

Statistic Assignment

Q. 2. Explain the different types of data
 (Qualitative and Quantitative) and provide examples of each. Discuss nominal, ordinal, interval, and ratio scales.

⇒ Qualitative data (categorical data)

Qualitative data refers to non-numerical information that describes characteristics or qualities. In other words, qualitative data is categorical data.

→ Qualitative data is often categorical, categorized into different groups or labels, which can be analyzed for patterns and insights.

Key features of Qualitative data:

1. Non-numerical nature: Qualitative data consists of description or attributes, such as color, type or preference. For example, in a survey colors favorite colors, responses like "red", "blue", or "green" are qualitative.

2. Categorical: This data is usually grouped into categories (hence it's also known as "categorical data"). These categories can either be nominal (no inherent order, like gender or nationality) or

ordinal (where there's a natural order, like rating something as "poor", "good", or "excellent").

Example

① Nominal and Ordinal

② Nominal

No order in categorical data: nationality, blood group, gender, color, etc.

③ Ordinal → order in data

Rank, good, better, best, grades, salary, high, medium, and low.

* Statistical Analysis

frequency counts

Percentages

Chi-square test

* Visual Representation

Bar charts

Pie charts

④ Quantitative data (numerical data)

Quantitative data is data that can be counted or measured and is represented as numbers.

→ data about numeric variables (e.g. how many, how much or how often)

Qualitative = Quality

allowing for various mathematical operations like (+, -, ×, ÷). Quantitative data provides concrete, measurable information that can be analyzed to identify trends, make predictions or evaluate relationships between variables.

* Types of Quantitative data

Discrete data:

This type of data consists of countable values, often integers. For example, the number of students in a classroom or the number of cars in parking lot.

- Discrete data cannot be subdivided into finer increments within the range (e.g. You can't have half a student).
- No of children.
- No of mobile phones.
- No of bank accounts.
- total numbers of employees in firm.
- Bathrooms, bedrooms, car numbers.

* Continuous: Continuous data can take any value within a given range and can be subdivided.

- Real numbers.
- Height, weight, speed, temperatures, gas, movie duration, length, distance, salary, price, etc.

Example:- Discrete Cens (Continues).

Analysis of Quantitative Data

1. Descriptive Statistics

• measures of Central Tendencies

- Mean

- Median

- mode

2. Measures of Dispersion

- Range

- Variance

- Standard Deviation

3. Inferential Statistics: Techniques used to make predictions or inferences about a population based on a sample data.

- Hypothesis Testing

- Regression Analysis

- Correlation

4. Interval Data (Measurement Scale)

Interval data is a type of quantitative data where the difference between values are meaningful, but there is no true zero point.

→ this type of data allows for the measurement of the difference between numbers, but ratio may not be meaningful.

Exam Pic:

Temperature in Celsius or Farenheit

Dates on a calendar

Time on a clock

* Ratio scale

The ratio scale is the highest level of measurement scale for quantitative data providing the most information and allowing for a wide range of mathematical operations. A ratio scale has all the properties of an interval scale but with an additional feature; it is true zero pointing to nothing. Zero represents the complete absence of the quantity being measured.

1. True zero point: A ratio scale has an absolute zero which means that zero indicates a total absence of the quantity being measured.

for example: In the context of weight, zero grams means no weight at all.

2. Equal intervals: like an interval scale, the difference between values on a ratio scale is meaningful and consistent.

→ The diff. between 10 and 20 kg is same as between 20 and 30 kg.

3. Ratios are meaningful: there is a true zero point) you can meaningfully compare values in terms of ratios.
For Example: If one object weight 20 kg and another object weighs 10 kg, you can say that the first object is twice as heavy as the second.

Example:

Height, weight, Time, Distance, income.

Q. 2 What are the measures of central tendency and when should you use each?
Discuss the mean, median, and mode with examples and situations where each is appropriate.

⇒ Measures of central tendency are statistical tools used to identify the central value within a dataset.
These measures summarize a dataset with a single number that represents the center of the data distribution.
The three main measures of central tendency are the mean, median, and mode.

1. mean (Arithmetic Average)

The mean is the sum of all data values divided by the number of values in the dataset. It is the most common measure of central tendency.

$$\text{mean} = \frac{\sum x_i}{N}$$

* When to use the mean:

1) Symmetric data

2) Continuous and discrete data

• Mathematical Analysis

Example:

If the test scores of five students are 70, 75, 80, 85 and 90, the mean is:

$$\text{Mean} = \frac{70 + 75 + 80 + 85 + 90}{5} = 80$$

2). Median (middle value)

The median is the middle value in a set of data when all the values are arranged in ascending or descending order. If there is an even number of values, the median is the average of the two middle numbers and number is.

odd then most middle element is the median.

* When to use median

• Skewed data

• Ordinal data

• Non-normal Distributions

Example

4, 5, 2, 3, 1, 2, 1

Sort

2, 1, 2, 1, 3, 1, 4, 5

\Rightarrow Count & the number of Element = 8
 if Count is even

1 2 3 4 5 6 7 8

median = average of first middle most element

2

$$= \frac{2+3}{2} = 2.5$$

4, 1, 5, 1, 2, 1, 3, 1, 2

Sort: 2, 1, 2, 1, 3, 1, 4, 1, 5

Count = 8

↓ odd

median = the middle most element

= 3

* Mode (most frequent value)

Mode The mode is the value that appears most frequently in the dataset. A dataset may have no mode, one mode, or multiple modes (bimodal, multimodal).

 \rightarrow When to use

- Categorical data
- Multimodal distributions
- Non-numerical data

Example: 2, 2, 1, 3, 1, 2, 4, 1, 4, 1, 3, 4, 2, 1

mode = 4

Q. 3. Explain the concept of dispersion. How do variance and standard deviation measure the spread of data?

→ The measure of dispersion help to interpret the variability of data. i.e. to know how much homogeneous or heterogeneous the data is in simple terms, it shows how squeezed or scattered the variable is.

Common measures of dispersion

- Range
- Interquartile Range (IQR)
- Variance
- Standard Deviation.

Variance

Variance measures the average squared deviation from the mean. It gives us an idea of how much the data points are spread out from the mean in a dataset. It quantifies the overall variability of the data.

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

Where,

- σ^2 is the variance
- x_i is each individual data point
- \bar{x} is the mean of the data set
- n is the number of data points

Interpretation

- Large Variance: Indicates that the data points are widely spread out around the mean.
- Small Variance suggests that the data points are closely clustered around the mean.

Standard Deviation

The standard deviation is the square root of the variance providing of a measure of dispersion in the same units as the data.

$$\sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}}$$

where:

- σ is the standard deviation.
- x_i is each data point.
- \bar{x} is the mean.
- N is the number of data points.

→ Both Variance and Standard deviation quantify how much individual data points deviate from the mean.

1. Variance: gives the average squared deviation from the mean which helps in understanding the overall dispersion but can be less intuitive due to its squared units.

2. Standard Deviation directly shows the average distance of data points

Date _____
Page _____

from the mean in the same units as the data, making it more interpretable for practical applications.

Examples

Dataset

Dataset 1: 5, 6, 7, 8, 9

Dataset 2: 1, 5, 7, 9, 13

Both datasets have the same mean 7.

Variance of Dataset 1 is small because the values are close to mean.

Variance of Dataset 2 is larger because the values are more spread out from the mean.

→ The standard deviation could reflect these differences, with Dataset 2 having a larger standard deviation than Dataset 1.

Q. 4. what is a box plot and what can it tell you about the distribution of data?

⇒ A box plot (also called a box-and-whisker plot) is a graphical representation of a dataset that visually displays its distribution, spread and variability. It is particularly

useful for identifying key features of the data such as the median, quartiles, range, and potential outliers. Box plots provide a summary of the distribution without making assumptions about the underlying statistical distribution of the data.

key components of a Box Plot

1. median (Q2): The line inside the box represents the median (the middle value of the dataset when ordered).
2. Box: The box represents the interquartile range (IQR), which is the range between the First Quartile (Q_1) and the Third Quartile (Q_3). The edges of the box are drawn at Q_1 (25th percentile) and Q_3 (75th percentile).
3. Whiskers: The "whiskers" extend from the edges of the box to the smallest and largest data points within a certain range. Usually, they are drawn up to 1.5 times the IQR from the quartiles.
4. Outliers: These are individual data points that fall outside the whiskers.

Breakdown of Components:

- Minimum: The smallest data point within 1.5 times the IQR of Q_1 .
- First Quartile (Q_1): The 25th Percentile, where 25% of the data lies below this value.
- Median/Median (Q_2): The 50th Percentile, dividing the dataset in half.
- Third Quartile (Q_3): The 75th Percentile, where 75% of the data lies below this value.
- Maximum: The largest data point within 1.5 times the IQR of Q_3 .
- Outliers: Data points that lie beyond the whiskers, indicating unusual or extreme values in the dataset.

Example of a box plot:
dataset of exam scores:

45, 50, 53, 55, 58, 60, 63, 65, 70, 75,
80, 85, 90, 95, 100

$$Q_1 = \frac{n+1}{4}^{\text{th}} = \frac{15+1}{4}^{\text{th}} = 4^{\text{th}} \text{ no.} = Q_1 = 55$$

$$Q_3 = \frac{3(n+1)}{4}^{\text{th}} = \frac{3(15+1)}{4}^{\text{th}} = 12^{\text{th}} \text{ no.} = Q_3 = 85$$

$$M_d = \frac{(n+1)}{2}^{\text{th}} = \frac{15+1}{2}^{\text{th}} = 8^{\text{th}} = 65$$

- Whiskers: Extend from Q1 to Q3 values.
- No obvious outliers in this case.

* What a box plot tells you about the data

1. Center (Median): The line inside the box shows the median which is a measure of central tendency. It gives an idea of where the middle of the data lies.
2. Spread (IQR): The width of the box represents the interquartile range (IQR), which indicates the spread at the middle 50% of the data.
3. Symmetry / Skewness:
 - If the median is roughly centered inside the box, the data is symmetrically distributed.
 - Right-skewed (positively skewed): The median is closer to Q3, and the right whisker is longer, indicating a longer tail on the right.
 - Left-skewed (negatively skewed): The median is closer to Q1, and the left whisker is longer, indicating a longer tail on the left.
4. Range: The distance between the whiskers gives an idea of the

Date _____
Page _____

total range of the data. A large range indicates more variability in the dataset.

5. Outliers: Outliers are shown as individual points - the whiskers.

Q. 5. Discuss the role of random sampling in making inferences about populations.

⇒ Random Sampling is fundamental to making accurate inferences about populations in statistics.

Creating a Representative Sample

Random Sampling ensures that every individual in a population has an equal chance of being selected. This results in a sample that accurately reflects the population's characteristics.

2. Facilitating Generalization

The primary goal of random sampling is to allow researchers to generalize their findings from the sample to the entire population.

3. Minimizing Bias

Random Sampling reduces the risk of Selection Bias which occurs when the sample does not accurately represent the population.

U. Supporting Statistical Validity

Most statistical methods, such as confidence intervals and hypothesis testing, assume that data is collected through random sampling.

5. Error Estimation

With random sampling, it's safe to estimate sampling error - the difference between the sample statistic and the true population parameter.

6. Cost-Efficiency

Random sampling allows researchers to collect data from a subset of the population rather than the entire population, which can be expensive and time-consuming.

→ Random sampling is crucial for obtaining a representative sample, minimizing bias, and ensuring the statistical validity of inferences.

It allows researchers to generalize their results to the entire population while accounting for errors and ensuring accuracy.

Q. Explain the Skewness and its types. How does skewness affect the interpretation of data?

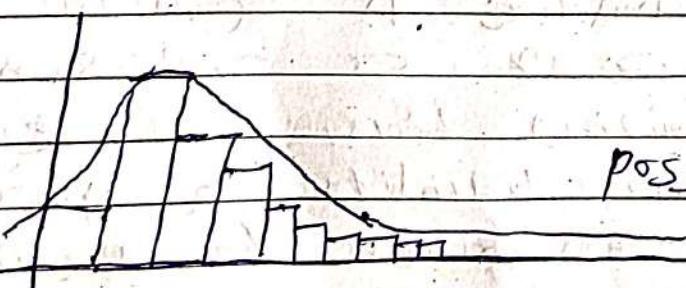
→ Skewness is a statistical measure that describes the asymmetry of a distribution of values around its mean. It indicates the direction and extent to which the values deviate from a normal distribution. Skewness can be positive, negative or zero, depending on the shape of the distribution.

TYPES OF SKEWNESS

1. Positive skewness (right-skewed)

In a positively skewed distribution, the tail on the right side (higher values) is longer or fatter than the left side.

The mean is usually greater than the median, and most of the data values cluster toward the lower end.



positive skewness distribution

2. Negative skewness (left-skewed)

In a negatively skewed distribution, the tail on the left side (lower values) is longer or fatter than the right side.

The mean is usually less than the median, and most of the data values cluster toward the higher end.

cluster toward the higher end.

3. Zero Skewness (Symmetric)

- In a Symmetric distribution, the data is evenly distributed around the mean, with no skewness.
- the mean, median and mode are all equal.

Skewness interpretation and effects on data

1. mean, median, and mode
 - in positively skewed data the mean is greater than the median and the median is greater than the mode ($\text{Mean} > \text{Median} > \text{Mode}$).
 - in negatively skewed data the mean is less than the median and the median is less than the mode ($\text{Mean} < \text{Median} < \text{Mode}$).

2 outliers

Skewed distributions are often driven by outliers. A right-skewed distribution may be influenced by unusually high values, while a left-skewed distribution may contain unusually low values. These outliers can distort the interpretation of summary statistics, particularly the mean.

4. Data Transformation
When data is highly skewed transformations like logarithmic or square root transformation can reduce skewness and make the data more symmetrical.

F. What is the interquartile range (IQR), and how is it used to detect outliers?
The interquartile range (IQR) is a measure of statistical dispersion that describes the spread of the middle 50% of a dataset. It is the difference between the third quartile (Q_3) and the first quartile (Q_1), which are the 75th and 25th percentiles of the data, respectively.

$$IQR = Q_3 - Q_1$$

- Q_1 (First quartile): The value below which 25% of the data fall.
- Q_3 (Third quartile): The value below which 75% of the data fall.

Purpose of IQR

IQR is used to measure the spread of the middle portion of data, filtering out extreme values which makes it less sensitive to outliers compared to other measures like range or standard deviation.

Date _____
Page _____

How IQR is used to Detect outliers
 Outliers are data points that are significantly different from the rest of the data. The IQR is often used to identify potential outliers through the following

Calculate the IQR

- Subtract the first quartile from the third quartile to find the IQR.
- A common rule of thumb is to consider any data points that fall below $Q_1 - 1.5 \times IQR$ or above $Q_3 + 1.5 \times IQR$
- Lower fence: $Q_1 - 1.5 \times IQR$
- Upper fence: $Q_3 + 1.5 \times IQR$
- Any data points below the lower fence or above the upper fence are considered potential outliers.

Example

10, 12, 13, 15, 16, 18, 19, 21, 24, 29,
U.S

n = 11

number is odd

then $Q_1 = \frac{n+1}{4} th = \frac{11+1}{4} th = \frac{12}{4} = 3rd m$

$$Q_1 = 13$$

$$Q_3 = \frac{3(n+1)}{4} th = \frac{3(11+1)}{4} th = \frac{36}{4} th = 9th \\ \text{no. } Q_3 = 24$$

Date _____
Page _____

$$IQR = Q_3 - Q_1 = 24 - 13 = 11$$

Lower Fence: $Q_1 - 1.5 \times IQR = 13 - 1.5 \times 11 = -3.5$

Upper Fence: $Q_3 + 1.5 \times IQR = 24 + 1.5 \times 11 = 40.5$

"Any value" below -3.5 or above 40.5 are considered outliers in this case. 45 is an outlier because it is greater than 40.5 .

Q. 8. Discuss the conditions under which the binomial distribution is used.

\Rightarrow The binomial distribution is used to model the number of successes in a fixed number of independent trials of a binary (yes/no or success/failure) experiment. For the binomial distribution to apply, certain conditions must be met.

Conditions for using the binomial distribution

1. Fixed number of trials
 - The experiment or process must be repeated a fixed number of times. Each repetition is called a trial.
 - For example, flipping a coin 10 times or rolling a die 15 times.
2. Two possible outcomes per trial
 - Each trial must result in

Has exactly two possible outcomes, commonly referred to as "success" and "failure".

- Example: flipping a coin results in either heads or tails, or a medical test is either positive or negative.

3- independence of trials.

- The outcome of each trial must be independent of the outcomes of other trials. This means that the result of one trial shall not affect the result of another.

4. Constant probability of success (P)

- The probability of success, denoted as P , must remain the same for each trial. Likewise, the probability of failure, $1-P$ must also remain constant.

Binomial Distribution formula

The

$$P(X=k) = \binom{n}{k} P^k (1-P)^{n-k}$$

- $P(X=k)$ is the probability of k successes.

- $\binom{n}{k}$ is the binomial coefficient $\frac{n!}{k!(n-k)!}$, representing the number of ways to choose k successes from

n trials.

p^n represents the probability of k successes.

$(1-p)^{n-k}$ represents the probability of $n-k$ failures.

Examples of Binomial distribution in use

1. Coin tossing

If a fair coin is tossed 5 times, the binomial distribution can be used to calculate the probability of getting exactly 3 heads out of 5 trials.

2. Quality Control

In manufacturing process, if the probability of a defective product is 0.02, and 300 products are tested, the binomial distribution can be used to determine the probability of finding exactly 2 defective products.

3. medical testing

If a diagnostic test has a 90% probability of accurately detecting a disease and 20 patients are tested, the binomial distribution can be used to find the probability of correctly diagnosis disease exactly 18 patients.

Q. 9. Explain the Properties of the Normal distribution and the Empirical rule (68-95-99.7 rule).

The Normal distribution, also known as the Gaussian distribution, is a continuous probability distribution characterized by a symmetrical bell-shaped curve. It is one of the most important distributions in statistics due to its wide applicability in real-world scenarios, such as heights, test scores, measurements errors, etc.

Properties of the Normal Distribution:

1. Symmetry:

The normal distribution is symmetric around its mean (μ). This means that the left and right sides of the curve are mirror images.

The mean, median and mode of the distribution are equal and located at the center of the curve.

2. Bell-shaped Curve:

The distribution has a characteristic bell-shaped curve that peaks at the mean and tapers off symmetrically in both directions.

Most of the data points are clustered around the mean, and the probability decreases as you move further from

From the mean.

3. Mean and Standard Deviation

Mean (μ): Determines the location or the center of the distribution.

Standard deviation (σ): Determines the spread or width of the distribution. A larger standard deviation means the data is more spread out, while a smaller one means the data is closely clustered around the mean.

4. Asymptotic:

The tails of the normal distribution extend indefinitely approaching but never touching the horizontal axis.

5. Total area under the curve.

The total area under the normal distribution curve equals 1, representing the entire probability space (or 100% of the probability).

6. Empirical Rule (68-95-99.7 Rule)

The Empirical Rule, also known as the 68-95-99.7 rule, applies to normally distributed data.

The Empirical Rule (68-95-99.7 Rule)

→ The Empirical rule provides insight into how data is distributed around the mean in a normal distribution. It helps to quickly estimate probabilities and understand the spread of the data.

1. 68% of the data lies within one standard deviation ($\mu \pm \sigma$) of the mean:
 - approximately 68% of the observations fall within one standard deviation (between $\mu - \sigma$ and $\mu + \sigma$).
 - this means that if the mean of dataset is 100, and the standard deviation is 20, then about 68% of the data will fall between 80 and 120.
2. 95% of the data lies within two standard deviations ($\mu \pm 2\sigma$) of the mean:
 - About 95% of the data falls within two standard deviations between $\mu - 2\sigma$ and $\mu + 2\sigma$.
 - Using the same example, about 95% of the data will fall between 80 and 120.
3. 99.7% of the data within three standard deviations ($\mu \pm 3\sigma$) of the mean:
 - Nearly all the data (99.7%) falls within three standard deviations. (between $\mu - 3\sigma$ and $\mu + 3\sigma$)
 - in this case, 99.7% of the data falls between 80 and 120.

Q. 10- Provide a real life example of a Poisson process and calculate the probability for a specific event.

⇒ Customer arrivals at a bank Suppose a bank receives, on average, 5 customer every 10 minutes. This scenario can be modeled as a Poisson process because:

1. Customer arrivals are independent of each other.
2. The rate of arrival (5 customer per 10 minutes) is constant over time.
3. Two customers cannot arrive at the exact same time.

Now let's calculate the probability that exactly 7 customers will arrive in the next 10 minutes.

Poisson Distribution formula

The probability of observing exactly k events (here 7 customers) in a fixed period, given the average rate of 0.5 customers per 10 minutes, is given by the Poisson Probability formula.

$$P(X=k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

- * $P(X=k)$ is the probability of observing k events (7 customers).
- * λ is the average rate of occurrence (5 customers in 10 minutes).

- k is the number of events (customers) for which we care calculating the Probability (P customers).
- e is Euler's number (ex 2.718).

Calculation

$$\lambda = 5 \quad k = 7$$

$$P(X=7) = \frac{5^7}{7!} e^{-5}$$

$$5^7 = 78125$$

$$e^{-5} = 0.006737 \quad ; \quad 7! = 5040$$

$$P(X=7) = \frac{78125 \times 0.006737}{5040} = 525.125$$

$$= 0.1049$$

The Probability that exactly 7 customers will arrives at the bank in the next 10 minutes is approximately 0.1049.

Q. 21. Explain what a random Variable is
Ans Differentiate between discrete and continuous random variables.

→ A random Variable is a numerical value that represents the outcomes of a random phenomenon or experiment. It assigns a real number to each possible outcome in a sample space, allowing random events to be analyzed mathematically.

tically. Random Variables are essential in probability and statistics as they allow us to quantify uncertainty.

types of random variables

1. Discrete Random Variables
2. Continuous Random Variables.

1. Discrete Random Variable

Discrete Random

Values Countable, finite or infinite set of values

Example Number of heads in coin tosses, dice rolls.

Probability distribution function (PMF)

Probabilities non zero probabilities single for specific values

Calculation sum of Probabilities of prob for individual outcomes

$$P(CD = c)$$

Continuous Random

→ uncountable, infinite number of values in a range.

Heights, weights, time distance

Probability density function (PDF)

Probability of any specific value is 0.

Integration of the probability density over a range

$$P(CD \in [a, b]) = \int_a^b f(x) dx$$

Q. 12. Provide an example dataset, calculate both Covariance and Correlation and interpret the results.

→ Suppose we have the following data for two variables, x (Number of study hours) and y (Exam Scores), for 5 students.

Date _____
Page _____

Student	1	2	3	4	5
Study hours (x_i)	2	3	5	7	9
Exam score (y_i)	50	60	80	85	95

Mean of $x = \bar{x} = \frac{2+3+5+7+9}{5} = \frac{26}{5} = 5.2$

mean of $y = \bar{y} = \frac{50+60+80+85+95}{5} = \frac{370}{5} = 74$

* Calculate Covariance

$$\text{Cov}(x_i y_j) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Student	x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
1	2	50	-3.2	-24	76.8
2	3	60	-2.2	-14	30.8
3	5	80	-0.2	6	-1.2
4	7	85	1.8	11	19.8
5	9	95	3.8	21	78.8

$$\text{Cov}(x_i y_j) = \frac{206}{5} = 41.2$$

* Calculate Correlation

$$P(x_i y_j) = \frac{\text{Cov}(x_i y_j)}{\sigma_x \sigma_y}$$

$$\sigma_x = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

Standard deviation of \bar{x}

Student	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
1	-3.2	10.24
2	-2.2	4.84
3	-0.2	0.04
4	-1.8	3.24
5	3.2	14.44

$$\sigma_{\bar{x}} = \sqrt{\frac{32.8}{5}} = \sqrt{6.56} \approx 2.56$$

Standard deviation of y

Student	$y_i - \bar{y}$	$(y_i - \bar{y})^2$
1	-24	576
2	-14	196
3	6	36
4	11	121
5	21	441

$$\sigma_y = \sqrt{\frac{1370}{5}} = \sqrt{274} \approx 16.5$$

Correlation Coefficient.

$$P(x,y) = \frac{41.2}{256 \times 16.55} = \frac{41.2}{42.38} \approx 0.97$$

Interpretation of Results.

1. Covariance: The covariance of 41.2 indicates that there is a positive relationship between the number of study hours (x) and exam scores (y). Since the covariance is positive, we know that as Study

hours increase, exam scores tend to increase as well.

2 Correlation: The correlation coefficient of 0.972 suggests a strong positive linear relationship between the two variables. This means that more study hours are strongly associated with higher exam scores.