

## Lecture (2)

### Sampling Distribution

First we should understand the difference between Population & Sample

Population : Collection of all items of interest  
[set of all elements under discussion or investigation]

+ Values we obtain from population are called Parameters

Sample subset of population

+ values we obtain from sample called Statistics

Ex to study IQ level of 700,000 student

Population  
every student (in campus, at home, part time, ...)

- ⇒ hard to define
- ⇒ hard to observe in real life
- ⇒ hard to contact

Sample  
any subset of 100 student

- ⇒ Much easier to gather
- ⇒ less time consuming
- ⇒ less cost

But choice of sample is  
⇒ really a big deal

①

That's why most of statistics studies work on sample of data

Sample must have two characteristics:

Random (Indep)

each member of sample is chosen from population by chance (unbiased)

representative

Sample should reflect accurately the member of entire population

Ex Survey on student for the entire school  
You contact your class mates

$\therefore$  the group you study called Sample, the values you get called Statistics

Random Sample

Ex measuring IQ level of 700,000 student  
Choose a sample of 100 student

Sample 1 :  $X_1, X_2, \dots, X_n$  Random variable  $n=100$   
Sample Size

$x_1, x_2, \dots, x_n$   
↓  
observations  
readings  
values

observed values of a random sample

Sample 2 :  $X_1, X_2, \dots, X_n$   
 $x_1, x_2, \dots, x_n$

different observed values

(7)



# Conditions of a Random sample

1)  $X_1, X_2, \dots, X_n$  are indep. R.V.

$$f(x_1, x_2, \dots, x_n) = f(x_1) f(x_2) \dots f(x_n)$$

joint pdf                      marginal pdfs

this means that selection<sup>#</sup>  $i+1$  does not depend on selection  $\#i$   
(1<sup>st</sup> selection does not affect 2<sup>nd</sup>, 3<sup>rd</sup>, ...)

2)  $X_1, X_2, X_3, \dots, X_n$  have the same distribution (identically distributed)

$$E(X_1) = E(X_2) \dots = E(X_n) = \mu \rightarrow \text{Population mean}$$
$$V(X_1) = V(X_2) \dots = V(X_n) = \sigma^2 \rightarrow \text{Population Variance}$$

$X_1$ : R.V. can take any value in each R.S.

$$\mu_{\text{population mean}} = \frac{\sum (X_i)}{N}$$

$\leftarrow$  all i<sup>th</sup> level of all 700,000 students  
 $N \rightarrow 700,000$

In most cases  $\mu$  is unknown  
So we use estimate for  $\mu$  called  $\hat{\mu}$   
(ch. 2, ch. 3)

But  $E(X_i) = \mu$                       identically distributed  
 $V(X_i) = \sigma^2$

From [1], [2] Any random sample

S:  $X_1, X_2, \dots, X_n$  are i.i.d  
independent identically distributed

Statistics It is a function of observable  
R.v.'s which itself an observable R.v.  
and does not contain any unknown parameter

For example, If  $X_1, X_2, \dots, X_n$  R.S.

$$\text{then } \overline{X} = \frac{1}{n} \sum_{i=1}^n x_i$$

Sample mean                      ↓  
sample size

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

Sample Variance

$\overline{X}, S^2$  are R.v.'s too

$\overline{X}, S^2$  are statistics

but  $\mu, \sigma^2$  are parameters  
(constants)



## Note

$\mu \neq \bar{X}$   
population mean  $\neq$  sample mean

but if we take a good sample we can converge to  $\mu$

$\hat{\mu}$  = value Point estimation ch. 2  
 $a < \hat{\mu} < b$  Confidence interval ch. 3

→ Each sample of same size  $n$   
Each sample has  $\bar{X}$ ,  $S^2$   
 $\bar{X}$ ,  $S^2$  are R.v.'s too

$$E(\bar{X}) = E\left(\frac{\sum x_i}{n}\right)$$

$$= \frac{1}{n} E(x_1 + x_2 + \dots + x_n)$$

$$= \frac{1}{n} \left[ E(x_1) + E(x_2) + \dots + E(x_n) \right]$$

identically distributed

$$= \frac{1}{n} \times n \mu = \mu$$

average mean of sample mean

$$\therefore \boxed{E(\bar{X}) = \mu}$$

= population mean

$$V(\bar{X}) = V\left(\frac{\sum x_i}{n}\right)$$

$$= \frac{1}{n^2} V(x_1 + x_2 + \dots + x_n)$$

identically distr

$$= \frac{1}{n^2} [n \sigma^2]$$

$$V(\bar{X}) = \frac{\sigma^2}{n}$$

variance of Sample  
mean =  $\frac{\text{Population var}}{\text{sample size}}$

## The distribution of $\bar{X}$

If  $X_1, X_2, \dots, X_n$  a R.S. from a normal  
Population with mean  $\mu$ , variance  $\sigma^2$   
if Population  $\sim \text{Normal}(\mu, \sigma^2)$

then  $\bar{X} \sim \text{Normal}(\mu, \frac{\sigma^2}{n})$

So 
$$\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = Z \sim N(0, 1)$$

this means that any R.S. from normal  
Population  $\sim \text{Normal}(\mu, \sigma^2)$  have sample mean  $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$   
where  $n$  is sample size

what happen if population is not normal

$\Rightarrow$  Central Limit theorem

If  $X_1, X_2, \dots, X_n$  a R.S. from any  
Population with mean  $\mu$  and variance  
 $\sigma^2$  then for large enough  $n$  ( $n \geq 30$ )

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

$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1) \quad \text{as } n \geq 30$$



Ex Certain tubes manufacture by a company have a mean lifetime of 900 hrs and standard deviation of 50 hrs

Find the probability that a random sample of 64 tubes taken from the group will have mean lifetime between 895, 910 hrs

Population Parameters  $\left\{ \begin{array}{l} \mu = 900 \text{ hrs} \\ \sigma = 50 \text{ hr} \end{array} \right.$

Since  $n = 64 \gg 30$  large enough by Central Limit theorem

$$\bar{X} \sim N\left(\mu_{\bar{X}} = \mu = 900, \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} = \frac{50}{8}\right)$$
$$\bar{X} \sim N\left(900, \frac{50}{8}\right)$$

$Pr(895 < \bar{X} < 910)$  by standardization  
Convert  $\bar{X}$  to  $Z$  (standard normal)

$$Pr\left(\frac{895-900}{\frac{50}{8}} < \frac{\bar{X}-\mu}{\frac{\sigma}{\sqrt{n}}} < \frac{910-900}{\frac{50}{8}}\right)$$

$$P(-0.8 < Z < 1.6)$$

$$= \Phi(1.6) - \Phi(-0.8)$$
$$= \Phi(1.6) - [1 - \Phi(0.8)]$$

$$= \Phi(1.6) + \Phi(0.8) - 1$$

As  $n \gg \mu \rightarrow 900$

$$= 0.733 \text{ using table or any package}$$

⑦