

# Statistical Measures

## 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.

- mean
- median
- mode.

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

$$\begin{aligned} \text{data} &= [1, 2, 3, 4, 5] & x_1 \ x_2 \ x_3 \ x_4 \ x_5 \\ & & = 1 + 2 + 3 + 4 + 5 \\ & & 5 = \cancel{5} \\ & & = \frac{15}{5} = 3 \end{aligned}$$

Sample Data  
(n)

Sample mean  
( $\bar{x}$ )

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

Population Data.  
(N)

Population mean  
( $\mu$ )

$$\mu = \frac{\sum_{i=1}^N x_i}{N} = \frac{x_1 + x_2 + x_3 + x_4 + x_5}{5}$$

- Median - The median is a middle value in the dataset when the data is arranged in order.

$$\text{data} = [6, 4, 3, 1, 5, 2]$$

$$\hookrightarrow \text{sort} = [1, 2, 3, 4, 5, 6]$$

$$\text{median} = \frac{3+4}{2} = 3.5$$

$$= [1, 2, 3, 4, 5]$$

median

- median is always Robust to outlier.

ex - classroom Age in BCA class Final year.

$$\text{Age} = [19, 21, 20, 19, 21, 22, 20, 21, 19, 20]$$

$$\text{mean} = \frac{19+21+20+19+21+22+20+21+19+20}{10}$$

$$= \frac{202}{10} \quad \boxed{\text{mean} = 20.2}$$

$$\text{median} = 19, 19, 19, 20, 20, 21, 21, 21, 22$$

$$= \frac{20+20}{2} = 20 \quad \boxed{\text{median} = 20}$$

Age = [19, 21, 20, 19, 21, 22, 20, 21, 19, 2020]

$$\text{mean} = \frac{182 + 2020}{10} = \frac{2202}{10}$$

$$\text{mean} = 220.2$$

outlier

median = 19, 19, 19, 20, 20, 21, 21, 21, 22, 2020

$$= \frac{20 + 21}{2} = 20.5$$

— Mode — The mode is the value that appears most frequently in dataset.

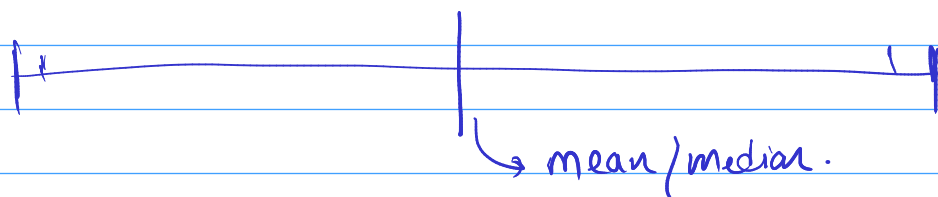
→ mode generally use for Categorical Data

data - Gender = [ male  
female  
male

Male = 60  
= Female = 40  
mode

Null  
missing  
100 student  
Female

## 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.

## Variance

**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.

Population ( $\sigma^2$ ) → sigma

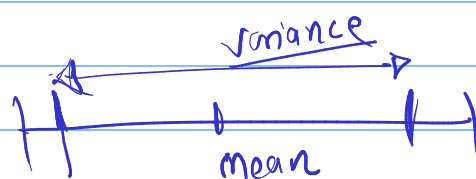
$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

Sample ( $s^2$ )

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

$$\text{km} = \text{km}^2$$

$$\sigma = \text{km}$$

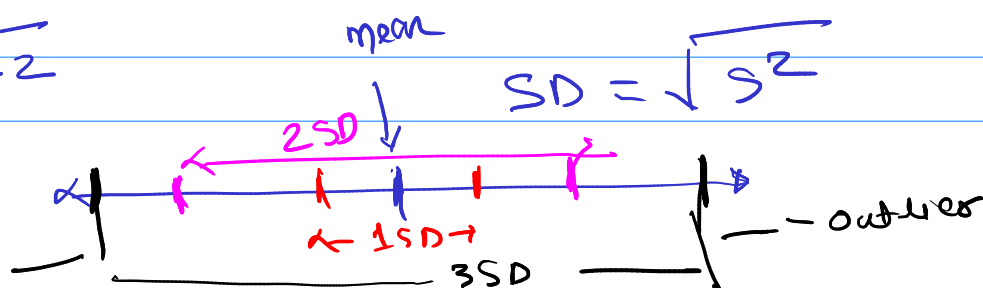


## Standard Deviation

**Standard Deviation:** The standard deviation is the square root of the variance. It is a widely used measure of dispersion that is useful in describing the shape of a distribution.

$$\sigma = \sqrt{\sigma^2}$$

$$SD = \sqrt{s^2}$$



## - Coefficient of Variation CV

Salary	Experience.
↑	↑
↓	↓

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 fields such as biology, chemistry, and engineering.

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

## Graph

Categorical

Numerical.

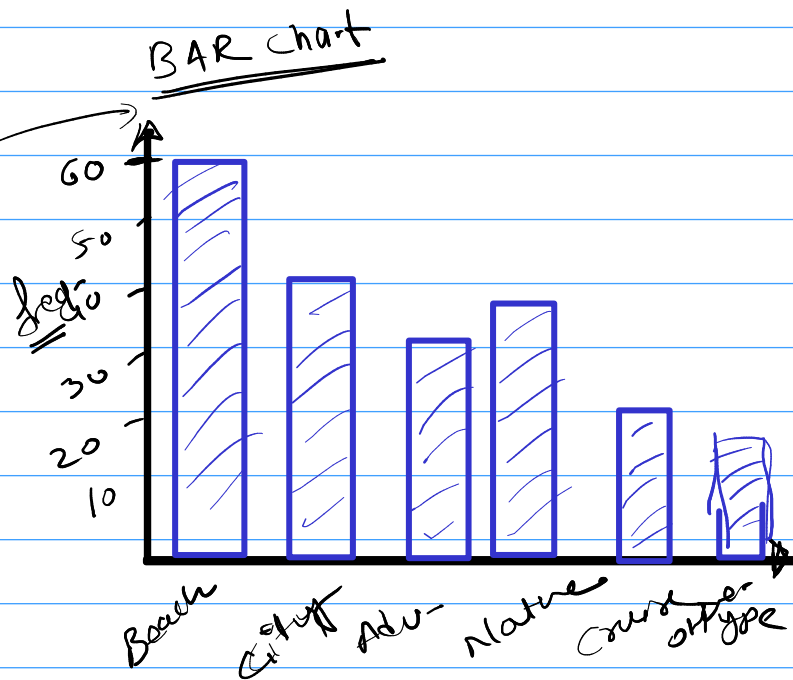
- ① Univariate Analysis - (one feature)
- ② Bivariate Analysis (2 feature)
- ③ Multivariate Analysis (All)

# Categorical Data

## Frequency Distribution Table.

ex = 200 people vacation

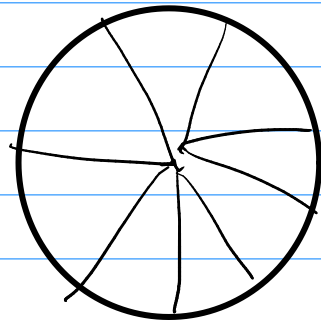
Type of vac.	Frequency
Beach	60
City	40
Adventure	30
Nature	35
Cruise	20
others	15
200	



Type of vac.	Frequency	Relative freq
Beach	60	0.3
City	40	0.2
Adventure	30	0.15
Nature	35	0.175
Cruise	20	0.1
others	15	0.075

$60/200 = 30\%$   
 $40/200 = 20\%$   
 $30/200 = 15\%$   
 $35/200 = 17.5\%$   
 $20/200 = 10\%$   
 $15/200 = 7.5\%$   
100%

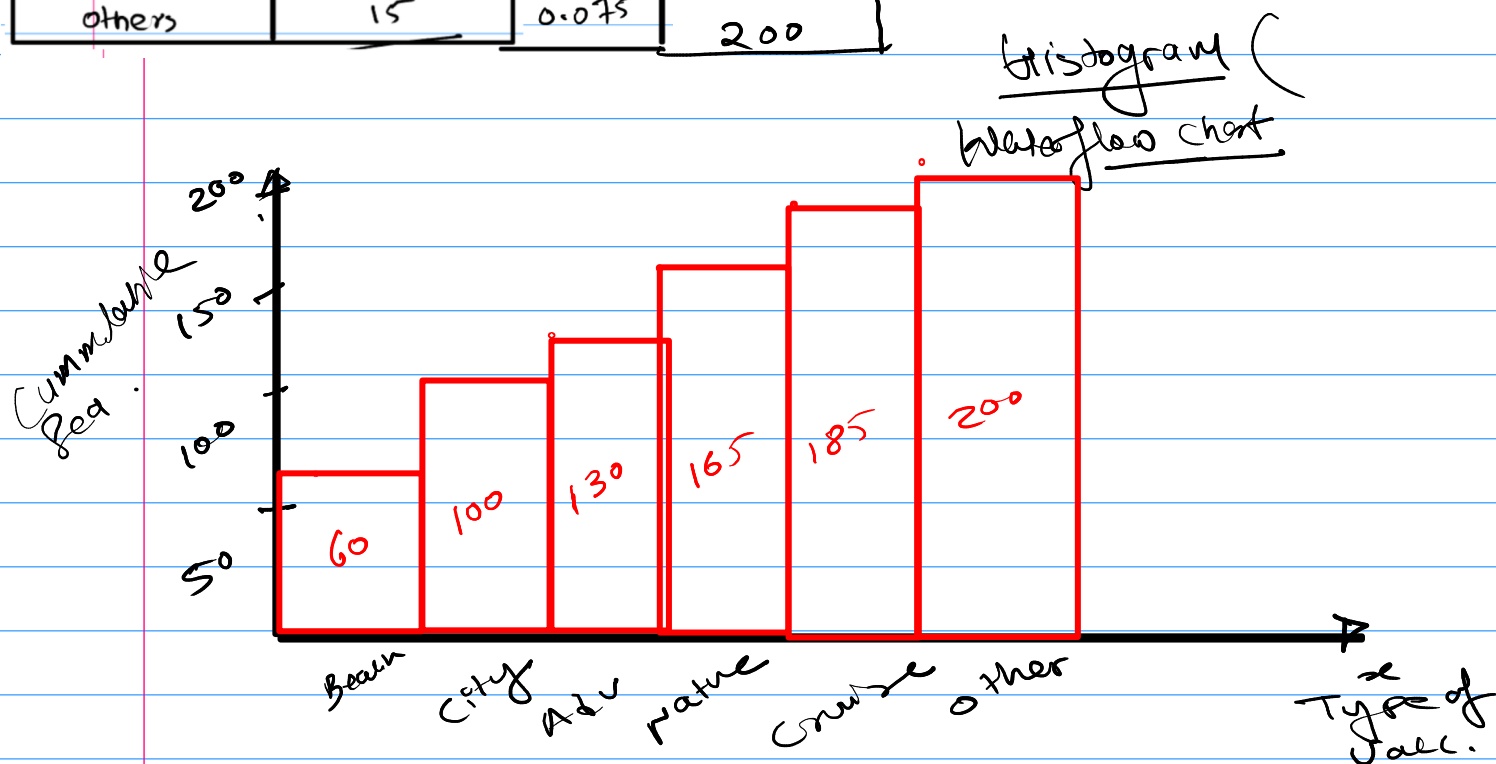
Pie



**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 observations in the dataset or sample.

**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.

Type of Jacc.	Frequency	Relative Freq	Cumulative
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
others	15	0.075	200



### Numerical Data

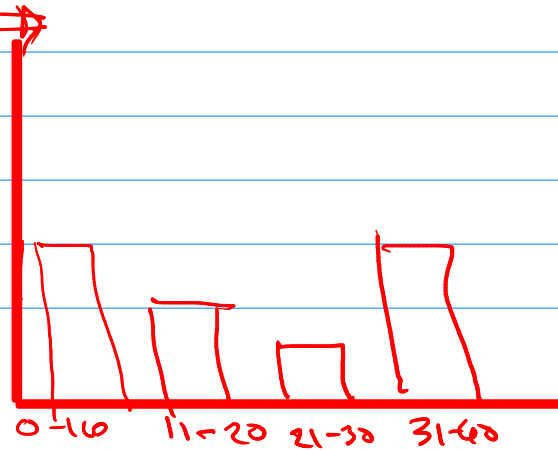
Hospital Patient Age

- 11
- 2
- 20
- 25
- 60
- 55

Age-group	Freq
0-10	—
11-20	—
21-30	—
31-40	—
41-50	—

Bins

Bar chart





Category + Category

Titantic Data = Survived Pclass

Death — 0 1

Survived — 1 2 3

Cross-Table

	class		
	1	2	3
Survived			
0	42	31	63
1	71	181	13



Numerical vs Numerical

Scatter plot

