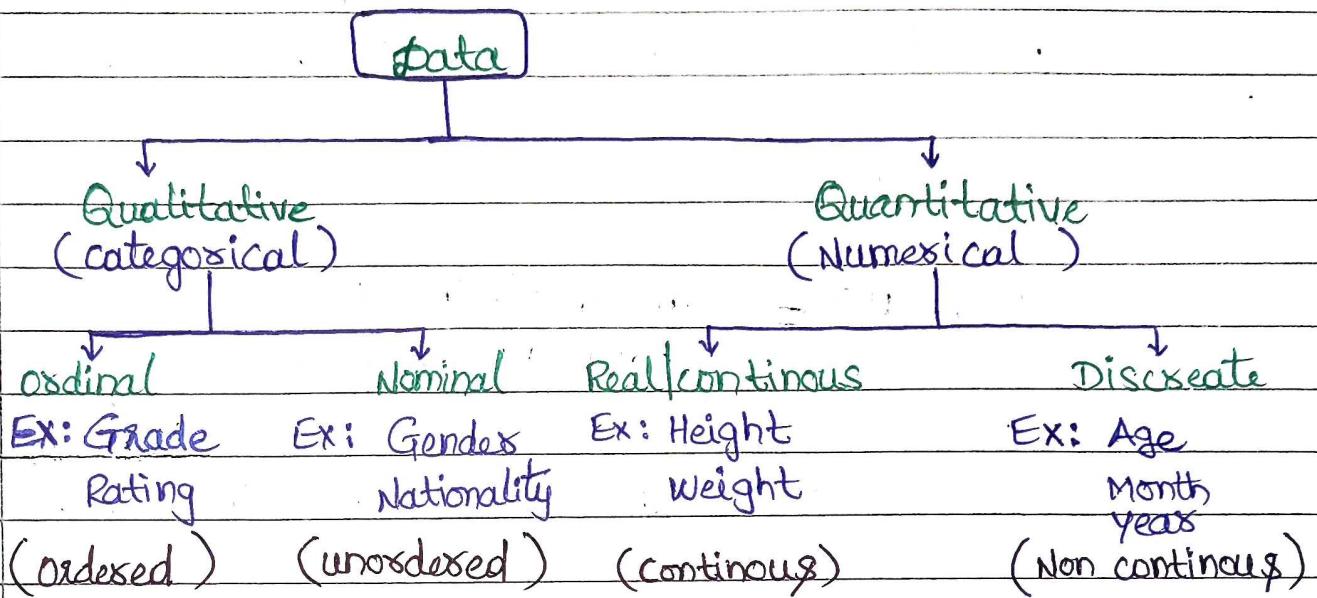


Descriptive Statistics

1 / 1

- Describe and summarizing the Data.



- * Descriptive statistics allows us to do non-visual data analysis technique for data analysis.

Case Study :

let us assume vijetha Super Market hires you to Analyze their data ?

(Why? → so that they can improve their Business .

* No matter in what format the data comes .

For Ex :- • csv , • xlsx , MySQL , PostgreSQL , MongoDB , Oracle , • json , • xml , • txt , • jpg , • png , Google sheets , AWS , GCP , Azure etc...)

* you should be able to LOAD the data in your favorite tool .

* you should be able to perform TRANSFORMATION of data so that it is ready for Analysis .

* perform Data Analysis .

Measure of central Tendency:

* These are the statistics that helps us measure / computes the central point in the data.

① Mean (Average):

$$\text{Mean} (\mu) = \frac{\text{sum of all the values}}{\text{total no. of values}}$$

② Median: It is the central value in your given dataset.

$$\text{Ex: } 10, 20, 30, 40, 50 \rightarrow \text{mean} = \frac{10+20+30+40+50}{5} = 30$$

→ Median → 10 20 30 40 50

computing the median:

① Sort the data

② Check the count (data) = n
values

- ②(a) 'n' is odd:

$$\text{median} = \left(\frac{n+1}{2} \right)^{\text{th}} \text{ term}$$

- ②(b) 'n' is Even:

$$\text{median} = \frac{\left(\frac{n}{2} \right)^{\text{th}} + \left(\frac{n}{2} + 1 \right)^{\text{th}}}{2}$$

Mean & Median allows us to find the central tendency of the data.

$$\text{Ex: } 10 \ 20 \ 40 \ 50 \ 300$$

$$\text{mean} = \frac{10+20+40+50+300}{5} = 84$$

median \Rightarrow sort data $\Rightarrow 10 \ 20 \ 40 \ 50 \ 300$

$$\text{median} = \frac{n+1}{2} = 40$$

Note : Mean is badly affected by outliers, whereas Median is somewhat it manage.

Outliers : Any extremely large or extremely small values when compared to rest of the data.

③ Mode : Applied on Categorical or Discrete Numerical columns.

↳ Most frequent value in the data.

- unimodal : Ex: 2 1 2 3 1 2 3 2 2
mode = 2

- Bimodal : Ex: 2 3 2 3 1 2 3 2 3
mode = 2, 3.

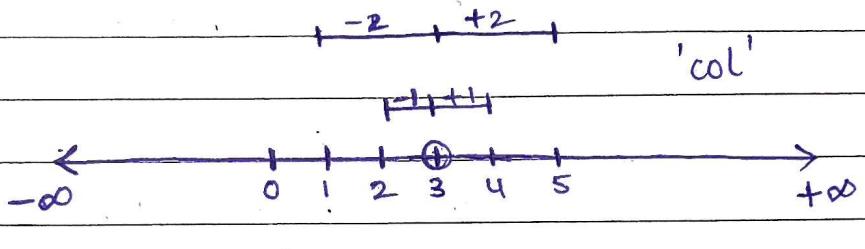
Measure of spread / Dispersion :

① Range : Range = maximum - minimum.

② Variance : It helps us measure the spread of the given data Across Mean.

Ex: col

5
4
3
2
1



* variance measures how the data is spread across mean.

- In simple terms ↴

variance measures how far away are the given data points from center of your data.

		$x_i - \text{mean}$	distance of x_i from mean
x_1	5	+2	
x_2	4	+1	
x_3	3	0	
x_4	2	-1	
x_5	1	-2	

$\sum_{i=1}^n (x_i - \text{mean})^2$ → problem with this formula
these are distances with
+ve and -ve magnitude &
they are cancelling each other.

solution :-

$$\frac{\sum \text{abs}(x_i - \text{mean})}{n}$$

- This was the not
solution

$$\begin{aligned}\sigma^2 &= \frac{\sum (x_i - \text{mean})^2}{n} \\ &= \frac{2^2 + 1^2 + (-2)^2 + (-1)^2 + 0^2}{5} \\ &= \sqrt{2} \\ \text{std} &= \sqrt{\text{Variance}}\end{aligned}$$

$x_i - \text{mean}$ → How far away is ' x_i ' from mean
In order to aggregate all distances ↴

$$\text{Variance} = \sigma^2 = \frac{\sum_{i=1}^n (x_i - \text{mean})^2}{n}$$

TO Negate the effect of negative value.

→ Standard deviation = $\sigma = \sqrt{\text{Variance}}$]

Percentiles:

- use cases : ① It helps us know the position / Rank.
② It helps you Break the data into Equal parts.

Example :

* you scored 95 percentage :

It tells how much you scored. It won't help you compare your performance with others.

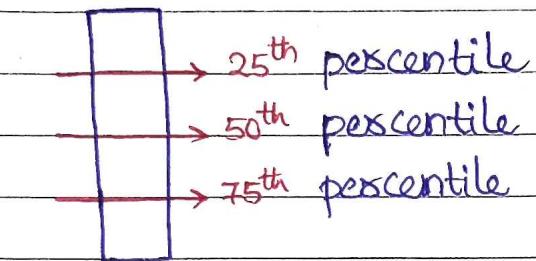
* you scored 95 percentile :

It tells out of all the participants you scored Better than 95% of participants.

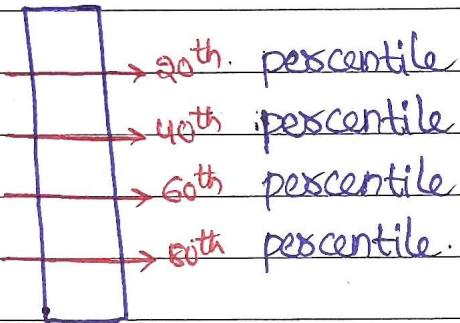
Rank / position in the data (i.e; out of all the quiz participants)

- Some popular percentiles are :

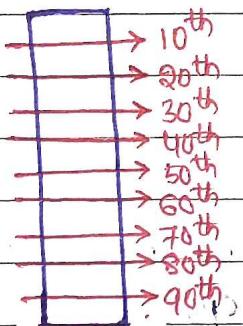
(A) Quantile (Divides the data into 4 equal parts)



(B) Quintile (Divides the data into 5 equal parts)



(C) Decile (Divides the data into 10 equal parts)



— / /

③ IQR (Inter Quartile Range)

↗ max - min
 ↗ 25th, 50th, 75th percentiles.

$$\boxed{IQR = 75^{\text{th}} \text{ percentile} - 25^{\text{th}} \text{ percentile}}$$

④ Mean Absolute Deviation (spread across mean)

$$x_1 - \mu_x$$

$$x_2 - \mu_x$$

$$x_3 - \mu_x$$

$$\vdots$$

$$x_n - \mu_x$$

Mean Aggregation

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

$$\frac{n}{\sum_{i=1}^n \text{abs}(x_i - \mu_x)}$$

Mean Square Deviation

②

Variance (σ^2)

Mean Absolute Deviation.

Note: For Example your manager asked you to calculate 'MAD' of the data.

you should always get a clarity, MAD means do you mean Median Absolute Deviation ②
Mean Absolute Deviation.

Measure of Relationship (Bi-variate Analysis)

① co-variance :

$$\text{Variance, } \sigma^2(x, x) = \frac{\sum_{i=1}^n (x_i - \mu_x)(x_i - \mu_x)}$$

$$\text{co-variance}(x, y) = \frac{\sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)}$$

Question : How do we Interpret from $\text{Cov}(x, y)$?

②

What does $\text{cov}(x, y)$ implies ?

Answer : $\text{Cov}(x, y)$ helps us understand the Relationship between the two columns.

- what kind of a Relationship ?

$x \propto y \rightarrow$ whenever $\text{cov}(x, y)$ is +ve

$x \propto \frac{1}{y} \rightarrow$ whenever $\text{cov}(x, y)$ is -ve

- Range of covariance (x, y) \rightarrow

$$-\infty \leq \text{cov}(x, y) \leq +\infty$$

- what is the meaning of $x \propto y$ ③ $x \propto \frac{1}{y}$

$x \propto y \rightarrow$ Means if 'x' increases, 'y' also increases
i.e; There is a directly proportional Relationship
between x & y .

$$0 < \text{cov}(x, y) < +\infty \rightarrow \text{cov}(x, y) \text{ is +ve}$$

$x \propto \frac{1}{y}$ → Means it 'x' increases; 'y' decreases
 i.e.; These is an inversely proportional relationship between x & y.

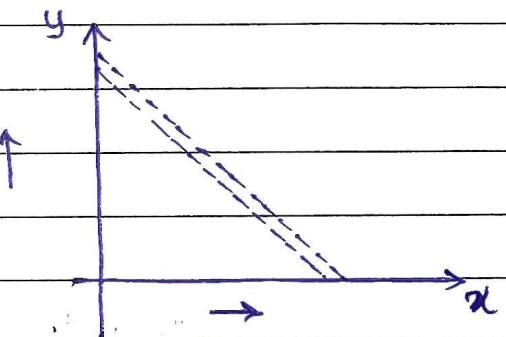
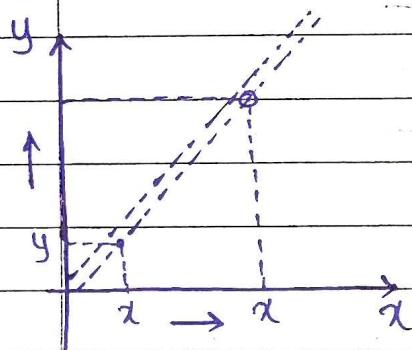
$-\infty < \text{cov}(x,y) < 0$ → cov(x,y) is Negative.

Bivariate
 (Num vs Num)

viz. Analysis

$$x \propto y$$

$$x \propto \frac{1}{y}$$



② Pearson correlation coefficient ($\rho_{x,y}$) :

$$\rho_{x,y} = \frac{\text{cov}(x,y)}{\sigma(x) * \sigma(y)}$$

$$-1 \leq \rho_{x,y} \leq +1$$

Note : { Earlier cov(x,y) was able to determine
 it $x \propto y$ or $x \propto \frac{1}{y}$ (direction of Relationship)

$r_{x,y}$ can determine even the strength of a relationship along with it $x \propto y$ or $x \propto \frac{1}{y}$

- Direction: if $r_{x,y} < 0 \Rightarrow x \propto \frac{1}{y}$
- else $r_{x,y} > 0 \Rightarrow x \propto y$
- else $r_{x,y} = 0 \Rightarrow$ No Relation

Strength of Relationship:

$0.5 \leq r_{x,y} \leq 1 \rightarrow$ strong $\leftarrow -1 \leq r_{x,y} \leq -0.5$

$-0.5 \leq r_{x,y} \leq 0.5 \rightarrow$ weak

$r_{x,y} = 0 \rightarrow$ No Relationship

0 \rightarrow No Relation
0 to 0.5 \rightarrow weak
0.5 to 1 \rightarrow strong

Note: Pearson correlation coefficient can only determine Strength & Direction of a Linear Relationship.