

1. Explain the different types of data (qualitative and quantitative) and provide examples of each. Discuss nominal, ordinal, interval, and ratio scales.

answer-Data can be classified into two primary types: qualitative and quantitative. Understanding these types helps in selecting appropriate statistical methods for analysis.

Types of Data

1. Qualitative Data (Categorical Data):

- This type of data describes characteristics or qualities and cannot be measured numerically. It is often used to categorize or label items.
- Examples:
 - Nominal Data: Categories without a specific order. Examples include:
 - Gender (male, female)
 - Colors (red, blue, green)
 - Types of cuisine (Italian, Mexican, Chinese)
 - Ordinal Data: Categories with a defined order but without a consistent difference between them. Examples include:
 - Education level (high school, bachelor's, master's, doctorate)
 - Satisfaction ratings (poor, fair, good, excellent)

2. Quantitative Data (Numerical Data):

- This type of data is numerical and can be measured or counted. It can be further divided into two subtypes:
- Examples:
 - Interval Data: Numerical data where the difference between values is meaningful, but there is no true zero point. Examples include:
 - Temperature (Celsius or Fahrenheit)
 - Dates (years)
 - Ratio Data: Numerical data with a true zero point, allowing for the comparison of absolute magnitudes. Examples include:
 - Height (inches or centimeters)
 - Weight (pounds or kilograms)
 - Time (seconds, minutes)

Scales of Measurement

1. Nominal Scale:

- The simplest form of measurement. Data is categorized without a specific order.

Example: Types of pets (dog, cat, fish). There is no ranking among the categories.

2. Ordinal Scale:

- Data is categorized and ordered, but the differences between categories are not quantifiable.
- Example: Movie ratings (1 star, 2 stars, 3 stars, etc.). The ratings indicate an order, but the difference between 1 and 2 stars isn't necessarily the same as between 2 and 3 stars.

3. Interval Scale:

- Data is ordered and the differences between values are meaningful. However, there is no true zero point.
- Example: Temperature in Celsius. The difference between 10°C and 20°C is the same as between 20°C and 30°C, but 0°C does not represent the absence of temperature.

4. Ratio Scale:

- Data is ordered, differences are meaningful, and there is a true zero point, allowing for a full range of mathematical operations.
- Example: Weight. A weight of 0 kg means no weight, and a weight of 10 kg is twice that of a weight

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 metrics that describe the center or typical value of a dataset. The three most common measures are the mean, median, and mode. Each measure has its own advantages and is suitable for different types of data and situations.

1. Mean

- Definition: The mean is the average of a set of values, calculated by summing all the values and dividing by the number of values.

$$\text{Mean} = \frac{\sum X_i}{N} \quad \text{Mean} = \frac{\sum X_i}{N}$$

- Example: For the dataset {3, 5, 7, 9}:

$$\text{Mean} = \frac{3 + 5 + 7 + 9}{4} = \frac{24}{4} = 6$$

- When to Use:

- The mean is appropriate for quantitative data that is normally distributed and when there are no extreme outliers, as outliers can skew the mean significantly.
- Example Situations:
 - Calculating the average score of students in a class.

- Determining average monthly income in a region (if income levels are fairly uniform).

2. Median

- Definition: The median is the middle value of a dataset when it is ordered from smallest to largest. If there is an even number of values, the median is the average of the two middle values.
- Example: For the dataset $\{3, 5, 7, 9\}$, the median is 666 (the average of 5 and 7, the two middle values). For the dataset $\{3, 5, 7\}$, the median is 555.
- When to Use:
 - The median is useful for skewed distributions or when there are outliers, as it is not affected by extreme values.
 - Example Situations:
 - Reporting household incomes in a region where a few households earn significantly more than others, skewing the mean.
 - Analyzing test scores where a few students score extremely low or high.

3. Mode

- Definition: The mode is the value that appears most frequently in a dataset. A dataset can have one mode (unimodal), more than one mode (bimodal or multimodal), or no mode at all if all values are unique.
- Example: For the dataset $\{1, 2, 2, 3, 4\}$, the mode is 222. For the dataset $\{1, 1, 2, 2, 3\}$, both 111 and 222 are modes (bimodal).
- When to Use:
 - The mode is particularly useful for categorical data, where we want to identify the most common category. It can also be useful for quantitative data with repeated values.
 - Example Situations:
 - Determining the most popular color of cars sold in a dealership.
 - Analyzing survey responses where one answer is significantly more frequent than others.

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

Dispersion, also known as variability or spread, refers to the extent to which data points in a dataset differ from the average (mean) or from each other. Understanding dispersion is crucial because it provides insight into the consistency, reliability, and overall distribution of the data. Two of the most commonly used measures of dispersion are variance and standard deviation.

Variance

- Definition: Variance measures the average squared deviation of each data point from the mean. It quantifies how much the values in a dataset spread out from their mean.
- Formula: For a population:

$$\sigma^2 = \frac{\sum (X_i - \mu)^2}{N}$$
For a sample:

$$s^2 = \frac{\sum (X_i - \bar{X})^2}{n - 1}$$
where:
 - X_i = each data point
 - μ = population mean
 - \bar{X} = sample mean
 - N = number of data points in the population
 - n = number of data points in the sample
- Interpretation: A higher variance indicates that the data points are spread out over a larger range of values, while a lower variance indicates that the data points tend to be closer to the mean.

Standard Deviation

- Definition: Standard deviation is the square root of the variance. It provides a measure of dispersion in the same units as the original data, making it more interpretable.
- Formula: For a population:

$$\sigma = \sqrt{\sigma^2}$$
For a sample:

$$s = \sqrt{s^2}$$
- Interpretation: Like variance, a higher standard deviation indicates a greater spread of data points around the mean. Since it is expressed in the same units as the data, it allows for more intuitive understanding and comparison of dispersion.

Example

Consider the following dataset: {4,6,8,10}.

1. Calculate the Mean:

$$\text{Mean} = \frac{4 + 6 + 8 + 10}{4} = 7$$

$$7 \times 4 = 28$$
2. Calculate the Variance:
 - Deviations from the mean: $(4-7)^2, (6-7)^2, (8-7)^2, (10-7)^2 \rightarrow 9, 1, 1, 9$
 - Sum of squared deviations: $9 + 1 + 1 + 9 = 20$
 - For a sample:
3. $s^2 = \frac{20}{4-1} = \frac{20}{3} \approx 6.67$
 $s = \sqrt{6.67} \approx 2.58$
4. Calculate the Standard Deviation:
 $s = 2.58$

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

A box plot, also known as a box-and-whisker plot, is a graphical representation of the distribution of a dataset that highlights its central tendency, variability, and potential outliers. It provides a visual summary of key statistical measures, making it easier to compare different datasets.

Components of a Box Plot

1. **Box:** The central box represents the interquartile range (IQR), which is the range between the first quartile (Q1) and the third quartile (Q3). This box contains the middle 50% of the data.
 - Q1 (First Quartile): The median of the lower half of the data (25th percentile).
 - Q3 (Third Quartile): The median of the upper half of the data (75th percentile).
 - The box extends from Q1 to Q3, with a line inside the box representing the median (Q2).
2. **Whiskers:** The lines extending from the box (the whiskers) indicate the range of the data outside the interquartile range. Typically, whiskers extend to the smallest and largest values within 1.5 times the IQR from the quartiles.
3. **Outliers:** Data points that fall outside the whiskers are considered outliers and are often plotted as individual points. These points may represent unusual observations that warrant further investigation.

What a Box Plot Can Tell You

1. Central Tendency: The position of the median line within the box gives a quick sense of the dataset's center.
2. Spread of the Data: The size of the box (IQR) indicates the spread of the middle 50% of the data. A larger box suggests greater variability, while a smaller box indicates more consistency.
3. Skewness: The relative lengths of the whiskers and the position of the median line can indicate skewness:
 - If the median is closer to Q1, the data may be right-skewed (more values on the lower end).
 - If the median is closer to Q3, the data may be left-skewed (more values on the higher end).
4. Outliers: The presence of outliers, marked as individual points beyond the whiskers, can highlight unusual observations that may require further investigation.
5. Comparison of Datasets: Box plots are particularly useful for comparing multiple groups side by side. By placing multiple box plots on the same graph, you can easily compare medians, spreads, and outlier presence across different categories.

Example

Consider a dataset of test scores from two classes:

- Class A: 55, 60, 65, 70, 75, 80, 85, 90
- Class B: 40, 55, 65, 70, 75, 80, 95, 100

A box plot for each class would reveal:

- Class A: A relatively narrow box with a median around 75, indicating that most students scored within a consistent range.
- Class B: A wider box with the presence of outliers (e.g., 40 and 100), suggesting greater variability in scores and potentially a less consistent performance.

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

Random sampling is a fundamental technique in statistics that plays a crucial role in making inferences about populations. It involves selecting a subset of individuals or observations from a larger population in such a way that every member of the population has an equal chance of being chosen. This method is vital for ensuring that the sample is representative of the population, which is essential for drawing valid conclusions.

Importance of Random Sampling

1. **Representation:**
 - Random sampling helps to ensure that the sample reflects the characteristics of the overall population. This reduces the likelihood of biases that could arise from selecting specific individuals or groups, making the results more generalizable.
2. **Minimizing Bias:**
 - By giving each member of the population an equal chance of being included, random sampling minimizes selection bias. This is crucial for the validity of the statistical inferences made from the sample data.
3. **Facilitating Statistical Inference:**
 - Random samples allow for the application of statistical methods and theorems that are based on probability. For example, the Central Limit Theorem states that the distribution of sample means will approximate a normal distribution as the sample size increases, regardless of the population's distribution, provided that the samples are random.
4. **Estimating Population Parameters:**
 - Random sampling enables researchers to estimate population parameters (such as the mean, variance, and proportions) and calculate confidence intervals around these estimates. This provides a quantifiable measure of uncertainty regarding the estimates.
5. **Hypothesis Testing:**
 - In hypothesis testing, random samples are essential for assessing the validity of null and alternative hypotheses. The results from random samples can help determine whether observed effects or differences are statistically significant.

Types of Random Sampling

1. **Simple Random Sampling:**
 - Every member of the population has an equal chance of being selected. This can be achieved through methods like drawing names from a hat or using random number generators.
2. **Stratified Sampling:**
 - The population is divided into subgroups (strata) based on specific characteristics (e.g., age, gender), and random samples are taken from each stratum. This ensures representation from each subgroup.
3. **Systematic Sampling:**
 - Members of the population are selected at regular intervals (e.g., every 10th person on a list). This method can be random if the starting point is chosen randomly.
4. **Cluster Sampling:**
 - The population is divided into clusters (e.g., geographical areas), and entire clusters are randomly selected for inclusion in the sample. This is often used when populations are too large or dispersed for simple random sampling.

6. Explain the concept of skewness and its types. How does skewness affect the interpretation of data?

Skewness is a statistical measure that describes the asymmetry of the distribution of data points in a dataset. It indicates whether the data are skewed to the left or right of the mean. Understanding skewness is important for interpreting data because it can affect the choice of statistical methods and the conclusions drawn from the data.

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. This means that a majority of the data points are concentrated on the lower end of the scale, with a few high outliers pulling the mean to the right of the median.
- **Characteristics:**
 - $\text{Mean} > \text{Median} > \text{Mode}$
- **Example:** Income distribution in many economies, where a small number of individuals have very high incomes compared to the majority.

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. Most data points are concentrated on the higher end of the scale, with a few low outliers pulling the mean to the left of the median.
- **Characteristics:**
 - $\text{Mean} < \text{Median} < \text{Mode}$
- **Example:** Age at retirement, where most people retire at a similar age but a few retire much earlier.

3. Zero Skewness (Symmetric Distribution):

- In a symmetric distribution, the data is evenly distributed around the mean. The left and right sides of the distribution mirror each other, meaning that the mean, median, and mode are all equal.
- **Example:** Normally distributed data, such as heights of adult men in a specific population.

Impact of Skewness on Data Interpretation

1. Mean vs. Median:

- Skewness affects the relationship between the mean and median. In skewed distributions, the mean can be significantly influenced by extreme values (outliers), while the median provides a better measure of central tendency.
- For example, in a right-skewed distribution, the mean may be higher than the median, suggesting that the average is influenced by a few high values.

2. Choice of Statistical Methods:

- Many statistical methods assume normality (symmetric distributions). When data are skewed, parametric tests (e.g., t-tests) may not be appropriate, and non-parametric tests (which do not assume a normal distribution) might be more suitable.
3. **Understanding Variability:**
 - Skewness provides insight into the variability and potential outliers in the data. In decision-making contexts, recognizing skewness can help identify risks associated with extreme values.
 4. **Visual Representation:**
 - Skewness can often be identified visually using histograms or box plots. Understanding the skewness helps in interpreting these visual representations more accurately.

7. 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 represents the range within which the central 50% of a dataset lies. It is calculated as the difference between the first quartile (Q1) and the third quartile (Q3):

$$\text{IQR} = Q3 - Q1$$

Calculation of IQR

1. Order the Data: Arrange the data points in ascending order.
2. Determine Q1 and Q3:
 - Q1 (the first quartile) is the median of the lower half of the data (25th percentile).
 - Q3 (the third quartile) is the median of the upper half of the data (75th percentile).
3. Calculate the IQR:
 - Subtract Q1 from Q3.

Example

Consider the following dataset: {3,7,8,12,14,18,21}

1. Order the Data: Already ordered.
2. Calculate Q1 and Q3:
 - Lower half: {3,7,8} → Q1 = 7
 - Upper half: {12,14,18,21} → Q3 = 18
3. Calculate the IQR: $\text{IQR} = 18 - 7 = 11$

Using IQR to Detect Outliers

The IQR is commonly used to identify outliers in a dataset through the following steps:

1. Determine the Lower and Upper Bound:
 - Lower Bound: $Q1 - 1.5 \times \text{IQR}$
 - Upper Bound: $Q3 + 1.5 \times \text{IQR}$
2. Identify Outliers:
 - Any data points below the lower bound or above the upper bound are considered outliers.

Example of Outlier Detection

Continuing with the previous example:

1. Calculate the bounds:
 - $\text{IQR} = 11, Q1 = 7, Q3 = 18$
 - Lower Bound: $7 - 1.5 \times 11 = 7 - 16.5 = -9.5$
 - Upper Bound: $18 + 1.5 \times 11 = 18 + 16.5 = 34.5$
2. Identify Outliers:
 - Any data points below -9.5 or above 34.5 are considered outliers.
 - Since all values in the dataset $\{3, 7, 8, 12, 14, 18, 21\}$ fall within this range, there are no outliers.

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

The binomial distribution is a discrete probability distribution that models the number of successes in a fixed number of independent Bernoulli trials, each with the same probability of success. It is used under specific conditions, which are essential for the validity of the model. Here are the key conditions:

Conditions for Binomial Distribution

1. Fixed Number of Trials (n):
 - The number of trials must be predetermined and constant. Each trial is conducted under the same conditions.
 - Example: Flipping a coin 10 times.
2. Two Possible Outcomes:
 - Each trial results in one of two outcomes, commonly referred to as "success" (e.g., heads in a coin flip) and "failure" (e.g., tails).

- Example: In a medical test, a patient can either test positive (success) or negative (failure).
- 3. Constant Probability of Success (p):
 - The probability of success remains the same for each trial. This means that the trials are independent, and the outcome of one trial does not affect the others.
 - Example: The probability of rolling a 3 on a fair six-sided die is always $\frac{1}{6}$ for each roll.
- 4. Independence of Trials:
 - The trials must be independent, meaning the outcome of one trial does not influence the outcome of another.
 - Example: Each flip of a coin does not affect the results of subsequent flips.

Binomial Distribution Formula

If all conditions are met, the probability of getting exactly k successes in n trials is given by the binomial probability formula:

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

where:

- $\binom{n}{k}$ is the binomial coefficient (the number of ways to choose k successes from n trials),
- p is the probability of success,
- $(1-p)$ is the probability of failure,
- n is the total number of trials,
- k is the number of successes.

Examples of Binomial Distribution Applications

1. Coin Flipping:
 - Flipping a coin 10 times and counting the number of heads (successes).
2. Quality Control:
 - Testing a batch of products to see how many are defective (e.g., testing 100 light bulbs to see how many are non-functioning).
3. Marketing Campaigns:
 - Sending out 500 email invitations and measuring how many recipients respond positively.

Summary

The binomial distribution is appropriate when there is a fixed number of independent trials, each with two possible outcomes, and a constant probability of success. Understanding these conditions

9. Explain the properties of the normal distribution and the empirical rule (68-95-99.7 rule).

The normal distribution is a fundamental probability distribution in statistics, characterized by its symmetric, bell-shaped curve. It is often referred to as a Gaussian distribution. Understanding its properties and the empirical rule is essential for interpreting data that follows this distribution.

Properties of the Normal Distribution

1. Symmetry:
 - The normal distribution is symmetric about its mean. This means that the left and right sides of the curve are mirror images.
2. Mean, Median, and Mode:
 - In a normal distribution, the mean, median, and mode are all equal and located at the center of the distribution. This central point is also the highest point of the curve.
3. Bell-Shaped Curve:
 - The shape of the distribution is bell-shaped, tapering off equally on both ends. This indicates that values further away from the mean are less likely to occur.
4. Asymptotic:
 - The tails of the normal distribution approach but never touch the horizontal axis. This means that extreme values (far from the mean) are theoretically possible, although very unlikely.
5. Defined by Two Parameters:
 - The normal distribution is completely defined by its mean (μ) and standard deviation (σ). The mean determines the location of the center of the graph, while the standard deviation determines the width of the graph.
6. Area Under the Curve:
 - The total area under the normal distribution curve is equal to 1, representing 100% of the data.

The Empirical Rule (68-95-99.7 Rule)

The empirical rule describes how data is distributed in a normal distribution in terms of standard deviations from the mean. It states:

1. 68% of the data falls within one standard deviation of the mean:
 - This means that approximately 68% of the values in a normal distribution lie between $\mu - \sigma$ and $\mu + \sigma$.
2. 95% of the data falls within two standard deviations of the mean:
 - About 95% of the values lie between $\mu - 2\sigma$ and $\mu + 2\sigma$.
3. 99.7% of the data falls within three standard deviations of the mean:
 - Roughly 99.7% of the values lie between $\mu - 3\sigma$ and $\mu + 3\sigma$.

Visualization of the Empirical Rule

- One Standard Deviation (σ):
 - Covers the range from $\mu - \sigma$ to $\mu + \sigma$ (approximately 68% of the data).
- Two Standard Deviations (2σ):
 - Covers the range from $\mu - 2\sigma$ to $\mu + 2\sigma$ (approximately 95% of the data).
- Three Standard Deviations (3σ):
 - Covers the range from $\mu - 3\sigma$ to $\mu + 3\sigma$ (approximately 99.7% of the data).

Applications of the Normal Distribution and Empirical Rule

The properties of the normal distribution and the empirical rule are widely used in various fields, including:

- Quality Control: To assess product consistency and identify defects.
- Psychometrics: In testing and measurement, where many traits are normally distributed.
- Finance: For modeling returns on investments and assessing risk.
- Social Sciences: To analyze survey results and population characteristics.

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

A Poisson process is a statistical model that describes events occurring randomly over a fixed interval of time or space. It is characterized by a constant average rate (λ) of occurrence and assumes that events happen independently of each other.

Real-Life Example: Customer Arrivals at a Coffee Shop

Consider a coffee shop that experiences an average of 3 customers arriving every 10 minutes. We can model the number of customer arrivals in a 10-minute interval as a Poisson process, where the average rate (λ) is 3.

Problem: Calculate the Probability of a Specific Event

Let's calculate the probability that exactly 5 customers arrive at the coffee shop in a 10-minute interval.

Poisson Probability Formula

The probability of observing k events (customers arriving) in a Poisson process is given by the formula:

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

where:

- $P(X=k)$ is the probability of observing k events.
- e is the base of the natural logarithm (approximately equal to 2.71828).
- λ is the average rate of occurrence (in this case, 3).
- k is the number of occurrences we want to find the probability for (in this case, 5).

Given Values

- $\lambda = 3$
- $k = 5$

Calculation

1. Calculate $e^{-\lambda}$:
 $e^{-3} \approx 0.0498$
2. Calculate λ^k :
 $3^5 = 243$
3. Calculate $k!$:
 $5! = 120$
4. Substitute into the formula:
 $P(X=5) = \frac{e^{-3} \cdot 3^5}{5!} = \frac{0.0498 \cdot 243}{120}$
5. Calculate the probability:
 $P(X=5) \approx 0.1008$

Conclusion

The probability that exactly 5 customers arrive at the coffee shop in a 10-minute interval is approximately 0.1008, or 10.08%. This example illustrates how a Poisson process can be used to model real-world events, providing insights into customer behavior and helping businesses make

11. Explain what a random variable is and differentiate between discrete and continuous random variables.

A random variable is a numerical outcome of a random phenomenon. It assigns a numerical value to each possible outcome in a sample space of a stochastic (random) process, allowing for quantitative analysis of uncertainty. Random variables are typically categorized into two main types: discrete and continuous.

Types of Random Variables

1. Discrete Random Variables

- **Definition:** A discrete random variable can take on a countable number of distinct values. These values are typically whole numbers or integers, and there are gaps between possible values.
- **Examples:**
 - The number of heads when flipping a coin three times (possible values: 0, 1, 2, 3).
 - The number of customers arriving at a store in an hour (possible values: 0, 1, 2, ...).
 - The outcome of rolling a die (possible values: 1, 2, 3, 4, 5, 6).
- **Probability Distribution:** Discrete random variables are described by a probability mass function (PMF), which gives the probability of each possible value.

2. Continuous Random Variables

- **Definition:** A continuous random variable can take on an infinite number of values within a given range. These values can be any real number, including fractions and decimals.
- **Examples:**
 - The height of students in a classroom (any value within a realistic range, e.g., 150.5 cm to 190.2 cm).
 - The time it takes for a computer to complete a task (could be any positive real number).
 - The temperature in a city over a day (can vary continuously).
- **Probability Distribution:** Continuous random variables are described by a probability density function (PDF), which provides the probabilities of the variable falling within a particular range of values. The total area under the PDF curve equals 1.

Key Differences

Feature	Discrete Random Variables	Continuous Random Variables
Nature of Values	Countable (e.g., integers)	Uncountable (e.g., real numbers)

Examples	Number of students, dice rolls	Height, weight, temperature
Probability Distribution	Probability mass function (PMF)	Probability density function (PDF)
Probability of Specific Value	Probability of a specific value is non-zero	Probability of a specific value is zero (area under a point is zero)
Range of Values	Finite or countably infinite	Infinite within a range

Summary

Random variables are fundamental concepts in probability and statistics, allowing us to quantify and analyze random phenomena. Discrete random variables take on specific, countable values, while continuous random variables can assume any value within a range. Understanding the difference

12. Provide an example dataset, calculate both covariance and correlation, and interpret the results.

Let's consider a simple example dataset to calculate both covariance and correlation. We'll use two variables: X (the number of hours studied) and Y (the scores obtained on a test).

Example Dataset

Student	Hours Studied (X)	Test Score (Y)
1	2	65
2	3	70

3	5	80
4	7	85
5	8	90

Step 1: Calculate the Means

1. Mean of X:

$$\bar{X} = \frac{2+3+5+7+8}{5} = \frac{25}{5} = 5$$

2. Mean of Y:

$$\bar{Y} = \frac{65+70+80+85+90}{5} = \frac{390}{5} = 78$$

Step 2: Calculate Covariance

Covariance is calculated using the formula:

$$\text{Cov}(X, Y) = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{n - 1}$$

where n is the number of data points.

Calculating Each Term

Student	X_i	Y_i	$X_i - \bar{X}$	$Y_i - \bar{Y}$	$(X_i - \bar{X})(Y_i - \bar{Y})$
1	2	65	$2 - 5 = -3$	$65 - 78 = -13$	39
2	3	70	$3 - 5 = -2$	$70 - 78 = -8$	16
3	5	80	$5 - 5 = 0$	$80 - 78 = 2$	0

4	7	85	$7 - 5 = 2$	$85 - 78 = 7$	14
---	---	----	-------------	---------------	----

5	8	90	$8 - 5 = 3$	$90 - 78 = 12$	36
---	---	----	-------------	----------------	----

Sum of Products

$$\sum (X_i - \bar{X})(Y_i - \bar{Y}) = 39 + 16 + 0 + 14 + 36 = 105$$

$$\sum (X_i - \bar{X})(Y_i - \bar{Y}) = 39 + 16 + 0 + 14 + 36 = 105$$

Covariance Calculation

$$\text{Cov}(X, Y) = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{n - 1} = \frac{105}{5 - 1} = \frac{105}{4} = 26.25$$

Step 3: Calculate Correlation

Correlation is calculated using the formula:

$$r = \frac{\text{Cov}(X, Y)}{s_X s_Y}$$

where s_X and s_Y are the standard deviations of X and Y .

Step 3.1: Calculate Standard Deviations

1. Standard Deviation of X :

$$s_X = \sqrt{\frac{\sum (X_i - \bar{X})^2}{n - 1}}$$

$$\sum (X_i - \bar{X})^2 = (-3)^2 + (-2)^2 + (0)^2 + (2)^2 + (3)^2 = 9 + 4 + 0 + 4 + 9 = 26$$

$$s_X = \sqrt{\frac{26}{4}} = \sqrt{6.5} \approx 2.55$$

2. Standard Deviation of Y :

$$s_Y = \sqrt{\frac{\sum (Y_i - \bar{Y})^2}{n - 1}}$$

$$\sum (Y_i - \bar{Y})^2 = (-13)^2 + (-8)^2 + (2)^2 + (7)^2 + (12)^2 = 169 + 64 + 4 + 49 + 144 = 430$$

$$s_Y = \sqrt{\frac{430}{4}} = \sqrt{107.5} \approx 10.37$$

Step 3.2: Calculate Correlation

$$r = \frac{\text{Cov}(X, Y)}{s_X s_Y} = \frac{26.25}{(2.55)(10.37)} \approx \frac{26.25}{26.49} \approx 0.99$$

Interpretation of Results

- Covariance (26.25): This positive covariance indicates a strong positive relationship between the number of hours studied and the test scores. As one variable increases, the other tends to increase as well.
- Correlation (0.99): The correlation coefficient is very close to 1, suggesting a nearly perfect positive linear relationship between the two variables. This means that as students study more hours, their test scores increase almost proportionally.

Conclusion

In this dataset, both covariance and correlation reveal a strong positive relationship between hours studied and test scores. This indicates that increased study time is associated with higher scores, which could inform strategies for academic performance improvement.