⊙ EDA

1. Measure of Central Tendency
    a. Mean
    b. Median
    c. Mode

Central Tendency: The property of data being concentrated in the centre

    a 1.1 km
    b 1.8 km
    c. 1.5 km
    d. 1.1 km
    e. 15 km

(here the ⑮ is outliers we exclude to avoid wrong mean)

∴ If there is no outliers then we use mean

Mean: It is the average of all numbers and is sometimes called as Arithmetic mean

Median: The statistical median is the middle number in a sequence of numbers.

Mode: The mode is the number that occur more with a set of number
    • Discrete numerical data is go with mode.
    • Used for categorical data mostly

— Mode Example:

$$3 \quad 7 \quad \underline{5} \quad \underline{25}$$

Mode = 5
(Most Frequent)

---

⊙ Measure of Spread / Data Variability

- Range.
  - The difference b/w the highest and lowest value within a set of numbers

- Interquartile range (IQR)
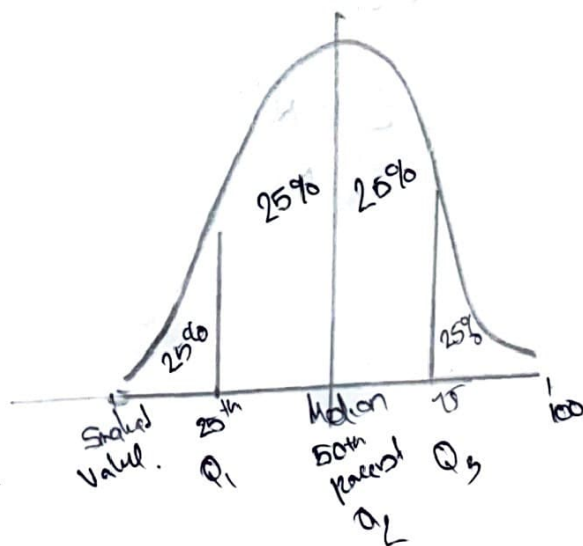  - The interquartile data is the middle half of the data

$Q_1 - 25\%$

$Q_2 - 50\%$

$Q_3 - 75\%$

∴ $Q_2$ is the median

$$1 \; 2 \; 3 \; 4 \; \boxed{5} \; 6 \; 7 \; 8 \; 9$$

2.5 → $Q_2$ → 7.5

$Q_1$ → $Q_3$

Figure:



$IQR = Q_3 - Q_1$

$= 7.5 - 2.5$

$= \dfrac{5}{2}$

∴ $Q_2$ is not always equal to IQR

25% 25%

25%

Smallest Value.    25th
                   $Q_1$    50th
                          percent    $Q_3$
                   Median    $Q_2$    75
                   
100

another example:

1  2  3  4  5  6  7  8  9  10
   └─┘           └─┘        └─┘
   2·5           6·5        8·5

$$IQR = 8.5 - 2.5 = \frac{6}{2}$$

## – Standard Deviation ($\sigma$)

The standard deviation is the measure that is used to quantify the amount of variations or dispersa of a set of data values

standard deviation = $\sqrt{\text{Variance}}$
($\sigma$)

· square root of variance is "SD"

Population

$$\sigma = \sqrt{\frac{\Sigma(x-\mu)^2}{N}}$$

Sample.

$$S = \sqrt{\frac{\Sigma(x-\bar{x})^2}{n-1}}$$

## – Variance

The variance is a measure of how far a set of data are dispersed out from their mean/average value

· The '$\sigma^2$' is Variance.

$$\sigma^2 = \frac{\Sigma(x-\mu)^2}{N} \qquad S^2 = \frac{\Sigma(x-\mu)^2}{N-1}$$

# Example

| $x$ | $\begin{array}{c}x-M\\ \text{or}\\ x-\bar{x}\end{array}$ | $(x-M)^2$ |
|---|---|---|
| 1 | $1-3=-2$ | 4 |
| 2 | $2-3=-1$ | 1 |
| 3 | $0$ | $0$ |
| 4 | $1$ | 1 |
| 5 | $2$ | 4 |

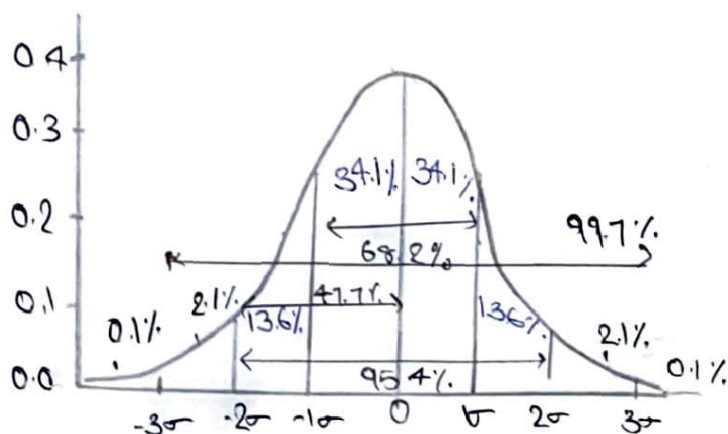$$\text{Variance} = \frac{\sum(x-M)^2}{M} = \frac{10}{5} = 2$$

$$\bar{x} = \frac{1+2+3+4+5}{5} = \frac{15}{5} = \frac{3}{2}$$

Variance $= 2$

and

$$SD = \sqrt{\text{Variance}}$$
$$= \sqrt{2} = \frac{1.414}{2}$$
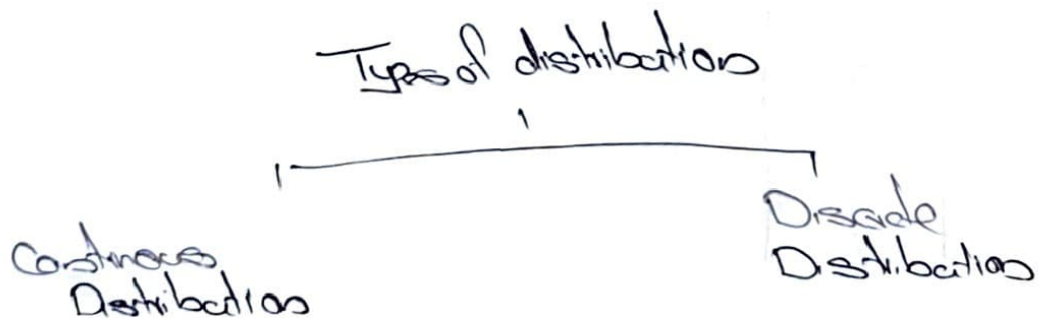
- Percentile / "Centile"

Fig.



A percentile is a measure used in statistics indicating the value below which a given percentage of observations in a group of observations fall.

⊙ Distributions

The graphical representation of all observations is also known as distribution.

Types of distribution

Continuous Distribution

Discrete Distribution

- Normal Distribution

Normal distribution, also known as the Gaussian df. It is a probability distribution that is symmetric about the mean.

· Which follow $3\sigma$ rule
· Most prefered in the pipeline
· Bell shaped Curve.

- Properties of Normal/Gausion dist
1 · Empirical rule
2 · Distortion is Normal dist
3 · Central limit theorem.
4 · Standard Normal dist
5 · Outliers
6 · QQ plot
7 · Log, sqrt, Boxcox transformation

# 1. Empirical rule

- Figure:

<u>Same figures as the previous percentile</u>

. The empirical rule states that for a normal distribution nearly all of the data will fall within three standard deviations of the mean. The empirical rule can be broken down into three parts:

- 68% data fall within the 1st SD from the mean
- 95% fall within 2 SD
- 99.7% fall within 3 SD
- Any point lying after 3 sigma is outlier.

# 2. Distortion in Normal Distribution

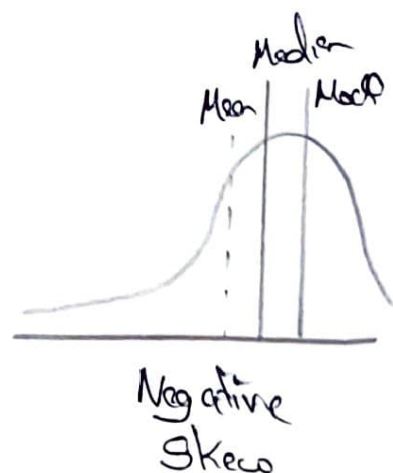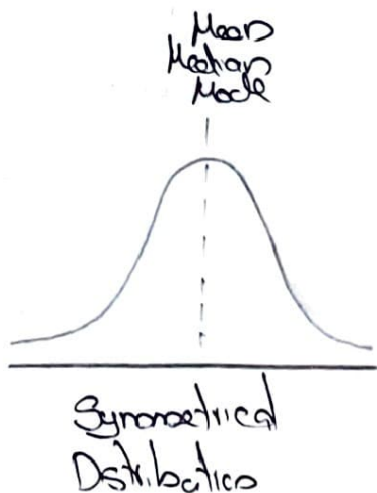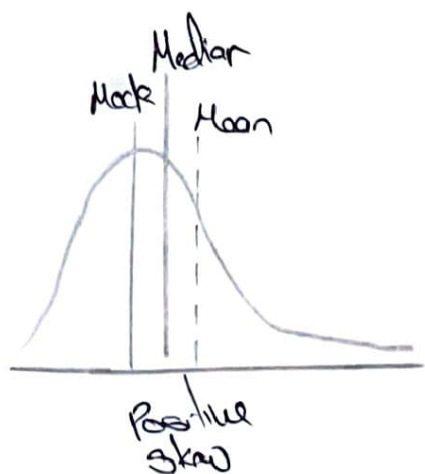The distortion is normally distributed curves can be quantified in 2 ways.
1. Skewness
2. kurtosis

# 1. Skewness in Normal Distribution

Skewness is asymmetry in a statistical distribution, in which the curve appears distorted or skewed either to the left or to the right. Skewness can be quantified to define the extent to which a distribution differs from a normal distribution



Positive Skew

Symmetrical Distribution

Negative Skew

- How much skewness and kurtosis
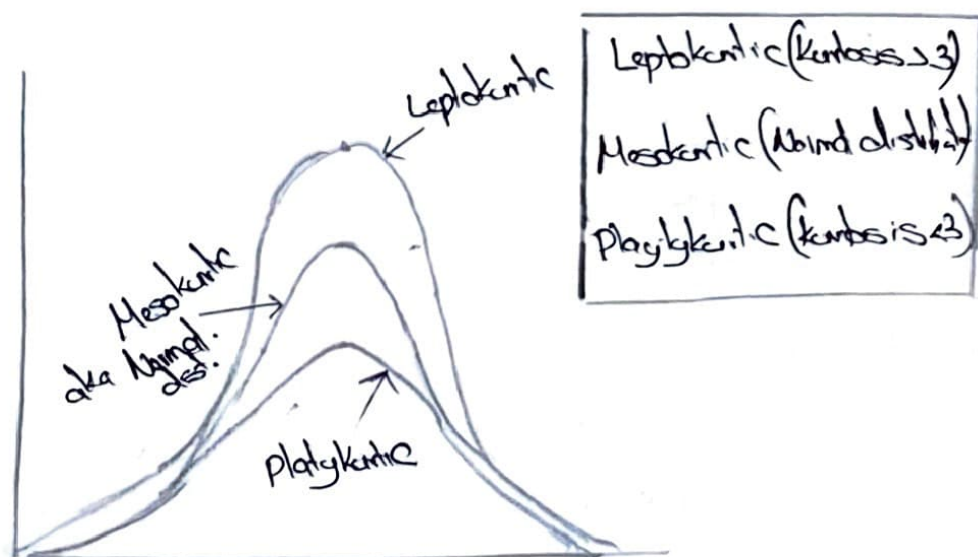    - If the skewness is between -0.5 and 0.5 → fairly Symm
    - If -1 and -0.5 or 0.5 and 1 → moderately skewed
    - If skewness is > 1 or < -1 → highly skewed
    - A standard normal distribution has 3 kurtosis.

- Kurtosis.

    In probability theory and statistics, kurtosis is a measure of the "peakedness" of the Probability distribution of a real-valued random variable



| Leptokurtic (kurtosis > 3) |
| Mesokurtic (Normal distribution) |
| Platykurtic (kurtosis < 3) |

- Which is Best - the Mean, Median, or Mode?
    When you have a symmetrical distribution for continuous data, the mean, median and mode are equal. In this case, analysts tend to use the mean because it includes all of the data in the calculations. However, if you have a skewed distribution, the median is often the best measure of central tendency

· When you have categorical or discrete data, the median or mode is usually the best choice

For categorical data, you have to use the mode

Preferable:

Normal dist ⟶ Mean

Skewed dist ⟶ Median

Categorical dist ⟶ Mode

- Corrected Distortion is Normal Distribution

Transformation is nothing but taking a mathematical function and applying it to the data.

Log Transformation [Each data point is replaced with $\log(x)$ to obtain ND]

Square-root Transformation [Each data point is replaced by it square root]

Reciprocal Transformation [It takes the reverse of x ie, $1/x$]

Box-Cox Transformation [Transformation of non-normal dependent Variable to normal shape]

Reason: To transform the data to either reduce the skewness or to normalize the data or simply make the data easier to understand
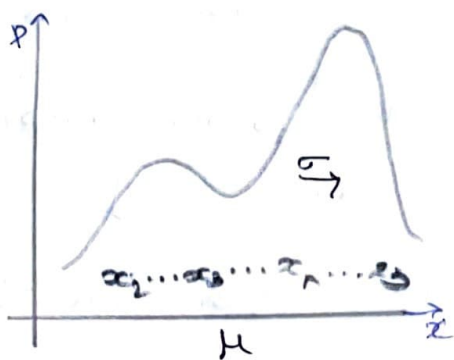
- Outliers.

An outlier is an observation point that is distant from other observations. An outlier may be due to variability in the measurement or it may indicate experimental error; the data are sometimes excluded from the dataset.

⊙ Central limit theorem

The central theorem states that the distribution of sample means approximates a normal distribution as the sample size get larger, regardless of population distribution shape
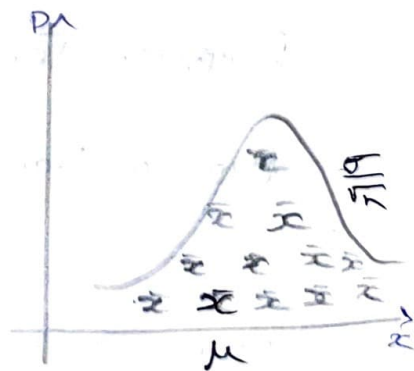
• CLT in one sentence "Even if it's not normal, the avg is normal"

• When we are collecting mean's the samples from any dist the no of samples taken for calculating the mean should be greater or equal to 30

• Rule of thumb: the sample should be bigger than 25 observations or ("30")



Population distribution



Samples of size n
x̄



sample distribution of the mean.