

(10)

MODE: The value that appears most frequently in a data set.

A set of data may have one mode, more than one mode, or no mode at all.

NOTE:

* The mode can be the same value as the mean and/or median, but this is usually not the case.

Mode is Unimodal and Bimodal.

MEASURE OF SPREAD: Also called as measure of dispersion / UNIVARIATE ANALYSIS

It is used to describe the variability in a sample or population. It is usually used in conjunction with a measure of central tendency,

such as ^{the} mean or median, to provide

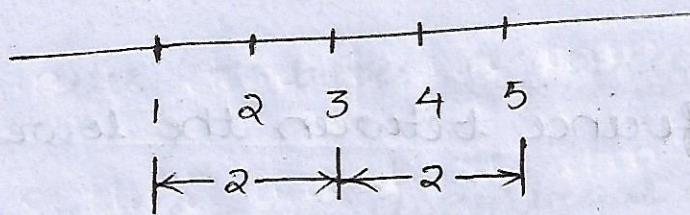
(11) an overall description of a set of data

VARIANCE: The measurement of the spread between numbers in data set.

(or)

How far the data points are from the mean.

Example:



x
$x_1 = 5$
$x_2 = 2$
$x_3 = 3$
$x_4 = 4$
$x_5 = 1$

$\rightarrow \mu_x = 3$

$$\rightarrow \mu_x - x_1 = 3 - 5 = -2$$

$$\rightarrow \mu_x - x_5 = 3 - 1 = 2$$

As the variance needs to be positive, we use:

Average.

$$\left[\frac{1}{n} \sum_{i=1}^n \right] |\mu_x - x_i| = \frac{1}{n} \sum_{i=1}^n (\mu_x - x_i)^2$$

Variance is average square distance b/w mean & data points
As the modulus value increase the

magnitude of the value, we use standard deviation.

So, The measure of spread is standard deviation. 12

STANDARD DEVIATION: The statistic that measures the dispersion of a dataset relative to its mean and is calculated as the square root of the variance.

$$\text{Std dev} = \sqrt{\text{var}}$$

RANGE: The difference between the lowest and highest values.

Example :

5	2	3	4	1
---	---	---	---	---

$$\text{Range} = 5 - 1 = 4$$

PERCENTILES: Also called as centile.

A type of quantile which divides the given probability distribution, or sample into 100 equal-sized intervals.

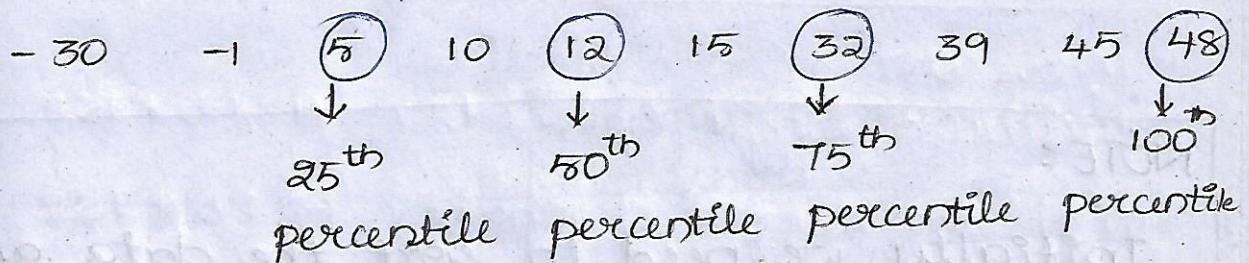
↳ This allows the data to be analyzed in the

terms of percentages. (B)

Example:

I SCORED $\underbrace{48/100}$ IN MATHEMATICS

↓
100 percentile - gives the rank.



There are 3 types of percentiles. They are

1. Quartile - It describes a division of observations into four defined intervals based on the values of the data and how they compare to the entire set of observations. (25%, 50%, 75%)

2. Quintile - It describes a division of observations into five defined intervals based on the values of the data as

20%, 40%, 60%, 80%.

3. Decile - It describes a division of \textcircled{W} observations into ten defined intervals based on the values of the data as 10%, 20%, 30%, 40%, 50%, 60%, 70%, 80%, 90%.

NOTE:

Initially we need to sort the data and then the percentile is taken.

IQR : INTERQUARTILE RANGE

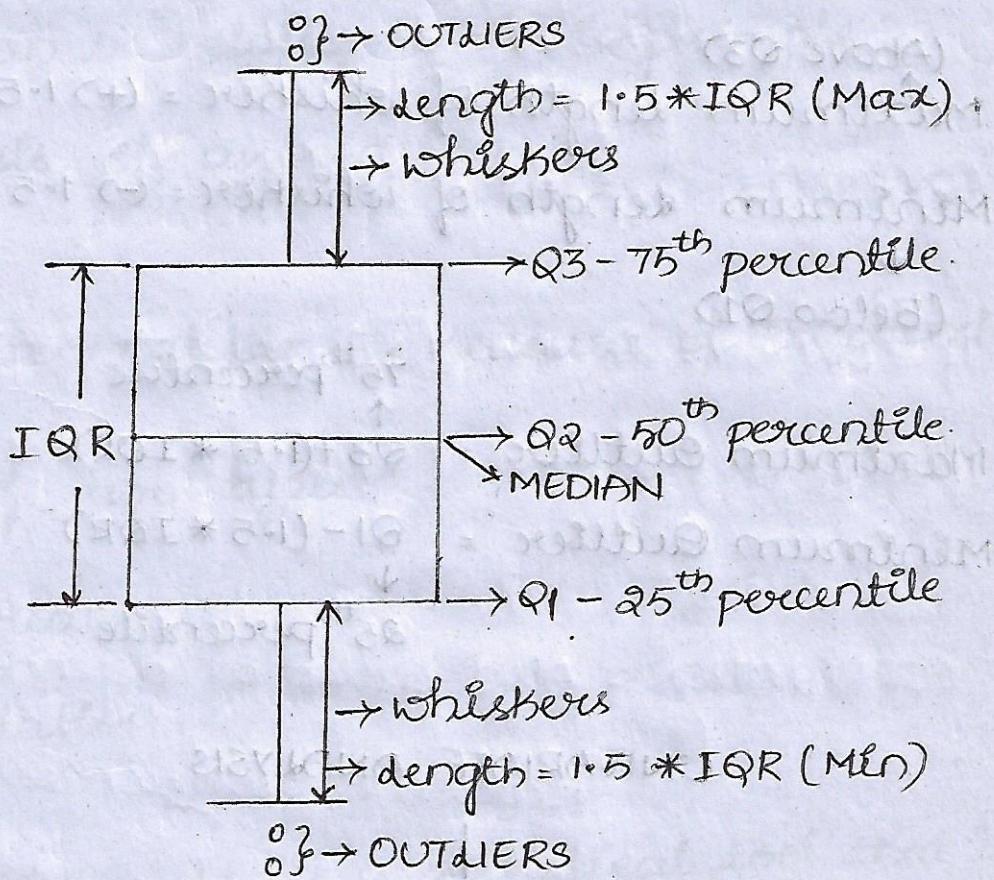
It is a measure of variability, based on dividing a data set into quartiles.

↳ The quartiles divide a rank-ordered data set into four equal parts.

↳ The values that divide each part are called the first, second and third quartiles; and they are denoted by Q_1, Q_2, Q_3 respectively.

(15)

BOX PLOT : Also known as box and whisker plot. It is a type of chart often used in EDA to visually show the distribution of numerical data and skewness through displaying the data quartiles (percentiles) and averages.



In statistical theory;

$$\text{length of whiskers} = 1.5 * \text{IQR}$$

OUTLIERS : It is defined as a data point that is located outside the whiskers of the box plot.

→ An outlier is an observation that is numerically distant from the rest of the data. (16)

Outliers → Extremely small values or large values.

(Above Q3)

Maximum length of whisker = (+) $1.5 * IQR$

Minimum length of whisker = (-) $1.5 * IQR$

(Below Q1)

75th percentile

↑
Maximum Outlier = $Q3 + (1.5 * IQR)$

↓
Minimum Outlier = $Q1 - (1.5 * IQR)$

25th percentile.

UNIVARIATE ANALYSIS

MEASURE OF SPREAD

- VARIANCE - Impacts with outliers
- STD. DEVIATION - "
- RANGE - Impacts with outliers
- IQR - Doesn't impacts with outliers

MEASURE OF CENTRAL TENDENCY

- MEAN - Impacts with outliers
- MEDIAN - Doesn't Impacts with outliers
- MODE - "

(17) **BIVARIATE ANALYSIS**: It is one of the simplest form of quantitative analysis. It involves the analysis of two variables, for the purpose of determining the empirical relationship between them.

↳ This can be helpful in testing simple hypothesis of association.

Bivariate Analysis → Measure of Relationships

It is of two types.

1. Covariance.

2. Correlation.

1. COVARIANCE: It is a statistical tool that is used to determine the relationship between the movement of two variables.

↳ When two stocks tend to move together, they are seen as having a positive covariance.

i.e., if $x \& y$ is taken then $x \propto y$

→ When the two stocks tend to move inversely, then the covariance is negative i.e., if $X \propto Y$ are taken then $X \propto \frac{1}{Y}$

SIGN MATTERS for covariance. \Rightarrow



$\text{cov}(X, Y) = \pm \text{values.}$

$$\text{cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (\mu_X - X_i)(\mu_Y - Y_i)$$

If $\text{cov}(X, Y)$ is +ve. ($X \uparrow, Y \uparrow$) $\rightarrow X \propto Y$

If $\text{cov}(X, Y)$ is -ve ($X \uparrow, Y \downarrow$) $\rightarrow X \propto \frac{1}{Y}$

2. CORRELATION : Also called as dependence or
(β)

Pearson correlation coefficient.

In statistics, correlation is any statistical relationship, whether causal or not, between two random variables or bivariate data.

→ When the two sets of data are strongly linked together we say that they have a

HIGH CORRELATION.

(19)

↳ When the correlation is positive, when the values increases together.

↳ The correlation is negative, when one value decreases as the other increases.

The word correlation is made of :

co-(together) and relation.

$$\text{correlation } \rho_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x * \sigma_y} \text{ and lies between } -1 \leq \rho_{xy} \leq 1$$

If $\rho = -1$ ($x\uparrow, y\downarrow$) - All points lies on same line.

If $\rho = 0$ - All points lies in a random manner

If $\rho = 1$ ($x\uparrow, y\uparrow$) - All points lies on same line.

$$-1 < \rho < 0$$

Negative Relationship

and variance is a lot
between the points

$$0 < \rho < 1$$

Positive Relationship

and variance is a
lot between the points

NOTE:

(20)

correlation can only detect the linear relationships.

DRAWBACK : Can't detect the patterns apart from the linear relationships