

## Unit - 1

### 1.) Data

#### and information

\* Data is a collection of different facts, figures, object, symbols and events.

\* Data can be also defined as collection of observations.

\* Data can be of two types

They are Qualitative data and Quantitative data.

#### Qualitative data:

\* Qualitative data is a descriptive information.

\* for example, the most common names india, etc.

#### Quantitative data:

\* Quantitative data is a numerical information.

\* for example : height, weight,  
customers in shop, leaves on tree  
and other measurement.

\* Quantitative data can be of  
two types they are continuous  
and Discrete.

## Information:

\* Information is defined as  
classified or organized data that  
has some meaningful value for  
the user.

\* Information is used to make  
decisions and take actions.

\* A Data is said to be  
informative if it has the following  
significant.

\* Accuracy

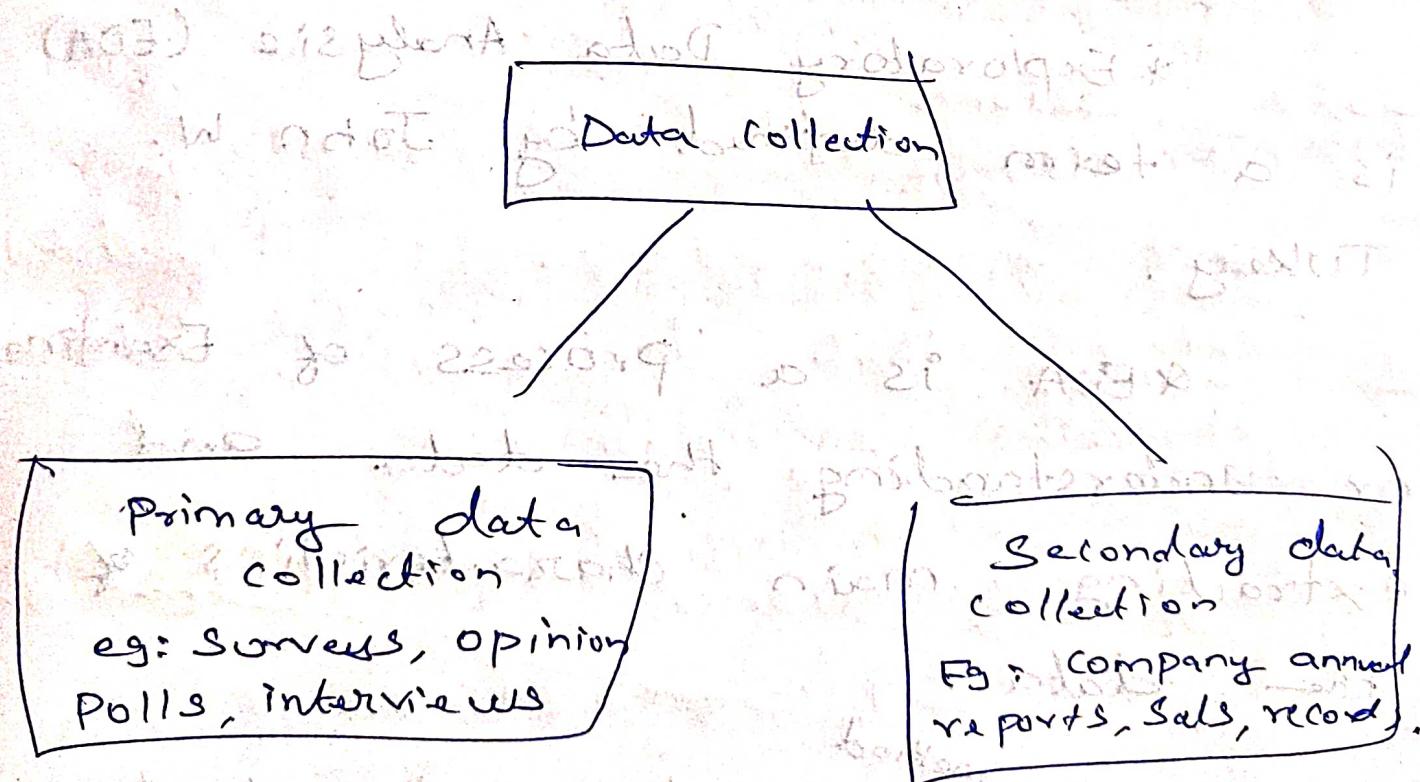
\* Completeness

\* Timeliness

## Data collections:

\* Data collection is defined as a method of collecting, analyzing the data for the purpose of validation and ~~don~~ research using some techniques.

\* Data Collection is classified into two types they are Primary and Secondary.



## Problems in data:

- \* Inconsistent data
- \* Ambiguous data - errors in data.
- \* Duplicate data -
- \* Too much data
- \* Inaccurate data
- \* Hidden data

## Q) EDA [Exploratory Data Analysis]

- \* Exploratory Data Analysis (EDA) is a term defined by John W. Tukey.
- & EDA is a process of Examining or Understanding the data and extracting main characteristics of the data.
- & EDA <sup>method</sup> is classified into two types they are graphical analysis and non-graphical analysis.

\* EDA is the first critical step in analyzing the data.

\* It is also known as visual analytics or descriptive statistics

\* It is the practice of inspecting and exploring the data before fitting predictors and other inferential goals.

### Statistical graphics:

\* A statistical graph or chart is a pictorial representation of statistical data in graphical form.

e.g.: Histograms, box plots, scatter plots

### Non-graphical analysis:

\* It uses statistical techniques to explore a single variable

## Reasons to use EDA:

- 1) Detection of mistakes
- 2) Examine the data distribution
- 3) Handling missing values.
- 4) Handling outliers
- 5) Removing duplicating
- 6) Normalizing and scaling
- 7) Determining relationships

## 3) Data Science:

\* Data science is the domain of study that deals with data in large volume.

\* Data science is combination of math and statistics.

## life cycle of data science.

Data science's are classified into  
five stages:

- 1) capture - data collection, Data acquisition, Data extraction
- 2) maintain:- Data warehousing, data cleansing, data processing.
- 3) process - Data mining, clustering, modeling.
- 4) Analyze - Predictive Analysis, regression, text mining.
- 5) Communicate. - Data reporting, Data visualization, Business intelligence, Decision making

Data science tools:

1) Data Analysis:

SAS, Jupyter, MATLAB, Excel.

## \* Data Warehousing:

Informatica, AWS Redshift

## 3) Data Visualization:

Jupyter, Tableau, Cognos

PowerBI

## 4) Machine Learning:

Spark, MKLlib, Azure ML

Studio, Mahout.

## Application of Data Science:

\* Healthcare

\* Streaming - to next level.

\* Image recognition

\* Recommendation

\* Logistics

\* Fraud detection

\* Internet search

\* Speech recognition

\* Finance

- 3) Types of EDA
- \* There are four Primary types of EDA, Namely
    - \* Univariate non-graphical
    - \* Univariate graphical
    - \* Multivariate non-graphical
    - \* Multivariate graphical.
  - \* Univariate non-graphical:
    - \* This is the Simplest form of data Analysis.
    - \* The data consists of 2 things, since it is an one variable.
    - \* Single variable.
    - \* It does not deals with relationships.
    - \* The main purpose of Univariate analysis is to describe the data and finding patterns.

## Univariate graphical:

\* Univariate non-graphical

methods do not provide a full picture of the data.

\* Hence graphical methods -

required. Shape of distribution,

frequency, mean, median, mode etc..

\* Steam-leaf plots, which shows all data values, and the shape of distribution.

\* Histograms, Bar plots : represents the frequency.

\* Box, Plots : mean, median, mode.

## Multivariate non graphical:

\* Multivariate data arises

from more than one variable

\* Multivariate non-graphical

EDA techniques generally shows

relationship between two or  
the more variables.

\* Multivariate non-graphical  
represents data through cross-  
tabulation or statics.

### \* Multivariate graphical:

\* Multivariate data uses  
graphics to display the relationships  
between two or more sets of  
the data.

\* The most used graphic  
is bar plot or bar plot.

\* Other common types of  
multivariate graphics include

\* scatter plot

\* Multivariate chart

\* Run chart

\* Bubble chart

\* Heat map

#### 4) Sense of Data:

\* In today's digital era data is very powerful tool.

\* Just having possession

of data is not sufficient. one should be able to analyze the data in proper manner to get insight of data.

\* making sense of the data means finding meanings from data generating inference from data.

\* Data analysis the process of making data meaning.

\* consider the process of creating something out of a set of raw materials.

for example: given a bolt of fabric and some notions.

Different people make different kinds of creation. one person might make a baby dress, another one make dance costume and another one make diwali lamp. But what they come up with is raw materials. It would be awfully difficult to make a working of car using bolt of fabric.

\* Basic steps to make data sense

- 1) managing data.
- 2) Getting familiar with data.
- 3) making sense of data.

\* Managing data

\* It is an initial process of analyzing qualitative data.

\* First task is arrange or organize data. So it will make sense of it.

\* Separating different types of data. This will help to access the data easily and also keep the data in chronological order makes data access fast.

\* organizing data in topic wise or document wise

### Getting familiar with data:

\* Based on data collection and its size, it might take time to understand about data. Getting familiar with data plays a vital role in output or result prediction.

### Making sense of data:

\* Arranging the data in certain presentable or viewable format.

\* Data is best seen using tables, charts, graphs and patterns.

## 6) Software for EDA:

### 1. R -

\* It is Open Source Programming language.

\* It is graphics Supported

\* R is foundation for Statistical

Computing .

\* The R language is widely

used for Statistical Observations  
and data Analysis

### 2. Python :

\* An interpreted, object oriented  
programming language .

\* It is high level and  
has built-in data structures .

\* It is rapidly used for  
application development and as well  
as for scripting .

### 3. Excel / spreadsheet:

- \* Excel tool is considered to indispensable part of analytics industry.
- \* It supports all the features like summarizing, visualizing data, data wrangling etc.
- \* It is managed by Microsoft.
- \* Microsoft excel is paid.

### 4. Weka:

- \* weka is an machine learning tool,
- \* Easy to learn
- \* having an intuitive interface.
- \* It provides options for data - pre processing, classification, clustering, association rule, and visualization.
- \* It is built in Java.

## 5. Tableau :

\* Tableau is data visualization Software.

\* It is a fast visualization which allows exploring of data.

## 7. Visual Aids For EDA:

### Univariate plots:

\* It shows the frequency or the distribution shape of variable.

#### i) Histograms:

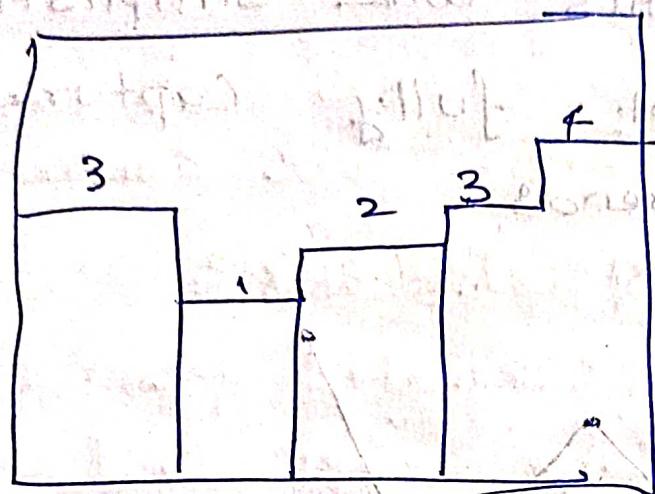
\* Histograms are two-dimensional plots which contains x-axis and x-axis.

\* x-axis shows time intervals

\* x-axis shows the frequency values.

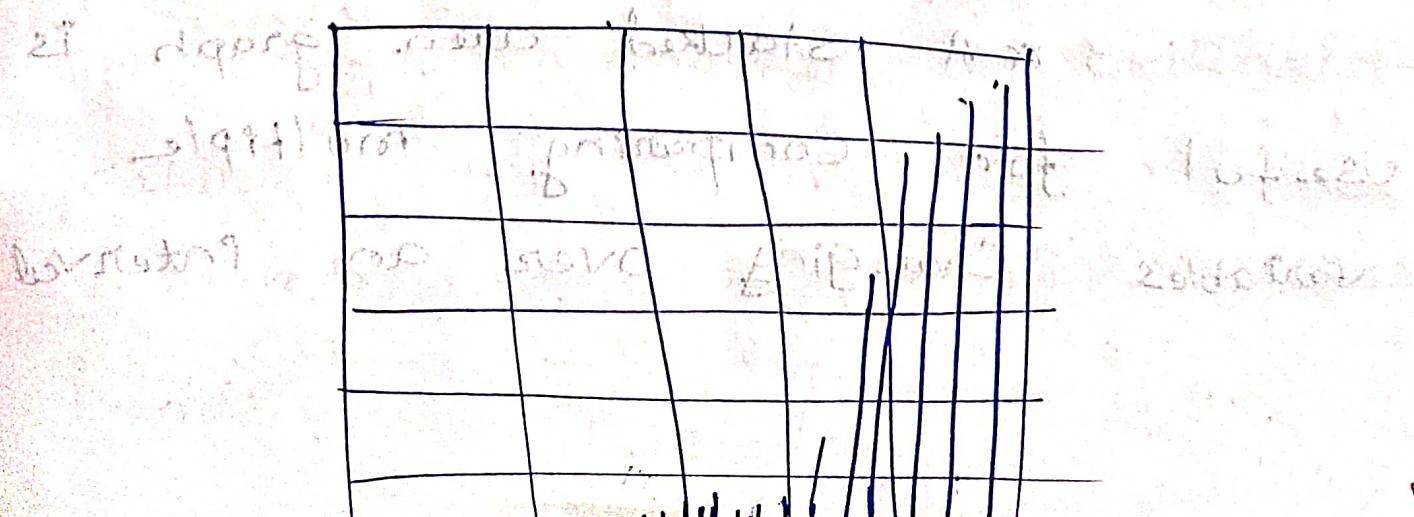
\* It is similar to Bar graph, but histogram does not contain gaps between bars.

e.g:



## ii) Dist plot:

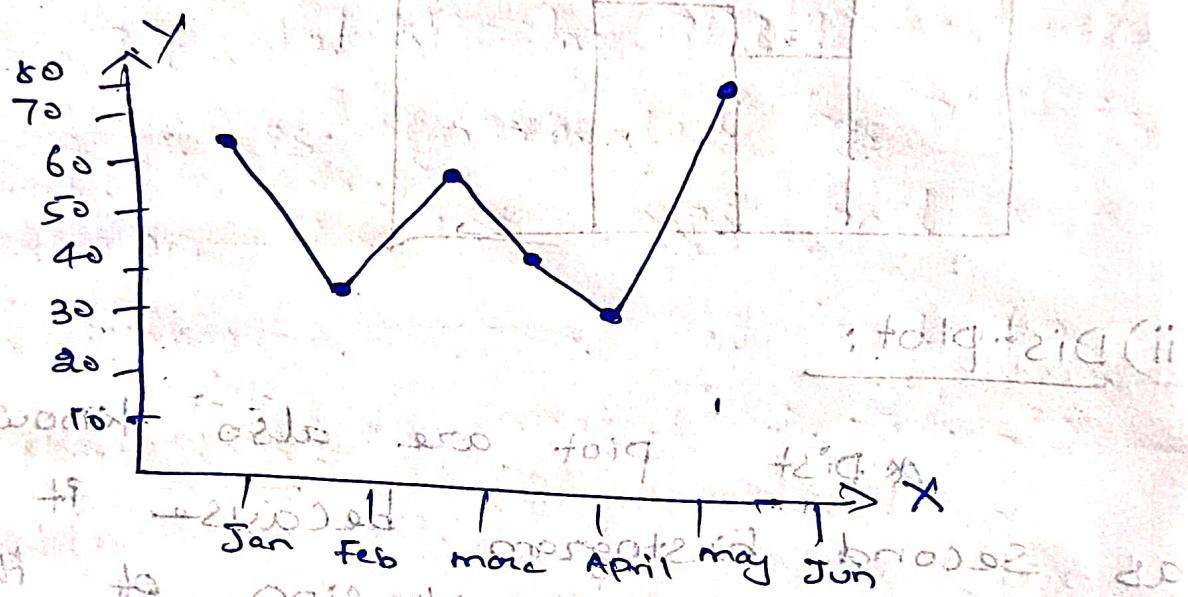
\* Dist plot are also known as second histogram because it is a slight improved version of the histogram.  
\* Distplot gives us as KDE (Kernel Density Estimation).



### iii) Line chart:

\* A line chart is a graphical representation of an asset's historical price.

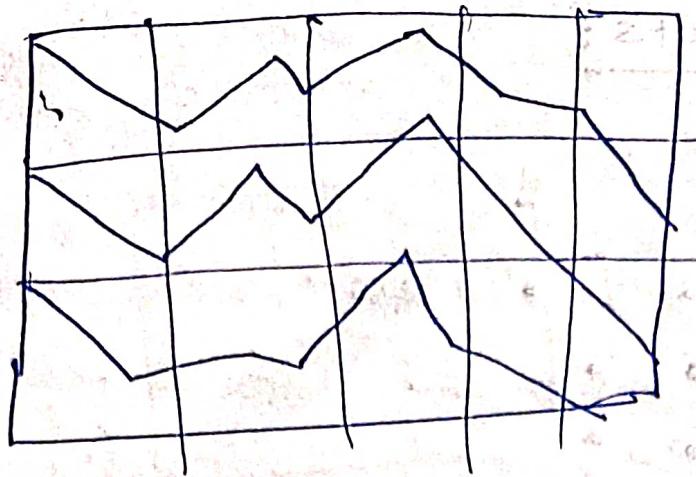
\* Line charts are simplistic and may not fully capture patterns or trends.



### iv) Stacked area plot/chart:

\* A chart combines the line chart and bar chart elements.

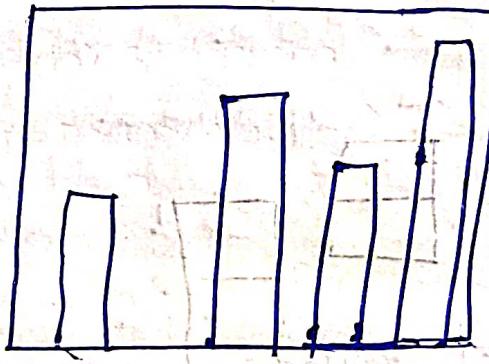
\* A stacked area graph is useful for comparing multiple variables changing over an interval.



## 2. Bivariate plots:

Bivariate plots display the relationship between two variables in exploratory data analysis.

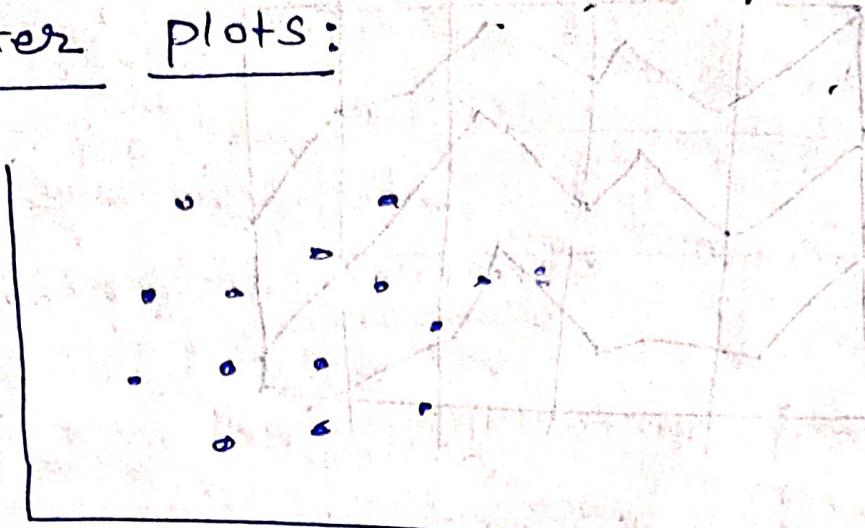
### i) Bar graphs:



& Bar charts can be used to compare nominal or ordinal data.

& They are helpful for recognizing trends.

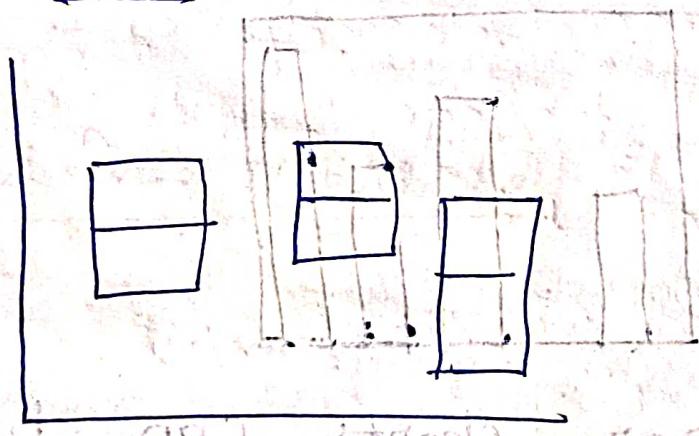
## ii) Scatter plots:



\* Scatter plots are commonly used in statistical analysis.

\* It is correlated by plotting them on the X-axis and Y-axis.

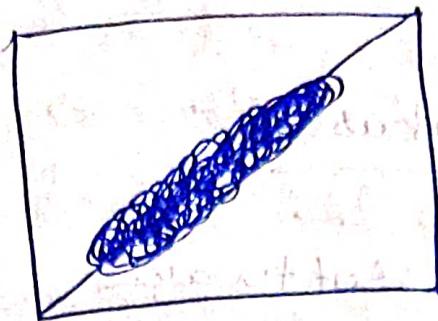
## iii) Box plots:



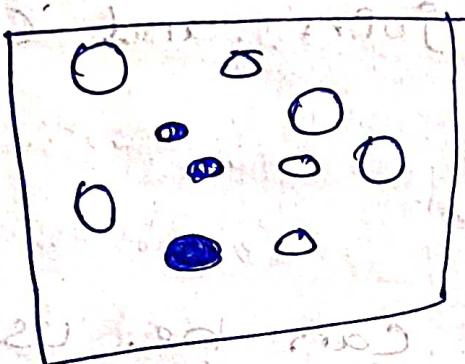
\* Box plots are graphical representation of two boxical relationships.

\* Box plots are suitable for identifying outliers.

#### iv) Correlation plots:



#### v) Cluster Map:



\* It is also known as a dendrogram.

\* It is also known as Heat map.

#### 3) Multivariate plots:

\* It deals with more than 3 variables.

\* It is also known as three dimensional plots.

\* It is used in very large dimensional.

Text data:

Image data:

## 8) Merging Database:

pandas Software in Python programming language for manipulation and analysis

\* merge, join and concatenate.

Pandas merge:

\* merge() can be used to a

database's join operation. It is

the most flexible operation

that can be applied to data.

\* it is used to combine

data using one or more keys.

\* Both Many-to-one and

Many-to-Many joins with merge().

\* merge() function is used

for merging datasets.