

Mathematics

What is statistics?

Statistics is a branch of mathematics that involves collecting, analysing, interpreting and presenting data. It provides tools and methods to understand and make sense of large amounts of data and to draw conclusions and make decisions based on the data.

In practice, statistics is used in a wide range of fields, such as business, economics, social science, medicine and engineering. It is used to conduct research studies, analyse market trends, evaluate the effectiveness of treatments and interventions, and make forecasts and predictions.

Examples

1. Business - Data Analysis (Identifying customer behavior) and Demand Forecasting.
2. Medical - Identify efficacy of new medicines (clinical trials); Identifying risk factors for diseases (Epidemiology)
3. Government & Politics - Conducting surveys, polling
4. Environmental Science - Climate research

Statistics

Descriptive

Inferential

Data is given, and the descriptive (detail) ^{summary} are generated with the help of statistics.

Inferential statistics use to prediction.

Population VS Sample

Population refers to the entire group of individuals or objects that we are interested in studying. It is the complete set of observation that we want to make inferences about. For example, the population might be all the students in a particular school or all the cars in a particular city.

A sample, on the other hand, is a subset of the population. It is a smaller group of individuals or objects that we select from the population to study. Samples are used to estimate characteristics of the population such as the mean or the proportion with a certain attribute. For example, we might randomly select 100 students

Example

1. All cricket fans vs fans who were present in the stadium.
2. All student vs who visit college for lectures

Things to be careful about while creating samples

1. Sample Size
2. Random
3. Representative

Parameter Vs Statistics

A parameter is a characteristic of a population, while a statistics is a characteristics of a sample.

Parameters are generally unknown and are estimated using statistics. The goal of statistical inference is to use the information obtained from the sample to make inference about the population parameter.

Inferential Statistics

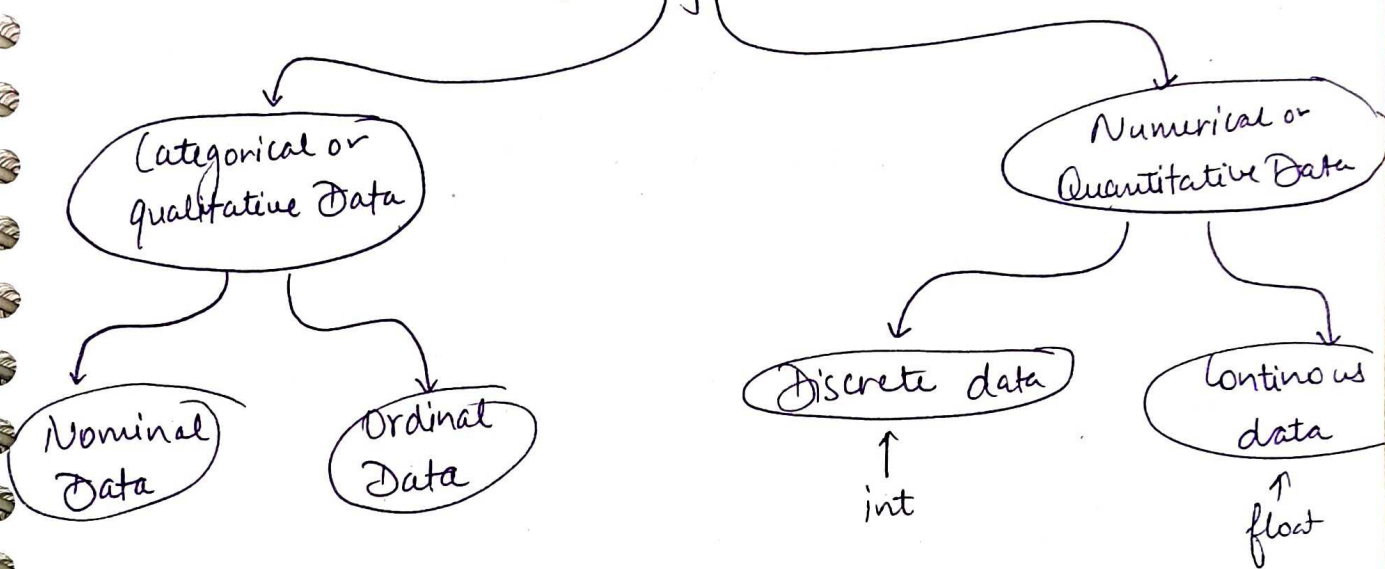
Inferential statistics is a branch of statistics that deals with making inferences or predictions about a larger population based on a sample of data. It involves using statistical techniques to test hypotheses and draw conclusion from data. Some of the topics that come under inferential statistics are:

1. Hypothesis testing: This involves testing a hypothesis about a population parameter based on a sample of data. For example, testing whether the mean height of a population is different from a given value.
2. Confidence Intervals: This involves estimating the range of values that a population parameter could take based on a sample of data. For example, estimating the population mean height within a given confidence level.
3. Analysis of Variance (ANOVA): This involves comparing means across multiple groups to determine if there are any significant differences. For example, comparing the mean height of individuals from different regions.
4. Regression analysis: This involves modelling the relationship between a dependent variable and one or more independent variables. For example, predicting the sales of a product based on advertising expenditure.
5. Chi-square tests: This involves testing the independence or association between two categorical variables. For example, testing whether gender and occupation are independent variables.

6. Sampling techniques: This involves ensuring that the sample of data is representation the population. For example, using random sampling to select individuals from a population.

7. Bayesian statistics

Types of data



Measure of Central Tendency

A measure of central tendency is a statistical measure that represents a typical or central value for a dataset. It provides a summary of the data by identifying a single value that is most representative of the dataset as a whole.

eg.

Age



Find the
Centralized

Central Tendency

Mean

Median

Mode

1. Mean

The mean is the sum of all values in the dataset divided by the number of values.

Population Mean

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

Sample Mean

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

Outlier \rightarrow create issue

2. Median

The median is the middle value in the dataset when the dataset is arranged.

3. Mode

The mode is the value that appearance most frequently in the dataset.

4. Weighted Mean

The weighted mean is the sum of the products of each value and its weight, divided by the sum of the weights. It is used to calculate a mean when the values in the datasets have different importance or frequency.

$$\boxed{\text{LR}} \quad 0.2 \rightarrow 10L$$

$$\boxed{\text{RF}} \quad 0.3 \rightarrow 15L$$

$$\boxed{\text{Xgboost}} \quad 0.5 \rightarrow 12L$$

$\underbrace{\quad\quad}_1$
weight

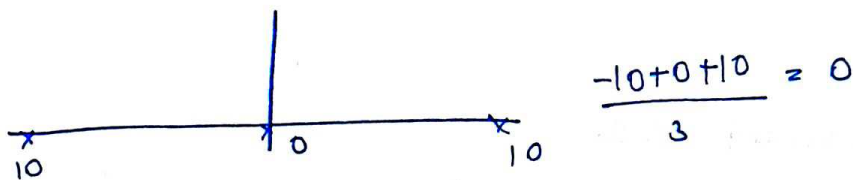
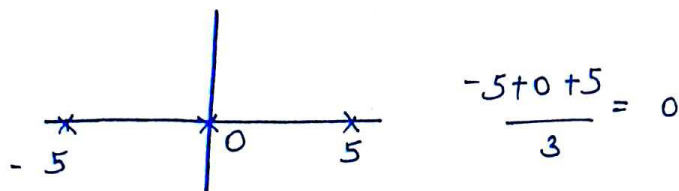
$$= \frac{0.2 \times 10 + 0.3 \times 15 + 0.5 \times 12}{0.2 + 0.3 + 0.5}$$

5. Trimmed Mean

A trimmed mean is calculated by removing a certain percentage of the smallest and largest values from the dataset and then taking the mean of the remaining values. The percentage of values removed is called the trimming percentage.

Measure of Dispersion

A measure of dispersion is a statistical measure that describes the spread or variability of a dataset. It provides information about how the data is distributed around the central tendency (mean, median or mode) of the dataset.



1. Range

The Range is the difference between the maximum and minimum values in the dataset. It is a simple measure of dispersion that is easy to calculate but can be affected by outliers.

$$\text{Range} = 5 - (-5) = 10 \text{ range}$$

$$10 - (-10) = 20$$

Range affected by outliers.

1. Range

The range is the difference between the maximum values in the dataset. It is a simple measure of dispersion that is easy to calculate but can be affected by outliers.

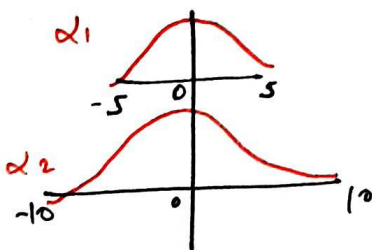
2. Variance

The variance is the average of the squared differences between each data point and the mean. It measures the average distance of each data point from the mean and is useful in comparing the dispersion of datasets with different means.

x	$x - \text{mean}$	$(x - \text{mean})^2$
3	$3 - 3$	0
2	$2 - 3$	1
1	$1 - 3$	4
5	$5 - 3$	4
4	$4 - 3$	1

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

Variance is not giving ~~exact~~ spread value but
variance is directly proportional to spread.



$$\sigma^2 = \frac{\sum (x - \mu)^2}{N} \quad \text{Population Variance}$$

$$s^2 = \frac{\sum (x - \bar{x})^2}{n - 1} \quad \text{Sample Variance}$$

we can also use (Mean Absolute Error)

$$MAD = \frac{\sum |x - \bar{x}|}{N}$$

↳ Not inference → pop data ko nikal variance nhi

↳ good for outlier

↳ Not use for finding variance

3. Standard Deviation

The standard deviation is the square roots of the variance. It is a widely used measured of dispersion that is useful in describing the shape of a distribution.

$$\sigma^2 = \frac{\sum (x - \mu)^2}{N} \quad \text{Population Variance}$$

$$s^2 = \frac{\sum (x - \bar{x})^2}{N} \quad \text{Sample Variance}$$

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}} \quad \text{Population SD}$$

$$s = \sqrt{\frac{\sum (x - \bar{x})^2}{n}} \quad \text{Sample SD}$$

Standard Deviation \rightarrow Same Unit bcz of under root

$$\text{Variance} = (5 \text{ lpa})^2 = 5 \text{ lpa}^2$$

$$\text{SD} = \sqrt{5 \text{ lpa}^2} = 5 \text{ lpa}$$

4. Coefficient of Variance

Coefficient of Variation (CV): The CV is the ratio of the standard deviation to the mean expressed as a percentage. It is used to compare the variability of datasets with different means and is commonly used in field such as biological, chemistry, and engineering.

The coefficient of variation (CV) is a statistical measure that expresses the amount of their variability in a dataset relative to the mean. It is an dimensionless quantity that is expressed as a percentage.

$$\text{CV} = \left(\frac{\text{"Standard deviation"}}{\text{mean}} \right) \times 100\%$$

Graph of Univariate Analysis

Categorical

Numerical

2. Categorical - Frequency Distribution Table & Cumulative Frequency

A frequency Distribution table is a table that summarizes the number of times (or frequency) that each value occurs in a dataset.

Let's say we have a survey of 200 people and we ask them about their favourite type of vacation, which could be one of six categories; Beach, City, Adventure, Nature, Cruise, or Other.

Types of Vacation

Frequency

Beach

60

City

40

Adventure

30

Nature

35

Cruise

20

Other

15

Bar chart

Relative Frequency is the proportion or percentage of a category in a dataset or sample. It is calculated by dividing the frequency of a category by the total number of observation in the dataset or sample.

<u>Types of Vacation</u>	<u>Frequency</u>	<u>Relative Frequency</u>
Beach	60	0.3
City	40	0.2
Adventure	30	0.15
Nature	35	0.175
Cruise	20	0.1
Other	15	0.075

Pie chart

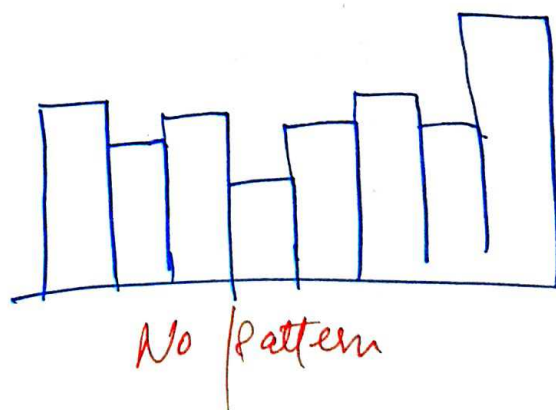
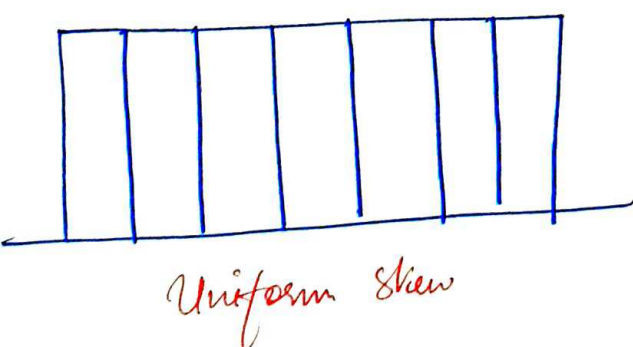
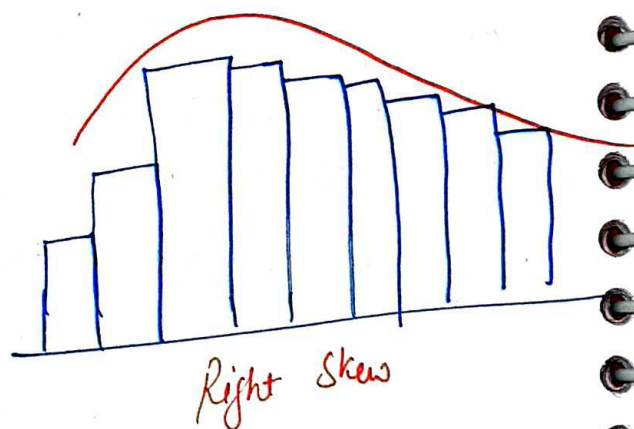
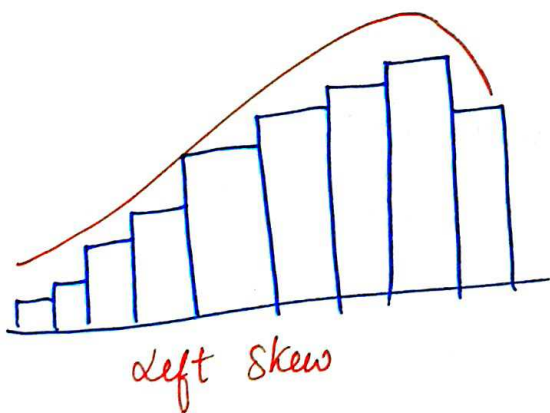
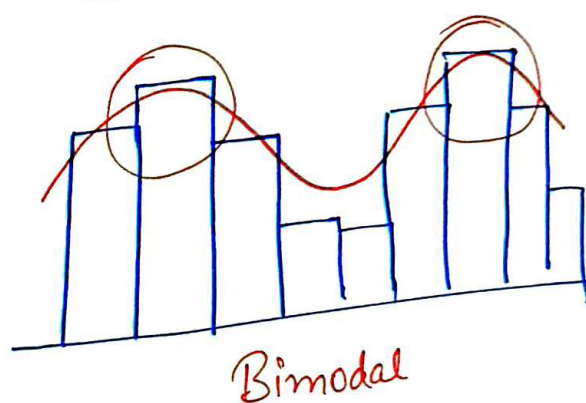
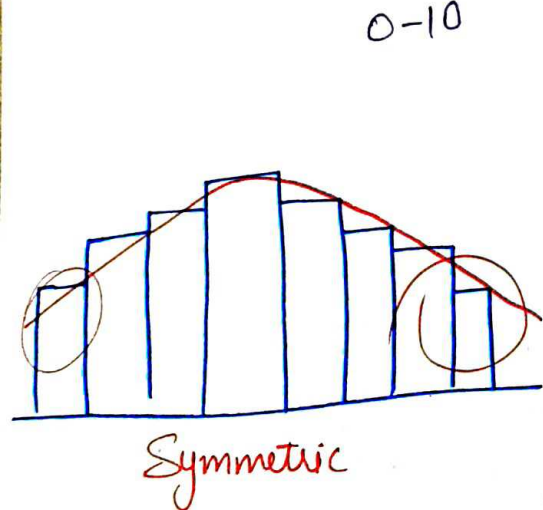
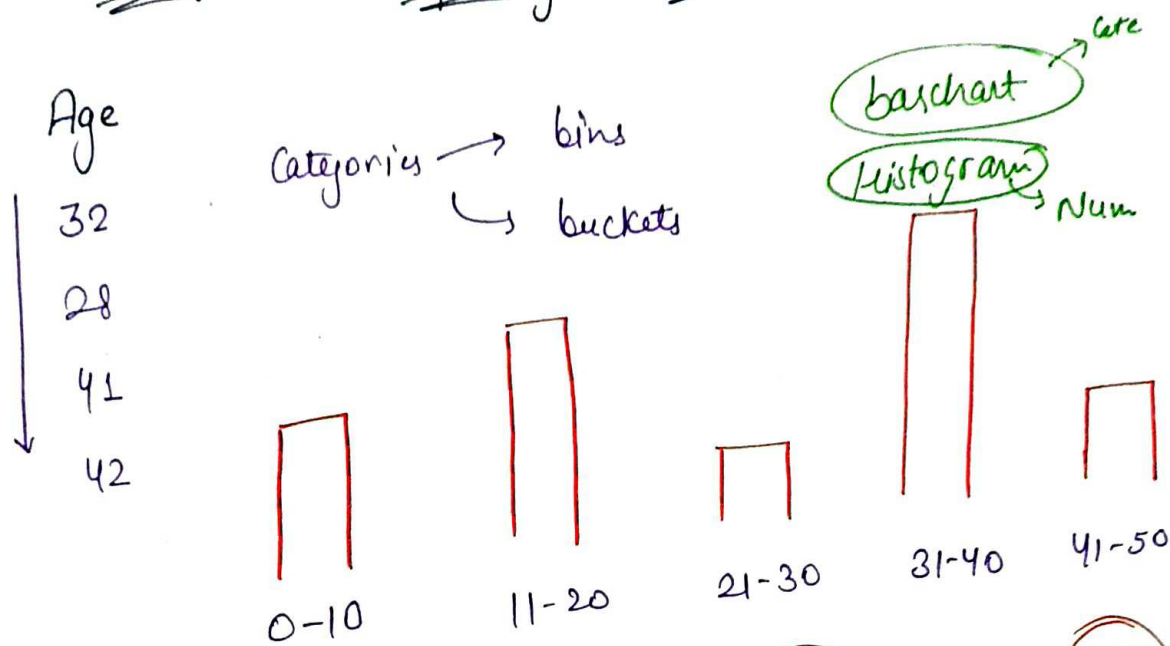
Percentage

Cumulative Frequency is the running total of frequencies of a variable or category in a dataset or sample. It is calculated by adding up the frequencies of the current category and all previous categories in the dataset or sample.

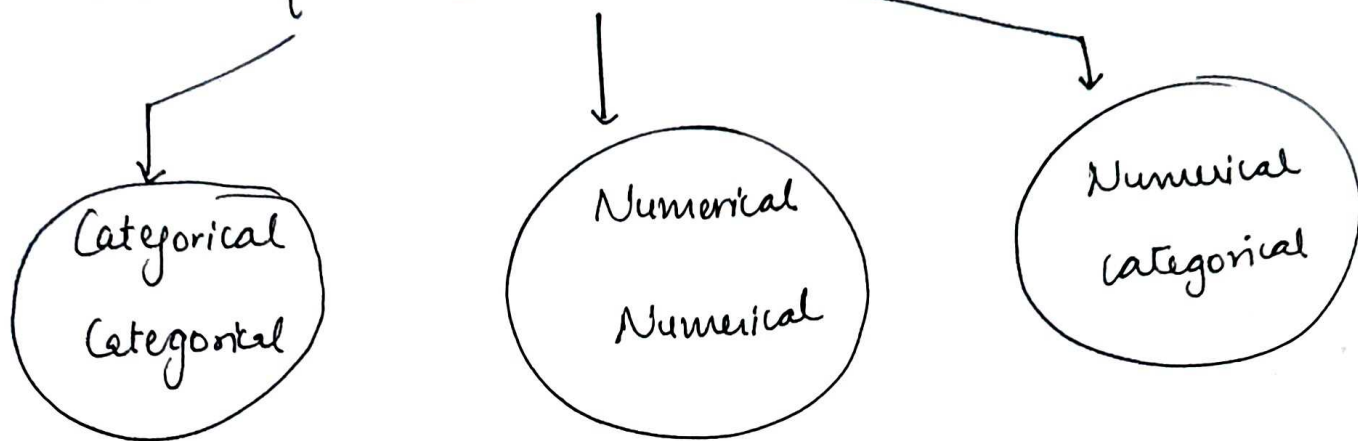
<u>Type of Vacation</u>	<u>Frequency</u>	<u>Relative Frequency</u>	<u>Cumulative Freq</u>
Beach	60	0.3	60
City	40	0.2	100
Adventure	30	0.15	130
Nature	35	0.175	165
Cruise	20	0.1	185
Other	15	0.075	200

line plot

2. Numerical - Frequency Distribution Table 2 Histogram



Graph for Bivariate Analysis



1. Categorical - Categorical

Contingency Table / Crosstab

A Contingency table, also known as a cross-tabulation or crosstab is a type of table used in statistics to summarize the relationship between two categorical variables. A contingency table displays the frequencies or relative frequencies of the observed values of the two variables, organized into rows and columns.

Survived	Pclass
0	1
1	2
:	3
:	:

⇒

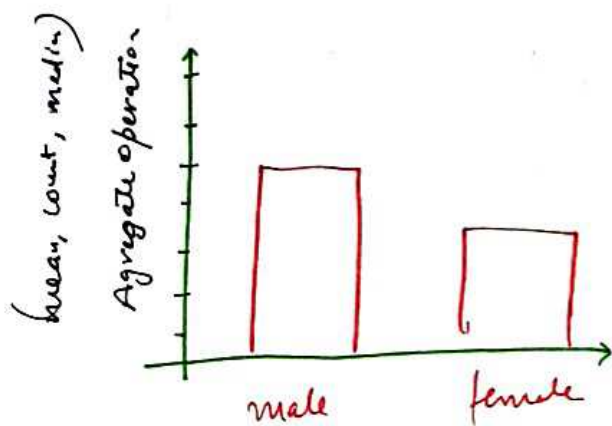
Survived	1	2	3
0	42	31	63
1	24	36	54

Contingency table

2. Numerical - Numerical

Scatter Plot

3. Categorical - Numerical



* also create Contingency table

	0-10	10-20	20-30
male	32	41	110
female	15	18	120