

## Q1. What is Statistics?

Statistics is the practice and science of collecting, analyzing, interpreting, and presenting numerical data. It involves using mathematical methods to gather information about a group or population and making inferences about that population based on the data collected from a smaller sample. Statistics plays a critical role in many fields, including business, science, social sciences, medicine, engineering, and many others. Some of the key concepts in statistics include probability theory, statistical inference, hypothesis testing, and regression analysis. By understanding statistical principles, researchers and decision-makers can make informed decisions based on data-driven insights.

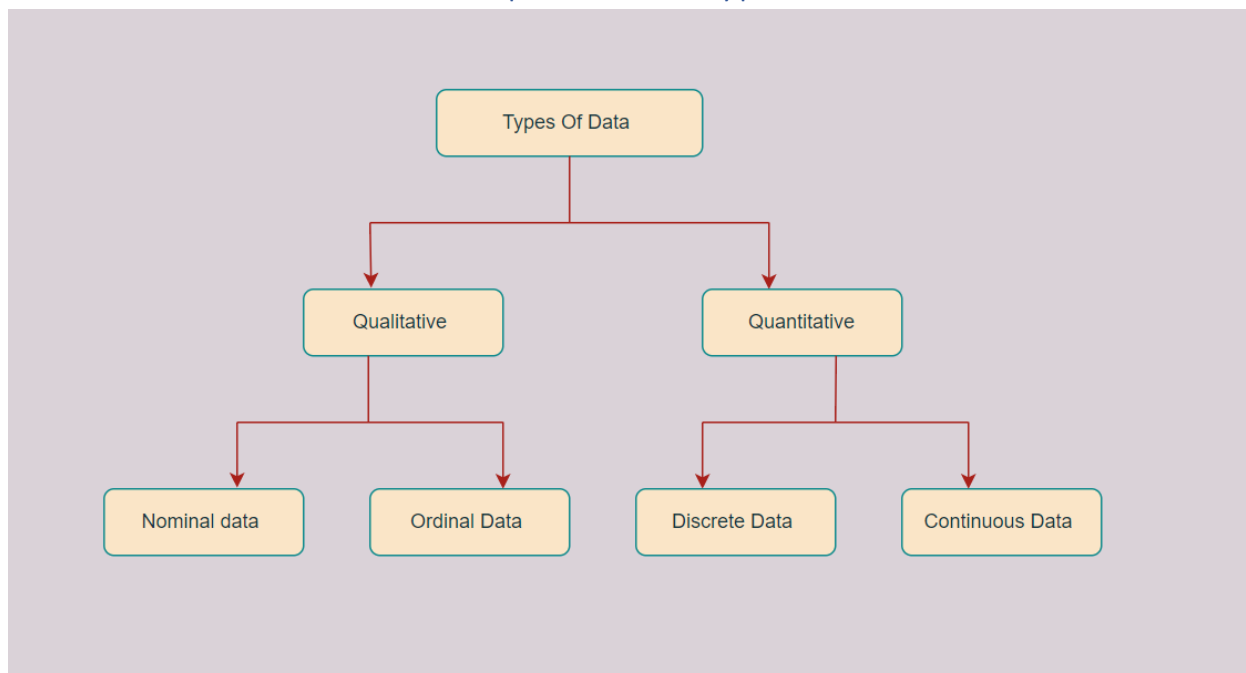
## Q2. Define the different types of statistics and give an example of when each type might be used.

There are two main types of statistics: descriptive statistics and inferential statistics.

**Descriptive statistics:** Descriptive statistics are used to summarize and describe the characteristics of a dataset. This type of statistics is used to summarize and describe the main features of a dataset, such as the mean, median, mode, variance, and standard deviation. Descriptive statistics can be used to analyze data in many fields, such as finance, healthcare, and social sciences. For example, a healthcare researcher may use descriptive statistics to summarize the demographic characteristics of patients in a clinical trial.

**Inferential statistics:** Inferential statistics are used to make inferences about a population based on a sample of data. This type of statistics is used to test hypotheses and make predictions based on the sample data. Inferential statistics are widely used in fields such as economics, psychology, and political science. For example, a political scientist might use inferential statistics to analyze survey data to predict voting behavior in an upcoming election.

Q3. What are the different types of data and how do they differ from each other? Provide an example of each type of data.



### 1. Qualitative or Categorical Data

Qualitative or Categorical Data is data that can't be measured or counted in the form of numbers. These types of data are sorted by category, not by number. That's why it is also known as Categorical Data. These data consist of audio, images, symbols, or text. The gender of a person, i.e., male, female, or others, is qualitative data.

Qualitative data talks about the perception of people. This data helps market researchers understand the customers' tastes and then design their ideas and strategies accordingly.

The other examples of qualitative data are:

- What language do you speak
- Favorite holiday destination
- Opinion on something (agree, disagree, or neutral)
- Colors

The Qualitative data are further classified into two parts:

#### 1.1. Nominal Data

Nominal Data is used to label variables without any order or quantitative value. The color of hair can be considered nominal data, as one color can't be compared with another color.

The name “nominal” comes from the Latin name “nomen,” which means “name.” With the help of nominal data, we can’t do any numerical tasks or can’t give any order to sort the data. These data don’t have any meaningful order; their values are distributed into distinct categories.

Examples of Nominal Data:

- Colour of hair (Blonde, red, Brown, Black, etc.)
- Marital status (Single, Widowed, Married)
- Nationality (Indian, German, American)
- Gender (Male, Female, Others)
- Eye Color (Black, Brown, etc.)

## **1.2. Ordinal Data**

Ordinal data have natural ordering where a number is present in some kind of order by their position on the scale. These data are used for observation like customer satisfaction, happiness, etc., but we can’t do any arithmetical tasks on them.

Ordinal data is qualitative data for which their values have some kind of relative position. These kinds of data can be considered “in-between” qualitative and quantitative data. The ordinal data only shows the sequences and cannot use for statistical analysis. Compared to nominal data, ordinal data have some kind of order that is not present in nominal data.

Examples of Ordinal Data:

- When companies ask for feedback, experience, or satisfaction on a scale of 1 to 10
- Letter grades in the exam (A, B, C, D, etc.)
- Ranking of people in a competition (First, Second, Third, etc.)
- Economic Status (High, Medium, and Low)
- Education Level (Higher, Secondary, Primary)

### Difference between Nominal and Ordinal Data

<b>Nominal Data</b>	<b>Ordinal Data</b>
Nominal data can't be quantified, neither they have any intrinsic ordering	Ordinal data gives some kind of sequential order by their position on the scale
Nominal data is qualitative data or categorical data	Ordinal data is said to be "in-between" qualitative data and quantitative data
They don't provide any quantitative value, neither can we perform any arithmetical operation	They provide sequence and can assign numbers to ordinal data but cannot perform the arithmetical operation
Nominal data cannot be used to compare with one another	Ordinal data can help to compare one item with another by ranking or ordering
<b>Examples:</b> Eye color, housing style, gender, hair color, religion, marital status, ethnicity, etc	<b>Examples:</b> Economic status, customer satisfaction, education level, letter grades, etc

## 2. Quantitative Data

Quantitative data can be expressed in numerical values, making it countable and including statistical data analysis. These kinds of data are also known as Numerical data. It answers the questions like "how much," "how many," and "how often." For example, the price of a phone, the computer's ram, the height or weight of a person, etc., falls under quantitative data.

Quantitative data can be used for statistical manipulation. These data can be represented on a wide variety of graphs and charts, such as bar graphs, histograms, scatter plots, boxplots, pie charts, line graphs, etc.

## **Examples of Quantitative Data :**

- Height or weight of a person or object
- Room Temperature
- Scores and Marks (Ex: 59, 80, 60, etc.)
- Time

The Quantitative data are further classified into two parts:

### **2.1. Discrete Data**

The term discrete means distinct or separate. The discrete data contain the values that fall under integers or whole numbers. The total number of students in a class is an example of discrete data. These data can't be broken into decimal or fraction values.

The discrete data are countable and have finite values; their subdivision is not possible. These data are represented mainly by a bar graph, number line, or frequency table.

Examples of Discrete Data:

- Total numbers of students present in a class
- Cost of a cell phone
- Numbers of employees in a company
- The total number of players who participated in a competition
- Days in a week

### **2.2. Continuous Data**

Continuous data are in the form of fractional numbers. It can be the version of an android phone, the height of a person, the length of an object, etc. Continuous data represents information that can be divided into smaller levels. The continuous variable can take any value within a range.

The key difference between discrete and continuous data is that discrete data contains the integer or whole number. Still, continuous data stores the fractional numbers to record different types of data such as temperature, height, width, time, speed, etc.

Examples of Continuous Data :

- Height of a person
- Speed of a vehicle
- "Time-taken" to finish the work
- Wi-Fi Frequency
- Market share price

Difference between Discrete and Continuous Data

Discrete Data	Continuous Data
Discrete data are countable and finite; they are whole numbers or integers	Continuous data are measurable; they are in the form of fractions or decimal
Discrete data are represented mainly by bar graphs	Continuous data are represented in the form of a histogram
The values cannot be divided into subdivisions into smaller pieces	The values can be divided into subdivisions into smaller pieces
Discrete data have spaces between the values	Continuous data are in the form of a continuous sequence
<b>Examples:</b> Total students in a class, number of days in a week, size of a shoe, etc	<b>Example:</b> Temperature of room, the weight of a person, length of an object, etc

Q4. Categorise the following datasets with respect to quantitative and qualitative data types:

(i) Grading in exam: A+, A, B+, B, C+, C, D, E

**Answer: Qualitative Data**

(ii) Colour of mangoes: yellow, green, orange, red

**Answer: Qualitative Data**

(iii) Height data of a class: [178.9, 179, 179.5, 176, 177.2, 178.3, 175.8,...]

**Answer: Quantitative Data**

(iv) Number of mangoes exported by a farm: [500, 600, 478, 672, ...]

**Answer: Quantitative Data**

Q5. Explain the concept of levels of measurement and give an example of a variable for each level.

Levels of measurement, also known as scales of measurement, are ways to categorize or classify variables in statistical analysis. There are four commonly recognized levels of measurement, which are:

**Nominal level:** This level of measurement is the lowest level of measurement and is used to categorize data into groups or classes. Variables measured at this level have no inherent order or numerical value. Examples of variables at the nominal level include gender (male or female), hair color (blonde, brunette, black, etc.), and marital status (single, married, divorced, etc.).

**Ordinal level:** This level of measurement is used to rank or order data. Variables measured at this level have a specific order but do not have an equal distance between them. Examples of variables at the ordinal level include educational levels (elementary, high school, college, etc.), socioeconomic status (low, middle, high), and performance ratings (poor, average, good, excellent).

**Interval level:** This level of measurement has equal intervals between values, but there is no true zero point. Examples of variables at the interval level include temperature (measured in degrees Celsius or Fahrenheit), calendar dates (January 1st, February 1st, March 1st, etc.), and IQ scores.

**Ratio level:** This level of measurement is the highest level of measurement and has all the properties of the previous levels. It has equal intervals between values, a true zero point, and the ability to compare and calculate ratios. Examples of variables at the ratio level include weight, height, income, and age.

Example variables for each level of measurement:

**Nominal level:** eye color (blue, green, brown, etc.), blood type (A, B, AB, O), political party (Republican, Democrat, Independent)

**Ordinal level:** academic achievement (first, second, third, etc.), military rank (private, sergeant, captain, etc.), survey responses (strongly disagree, disagree, neutral, agree, strongly agree)

**Interval level:** temperature (measured in Celsius or Fahrenheit), calendar dates (January 1st, February 1st, March 1st, etc.), IQ scores

**Ratio level:** weight, height, age, income, number of children.

Q6. Why is it important to understand the level of measurement when analyzing data? Provide an example to illustrate your answer.

It is important to understand the level of measurement when analyzing data because different statistical techniques are appropriate for different levels of measurement. Using the wrong statistical technique can result in incorrect or misleading conclusions.

For example, suppose we have data on the hair color of a group of people. Hair color is a nominal variable because it cannot be ordered or ranked. If we were to calculate the mean hair color of this group, it would be meaningless because the values of the variable cannot be numerically averaged. Instead, we could report the percentage of people in the group with each hair color or create a frequency distribution.



In contrast, if we have data on the heights of a group of people, height is a continuous variable measured at the ratio level. We can use statistical techniques such as mean, median, and standard deviation to analyze this data. If we were to use a technique appropriate for nominal data, such as a frequency distribution, we would lose important information about the variability of the heights in the group.

Understanding the level of measurement of a variable allows us to choose appropriate statistical techniques, interpret results correctly, and avoid making errors in our analysis.

### Q7. How nominal data type is different from ordinal data type.

Nominal and ordinal data types are both categorical data types, but they differ in terms of the level of measurement and the nature of the categories.

Nominal data is the lowest level of measurement and is used to categorize data into groups or classes. The categories at the nominal level have no inherent order or numerical value. Nominal data is characterized by the fact that categories cannot be arranged in any specific order or sequence. Examples of nominal data include gender (male or female), race (Asian, Black, White, etc.), and eye color (blue, green, brown, etc.).

In contrast, ordinal data is used to rank or order data. The categories in ordinal data have a specific order or rank, but they do not have an equal distance between them. Ordinal data is characterized by the fact that categories have a natural order. Examples of ordinal data include educational levels (elementary, high school, college, etc.), economic status (low, middle, high), and performance ratings (poor, fair, good, excellent).

The key difference between nominal and ordinal data is that ordinal data has a natural ordering or ranking of categories, while nominal data does not. Additionally, ordinal data allows for comparisons between categories that are ranked higher or lower, whereas nominal data does not allow for comparisons based on the order or rank of categories.

### Q8. Which type of plot can be used to display data in terms of range? A box plot, also known as a box-and-whisker plot, can be used to display data in terms of range.

A box plot displays the distribution of a dataset by showing the median, quartiles, and the range of the data. The box represents the middle 50% of the data, with the bottom of the box representing the first quartile (Q1) and the top of the box representing the third quartile (Q3). The line inside the box represents the median of the data. The "whiskers" or lines extending from the box represent the range of the data, with the maximum and minimum values displayed as points at the end of the whiskers.

By displaying the range of the data in a box plot, we can easily see the spread of the data and any outliers. Box plots are useful for comparing the distributions of different datasets or for visualizing changes in a single dataset over time.

Overall, box plots are an effective way to display data in terms of range because they provide a concise and clear visual representation of the spread and distribution of the data.

Q9. Describe the difference between descriptive and inferential statistics. Give an example of each type of statistics and explain how they are used.

Descriptive and inferential statistics are two types of statistical analysis used to understand and draw conclusions from data.

Descriptive statistics are used to summarize and describe the characteristics of a set of data. They provide a way to organize, visualize, and summarize data in a meaningful way. Examples of descriptive statistics include measures of central tendency such as the mean, median, and mode, and measures of variability such as the range, standard deviation, and variance.

For example, if we were analyzing the ages of a group of people, we might use descriptive statistics such as the mean age, median age, and standard deviation to summarize the data and understand the distribution of ages in the group.

Inferential statistics, on the other hand, are used to make inferences or predictions about a population based on a sample of data. They involve using statistical techniques to analyze the sample data and make generalizations about the larger population. Examples of inferential statistics include hypothesis testing, confidence intervals, and regression analysis.

For example, if we were conducting a study to test the effectiveness of a new drug, we might use inferential statistics to determine whether the drug has a statistically significant effect on the target population. We might use hypothesis testing to compare the outcomes of a group of patients who received the drug with a control group that did not receive the drug, and determine whether the difference in outcomes is statistically significant.

In summary, descriptive statistics are used to summarize and describe data, while inferential statistics are used to make inferences and predictions about a larger population based on a sample of data. Both types of statistics are important in data analysis and are used to draw meaningful conclusions from data.

Q10. What are some common measures of central tendency and variability used in statistics? Explain how each measure can be used to describe a dataset.

Measures of central tendency and variability are commonly used in statistics to describe a dataset.

Measures of central tendency describe the typical value of a dataset, or the value around which the data tends to cluster. The three most common measures of central tendency are:

**Mean:** The mean is the arithmetic average of a set of data. It is calculated by adding up all the values in the dataset and dividing by the number of observations. The mean is useful for describing a dataset with a roughly symmetrical distribution.

**Median:** The median is the middle value in a dataset when the values are arranged in numerical order. The median is useful for describing a dataset with outliers or skewed data, as it is less sensitive to extreme values than the mean.

**Mode:** The mode is the value that appears most frequently in a dataset. The mode is useful for describing a dataset with discrete or categorical data, where the most common value may be more informative than the mean or median.

Measures of variability describe how spread out the data is, or how much the individual data points differ from the central tendency. The three most common measures of variability are:

**Range:** The range is the difference between the largest and smallest values in a dataset. The range provides a quick and easy way to describe the spread of a dataset, but it can be sensitive to outliers.

**Standard deviation:** The standard deviation is a measure of the spread of the data around the mean. It indicates how much the individual data points differ from the mean. A smaller standard deviation indicates that the data is tightly clustered around the mean, while a larger standard deviation indicates that the data is more spread out.

**Variance:** The variance is a measure of the variability of the data, calculated by squaring the standard deviation. Like the standard deviation, a smaller variance indicates that the data is tightly clustered around the mean, while a larger variance indicates that the data is more spread out.

In summary, measures of central tendency and variability are used in statistics to describe a dataset. The choice of measure depends on the nature of the data and the research question being investigated. Together, measures of central tendency and variability provide a comprehensive summary of a dataset and can be used to draw meaningful conclusions from the data.