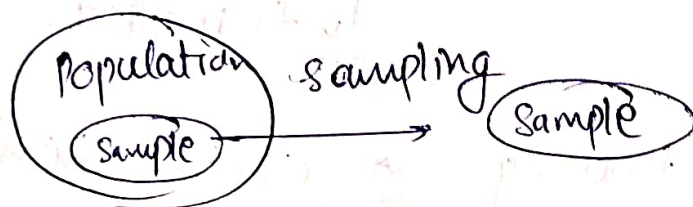


Data Science

Population and sample ÷ sample is subset of population.

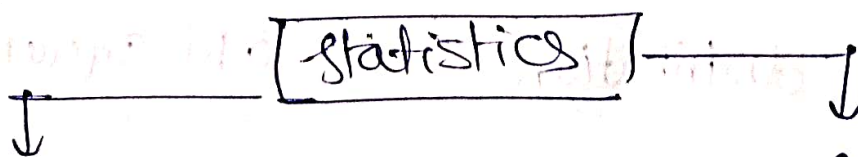
Population ÷ A population includes all the elements from a set of data. The population is whole set of values or individuals.

Sample ÷ The sample is a subset of population and is set of value you actually use in your estimation.



STATISTICS ÷

Statistics is a set of mathematical methods and tools that enables us to answer



Descriptive

Inferential.

Descriptive statistics ÷ this offers method to summarized data by transforming raw observation to meaningful information.

Inferential Statistics ÷ this offers method to study experiment done on small sample of chalk and the inference to entities.

Why Inferential?

- ↳ ① → suitable for large amount
- ↳ ② → sample data \approx inferential stats

Descriptive

- * measure of central tendency
- * measure of variability
- * measure of shapes
- * Percentiles
- * frequency distribution
- * Covariance and correlation

Inferential

- * Central Limit Theorem
- * hypothesis testing
 - z-test
 - t-test
- * Chi square test

measure of central tendency:

Mean / Avg \div $\frac{\text{sum of all data}}{\text{total no. of data units}}$

// to read a dataset \div

dataset = pd.read_csv("titanic.csv")

Median \div (name/data) \rightarrow sort \rightarrow middle values

Mode \div repeating itself (Most) \rightarrow Mode

Range \div Range is difference b/w maximum and minimum value in dataset.

$$\text{Range} = \text{Maximum} - \text{Minimum}$$

Mean absolute deviation:

The mean absolute deviation \div

$$\boxed{\text{MAD} = \frac{\sum |x_i - \bar{x}|}{n}}$$

\hookrightarrow Idea about variability in dataset

Standard deviation:

The standard deviation is a measure of the amount of variation or dispersion of a set values. A low standard deviation tend close to Mean

$$\sigma = \sqrt{\frac{\sum (x_i - M)^2}{N}}$$

variance:

variance is measure of

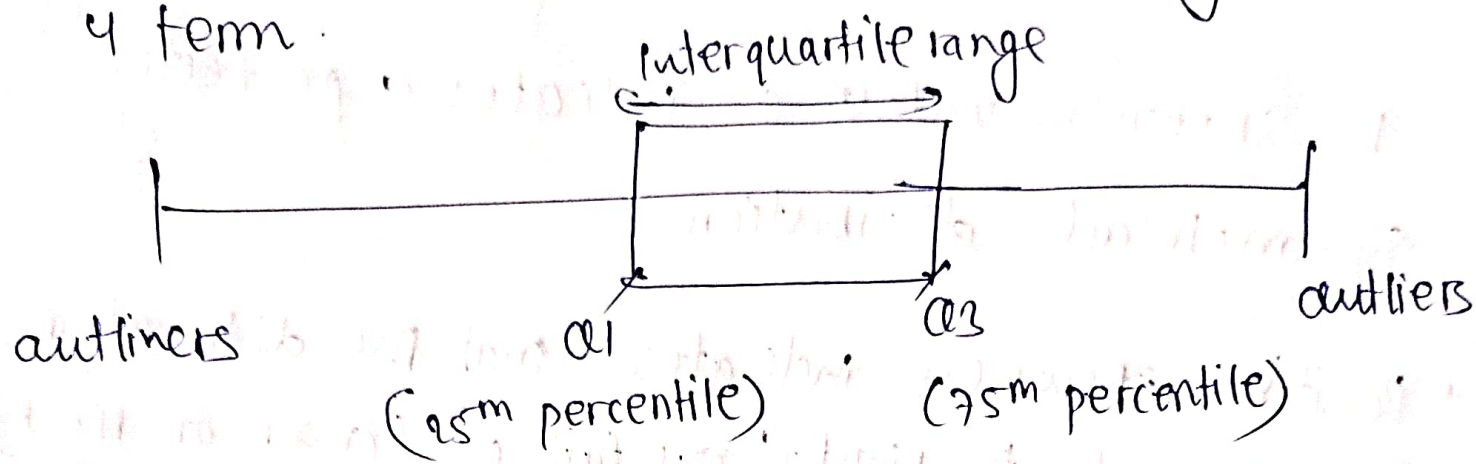
how data point diff from Mean.

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{N}$$

Acc to layman a variance is measure of how far a set of data are spread out.

Percentile ÷ Percentiles are used in statistics to gives you a number that describe the value that a given percent of the values are lower than.

Quartile \rightarrow statistic team dividing into 4 term.



Inter quartile range \div to find outliers.

$$\rightarrow IQR = Q_3 - Q_1$$

Q_1 (first quartile) \div value below which 25% data

Q_3 (third quartile) \div value below which 75% data

Outliers are defined data points fall below

$$Q_1 - 1.5 * IQR \text{ or above } Q_3 + 1.5 * IQR$$

Skewness \div how much data is spreaded.

↳ Skewness is how much asymmetric distribution

↳ A skewness value 0 indicates a perfectly symmetrical distribution

↳ Positive skewness indicates that the distribution is skewed to right (tail is longer on right)
negative skewness (tail is longer on left)

$$\text{Skewness} = \frac{\sum (x_i - \bar{x})^3}{(N-1)\sigma^3}$$

Frequency and cumulative Distribution,

A frequency distribution is a table that shows the numbers of occurrences of different value in data set.

A cumulative distribution shows the accumulation of frequency up to certain point.

SKWENESS

