

29/10/23

Statistics :

Introduction :

→ To check the quality of the data.

For this check we have to use 2 statistics.

1. Descriptive Statistics
2. Inferential Statistics.

1. Descriptive Statistics :

To make a statement or a conclusion.

2. Inferential Statistics :

It uses sample data to make inference or generalizations about a population.

What is Statistics ?

→ It is science of collecting, organizing and analyzing data. (for better decision making).

What is data ?

→ Data means facts or pieces of information that also can be measured.

→ Information contain data but data didn't contain information.

Example :

The IQ of a class Students.

{ 98, 97, 68, 57, 110 } \rightarrow Avg, Min, Max.

Descriptive Statistics :

It consist of organizing and Summarizing Data.

Inferential Statistics :

Technique where we used the data that we have measured to form conclusion.

(OR)

To make a statement / Conclusion on a Descriptive Statistics we use inferential Statistics.

Ex :

1. Are the Avg marks of the Java class Students is same as python class Students in the Besant?

A) The above statement is inferential Statistics.

2. What is the Avg Marks SQL Students?

a) Descriptive Statistics.

population (N) and Sample (n)

population (N) :

The entire group of the data we call it as a Population.

Ex: All people in India.

→ A Subset of a population we call it as a sample.

Ex:

= One lakh people from different region of India.

Key points:

- populations are larger than samples.
- Sample should be representative of the population.
- Samples allow for easier, faster and less costly data collection.

Types of Sampling techniques:

1. Simple random sampling.

Every member of a population has an equal chance of being selected for our sample.

Example:

- The Avg mileage of a bike
- What is Avg ratio of married people in Bangalore.

2. Stratified Sampling:

- Strata is nothing but a group.
- Where the population is split into non overlapping groups.

Example :

1. When the man is alive or dead.
2. The person is good ^{or} ~~but~~ bad. also.

3. Systematic Sampling :

From the population every n^{th} sample we will collect.

Example :

While doing Survey in the mall on ^{the topic of} modernization
Collecting information on every 5th person who is coming
out from the mall.

4. Convenience Sampling :

The sample is collected based on our convenience
from the particular domain experts.

30/10/25

Note : Sampling technique selection always depends on
problem statement.

Variable : A variable is the property that can take
on any value

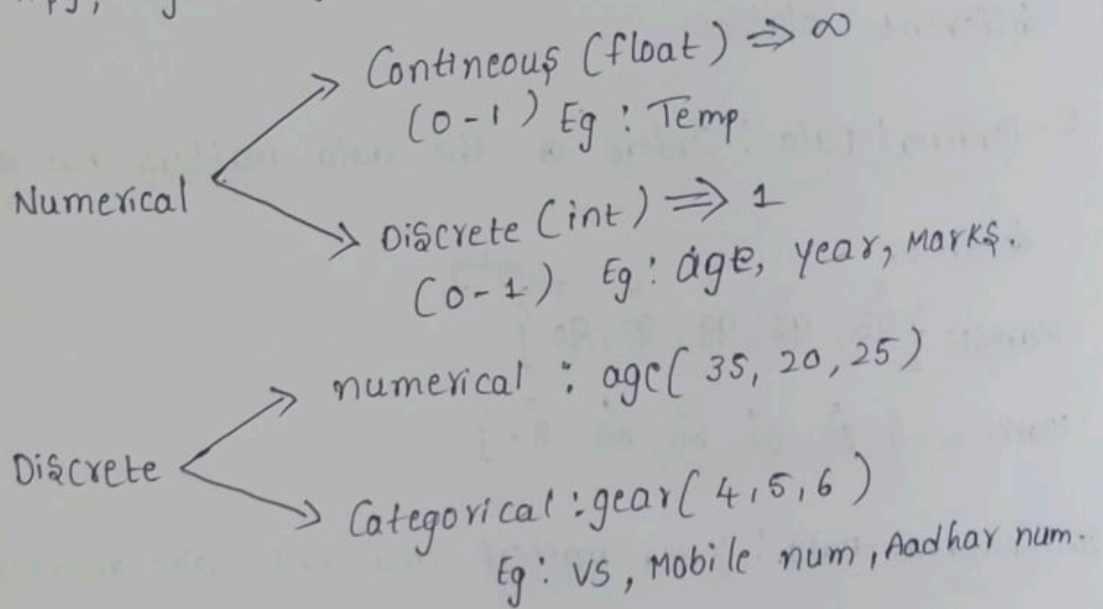
Two kinds of variables :

1. Quantitative (Numerical) variables.
2. Qualitative (Categorical) variables.

1. Quantitative variable:

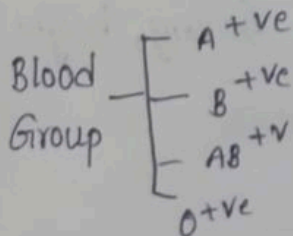
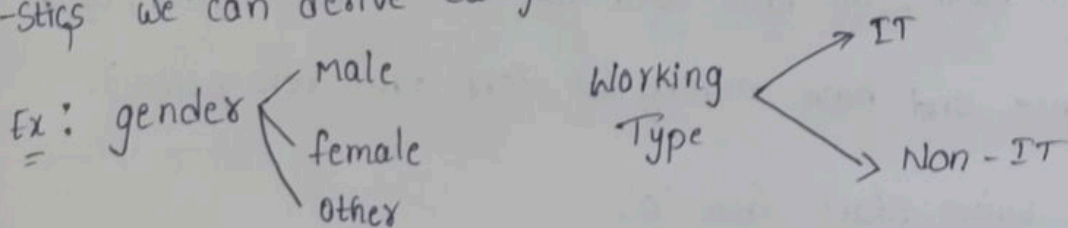
A value can be measured and we can perform mathematical operation like [Add, sub, mul, Div].

Ex: mpg, height, weight,



2. Qualitative variable:

Non-Measurable data and Based on some characteristics we can derive categorical variables.



Variable measurement scales:

4 types of measured variables.

1. Nominal Data: The categorical data which are having different classes.

2. Ordinal Data: Order of the data matters but value doesn't.

Marks₁ = { ³95, ²98, ¹99, ⁴92, ⁵90 } ^{Ranking}

Marks₂ = { ⁵80, ⁴82, ³85, ¹89, ²87 }

3. Interval Data: Order matters and value also matters but natural zero is not present.

Ex: Eye Sight, Temperature, Waves

4. Ratio Data: The ratio data can be measured, order, equal-distance and have meaningful zero (true zero point).

→ The value starts from '0'.

Ex: Speed, Rating\$, No. of Students, Salary, height, weight, age.

Descriptive Statistics:

1. Measure of Central tendency. (Mean, Median, Mode).

Mean:

1. population Mean (μ)

$$\mu = \frac{\sum_{i=1}^N x_i}{N}$$

2. Sample Mean (\bar{x})

$$\bar{x} = \frac{\sum_{i=1}^N x_i}{n}$$

Median:

* Sort the value either Ascending or descending order.

* Choose the mid value.

* If you get mid 2 value take avg of those 2 values.

→ Mean will be affected by outliers where as Median won't affect by outliers.

→ We use for null value Imputation. using Mean & Median.

Mode: Most repeatedly repetitive values.

Ex: [1, 2, 3, 1, 2, 2, 3, 4, 5]

More repeated value is 2.