

## Assignment Questions:

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

Ans: Data can be broadly categorized into qualitative and quantitative data.

1. Qualitative (Categorical) Data Qualitative data describes characteristics or attributes that cannot be measured numerically. It is used to classify or label elements within a dataset.

- Examples:
  - Eye color (Brown, Blue, Green)
  - Customer feedback (Positive, Neutral, Negative)
  - Car brands (Toyota, Ford, BMW)

#### Subtypes of Qualitative Data:

- Nominal Scale: Data that consists of categories without any order or ranking.
  - Example: Blood types (A, B, AB, O), types of pets (Dog, Cat, Bird).
- Ordinal Scale: Data that has a meaningful order but no consistent difference between categories.
  - Example: Education levels (Primary, Secondary, Tertiary), customer satisfaction ratings (Poor, Fair, Good, Excellent).

### 2. Quantitative (Numerical) Data

Quantitative data consists of numbers and can be measured or counted.

- Examples:
  - Age (25 years)
  - Salary (₹50,000)
  - Temperature (37°C)

#### Subtypes of Quantitative Data:

- Interval Scale: Data with ordered values and equal intervals, but no true zero point.
  - Example: Temperature in Celsius (0°C does not mean "no temperature"), IQ scores.
- Ratio Scale: Data with ordered values, equal intervals, and a true zero, meaning zero represents an absence of the quantity.
  - Example: Weight (0 kg means no weight), salary (₹0 means no income), height.

## **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.**

Ans: Measures of central tendency are statistical values used to describe the "middle" or typical value within a data set, with the three main measures being the mean (average), median (middle value), and mode (most frequent value); each is best used depending on the distribution of data and presence of outliers:.

Mean:

- Definition: The sum of all values in a data set divided by the number of values.
- When to use: When the data is fairly evenly distributed with no extreme outliers, as the mean represents the "center of mass" of the data.
- Example: A class of 20 students scores an average of 85 on a test, meaning the mean score is 85.

Median:

- Definition: The middle value in a data set when arranged in ascending order.
- When to use: When the data is skewed or has outliers, as the median is less affected by extreme values.
- Example: In a neighborhood where most houses are valued at \$200,000 but a few mansions are worth millions, the median house price might be a better representation of the "typical" price than the mean.

Mode:

- Definition: The value that appears most frequently in a data set.
- When to use: When you want to identify the most common category or value in a data set, especially for categorical data.
- Example: In a survey of favorite colors, if "blue" is chosen by the most people, then "blue" is the mode.

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

Ans: Dispersion in statistics refers to the degree to which data points are spread out or scattered around a central value, essentially measuring how varied the data is within a set; while variance and standard deviation are two key metrics used to quantify this spread, with variance representing the average squared deviation from the mean and standard deviation being the square root of variance, providing a more interpretable measure in the same units as the data itself.

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

Ans: Box plots are used to show distributions of numeric data values, especially when you want to compare them between multiple groups. They are built to provide high-level information at a glance, offering general information about a group of data's symmetry, skew, variance, and outliers.

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

Ans: Random sampling plays a crucial role in making inferences about populations by ensuring that a selected sample is representative of the larger population, minimizing bias and allowing researchers to draw reliable conclusions about the population characteristics based on the data collected from that sample, rather than having to study every individual within the population

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

Ans: Skewness refers to the degree of asymmetry in a data distribution, indicating whether the data is unevenly distributed around the mean, with one side of the distribution having a longer "tail" than the other; essentially, it measures how much a distribution deviates from a perfect bell curve symmetry, where the left and right sides are mirror images.

Types of Skewness:

- **Positive Skewness (Right Skewed):** When the tail of the distribution extends further to the right side, meaning there are more low values and fewer extreme high values. In this case, the mean is typically greater than the median and mode.
- **Negative Skewness (Left Skewed):** When the tail of the distribution extends further to the left side, indicating more high values and fewer extreme low values. Here, the mean is usually less than the median and mode.
- **Zero Skewness (Symmetrical):** When the distribution is perfectly balanced with no noticeable asymmetry, meaning the mean, median, and mode are equal.

How Skewness Affects Data Interpretation:

- **Misinterpretation of Central Tendency:** In a skewed distribution, relying solely on the mean can be misleading as it can be pulled towards the longer tail, not accurately representing the "typical" value. In such cases, the median is often a better indicator of central tendency.
- **Outlier Identification:** Skewness can highlight the presence of outliers, as extreme values on the longer tail can significantly impact the distribution.
- **Choice of Statistical Tests:** Certain statistical tests are designed for normally distributed data, so if data is significantly skewed, using such tests could lead to inaccurate results. In such cases, non-parametric statistical tests might be more appropriate.

## 7. What is the interquartile range (IQR), and how is it used to detect outliers?

Ans: Interquartile Range (IQR) & Outlier Detection

- IQR Definition: The Interquartile Range (IQR) is the difference between the third quartile (Q3) and the first quartile (Q1).  
$$\text{IQR} = Q3 - Q1$$
- Quartiles Explanation:
  - Q1 (25th percentile): The median of the lower half of the data.
  - Q3 (75th percentile): The median of the upper half of the data.
- Outlier Detection:
  - A data point is considered an outlier if it falls below or above the following thresholds:  
$$\text{Lower Bound} = Q1 - 1.5 \times \text{IQR}$$
$$\text{Upper Bound} = Q3 + 1.5 \times \text{IQR}$$
- Usage:
  - Identifies extreme values that might distort analysis.
  - Used in box plots to visualize data distribution and outliers.
  - Removes anomalies in financial analysis, machine learning, and statistics.

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

Ans: The binomial distribution is used when a situation involves a fixed number of independent trials, where each trial has only two possible outcomes (typically called "success" and "failure"), and the probability of success remains constant across all trials; in simpler terms, there must be a fixed number of trials, each trial must be independent of the others, and the probability of success must be the same for every trial.

Key conditions for using the binomial distribution:

- Fixed number of trials (n): The experiment must consist of a predetermined number of trials that do not change throughout the analysis.
- Independent trials: The outcome of one trial should not influence the outcome of any other trial.
- Two possible outcomes: Each trial can only result in one of two mutually exclusive outcomes, typically labeled as "success" and "failure".

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

Ans: The normal distribution, often visualized as a bell-shaped curve, is a probability distribution where most data points cluster around the mean, with values tapering off symmetrically on either side, and the "empirical rule" (or 68-95-99.7 rule) states that in a normal distribution, approximately 68% of data falls within one standard deviation of the mean, 95% falls within two standard deviations, and 99.7% falls within three standard deviations of the mean; essentially describing how data is spread around the average value.

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

Ans: A real-life example of a Poisson process is the number of phone calls received by a customer service center per hour, assuming the calls arrive independently at a constant average rate; let's say the average call rate is 5 calls per hour. To calculate the probability of receiving exactly 2 calls in a given hour, we would use the Poisson distribution formula with  $\lambda$  (average rate) = 5 and  $k$  (desired number of calls) = 2, resulting in the calculation:  $P(X=2) = (e^{-5} * 5^2) / 2! = 0.084$ , meaning there is an 8.4% chance of receiving exactly 2 calls in that hour.

Key points about this example:

- Poisson process characteristics:
- Events occur independently of each other.
- The average rate of events occurring is constant over time.
- The probability of an event occurring in a small time interval is proportional to the length of that interval.

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

Ans:

Type	Definition	Examples	Key Characteristics
Discrete Random Variable	Takes <b>finite</b> or <b>countable</b> values	Number of students in a class, Number of calls received in an hour	Values are <b>distinct</b> and <b>countable</b> (0,1,2,...)
Continuous Random Variable	Takes <b>infinite</b> values within a range	Height of a person, Time taken to finish a task, Temperature	Values are <b>measured</b> and can take <b>any value</b> in an interval (e.g., 2.34, 5.67)

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

Ans: | Student | Study Hours (X) | Exam Score (Y) |

|---|---|---|

| 1 | 2 | 60 |

| 2 | 5 | 75 |

| 3 | 3 | 65 |

| 4 | 8 | 90 |

| 5 | 1 | 50 |

Calculating Covariance:

Step 1: Find the mean of each variable:

Mean of Study Hours (X) =  $(2 + 5 + 3 + 8 + 1) / 5 = 3.8$

Mean of Exam Scores (Y) =  $(60 + 75 + 65 + 90 + 50) / 5 = 70$

Step 2: Calculate the deviation from the mean for each data point:

For Student 1:  $(2 - 3.8) * (60 - 70) = 18$

For Student 2:  $(5 - 3.8) * (75 - 70) = 6.2$

... and so on for all students.

Step 3: Multiply the deviations and sum them up:

Sum of (deviation of X \* deviation of Y) =  $18 + 6.2 + 0.2 - 12.2 + 18 = 30.4$

Step 4: Divide by the number of data points minus 1:

Covariance (Cov(X,Y)) =  $30.4 / (5 - 1) = 7.6$

Calculating Correlation:

Step 1: Calculate the standard deviation of each variable:

Standard Deviation of X = 2.45

Standard Deviation of Y = 14.14

Step 2: Divide the covariance by the product of standard deviations:

$$\text{Correlation (r)} = \text{Covariance} / (\text{SD}_X * \text{SD}_Y) = 7.6 / (2.45 * 14.14) = 0.22$$

Interpretation:

Positive Covariance:

The positive value of the covariance (7.6) indicates that when study hours increase, exam scores tend to increase as well, showing a positive relationship between the two variables.

Moderate Correlation:

The correlation coefficient (0.22) is positive but relatively small, suggesting that while there is a positive association between study hours and exam scores, it is not a very strong relationship.