

Assignment - 01

Part I → Descriptive analysis

(1) Mean → 30, 40, 45, 50, 200

$$\text{Sgl} \quad \frac{30+40+45+50+200}{5} \rightarrow \frac{365}{5} \rightarrow 73$$

2) Median → 30, 40, 45, 50, 200

→ 45 → median

(3) Mode → 30, 40, 45, 50, 200

Add the salary 40 to the above list
30, 40, 40, 45, 50, 200

→ Mode = 40

(4) interpretation →

~~Advantages~~

Mean → Advantage:- * Take into account all data
* Useful for further calculation

Limitation → * Sensitive to outliers

* May not accurately represent data in skewed distribution

Mode → Advantage → * Simple to understand and calculate

* Not affected by outliers

Limitation \rightarrow

- * Can be multiple mode in dataset
- * less informative for continuous data

Median \rightarrow

Advantage \rightarrow

- * Not affected by outliers
- * Represent the middle value, providing better central tendency

Limitation \rightarrow

- * Does not consider the magnitude of all data points

- * less useful in dataset with a small number of observations

Conclusion \rightarrow I would recommend using median to describe central tendency of the data as it is not affected by outliers and providing better central tendency.

(2) Standard deviation and variance

(1) Scores from a test \Rightarrow 5, 10, 10, 20, 30

$$\text{Mean} = \frac{5+10+10+20+30}{5} = \frac{75}{5} = 15$$

$$\text{variance} = \frac{(15-5)^2 + (15-10)^2 + (15-10)^2 + (15-20)^2 + (15-20)^2}{5}$$

$$\text{variance} = \frac{(10-15)^2 + (5-15)^2 + (5-15)^2 + (15-15)^2 + (20-15)^2}{5} \Rightarrow \frac{400}{5} \rightarrow 80$$

$$S.D = \sqrt{80} \rightarrow 8.94$$

Interpretation

Standard deviation is closed to the mean if mean data is constaince and low spread.

② Box plot outliers

The ages of a group of people $\geq 10, 12, 13, 15$
 $16, 18, 25, 35, 80$

$$Q_2 \Rightarrow 10, 12, 13, 15, 16, 18, 25, 35, 80$$

$$16 \leftarrow Q_2$$

$$Q_1 \Rightarrow 10, 12, 13, 15 \Rightarrow \frac{12+13}{2} \Rightarrow 12.5$$

$$Q_3 \Rightarrow 18, 25, 35, 80 \Rightarrow \frac{25+35}{2} \Rightarrow \frac{60}{2} \rightarrow 30$$

$$IQR = Q_3 - Q_1 \\ = 30 - 12.5 \Rightarrow 17.5$$

Ciii) IQR Rule

$$\text{lower bound} = Q_1 - 1.5 \times \text{IQR}$$

$$= 12.5 - 1.5 (17.5)$$

$$= 12.5 - 26.25 = -13.75$$

$$\text{Upper bound} = Q_3 + 1.5 \times \text{IQR}$$

$$= 30 + (1.5 \times 17.5)$$

$$= 30 + 26.25 = 56.25$$

Now any values < -13.75 and > 56.25 are outliers

Here 80 is an outlier

Part - 2

Inferential Statistics

Simple Linear Regression

Hours Studied \Rightarrow 1, 2, 3, 4, 5

Exam Scores \Rightarrow 50, 55, 65, 70, 75

(i) Equation of Simple linear regression

$$y = B_0 + B_1 x + \epsilon$$

where $B_0 \rightarrow B_0$ is the intercept

$B_1 \rightarrow B_1$ is the slope

$\epsilon \rightarrow$ epsilon is the error term

$$B_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

\Rightarrow

hours studied	Exam scores	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	50	(3-1)(63-50)	4
2	55	(3-2)(63-55)	1
3	65	(3-3)(63-65)	0
4	70	(3-4)(63-70)	1
5	75	(3-5)(63-75)	4
		$\sum (x_i - \bar{x})(y_i - \bar{y}) = 10$	$\sum (x_i - \bar{x})^2 = 10$

$$\bar{x} = \frac{1+2+3+4+5}{5} = \frac{15}{5} \Rightarrow 3$$

$$\bar{y} = \frac{50+55+65+70+75}{5} = \frac{325}{5} = 65$$

$$B_1 = \bar{x}(\bar{y} - \bar{y})$$

$$\approx (\bar{x}, \bar{y})^2$$

$$6.5 \\ 7.5$$

$$6.5$$

$$B_0 = \bar{y} - B_1 \cdot \bar{x}$$

$$= 63 - 6.5 \times 3 \\ = 63 - 19.5 = 43.5$$

Interpretation of Slope (6.5)

The slope of 6.5 means that for every additional hour studied, the exam score increases by 6.5 points on average.

This shows a positive linear regression between hours studied and exam performance.

(5). Hypothesis Testing T-Test

Given $\hat{\rightarrow}$ Avg. completion rate = 75%

company claim that a new interactive feature has increased the average completion rate = 80%

you collect a sample = 50 users

Avg. completion rate = 78%.

std. deviation = 5%.

- Q. \rightarrow
- $H_0: \bar{x} = 75\%$] \rightarrow hypothesis
 - $H_1: \bar{x} \neq 75\%$

(2.) Test statistic:

$$t = \frac{\bar{x} - u}{s / \sqrt{n}}$$

t = test statistics

\bar{x} = Average sample mean

u = hypothesized mean

s = std. deviation

n \rightarrow Sample size

$$t = \frac{78\% - 75\%}{5\% / \sqrt{50}}$$

$$= \frac{+3\%}{5\% / \sqrt{50}} \rightarrow \frac{+3\%}{5\% / \sqrt{2}}$$

$$= \frac{+3\% \times \sqrt{50}}{5\%}$$

$$= \frac{+3 \times 5 \times 1.414}{5} = +4.242$$

Sample size = 50

degrees of freedom (df) = $n-1 = 49$

Significance level (α) = 0.05

Test is one tailed

Critical value ≈ 1.677

$4.24 > 1.677 \rightarrow$ Null hypothesis rejected

⑥ Chi square test

	Like	Dislike	Row total
Male	40	60	100
Female	50	50	100
Total	90	110	200

H_0 = There is no relation with Gender and preferences

H_1 = There is relations with Gender and preferences

Formula For Expected Frequencies

$$= \frac{(\text{Row Total}) \times (\text{Column Total})}{\text{Grand total}}$$

$$\text{Male-Like} = \frac{(100 \times 90)}{200} = 45$$

$$\text{male dislike} = \frac{100 \times 10}{200} \rightarrow 5$$

$$\text{female-Like} = \frac{100 \times 90}{200} \rightarrow 45$$

$$\text{female-Dislike} = \frac{100 \times 10}{200} \rightarrow 5$$

	Like	Dislike	Total
Male	45	55	100
Female	45	55	100
	90	110	200

$$\chi^2 = \frac{\sum (O - E)^2}{E}$$

O = observed frequency

E = expected frequency

	Observed freq.	Expected freq.	$(O-E)^2/E$
male (Like)	40	45	0.556
male (Dislike)	60	55	0.455
female (Like)	50	45	0.556
female (Dislike)	50	55	0.455
			$\chi^2 = 2.02$

df → degree of freedom

$\text{df} = (\text{rows}-1) \times (\text{columns}-1)$

$$= (2-1)(2-1)$$

$$1 \times 1 \Rightarrow 1$$

$\alpha = 0.05$ for $\text{df} = 1$ is 3.841

Since $x^2 = 2.022 < 3.841$ hence we fail to reject the null hypothesis

ANOVA

→ Three different teaching methods result in different exam scores:

Group A \Rightarrow 70, 75, 80

Group B \Rightarrow 60, 65, 70

Group C \Rightarrow 85, 90, 95

Anova \Leftrightarrow Anova is a statistical method used to compare the means of three or more groups.

One-way - Anova

Compare the means of three or more independent groups based on one factor.

Hypotheses:- $H_0 \rightarrow$ All group means are equal
 $H_1 \rightarrow$ At least one group mean is different

Step-1

Calculate Group mean

$$\text{Group A mean} \rightarrow \frac{70+75+80}{3} = \frac{225}{3} \Rightarrow 75$$

$$\text{Group B mean} = \frac{60+65+70}{3} = \frac{195}{3} \Rightarrow 65$$

$$\text{Group C mean} = \frac{85+90+95}{3} = \frac{270}{3} \Rightarrow 90$$

Step-2

Calculate Grand mean

$$GM = \frac{70+75+80+60+65+70+85+90+95}{9}$$

$$= \frac{690}{9} \rightarrow \cancel{78.89} \cancel{76.67} \Rightarrow 76.67$$

Step-3 Sum of squares between (SSB)

$$SSB = n \sum (\bar{x}_i - GM)^2$$

$$= 3[(75-76.67)^2 + (65-76.67)^2 + (90-76.67)^2]$$

$$= 3 [2.79 + 136.19 + 172.69]$$

$$= (316.67) 3$$

$$= 950.0$$

Step-4

Sum of squares within (SSW)

$$SSW = \sum (x_{ij} - \bar{x}_j)^2$$

Group A \rightarrow mean = 75

$$(70-75)^2 + (75-75)^2 + (80-75)^2 \\ = 25 + 0 + 25 = 50$$

Group B

$$(60-65)^2 + (65-65)^2 + (70-65)^2 \\ 25+0+25 = 50$$

Group C

$$(85-90)^2 + (90-90)^2 + (95-90)^2 \\ 25+0+25 = 50$$

Group A + Group B + Group C = 150

Step-4 Calculate $\rightarrow F$ -ratio

$$df_{between} = K-1$$

$$df_{within} = N-K$$

K \rightarrow Total no. of groups

n \rightarrow Total no. of observations

$$df_{between} = K-1 = 3-1 = 2$$

$$df_{within} = N-K = 9-3 = 6$$

Mean Square between (MSB)

$$SSB = 950.01 \div 4 = 247.5005$$

df between

Mean Square within (MSW)

$$= SSW = 150 \div 25$$

df within 25

19.0002

$$\text{F ratio} \Rightarrow \frac{MSB}{MSW} = \frac{247.5005 \times 1}{19.0002 \div 1000} \\ = 19.003$$

Significance level of F at $\alpha = 0.05$

df₁ = 2, df₂ = 6

Significance F ≈ 5.14

$$F = 19.0002 > 5.14$$

we reject null hypothesis

Part - 3 Parametric vs. Non parametric Test

(i) \rightarrow Parametric and non-parametric tests are two categories of statistical tests. Parametric test rely on assumptions about the distribution of population data, typically requiring a normal distribution while non-parametric test does not require any assumption.

Difference between parametric and non-parametric tests

Parametric tests

These tests assume that the data comes from a population with a specific distribution, often a normal distribution. They also assumed data is measured.

- * They also assumed data is an interval or ratio scale.
- * The variance of the group being compared are roughly equal.

Non-parametric tests

These tests do not make assumption about the distribution of the population data.

- * They are suitable for data measured on ordinal or nominal scales.
- * They can handle outliers and non-normal distribution.

(ii) Examples

parametric test

- * T-test \Rightarrow Comparing the average height of student in two different classes.
- * Anova \Rightarrow Comparing the Average test scores of student across multiple teaching methods.

Non-parametric test

Mann - Whitney U test :-

Comparing the satisfaction levels (ranked on a scale) of patients receiving two different treatment

Krusal - Wallis test :-

Comparing the preferences ranking of different brands of coffee.

Non-normal data

If the data is not normally distributed, non-parametric tests are generally preferred. This is because parametric tests are sensitive to violations of the normality assumption and may produce inaccurate results.

Part - 4 . Test

(9) Z-test vs. T-test

① What is the difference between a Z-test and a T-test

Z-tests are statistical calculations that can be used to compare population means to samples.

T-test are calculation used to test a hypothesis but they are most useful when we need to determine if there is a statistically significant difference between two independent sample groups.

Bases	Z-test	T-test
Population Std. deviation	Known	Unknown
Sample size	Generally used for larger samples.	Generally used for smaller samples.
Distribution	Assume a normal distribution of the population	Assume a normal distribution of the population, but also accounts for the uncertainty in estimating the population std. deviation
Test-Statistics	Z-statistics	T-statistics

checklist

Page No. _____
Date _____

(c) P-value interpretation

A p-value is the probability of obtaining a test statistic at least as extreme as the one observed, assuming the null hypothesis is true. It quantifies the evidence against the null hypothesis.

A small p-value suggests that the observed data is unlikely under the null hypothesis, indicating stronger evidence against the null hypothesis.

A large p-value suggests the data is consistent with the null-hypothesis, providing less evidence against it.

(d) Interpreting a p-value of 0.03 at $\alpha = 0.05$

If you obtain a p-value of 0.03 and are using a significance level (α) of 0.05.

since $0.03 < 0.05$ you reject the null hypothesis.

That means there is statistically significant evidence at the 5% level to suggest that the observed effect is not due to random chance.