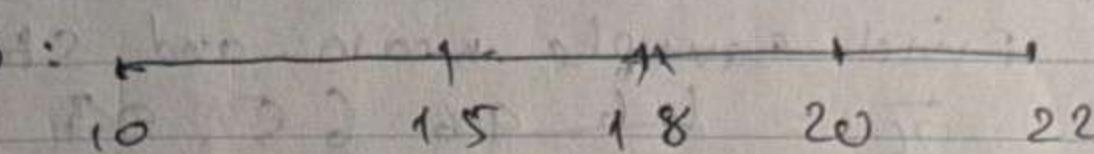
### ⑥ Quartiles



⑦ Interquartile range  $IQR = Q_3 - Q_1$

Ex: 55, 52, 52, 52, 49, 79, 67, 55

Find the sample quartiles and interquartile range

Ex 3: 

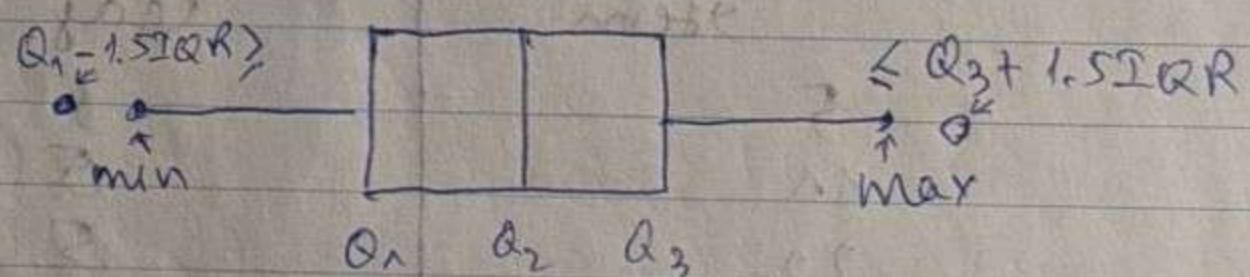
How many of the following statements are correct

i) About 50% of the data lies above 15

ii) IQR = 15

iii) The number 10 must be in the data

## 8. Box plot



Ex: 15, 20, 31, 31, 13, 2, 40, 41, 41, 41, 42, 43, 45, 45, 50, 70

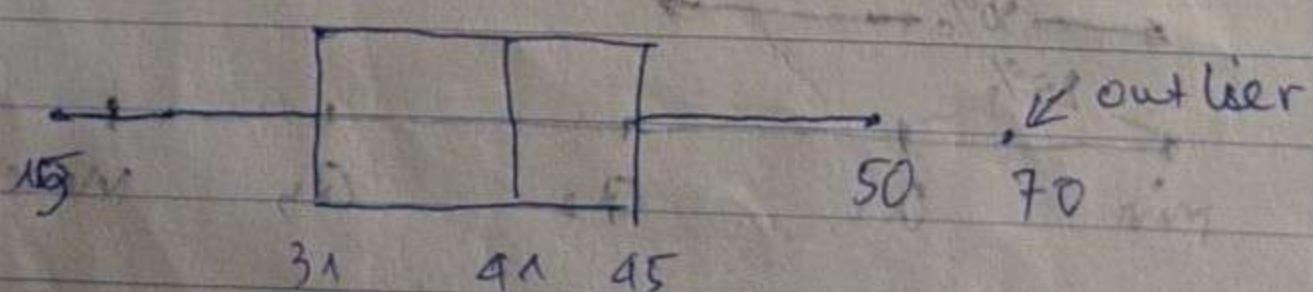
Draw a box plot for this data

$$\text{min} = 15, Q_1 = 31, Q_2 = 41, Q_3 = 45, \text{Max} = 70$$

$$\text{IQR} = Q_3 - Q_1 = 14$$

$$Q_1 - 1.5 \text{ IQR} = 10$$

$$Q_3 + 1.5 \text{ IQR} = 66$$



<sup>box plot</sup>  
Ex 2: Construct the ~~data set~~ for the data set  
 100, 113, 115, 116, 118, 123, 125, 136, 136,  
 138, 140, 142, 145, 146, 149, 150, 156, 162, 163

Ex 3:

stem		leaf
0		5 5 6
1		1 4 6 2 3
2		0 2 1

Ex 2:

$$Q_1 = 117, Q_2 = 137, Q_3 = 147.5$$

$$\text{IQR} = 30.5$$

$$\text{Min} = 117 - 1.5 \times 30.5 = 71.25$$

$$\text{Max} = 147 + 1.5 \times 30.5 = 183.25$$

No outlier

Ex 3:

$$Q_3 = 20$$

⑨ Frequency distribution (both ~~cate~~ and cont )

Ex: 1 1 2 2 2 3 3 4 5 6 7 7 8

x		Freq
1		2
2		3
3		2
4		1
5		1
6		1

10

15

20

KOKUYO

=> data : divide : 5-20 [ classes interval bins]

X	Freq
[1-2]	5
[3-4]	3
[5-6]	2
[7-8]	3

E x 2:

X	Freq	Cumulative Freq	Relative Freq	Cum Rel Freq
[0,2)	2	2	2/30	2/30
[2,4)	5	7	5/30	7/30
[4,6)	8	15	8/30	15/30
[6,8)	11	26	11/30	26/30
[8,10)	4	30	4/30	1

E x 3: The age of 25 member : take part in club

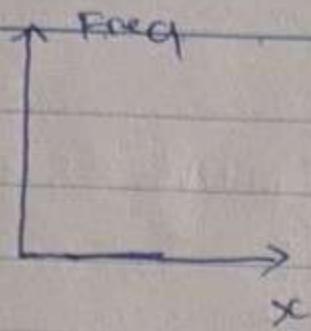
16 - 20	20%
21 - 30	28%
31 - 40	36%
41 - 50	16%

How many members in club that their age are not exceed 40?

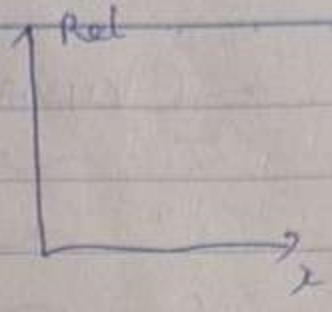
## ⑩ Histogram

Thứ  
Ngày

No.



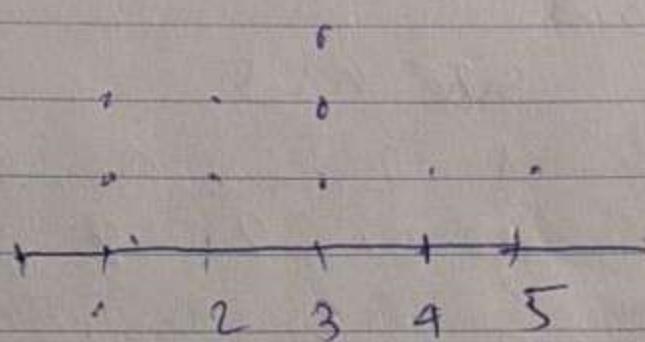
Histogram



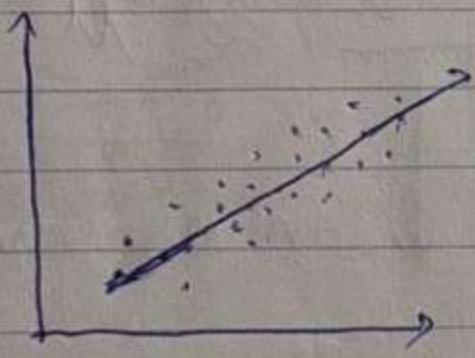
Relative Histogram

## (11) Dot plot

Ex: 1 1 1 2 2 3 3 3 4 5



## (12) Scatter plot

CHAPTER 7: Sampling distribution  
and point estimator

Big Picture :

\* Claim

$$\begin{aligned}\bar{x} &\rightsquigarrow \mu \\ s^2 &\rightsquigarrow \sigma^2 \\ \hat{P} &\rightsquigarrow P\end{aligned}$$

$\bar{x}$ : point estimator for  $\mu$   
 $s^2$ :  $\hat{s}^2$   
 $P$ :  $\hat{P}$

① Point estimator for  $\mu$

X   1 2 3 4 5 6
P(X)   $1/6$ $1/6$ $1/6$ $1/6$ $1/6$ $1/6$

$$\mu = 3.5$$

sample : size  $n=5$

$$c=1: 1 2 3 4 5 \rightarrow \bar{x}_1 = 3$$

$$c=2: 1 1 2 2 6 \rightarrow \bar{x}_2 = 2.4$$

$$\dots$$

$$\rightarrow \bar{x}_n = \dots$$

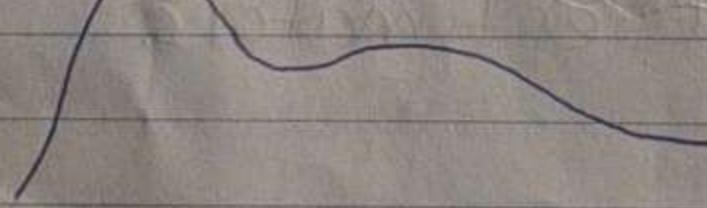
$$\frac{\bar{x}_1 + \bar{x}_2 + \dots + \bar{x}_n}{n} \approx \mu$$

Remark : The mean of sample :  $\mu_{\bar{x}} = \mu$

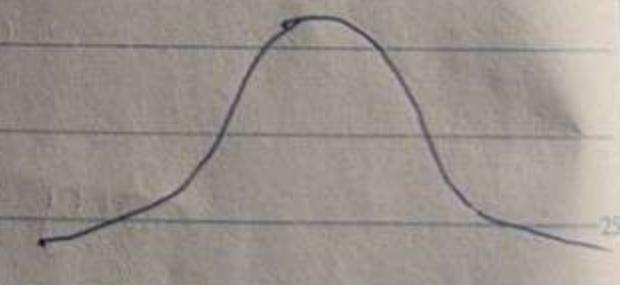
point estimator for  $\mu$   $\bar{x} \sim \mu$   
 unbiased estimator

$$(\mu, \sigma^2)$$

$$\bar{x} \sim N(\mu, \sigma^2)$$



sampling  $\rightarrow$



NORMAL

$\bar{x}$

Normal

$$\bar{x} \left\{ \begin{array}{l} \mu_{\bar{x}} = \mu \\ \sigma_{\bar{x}} = \sigma \end{array} \right.$$

10

15

20

$$\text{Standard error} = \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

\* Central limit theorem

$$[\text{If } X \text{ is normal and } n \geq 30 \rightarrow \bar{X} \sim N(\mu, \frac{\sigma^2}{n}) \xrightarrow{\text{sample}} \text{normal}]$$

$$Z \sim N(0, 1)$$

standard normal

- Theorem [CLT for two population]

If we have 2 independent populations with parameters  $(\mu_1, \sigma_1^2)$  and  $(\mu_2, \sigma_2^2)$ , and if  $\bar{X}_1$  and  $\bar{X}_2$  are the same means of 2 independent random samples of size  $n_1$  and  $n_2$  from these populations, then the sampling distributions of

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

is approximately standard normal for large  $n_1, n_2$ . It is exactly standard normal if the 2 populations are normal.

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

$$\bar{X}_1 - \bar{X}_2 \geq \mu_1 - \mu_2 + \frac{\sigma_1^2 + \sigma_2^2}{n_1 + n_2} =$$

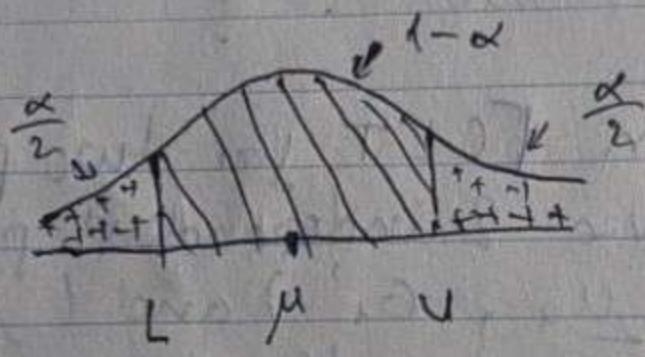
## CHAPTER 8: CONFIDENCE

### INTERVAL

*95% confident level*

$$\text{Ex: } \bar{x}_1 \xrightarrow{\mu} \xrightarrow{95\% \text{ confident level}}$$

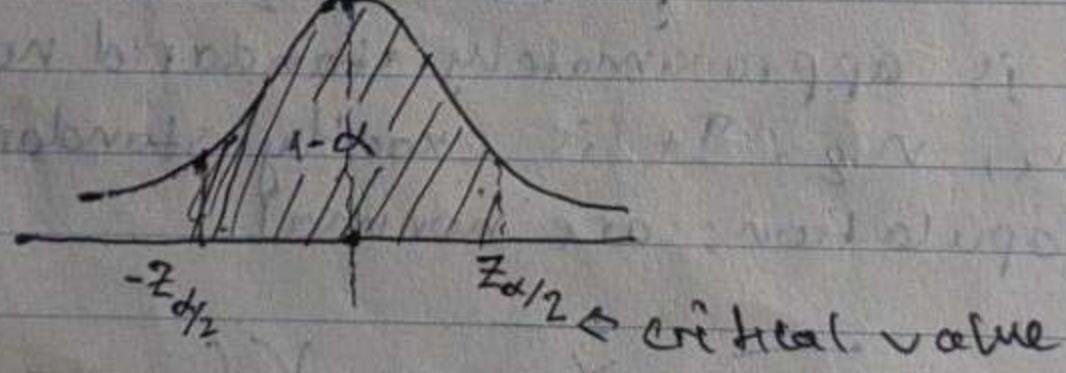
- ① Confidence level for mean ( $\sigma$ : known)



$$P(L \leq \mu \leq U) = 1 - \alpha$$

*confidence coefficient level*

Sample  $n$   $Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$



$$P(-z_{\alpha/2} < \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} < z_{\alpha/2}) = 1 - \alpha$$

$$= P(\bar{x} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

\* Margin of error

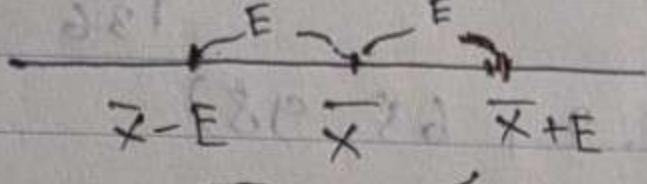
$$E = Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

$$\Rightarrow \bar{x} - E < \mu < \bar{x} + E$$

$\Rightarrow (1 - \alpha)$ , 100% confidence interval for mean

$$(\bar{x} - E, \bar{x} + E)$$

(2)



$$\text{width of confidence interval} = 2E$$

(3)

$$E \leq \alpha, \sigma, 1 - \alpha$$

$\Rightarrow$  How large of sample mean?  
 $\hookrightarrow$  find  $n = ?$

$$E = Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \Rightarrow n = \lceil \left( \frac{\sigma \cdot Z_{\alpha/2}}{E} \right)^2 \rceil$$

= critical value / standard error

$$Z_{\alpha/2} = -\text{NORM.S.INV} \left( \frac{\alpha}{2} \right) = \text{NORM.S.INV} \left( 1 - \frac{\alpha}{2} \right)$$

Confidence level

$Z_{\alpha/2}$

90%

1.64

95%

1.96

98%

2.33

99%

2.58

Ex 1:

$$n = 36$$

$$\bar{x} = 65.3, \sigma = 2.5$$

$$1 - \alpha = 0.9 \Rightarrow \frac{\alpha}{2} = 0.05 \Rightarrow z_{0.05} = 1.69$$

$\Rightarrow 90\%$  CI for mean is  $(\bar{x} \pm E)$

$$(65.3 - 1.69 \cdot \frac{2.5}{\sqrt{36}}, 65.3 + 1.69 \cdot \frac{2.5}{\sqrt{36}})$$

$$= (64.2, 65.98)$$

90% confidence that the true mean is between 64.2 and 65.98

Ex 2:

$$1 - \alpha = 0.98 \Rightarrow \frac{\alpha}{2} = 0.01 \Rightarrow z_{0.01} = 2.33$$

$$\sigma = 6, n = ?$$

$$n = \left\lceil \left( \frac{2.33 \cdot 6}{2} \right)^2 \right\rceil = 13$$

Ex 3:

$$\mu_{\text{MC}} = 21000, \sigma = 1000, z_{0.05} = 1.96, n = 69$$

$$2E = 2 \cdot z_{0.05} \frac{\sigma}{\sqrt{n}} = 2 \cdot 1.96 \cdot \frac{1000}{\sqrt{69}} = 490$$

Ex 4

$$z_{0.05} = 1.96$$

$$E \leq 500$$

$$\sigma = 1000$$

$$n = \left\lceil \left( \frac{4000 \cdot 1.96}{500} \right)^2 \right\rceil = 12$$

Ex5:

$$\begin{aligned} z_{0.005} &= 2.58 \\ E &= 8 \\ n &= 21 \end{aligned}$$

$$n = \left\lceil \left( \frac{21 \cdot 2.58}{8} \right)^2 \right\rceil = 9$$

\*  $E \downarrow \Rightarrow [ \bar{x} \pm z_{\alpha/2} \downarrow \Rightarrow \alpha/2 \uparrow \Rightarrow 1 - \alpha \downarrow ]$

\* One-sided confidence bounds

$(1 - \alpha) \cdot 100\%$  lower-confidence bound is  $\bar{x} - E$

$(1 - \alpha) \cdot 100\%$  upper-confidence bound is  $\bar{x} + E$

$E$  = critical value, standard error

$$E = z_{\alpha} \cdot \frac{s}{\sqrt{n}}$$

Ex9:  $\alpha = 0.01 \quad 1 - \alpha = 0.98 \Rightarrow \alpha = 0.02 \downarrow$ 

$$n = 15$$

$$\bar{x} = 1.5$$

$$z_{0.02} = 2.05$$

$$\rightarrow 98\% \text{ lower confidence bound: } \bar{x} - E = 1.5 - \frac{2.05 \cdot 0.01}{\sqrt{15}} \approx 1.495$$

Ex11:

$$\bar{x} = 82 \quad n = 30$$

$$\alpha = 0.02$$

$$\text{a)} \quad 1 - \alpha = 0.98 + \frac{0.1}{10} \rightarrow \alpha = 0.01 \rightarrow z_{0.01} = 2.58$$

$$90\% \text{ CI} \rightarrow \left[ \bar{x} - \frac{2.6}{\sqrt{n}}, \bar{x} + \frac{2.6}{\sqrt{n}} \right]$$

$$= (76.25, 87.7487)$$

$$\text{b)} \quad z_{0.02} = 2.05$$

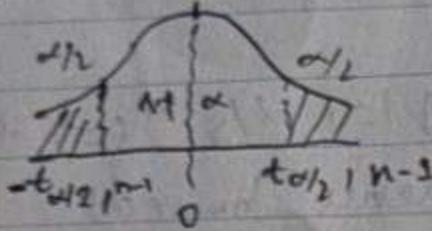
$\rightarrow 98\% \text{ upper-confidence bound}$

$$= \bar{x} + \frac{z_{0.02} \cdot s}{\sqrt{n}} \approx 86.56$$

KOKUYO

② Confidence interval for mean ( $\sigma$ : unknown)

$$\xrightarrow{\text{sample size } n} \begin{array}{c} \text{Normal distribution curve} \\ \text{Mean } \mu, \text{ Standard deviation } \sigma \\ \text{Confidence interval: } \bar{x} \pm t_{\alpha/2, n-1} \frac{s}{\sqrt{n}} \end{array}$$



$$s = (\text{sample standard deviation}) = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

unbiased estimator

degree of freedom  
 $df = n-1$

$$E = t_{\alpha/2, n-1} \cdot \frac{s}{\sqrt{n}}$$

Ex 5: Find the confidence interval of BP C.

$$n = 12$$

$$\bar{x} = 916.8$$

$$s = 15.2$$

$$1 - \alpha = 0.99 \rightarrow \alpha = 0.01 \Rightarrow t_{0.005, 11} = 3.106$$

$$(\bar{x} - E, \bar{x} + E) = (803, 121, 330, 428)$$

Ex 6:

$$\bar{x} = 2259, 912$$

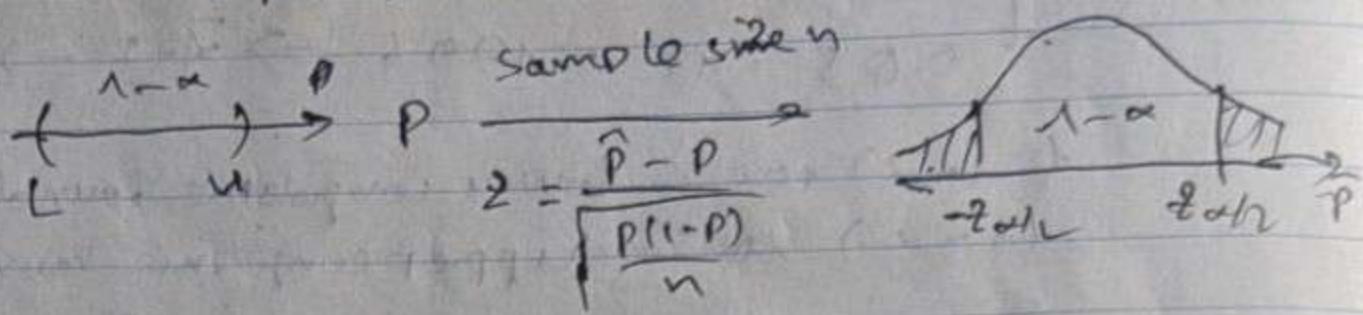
$$s = 35, 57$$

$$n = 12$$

$$\text{as } 1 - \alpha = 0.95 \rightarrow t_{0.025, 11} = 2.201$$

$$\Rightarrow (\bar{x} - E, \bar{x} + E) = (12239, 316, 2282, 517)$$

when  $| nP \geq 5 | nP(1-P) \geq 5 \Rightarrow \hat{P} \sim N(P, \frac{P(1-P)}{n})$



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

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

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

\* If  $P$  is unknown,  $\hat{P} = \hat{P}_{\text{max}} = 0.5$

\* One-sided confidence bound for proportion  
 $(1-\alpha) 100\%$  lower-confidence bound for proportion is  $\hat{P} - E$

$(1-\alpha) 100\%$  upper-confidence bound for proportion is  $\hat{P} + E$

$E$  = critical value standard error

$$E = z_{\alpha} \cdot \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}$$

② Confidence Variance

$$\text{for } S^2 \rightsquigarrow \sigma^2$$

$$\xrightarrow{\text{Sample size } n} \frac{x^2 = \frac{(n-1)s^2}{\sigma^2}}{\chi^2_{\alpha/2, n-1}} \xrightarrow{\text{Distribution}} \chi^2_{\alpha/2, n-1}$$

$$\rightarrow 1 - \alpha = P\left(\frac{(n-1)s^2}{\chi^2_{\alpha/2, n-1}} < \sigma^2 < \frac{(n-1)s^2}{\chi^2_{1-\alpha/2, n-1}}\right)$$

$$\Rightarrow (1 - \alpha) \cdot 100\% \text{ CI for variance is } \left( \frac{(n-1)s^2}{\chi^2_{\alpha/2, n-1}}, \frac{(n-1)s^2}{\chi^2_{1-\alpha/2, n-1}} \right) = N(A, B)$$

$(1 - \alpha) \cdot 100\%$  CI for standard deviation is  $N(\sqrt{A}, \sqrt{B})$

## CHAPTER 9: Test of hypothesis for a single sample

Ex: 95% CI for mean is  $(65.21, 66.79)$

95% confident that the true mean is between 65.21 and 66.79

① claim:  $\mu = 60$

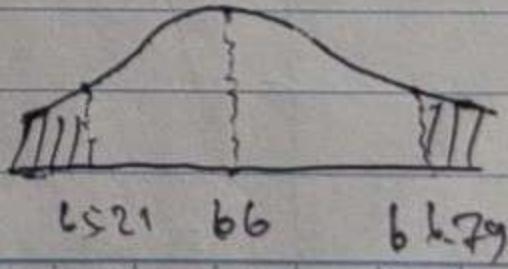
reject claim

② claim:  $\mu \neq 66$

reject claim

③ claim:  $\mu = 66$   
↓ maybe

Fail to reject claim



$\epsilon \square$ : reject

$\epsilon \square$ : fail to reject

KOKUYO

: acceptance region  
 : critical region

Def:

1, Statistical hypothesis is a statement about the parameter ( $\mu, \sigma^2, P$ ) of one or more population

3, claim  $\Rightarrow \begin{cases} H_0: \text{null hypothesis (contain "=")} \\ H_1: \text{alternative hypothesis (not contain "=")} \end{cases}$

Ex 1: 1) claim:  $\mu = 60 \rightarrow \begin{cases} H_0 = \mu = 60 \\ H_1 = \mu \neq 60 \end{cases}$

2) claim  $\mu \neq 66 \rightarrow \begin{cases} H_0 = \mu = 66 \\ H_1 = \mu \neq 66 \end{cases}$

3) claim  $\mu > 66 \rightarrow \begin{cases} H_0 = \mu \leq 66 \\ H_1 = \mu > 66 \end{cases}$

4) claim  $\mu < 66 \rightarrow \begin{cases} H_0 = \mu \geq 66 \\ H_1 = \mu < 66 \end{cases}$

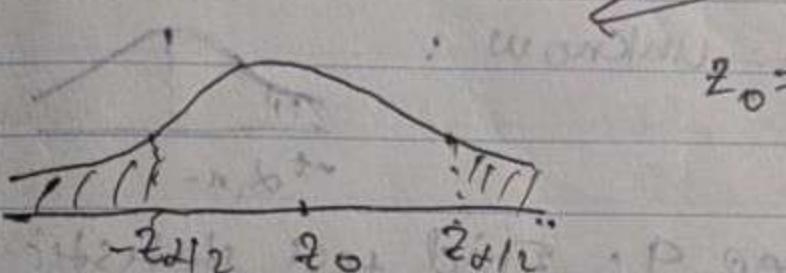
Ex: claim  $\mu = 66 \Rightarrow \begin{cases} H_0: \mu = 66 \\ H_1: \mu \neq 66 \end{cases}$   
 $(65, 21; 66, 29)$

Sample  $\Rightarrow \bar{x}$   
 Chap 8

$$\xrightarrow{\text{Chap 8}} \frac{1}{\bar{x}-E} \quad \frac{1}{\bar{x}+E} \Rightarrow (\bar{x}-E, \bar{x}+E)$$

$\rightarrow \begin{cases} \mu \in \square \\ \mu \notin \square \end{cases} \rightarrow \text{conclusion}$

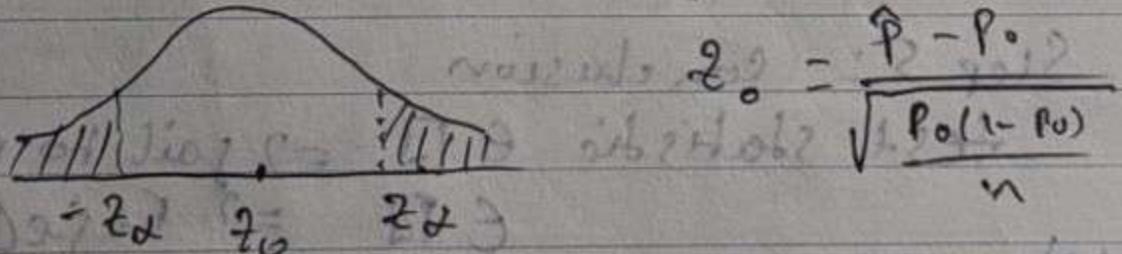
$$z_0 = \frac{\bar{x} - \mu}{G/\sqrt{n}}$$



$\rightarrow \begin{cases} z_0 \in \square : \text{fail to reject } H_0 \\ z_0 \in \square : \text{reject } H_0 \Rightarrow \text{reject claim} \end{cases}$

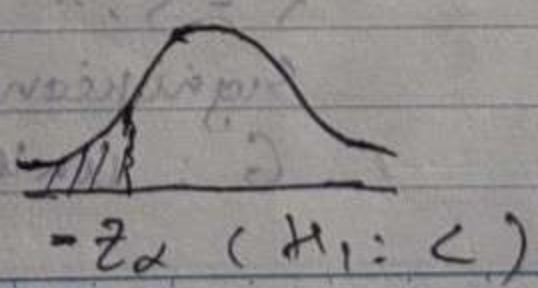
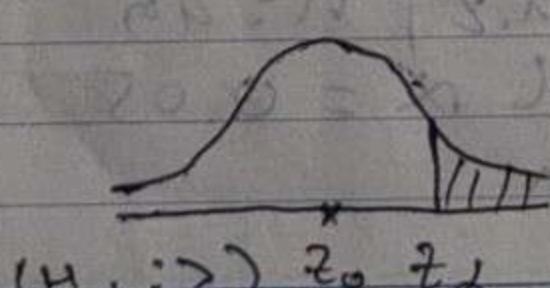
Ex: claim:  $p > 0.4$

$\Rightarrow \begin{cases} H_0: p \leq 0.4 \Leftrightarrow p = 0.4 \\ H_1: p > 0.4 \end{cases}$



$z_0 \in \square : \rightarrow \text{fail to reject } H_0 \Rightarrow \text{reject}$

$z_0 \in \square : \rightarrow \text{reject } H_0 \Rightarrow \text{fail to reject}$



$(H_1: >) z_0 z_d$

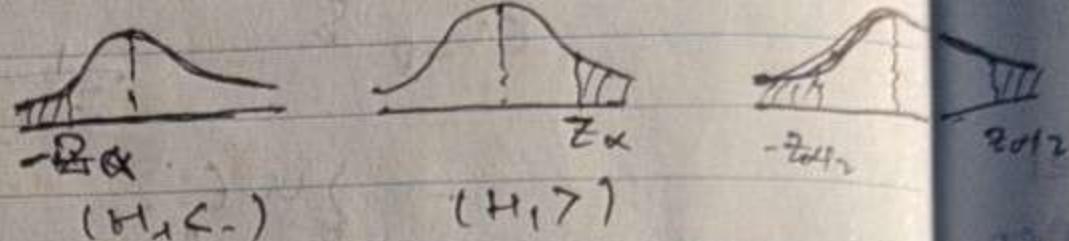
$-z_d (H_1: <)$

KOKUYO

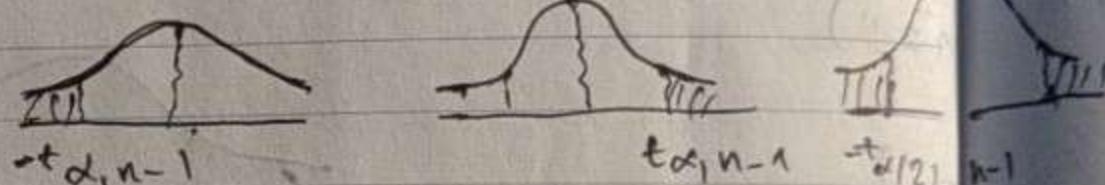
\* Test w/ a hypothesis (right tail test) (two + all test)  
 Step 1: claim  $\mu \leq \mu_0$   $H_0: \mu \leq \mu_0$   $H_1: \mu > \mu_0$   
 Step 2:  $H_0: \mu \geq \mu_0$   $H_0: \mu \leq \mu_0$   $H_0: \mu = \mu_0$   
 $H_1: \mu < \mu_0$   $H_1: \mu > \mu_0$   $H_1: \mu \neq \mu_0$

Step 3: Find critical value

6: Know:



6: unknown:



Step 4: Find test statistic

• For mean:  $\hat{\sigma}_x$  : unknown :  $z_0 = \frac{\bar{x} - \mu_0}{\hat{\sigma}_x / \sqrt{n}}$   
 [normal]  
 $n \geq 30$

$\hat{\sigma}_x$  : unknown :  $t_0 = \frac{\bar{x} - \mu_0}{s / \sqrt{n}}$

• For proportion

$$z_0 = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

Step 5: Conclusion

test statistic  $\in \text{I}$   $\Rightarrow$  fail to reject  $H_0$

$\in \text{II}$   $\Rightarrow$  reject  $H_0$

(?) Note

Ex 3:

$$\bar{x} = 5.1, s = 1.2, n = 99$$

$$\text{Significant level } \alpha = 0.05$$

6: unknown

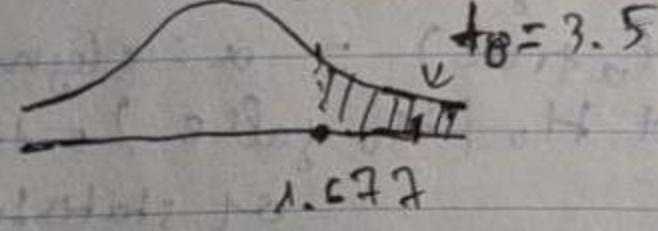
test)

$$\text{claim } \mu \geq 9.5 \rightarrow \begin{cases} H_0: \mu \leq 9.5 \\ H_1: \mu > 9.5 \end{cases}$$

(right-tailed test)

$$\text{critical value: } t_{\alpha/2, n-1} = t_{0.05, 48} = 1.677$$

$$\text{test statistic: } t_0 = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}} = \frac{5.1 - 9.5}{\frac{4.2}{\sqrt{7}}} = -3$$



$t_0 = 3.5$   
 $t_0 > 1.677 \rightarrow \text{reject } H_0$

$\rightarrow \text{fail to reject } C$

$$\textcircled{4} \quad \text{claim } p = 0.6 \quad \begin{cases} H_0: p = 0.6 \\ H_1: p \neq 0.6 \end{cases}$$

$$n = 190, \hat{p} = \frac{n}{n} = \frac{75}{190} = 0.53$$

$$\alpha = 0.1$$

$$\rightarrow z_{\alpha/2} = z_{0.05} = 1.695$$

$$\text{and } z_0 = \frac{p - p_0}{\sqrt{\frac{p(1-p)}{n}}} = -1.55 \in \square$$

$\rightarrow \text{fail to reject } H_0$

$$\text{CI } (65.21, 66.79)$$

↓

claim

↓

reject claim  $\rightarrow$  error

$H_0: \text{true } \rightarrow \text{error}$   
 $\text{rejection region}$

$H_0$  true       $H_0$  false  
 Fail to      ✓      Error type II ( $\beta$ )  
 reject  $H_0$       ✓      ( $H_0$  false)  
 Reject  $H_0$       Error type I ( $\alpha$ )      ✓       $(1-\beta)$

$1-\alpha$ : confidence level

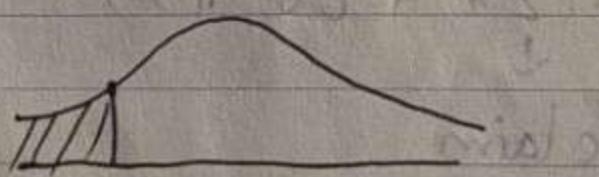
$P(\text{Reject } H_0 \mid H_0 \text{ true}) = \alpha$ : significant level  
 $P(\text{fail to reject } H_0 \mid H_0 \text{ false}) = 1-\beta$ : the power  
 of statistical test

Fix  $n \rightarrow \left\{ \begin{array}{l} \alpha : \text{fix} \\ \beta : \downarrow \end{array} \right.$

Ex5: Chon  $\mu < 150$   
 $\Rightarrow \left\{ \begin{array}{l} H_0: \mu \geq 150 \quad (\mu = 150) \\ H_1: \mu < 150 \end{array} \right.$

$$G = 16.2, \bar{X} = 153.36, \alpha = 0.01, n = 11$$

$$Z_0 = \frac{\bar{X} - \mu}{G/\sqrt{n}} = \frac{153.36 - 150}{16.2/\sqrt{11}} = 0.67$$



$$\alpha = 0.01$$

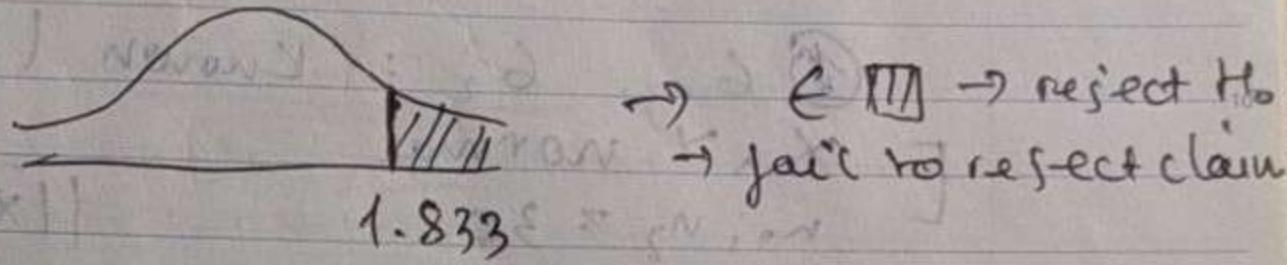
$$(H_1: \mu < 150)$$

$\rightarrow z_0 \in \square \rightarrow \text{fail to reject } H_0$   
 $\rightarrow \text{reject claim}$

$$\text{Ex 6: } \bar{x} = 537.1, \alpha = 0.05, t_{0.05, 19} = 1.833$$

$$\mu > 520 \rightarrow \begin{cases} H_0: \mu \leq 520 \\ H_1: \mu > 520 \end{cases}$$

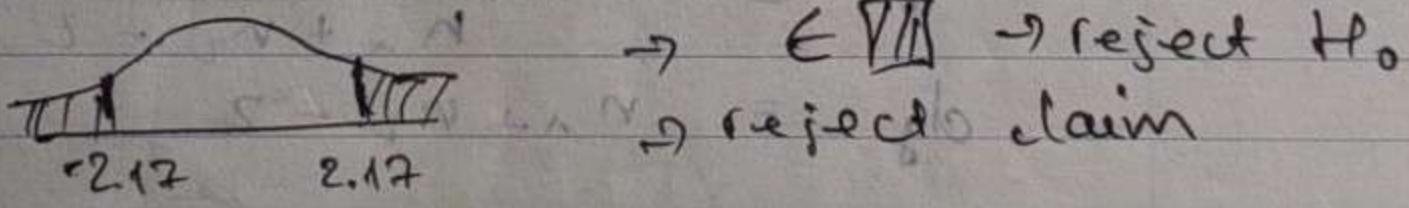
$$t_0 = \frac{537.1 - 520}{20.7 / \sqrt{19}} = 2.41$$



$$\text{Ex 7: } \hat{p} = 0.06, n = 200$$

$$\begin{cases} H_0: \mu = 3\% = \alpha \rightarrow 2\alpha_2 = 2.12 \\ H_1: \mu \neq 3\% \end{cases}$$

$$z_0 = \frac{0.06 - 0.03}{\sqrt{\frac{0.03(1-0.03)}{200}}} = \cancel{4.76} 2.487$$



$$\sigma^2 = \frac{pq}{n} = \frac{0.06 \cdot 0.94}{200} = 0.00141$$

$$\sigma = \sqrt{0.00141} \approx 0.0375$$

CHAPTER 10: Statistical Inference  
for two sample

For mean ( $\mu_1 - \mu_2$ )

$$\begin{cases} X, Y: \text{normal} \\ n_1, n_2 \geq 3 \end{cases} \Rightarrow \begin{cases} \bar{X} \sim N(\mu_1, \frac{\sigma_1^2}{n_1}) \\ \bar{Y} \sim N(\mu_2, \frac{\sigma_2^2}{n_2}) \end{cases}$$

$$\Rightarrow \bar{X} - \bar{Y} \sim (\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2})$$

Table for  $\mu_1 - \mu_2$

①  $\sigma_1, \sigma_2$ : known  $((\bar{X} - \bar{Y}) \pm E)$

$X, Y$  normal  $\stackrel{CI}{\Rightarrow} ((\bar{X} - \bar{Y}) \pm E)$

$n_1, n_2 \geq 30$

$\Rightarrow$  use Z  $E = \text{critical value} \cdot \text{Standard err or}$

$$E = z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad Z_0 = \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

②  $\sigma_1, \sigma_2$ : unknown

assume:  $\sigma_1^2 = \sigma_2^2$

$s_p^2 \Rightarrow$  pooled variance

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

$$\text{min. } df = n_1 + n_2 - 2$$

$((\bar{X} - \bar{Y}) \pm E)$

$$E = t_{\alpha/2, df} \sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}}$$

$$t_0 = \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}}}$$

$$1 < \frac{s_{\text{large}}^2}{s_{\text{small}}^2} < 2 \Rightarrow \sigma_1^2 = \sigma_2^2$$

b) we assume:  $\sigma_1^2 = \sigma_2^2$

~~df =~~

$$df = \frac{\left( \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^2}{\frac{(s_1^2/n_1)^2}{n_1-1} + \frac{(s_2^2/n_2)^2}{n_2-1}}$$

$$E = \text{F}_{12, df} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$t_0 = \frac{(\bar{x} - \bar{y}) - \Delta_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Ex:

$$\mu_1 - \mu_2 > 0$$

$$H_0: \mu_1 - \mu_2 \leq 0$$

$$(H_A: \mu_1 - \mu_2 > 0)$$

test statistic  $t_0 = \frac{\bar{x}_1 - \bar{x}_2 - \Delta_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$

$$= -0.1$$

$$S_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} \approx 6.97875$$

$$df = n_1 + n_2 - 2 = 24$$

$$t_0 = \frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = 0.0622$$

F1 F2 F3 F4 F5 F6 F7

@ # & ^ \_

Thứ  
Ngày

No.

$$d_f = 29$$

critical value:

$$t_{0.05, 29} = \dots$$

Ex 4:

$$H_0: \mu_1 - \mu_2 = 0 \quad b_1 = 0.526$$

$$H_1: \mu_1 - \mu_2 \neq 0 \quad b_2 = 0.776$$

$$\alpha = 0.05 \rightarrow \alpha/2 = 0.025 \quad n_1 = 20, n_2 = 17$$

$$\bar{x}_1 = 3.9, \bar{x}_2 = 3.2$$

$$z_{0.025} = 1.96$$

$$a) z_0 = \frac{3.4 - 3.2}{\sqrt{\frac{0.526^2}{20} + \frac{0.776^2}{17}}} = 0.877 \approx 0.88$$

$$\begin{aligned} p\text{-value} &= 2 P(Z > 0.88) \\ &\approx 2(1 - P(Z \leq 0.88)) \\ &= 2(1 - 0.81) = 0.38 \end{aligned}$$

$$b) E = 1.867 \cdot \sqrt{\frac{0.526^2}{20} + \frac{0.776^2}{17}} = 0.379$$

$$\text{Upper: } (\bar{x}_1 - \bar{x}_2) + E = 0.574$$

Table for  $P_1 - P_2$ :

$$\hat{P}_1 = \frac{x_1}{n_1}$$

$$\hat{P}_2 = \frac{x_2}{n_2}$$

$$\bar{P} = \frac{x_1 + x_2}{n_1 + n_2}$$

$$CI: (\hat{P}_1 - \hat{P}_2) \pm E$$

$$E = z_{\alpha/2} \sqrt{\frac{\hat{P}_1(1-\hat{P}_1)}{n_1} + \frac{\hat{P}_2(1-\hat{P}_2)}{n_2}}$$

$$H_0: P_1 - P_2 = 0$$

$$Test Statistic: z_0 = \frac{(\hat{P}_1 - \hat{P}_2) - 0}{\sqrt{\bar{P}(1-\bar{P})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

Ex 5:

$$\mu_1 = 7.6, \mu_2 = 6.0, 6_1 = 1.9, 6_2 = 1.7, \alpha = 0.01$$

$$n_1 = n_2 = 50, z_{0.01} = 2.33$$

$$H_0: \mu_1 - \mu_2 = 0 \rightarrow \text{two tails}$$

$$H_1: \mu_1 - \mu_2 \neq 0$$

$$\alpha = 0.02 \rightarrow z_{0.01} = 2.33$$

$$P\text{-value} = 2 P(Z > |z_0|)$$

$$z_0 = \frac{\mu_1 - \mu_2}{\sqrt{\frac{6_1^2}{n_1} + \frac{6_2^2}{n_2}}} = \frac{7.6 - 6.0}{\sqrt{\frac{1.9^2}{50} + \frac{1.7^2}{50}}}$$

$$= 2.297$$

KOKUYO

$\rightarrow z_0 \in \mathbb{D} \rightarrow$  fail to reject claim

7.

$$\bar{x}_1 = 12.0, \bar{x}_2 = 16.3$$

$$s_1 = 4.2, s_2 = 4.9$$

$$n_1 = 19, n_2 = 17$$

$$\alpha = 0.05$$

$$H_0: \mu_1 - \mu_2 \geq 0$$

$$H_1: \mu_1 - \mu_2 < 0 \rightarrow \text{left tail}$$

$$\rightarrow z_2 = 20.05, 29 = 170$$

$$t_0 = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s_1^2/n_1 + s_2^2/n_2}} = 0$$

$$1 < \frac{s_2^2}{s_1^2} < 1 \rightarrow s_1^2 = s_2^2$$

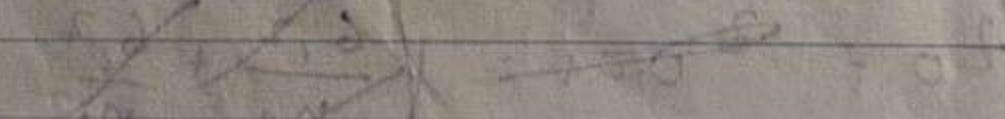
$$s_p = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}} \approx 18.59$$

$$t_0 = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = -2.185$$

$$\text{elbow out } 0.05, 29 = 1.7$$

$$\rightarrow t_0 \in \mathbb{D}$$

$$(1.81 < -2.185 < 2.77)$$



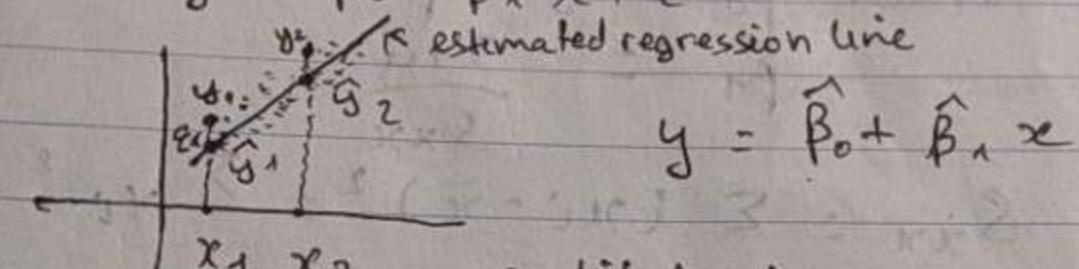
## CHAPTER 11: Simple linear regression and Correlation

\* Regression line:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

\* Simple linear regression model

$$y = \beta_0 + \beta_1 x + \epsilon$$



$$\epsilon_i = y_i - \hat{y}_i : \begin{cases} \text{predicted value} \\ \text{random error} \end{cases}$$

observed at  $x=x_i$

$\beta_0$ : y-intercept of population regression line

$\beta_1$ : slope of population regression line

estimated y-intercept

$$\hat{\beta}_0 \rightsquigarrow \beta_0$$

estimated slope

$$\hat{\beta}_1 \rightsquigarrow \beta_1$$

$$\hat{\sigma}^2 \rightsquigarrow \sigma^2$$

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

min

Def: The fitted or estimated regression line  
for data  $X, Y$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 \frac{x_i}{\bar{x}}, \quad \bar{x} = \frac{\sum x_i}{n}, \quad \bar{y} = \frac{\sum y_i}{n}$$

$$\hat{\beta}_1 = \frac{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} = \frac{s_{xy}}{s_{xx}}$$

Note:  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$

$$\hat{\beta}_1 = \frac{s_{xy}}{s_{xx}}$$

Sample variance of x

$$s_{xx} = \sum (x_i - \bar{x})^2 = \sum x_i^2 - \frac{(\sum x_i)^2}{n} = (n-1) s_x^2$$

$$s_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n}$$

\* Standard error of estimate

- Total variation (Total sum of square)

$$SS_T = s_{yy} = \sum (y_i - \bar{y})^2$$

- Error sum of square

$$SS_E = \sum (y_i - \hat{y}_i)^2 \quad (\text{unexplained variation})$$

- Regression sum of squares:

$$SS_R = \sum (\hat{y}_i - \bar{y})^2 = \hat{\beta}_1 s_{xy}$$

$$SS_T = SS_R + SS_E$$

Deg: unbiased estimator of standard error  
of estimate  $n$

$$\text{viet cho } \hat{s} = \sqrt{\frac{SS_E}{n-2}}$$

\* Coefficient of determination

$R^2 = \frac{SS_R}{SS_T}$

↑  
% variation of  $Y$  is explained  
by % variation of  $X$

% >>  $\Rightarrow$  cor relation: strong  
(correlation  $\in \mathbb{R}^2$ )

Note:  $0 \leq R^2 \leq 1$

$$+ R^2 = 1 \Leftrightarrow SS_R = SS_T$$
$$\Leftrightarrow SS_E = 0 \Leftrightarrow \sum (y_i - \hat{y})^2 = 0$$
$$\Leftrightarrow y_i = \hat{y}$$

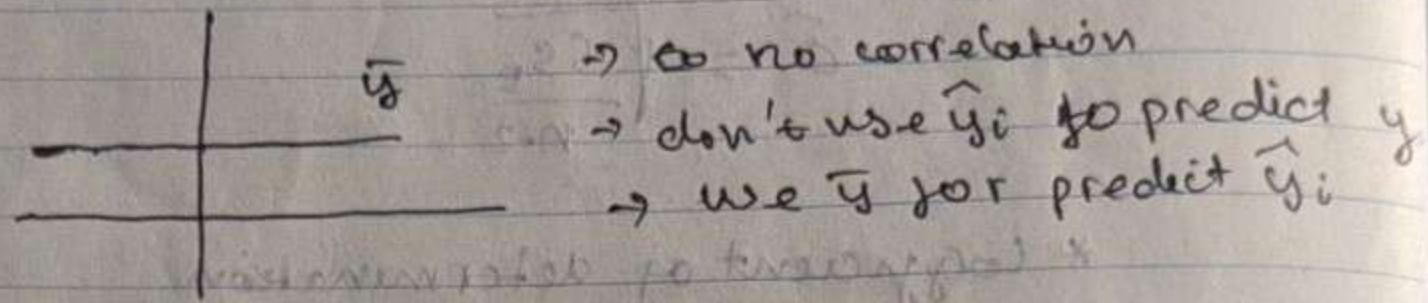
trong:  $100 \times 2$

(nhiều hơn)  $0 = 10 \cdot 10$

do đó trung bình  $E[10 \cdot 10] = 50$

$$+ R^2 = 0 \Rightarrow SS_R = 0 \Rightarrow \sum (\hat{y}_i - \bar{y})^2 = 0$$

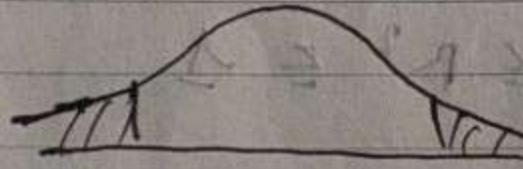
$$\Rightarrow \hat{y}_i = \bar{y}$$



1. Test of hypothesis about the slope  $\beta_1$

$$H_0: \beta_1 = \beta_{1,0} \quad (\text{Excel: } H_0: \beta_1 = 0) \quad H_1: \beta_1 \neq \beta_{1,0}$$

$$d_f = n - 2$$



Test statistic

$$t_0 = \frac{\hat{\beta}_1 - \beta_{1,0}}{se(\hat{\beta}_1)}$$

estimated Standard error of slope

$$se(\hat{\beta}_1) = \sqrt{\frac{S^2}{S_{xx}}}$$

Excel: Test

$$H_0: \beta_1 = 0 \quad (\text{no correlation})$$

$$H_1: \beta_1 \neq 0 \quad (\exists \text{ significant correlation})$$

Test  $H_0: \beta_1 = 0$   
 $H_1: \beta_1 \neq 0$

- Reject  $H_0$ :  $\exists$  a significant correlation between  $x_i$  and  $y$   
 $\Rightarrow$  the best predicted value of  $y$  at  $x = x_i$  is  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_i$
- Fail to reject  $H_0$ : no correlation  
 $\Rightarrow$  the best predicted value of  $y$  at  $x = x_i$  is  $\bar{y} = \frac{\sum y_i}{n}$

2. Test of hypothesis about the intercept coefficient  $\beta_0$ .

- y nh' slope  $\beta_1$

Khai: Tính chất test statistic

$$t_0 = \frac{\beta_0 - \beta_{0,0}}{se(\hat{\beta}_0)}$$

estimated standard error of intercept

$$se(\hat{\beta}_0) = \sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}$$

\* Continue:

④ Sample coefficient of correlation

$$R = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} \quad S_{xy} = \sum x_i y_i - n \bar{x} \bar{y}$$

$$S_{xx} = \sum x_i^2 - n \bar{x}^2$$

$$S_{yy} = \sum y_i^2 - n \bar{y}^2$$

$$- \text{with } S_{xx} = (n-1) S_{x^2}$$

$$S_{yy} = (n-1) S_{y^2}$$

$\square - F_1$  |  $\square + F_2$  |  $\times F_3$  |  $C F_4$  |  $\square F_5$  |  $\square F_6$  |  $\rightarrow F_7$  |  $\square F_8$  |  $\square F_9$  |  $\square F_{10}$  |  $\square F_{11}$  |  $\square F_{12}$

Thứ  
Ngày

№.

11)  $\hat{y} = 10225.802 + 69.964x$   
 $\hat{y}(30) = 10225.802 + 69.964 \cdot 30$

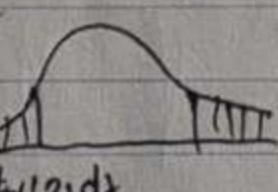
Review:

$$R = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} \quad \text{correlation coefficient}$$

Remark:  $R, \hat{B}_1$  same sign

$\rho$ :  $R \rightsquigarrow \rho$  unbiased estimator

sample coefficient population correlation coefficient



$\rightarrow$  test statistic

$$t_0 = \frac{R - \rho}{\sqrt{\frac{1-R^2}{n-2}}} \quad d\gamma = n-2$$

Ex 7: Test claim  $\rho = 0$ ,  $R = 0.833$ ,  $n = 12$ ,  $\alpha = 0.05$

$H_0: \rho = 0$

$H_1: \rho \neq 0$

$$t_0 = \frac{0.833 - 0}{\sqrt{\frac{1 - 0.833^2}{12 - 2}}} = 4.761$$

$$t_{\alpha/2, 10} = t_{0.025, 10} = +2.23$$

$t_0 \in \text{W} \rightarrow \text{reject } H_0$

$\rightarrow$  Therefore

Ex?? :  $x | 5 \quad 10 \quad 16 \quad 6 \quad 10 \quad 9$

$$y | 64 \quad 16 \quad 69 \quad 80 \quad 59 \quad 87$$

$$\alpha = 0.05, t_{0.025, 4} = 2.78,$$

a) Test claim there are regression linear correlation between  $x$  and  $y$ ?

$$\Rightarrow \rho \neq 0$$

$$H_0: \rho = 0 \quad (\text{no relation})$$

$$H_1: \rho \neq 0 \quad (\text{linear trend})$$

$$r = 0.2242$$

$$t_0 = \frac{0.2242 - 0}{\sqrt{\frac{1 - 0.2242^2}{4}}} = 0.96$$

$$\pm t_{0.025, 4} = \pm 2.78$$

$\rightarrow$  fail to reject  $H_0$

$\rightarrow$  don't use  $\hat{y}$  to predict  $y_i$  value at  $x_i$

$\rightarrow$  use  $\bar{y}$  to predict

Ex??  $x | 0 \quad 3 \quad 4 \quad 5 \quad 12$

$$y | 1 \quad 2 \quad 6 \quad 9 \quad 12$$

$$\alpha = 0.05$$

a) Test claim there are significant linear regression correlation between  $x$  and  $y$ ?

b) Find the best predicted value of  $y$  at  $x=0$