

Conclusion or Inferences about population data.

Date \_\_\_\_\_  
Page \_\_\_\_\_

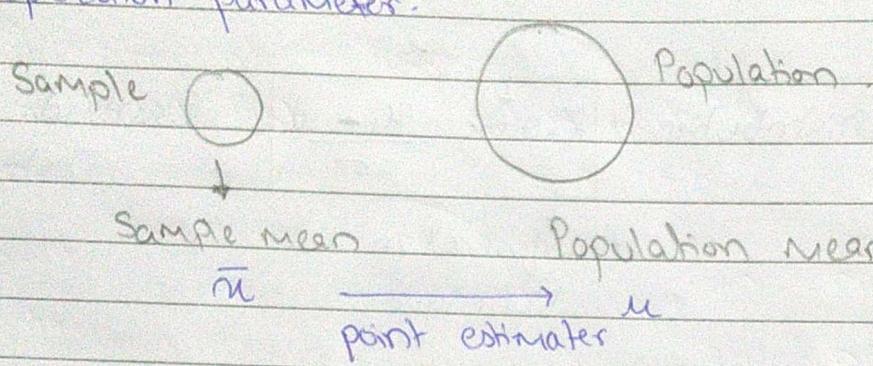
## INFERENTIAL STATS

### \* Estimate

It is an observed numerical value used to estimate an unknown population parameter.

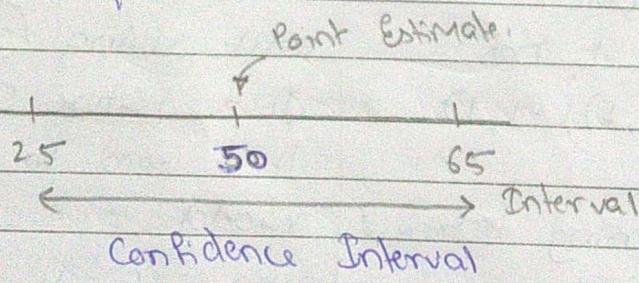
#### 1) Point Estimate

Single numerical value used to estimate the unknown population parameter.

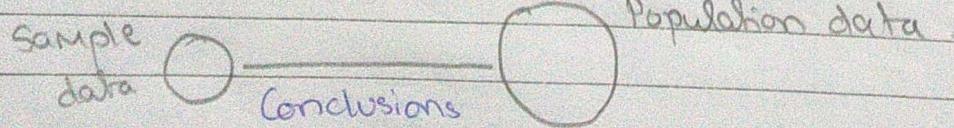


#### 2) Interval Estimate

Range of values used to estimate the unknown population parameter



### \* Hypothesis & Hypothesis Testing Mechanism



→ Made of hypothesis testing

what is given in question always goes to null hypothesis

Date \_\_\_\_\_  
Page \_\_\_\_\_

### 1) Null hypothesis ( $H_0$ )

- Assumption you are beginning with
- Ex: The person is not guilty

### 2) Alternate hypothesis ( $H_1$ )

- opposite of null hypothesis
- Ex: The person is guilty

z test,  
t test,  
Anova,  
chi-square

### 3) Evidences

- performing experiments  $\Rightarrow$  Statistical Analysis
- Ex: (fingure test, DNA, etc.. test)

### 4) we fail to reject the null hypothesis /

Reject the null hypothesis

#### \* P value

P value is a number, calculated from a statistical test, that describes how likely you are to have found a practical set of observations if null hypothesis were true.

- P values are used in hypothesis testing to help decide whether to reject the null hypothesis

→ Ex: Coin is fair or not

### Hypothesis Testing:

1) Null hypothesis  $H_0$  = Coin is Fair

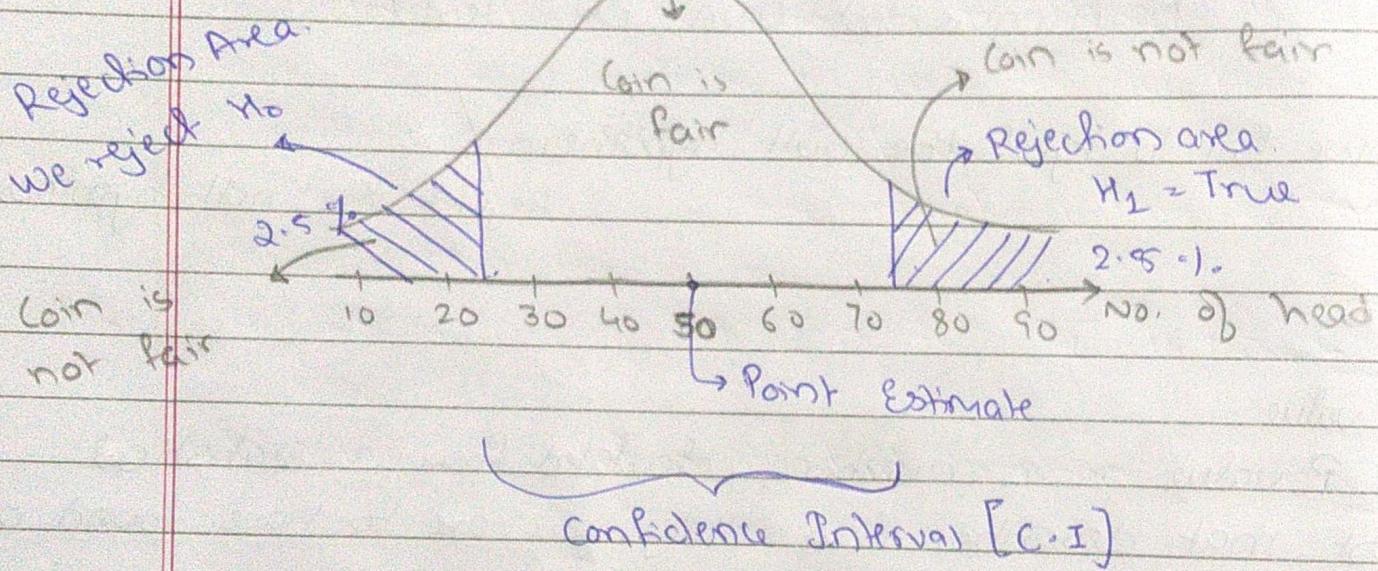
2) Alternate hypothesis  $H_1$  = Coin is not Fair

3) Experiment → Tossing

- 2 Tailed Test

Fail to reject  $H_0$

C.I = 95 %.



- Significance value ( $\alpha$ )

$$\alpha = 1 - C.I = 0.05$$

P value < Significance



we reject null hypothesis [Coin is not Fair]

else

We fail to reject null hypothesis [Coin is fair]

## \* Confidence Interval & Margin of Error

We construct a confidence interval to help estimate what is actual value of unknown population mean

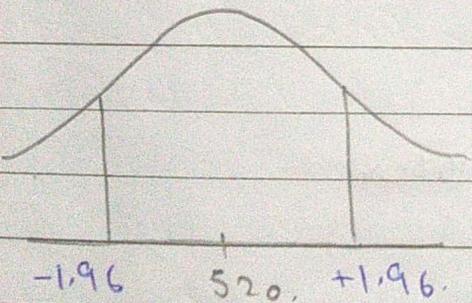
Point Estimate  $\pm$  Margin of Error

$$\bar{x} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

↑ Z-table  
↓ Significance value.

- Q. In verbal section of CAT Exam, standard deviation is known to be 100. A sample of 25 test takers has a mean of 520. Construct a 95% C.I about the mean?

Ans.  $\sigma = 100$ ,  $n = 25$ ,  $\bar{x} = 520$ , C.I = 0.95,  $\alpha = 0.05$



$$\bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$Z_{0.05/2} = Z_{0.025}$$

$$\begin{aligned} \downarrow \\ 480.8 \end{aligned} \quad \begin{aligned} \downarrow \\ 559.2 \end{aligned} \quad \begin{aligned} \downarrow \\ \text{Lower C.I} = 520 - (1.96) * \frac{100}{\sqrt{25}} = 480.8 \end{aligned}$$

$$\therefore \text{Higher C.I} = 520 + (1.96) * \frac{100}{\sqrt{25}} = 559.2$$

I am 95% confident/sure that Mean CAT score lies between 480.8 and 559.2

## → Hypothesis Testing & Statistical Analysis

- 1) Z test       $\downarrow \rightarrow$  Average
- 2) t test
- 3) Chi Square  $\rightarrow$  Categorical
- 4) Anova  $\Rightarrow$  variance

### • Z test :

- Q The average height of all residents in a city is 168 cm with  $\sigma = 3.9$ . A doctor believes the mean to be different. He measured the height of 36 individuals & found the avg height to be 169.5 cm.
- 1) State null & alternate hypothesis
  - 2) At a 95% C.I., is there enough evidence to reject null hypothesis

Ans.  $\mu = 168 \text{ cm}$ ,  $\sigma = 3.9$ ,  $n = 36$ ,  $\bar{x} = 169.5 \text{ cm}$

\* Whenever, population standard deviation is given used Z-test

- a) null hypothesis  $H_0: \mu = 168 \text{ cm}$
- b) Alternate hypothesis  $H_1: \mu \neq 168 \text{ cm}$
- c) C.I. = 0.95     $\alpha = 1 - 0.95 = 0.05$

### d) Statistical Analysis

$$Z\text{-test} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} \rightarrow \text{Standard Error}$$

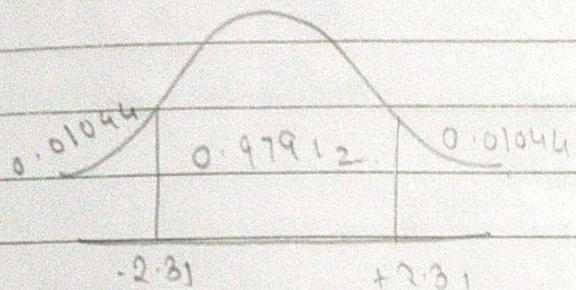
$$= \frac{169.5 - 168}{3.9 / \sqrt{36}} = 2.31$$

## Conclusion:

If Z-test value is less than -1.96 or greater than +1.96 we reject the null hypothesis.

$2.31 > 1.96 \rightarrow$  reject the null hypothesis

## • P-value



$$\begin{aligned}P\text{-value} &= 0.01044 + 0.01044 \\&= 0.02088\end{aligned}$$

If  $P\text{-value} < \text{significance value}$   
 $0.02088 < 0.05$

∴ Reject the null hypothesis.

When we have limited sample size (ie  $n < 30$ )

• T test → When we don't know population standard deviation

- Q. In the population the average IQ is 100. A team of researchers want to test a new medication to see if it has either a positive or negative effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140 with a standard deviation of 20. Did the medication affect intelligence? (I=95)

Ans

Given:

$$u = 100, n = 30, \bar{x} = 140, s = 20, C.I = 0.95, \alpha = 0.05$$

① Null hypothesis  $H_0 : u = 100$

② Alternative hypothesis  $H_1 : u \neq 100$  {2 tail test}

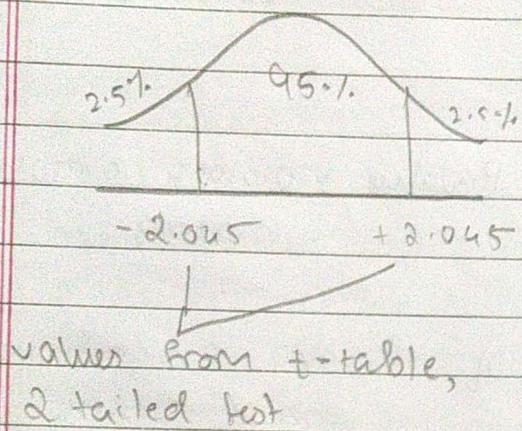
- 1 tailed test : sample mean = population mean / true

Date \_\_\_\_\_  
Page \_\_\_\_\_

③  $\alpha = 0.05$       C.I = 0.95

④ Degree of freedom [dof] =  $n - 1 = 30 - 1 = 29$

### ⑤ Decision boundary



Conclusion : If t-test is less than -2.045 & greater than 2.045, reject the null hypothesis

### ⑥ t test statistics

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{140 - 100}{20/\sqrt{30}} = 10.96$$

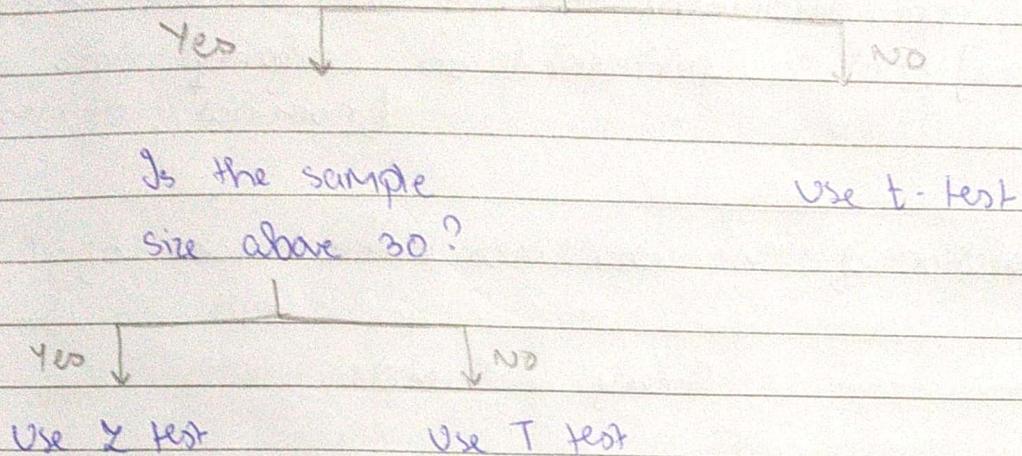
### ⑦ Conclusion

$10.96 > 2.045$  ... we reject null hypothesis

Here we don't know population standard deviation that's why we are doing t-Test

## \* When to use T-test vs Z-test

Do you know the population Std  $\sigma$ ?



## \* DISPERSION

Spread of the Data.

① Range : Upper limit - lower limit

10, 50, 9, 8, 25, 60

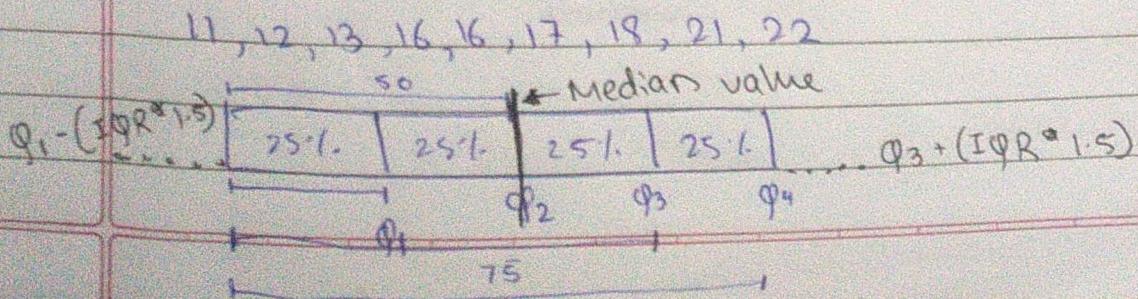
Sort  $\rightarrow$  8, 9, 10, 25, 50, 60  $\rightarrow$  upper  
low  $\downarrow$

$$60 - 8 = 52$$

② Percentile [Quantile, Decile]

$\rightarrow$  5 no. of summary  $\Rightarrow$  which describes the data

$$\begin{aligned} IQR &= q_3 - q_1 \\ &= 75\% - 25\% \\ &= 50\% \end{aligned}$$



## • CHI SQUARE TEST

The chi square test for goodness of fit claims about population proportions [categorical variables]

- It is a non parametric test [i.e will not use mean, median or mode] that is performed on categorical data.  
[nominal, ordinal]

Eg: Population of Male who likes different color of bikes

	Theory	Sample
Yellow Bike	1/3	22
orange Bike	1/3	17
Red Bike	1/3 ↓	59
talks more about proportions of distribution.	Theory categorical distribution [percentage]	observed categorical Distribution [count]

- Q. In a student class of 100 students, 30 are right handed. Does this class fit the theory 12% of people are right handed

	Observed	Expected	Theory
Right handed	30	12	
Left handed	70	88	
Observed			

→ Chi-Square : For Goodness of Fit

Q In 2010 census of the city, the weight of individuals in a small city were found to be following.

$< 50 \text{ kg}$	$50 - 75$	$> 75$	Expected
20%	30%	50%	

In 2020, weight of  $n=500$  individuals were sampled. Below are the results.

Confidence Interval	$< 50$	$50 - 75$	$> 75$	Observed
	95%	140	160	200

Using  $\alpha = 0.05$  would you conclude the population differences of weights has changed in last 10 years?

Ans.

Expected	$< 50$	$50 - 75$	$> 75$	
	$20 \times 500$	$30 \times 500$	$50 \times 500$	
	100	150	250	

- ① Null hypothesis:  $H_0$ : The data meets the expectation  
 ② Alternative hypothesis:  $H_1$ : The data does not meet the expectation

③  $\alpha = 0.05$ , CI = 95%.

④ Degree of freedom

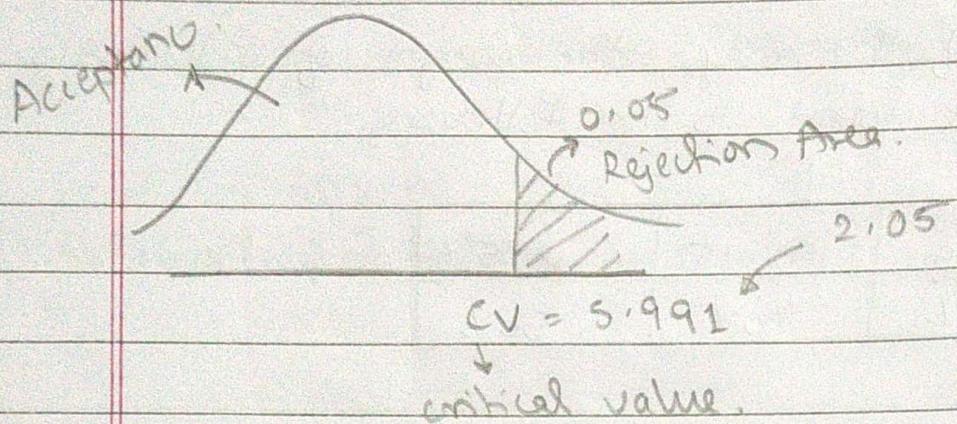
$$df = k-1 = 3-1 = 2$$

Date \_\_\_\_\_  
Page \_\_\_\_\_

Chi-square  
is always 1-tail table

## ⑤ Decision Boundary → CHI SQUARE Table

- Right skewed



If  $\chi^2$  is greater than 5.991, Reject  $H_0$   
else

we Fail to reject  $H_0$ .

## ⑥ Calculate Chi-Square Test Statistics

$$\chi^2 = \sum_{i=1}^k \frac{(O - E)^2}{E}$$

O: Observed  
E: Expected

$$= \frac{(140-100)^2}{100} + \frac{(160-150)^2}{150} + \frac{(200-250)^2}{250}$$

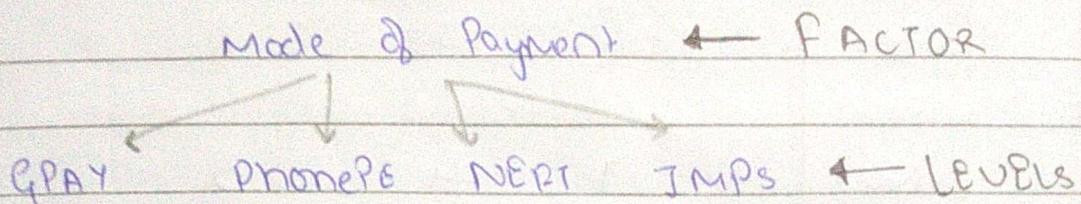
$$= 16 + 0.66 + 10$$

$$\boxed{\chi^2 = 26.66}$$

$26.66 > 5.99$ , Reject  $H_0$

## • ANOVA [Analysis of Variance]

It is a statistical method used to compare the mean of 2 or more groups



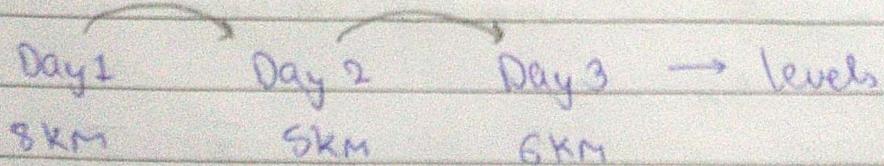
### → Types of Anova

1) one way Anova : one factor with at least 2 levels, these levels are independent

Eg: factor is brand of soda, & you collect data on Coke, Pepsi, sprite to find out if there is difference in price per 100 ml

2) Repeated Measures Anova : one factor with at least 2 levels, levels are dependent

Running → Factor



3) factorial Anova : Two or more factors [each of which with at least 2 levels], levels can be either dependent or independent

Running  $\leftarrow$  Factor

		Day 1	Day 2	Day 3	$\leftarrow$ dependent
		Male	8	5	6
Gender Factor	Male	7	4	3	
	Female				

### → Hypothesis Testing In Anova

- Null hypothesis  $H_0 : \mu_1 = \mu_2 = \mu_3 = \dots = \mu_K$
- Alternate hypothesis  $H_1 : \text{Atleast one of mean is not equal}$

### → F Test Statistics

$$F = \frac{\text{Variation between Samples}}{\text{Variation within Samples}}$$

$x_1 \quad x_2 \quad x_3$  Variance between sample.

1	6	5
2	7	6
4	3	3
5	2	2
3	1	4

$$H_0 : \bar{X}_1 = \bar{X}_2 = \bar{X}_3$$

$H_1 : \text{Atleast one sample mean is not equal.}$

$$\bar{X}_1 = 3 \quad \bar{X}_2 = 10/5 \quad \bar{X}_3 = 4$$

Q Doctors want to test a new medication which reduces headache. They split the participant into 3 condition [15mg, 30mg, 45mg]. Later on doctor ask the patient to rate the headache between [1-10]. Are there any differences between the 3 conditions using alpha = 0.05?

	15mg	30mg	45mg
N: Total number	9	7	4
a: No. of samples	8	6	3
	7	6	2
a: No. of categories	8	7	3
	8	8	4
n: number of samples in 1 column	9	7	3
n: number of samples in 1 column	8	6	2

① Null hypothesis  $H_0: \mu_{15} = \mu_{30} = \mu_{45}$

② Alternate hypothesis  $H_1:$  Atleast one mean is not equal.

③ State significance value

$$\alpha = 0.05 \Rightarrow C.I = 0.95$$

④ Calculate degree of freedom (df)

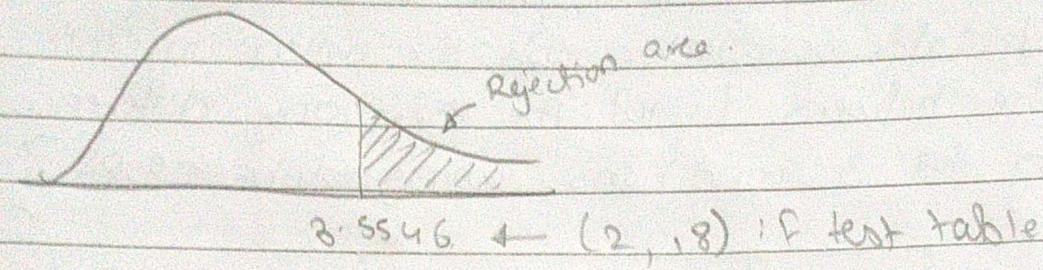
$$N = 21, a = 3, n = 7$$

$$df \text{ between} = a - 1 = 3 - 1 = 2$$

(2, 18)  $\Rightarrow$  F test table

$$df \text{ within} = N - a = 21 - 3 = 18$$

Decision Boundary  $\alpha=0.05$



### ⑤ Calculate F-test Statistics

	SS [sum of squares]	df	MS [Mean Square]	F
Between	98.67	2	49.34	
within	10.29	18	0.54	86.56
Total.	108.96	20		

$$\textcircled{1} \quad SS_{\text{Between}} = \frac{\sum (\sum a_i)^2}{n} - \frac{T^2}{N}$$

$$* 98.67/2 \\ * 10.29/18$$

$$15Mg = 9+8+7+8+8+9+8 = 57$$

$$30Mg = 7+6+6+7+8+7+6 = 47$$

$$45Mg = 4+3+2+3+4+3+2 = 21$$

$$\begin{aligned} \sum y^2 &= 9^2 + 8^2 + 7^2 + 8^2 + 8^2 + 9^2 + 8^2 + 7^2 \\ &\quad + 6^2 + 6^2 + 7^2 + \dots \\ &= 853 \end{aligned}$$

$$= \frac{57^2 + 47^2 + 21^2}{7} - \frac{[57 + 47 + 21]}{21}$$

$$= 98.67$$

$$\textcircled{2} \quad SS_{\text{within}} = \sum y^2 - \frac{\sum (\sum a_i)^2}{n}$$

$$= 853 - \left[ \frac{57^2 + 47^2 + 21^2}{7} \right]$$

$$= 10.29$$

$$F = \frac{\text{variation between samples}}{\text{variation within samples}} = \frac{MS_{\text{Between}}}{MS_{\text{within}}}$$

$$= \frac{49.34}{0.54} = 86.56$$

$86.56 > 3.556$ , we reject null hypothesis