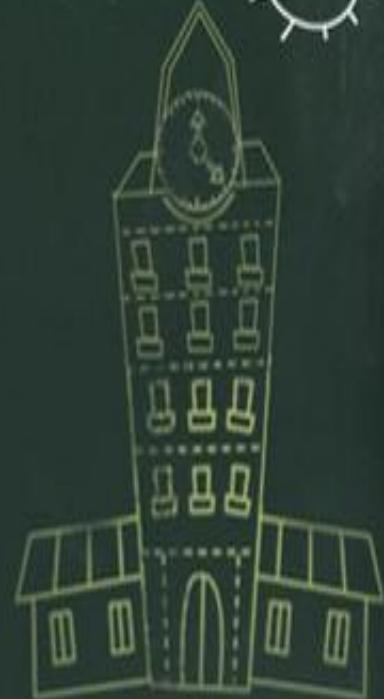


understanding of statistics for Data Science

Dr. Sheetal Dhande-Dandge

CONTENTS

- | | | | |
|----|---------------------------|----|----------------------|
| 1. | Introduction f statistics | 5. | Correlation |
| 2. | Sampling | 6. | Normal Distribustion |
| 3. | Central tendancy | 7. | Empirical Rule |
| 4. | Variation | 8. | Z-Score |



01

What is Statistics

Statistics is branch of mathematics that deals with collection, analyzing and interpreting large amount of data

01

Why Statistics is important

Statistics allows us to drive knowledge from large dataset and this knowledge can then used to make prediction ,decision classification,etc

01

Where Statistics is used

Statistics used in various field, some of them are

1. Medical Research
2. Stock Market
3. Sales Projection
4. Weather forecasting

02

Sampling

Sampling is the process of collecting data to perform analysis on

02

Sampling

Sample Vs Population

Sampling Frame is List From (Population) Which Sample is Selected.

02

Sampling Error

Sample Error is an error that leads to our Sample not accurately representing our population

02

Random Sampling

Random Sampling is process of selecting subset/sample from a population in such a way that every data point is equally likely to be included in sample

02

Systematic Sampling

Systematic Sampling is process of selecting your sample by picking every k th element in your population , you don't need to list this.

03

Central Tendency

Is used to indicate where does the middle or center of distribution of our data lies

- Mean
- Mode
- Median

03

Central Tendency: MODE

is used to indicate the most frequent data point , in other words , the one which occur most number of times.

03

Central Tendency: Median

is the middle of data , if the dataset arrange is asending order then the data element which occurs right at center , is the median

03

Central Tendency: Mean

is average of data ,in simpler terms its sum of values divided by total number of values , its represented by Greek letter 'sigma'

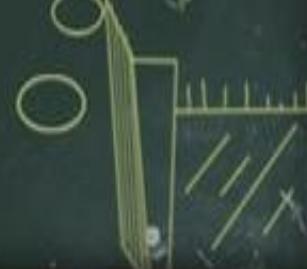
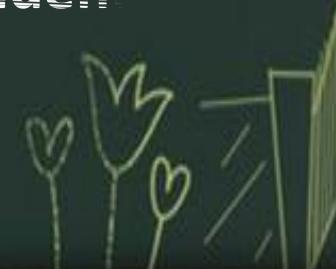
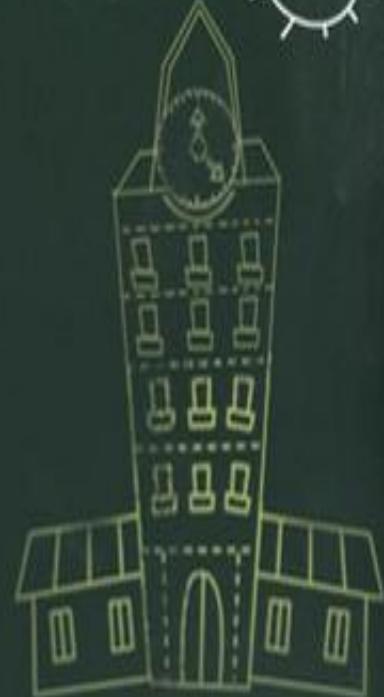
03

Central Tendency: Mean

Trimmed Mean : is used to deal with outliers by trimming or removing same data from both ends so as to get rid of outliers.

Weighted Mean: is used to certain values are supposed to count more in same context.

Eg. Calculating average grade of student based on their grade distribution.



04

Variation

variation in statistics is used to show how data is dispersed, or spread out , several measures of variation are used in statistics

- Range
- Quartile
- Variance

04

Variation: Range

Range : Range is the difference between the highest & lowest values in our dataset.
range tells us the distance between the lowest and highest value in data.

04

Variation: Percentile

Percentile: Percentile are score that are used to describe a value below which some observations fall .

If X is at 70th percentile it mean 70% of other data points from our sample are below X.

04

Variation: Quartile

Quartile : are used to break the data into four parts , so as to better find the spread of data in a way that is less influenced by outlier

Quartile are expressed in percentile, 1st Quartile is 25th percentile, 2nd quartile is 50th percentile(median) & 3rd Quartile is 75th percentile.

04

Variation: interQuartile Range (IQR)

InterQuartile Range (IQR) : is the difference between the lower and upper quartile. this gives us a better idea of the range of data.

Standard Variance

01

Standard variance measures how far a set of numbers are spread out from their average value

Standard Deviation

02

Standard Deviation is used to express the Magnitude by which the members of a group differ from the mean value for the group

Sample Variance

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

Sample Standard Deviation

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

correlation

01

Correlation is a term that is measure of the strength of linear relationship between two Quantitative Variable..

$$r = \frac{\sum(X-\bar{X})(Y-\bar{Y})}{\sqrt{\sum(X-\bar{X})^2} \sqrt{\sum(Y-\bar{Y})^2}}$$

Where, \bar{X} = mean of X variable
 \bar{Y} = mean of Y variable

Positive Correlation : is the term that is used a term that is used to describe a positive linear relationship between two quantitative variable

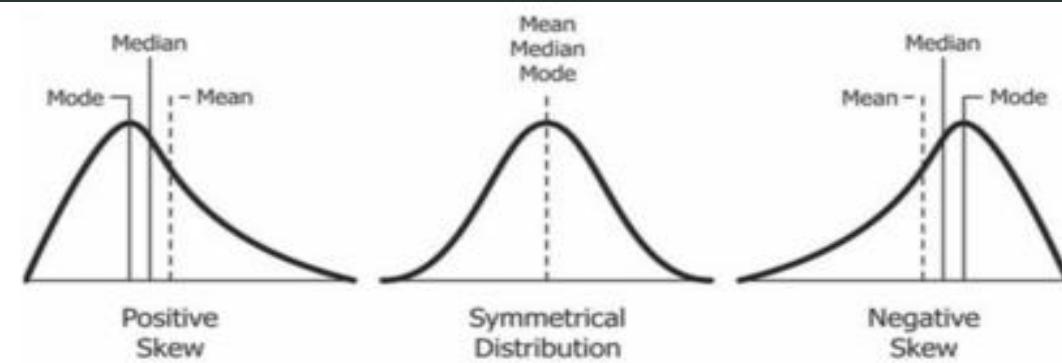
No Correlation : is the term that is used a term that is used to describe a No linear relationship between two quantitative variable

Negative Correlation : is the term that is used a term that is used to describe a Negative linear relationship between two quantitative variable

Normal Distribution

01

is the term used to describe a distribution which when plotted gives us a shape of bell curve. It has mean of zero and std deviation of 1.



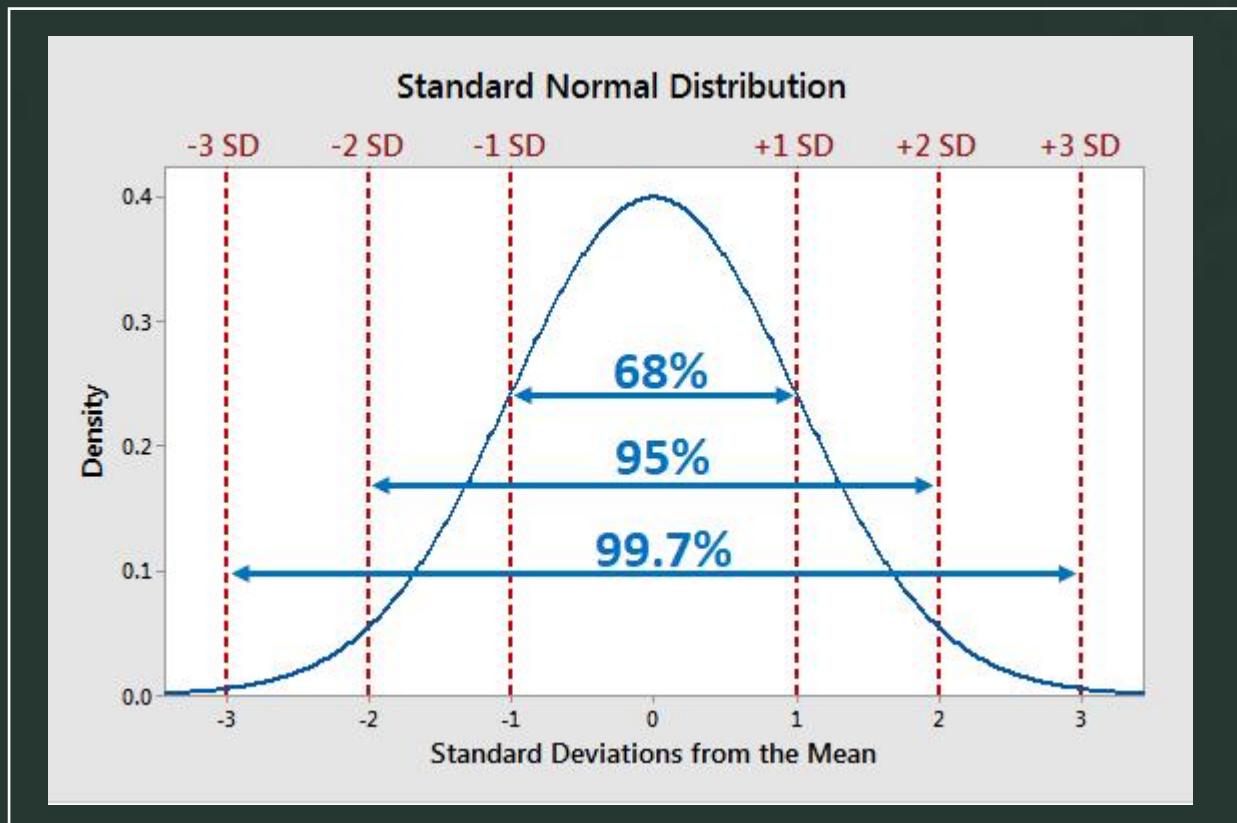
Empirical Rule

The empirical rule, also referred to as the three-sigma rule or 68-95-99.7 rule, is a statistical rule which states that for a normal distribution, almost all observed data will fall within three standard deviations (denoted by σ) of the mean or average (denoted by μ).

In particular, the empirical rule predicts that 68% of observations falls within the first standard deviation ($\mu \pm \sigma$), 95% within the first two standard deviations ($\mu \pm 2\sigma$), and 99.7% within the first three standard deviations ($\mu \pm 3\sigma$).

01

Is used to remember the Percentage of values that lie within a band around the mean in Normal Distribution with a width of two, four six stndard deviation



Z-score

Z-Score is Measure of how many standard deviation below or above the population mean a raw score

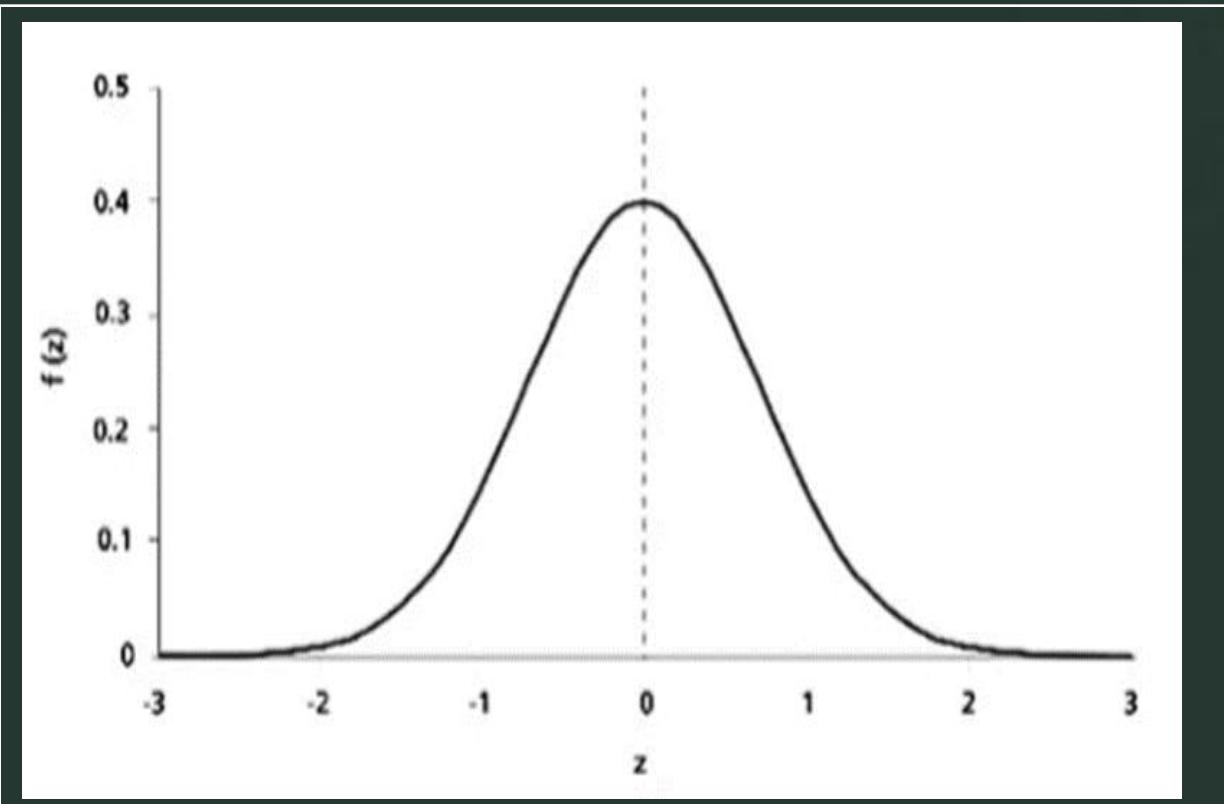
Z-score can be placed on a normal distribution curve.

$$z = \frac{x - \mu}{\sigma}$$

Score x minus Mean μ divided by SD σ

Z-score is a statistical measurement that describes a value's relationship to the mean of a group of values
Z-score is measured in terms of standard deviations from the mean. If a Z-score is 0, it indicates that the data point's score is identical to the mean score

01



THANK YOU

個日生

