

CS 3352 - Foundation of Data Science

Unit I - Introduction

Data science: benefits and uses - facets of data
Data science process: overview - Defining Research Goals
Retrieving Data - Data preparation - Exploratory Data
Analysis - Build the model - Presenting findings and
building applications - Data mining - Data warehousing
Basic statistical descriptions of data.

Big data:
It is a blanket term for any collection of data sets so large or complex that it becomes difficult to process them using traditional data management techniques such as RDBMS.

Data science:
It involves using methods to analyze massive amounts of data and extract the knowledge it contains.

→ Relationship between big data and data science is like the relationship between crude oil refinery

→ Data science & big data evolved from statistics and traditional data management.

The characteristics of big data:

- i) Volume - how much data?
- ii) Variety - how diverse are different types of data?
- iii) Velocity - At what speed is new data generated?

Additional characteristics:

- i) Veracity - how accurate?
- ii) Value - usefulness of data.

→ It often referred to as three Vs.

Benefits and uses of data science & big data.

→ Data science & big data are used everywhere in both commercial & non commercial purposes.

→ Commercial companies use data science and big data to understand in deep about their customers, processes, staff and products.

→ To offer customers a better user experience
↳ cross-sell - upsell products.

Eg: Google AdSense which collects data from internet users and recommends ads based on the history search.
maxpoint eg of real time personalized advertising.

People Analytix, HR professionals use people analytics data mining to screen candidates, monitor the mood of employees and study informal networks among workers.

Financial Institutions use data science & big data to predict stock markets, risks of losing money & learn to attract new clients.

Government organizations rely on data science to discover valuable information.
eg: Data.gov, ITR the use of us gov open data.

Non government organizations use it saves their funds. Used to increase the effectiveness of their fund raising efforts.

Online Learning Platforms.

Learn bay, courses, Udacity, Udemy & etc, mass open online courses.

Types of data: Different types of data (or) main categories of data

- 1) Structured
- 2) Unstructured
- 3) Natural language
- 4) Machine generated
- 5) Graph based
- 6) Audio, video and Images
- 7) Streaming

1. Structured data:

Structured data is data that depends on a fixed field within a record. It is easy to store & manage data in tables within database or Excel file.

eg. An excel table.

s.no Register No Student name mobile no email id.

2. Unstructured data

It is not easy to fit into a data model because its content is context specific or varying.

eg: Email.

3. Natural Language:

It is a special type of unstructured data. It is difficult to process because it requires knowledge of specific data science techniques and linguistics. eg: Email, word documents, Natural language processing, summarization, text completion, sentiment analysis.

4. Machine generated data:

It is information that is automatically created by a computer, process, application or other machine without human intervention.

eg: web server logs, call detail records, network event logs.

CSI PERF: Jca Commit T=313236 2016-11-28 11:36:13

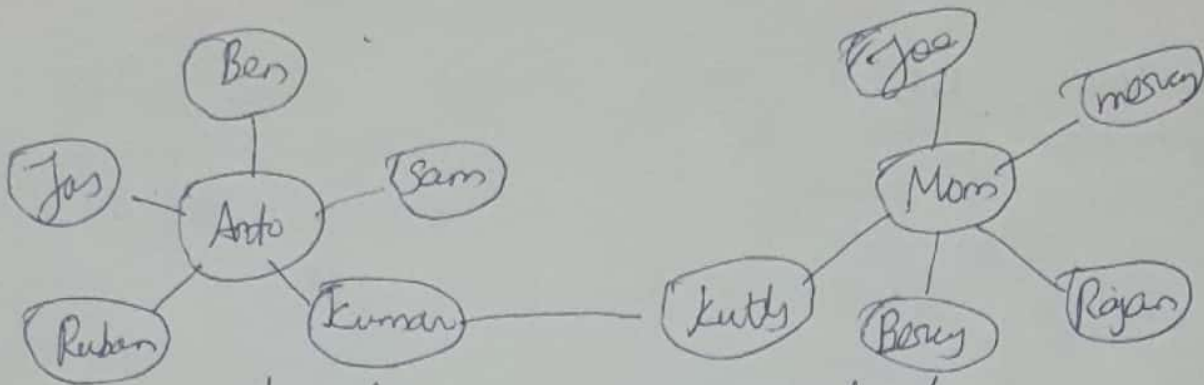
Info CSI 00000153 Creating NT transaction

Analysis of machine data relies on highly scalable tools due to its high volume & speed.

5. Graph based (or) Network data:

It is data that focuses on relationship between objects. It uses graph structure with nodes & edges to represent objects and their relationships.

eg: friends in social networks.



friends in a social network.

⑥ Audio, Image & Video:

multimedia data in the form of audio, video, images become an integral part of everyday life. Object recognition: challenges for computers. Deepmind: developed an algorithm which is - capable of learning how to play video games - able to interpret everything in the video screen via deep learning.

⑦ Streaming data:

It takes almost any of the previous forms but it means the data flows into the system. When an event happens instead of being loaded into a data store in a batch.

eg: what's trending on twitter, live sports, or music events.

1.3 The Data Science Process:

The data science process consists of 6 steps they are:

- ① Setting the research goal.
- ② Retrieving data
- ③ Data preparation
- ④ Data Exploration
- ⑤ Data modeling
- ⑥ Presentation and automation.

1. Setting the research goal

The purpose of this step is making sure all the stakeholders understand the what how and why of the project.

2. Retrieving data

This step includes providing suitable data & getting access to the data from its data owner. The result is data in its raw form.

3. Data preparation

This includes transforming the data from a raw form into usable data. This involves detection & correction of different kinds of error in the data. Combines data from different data sources & transform it.

4. Data Exploration

The step involves finding patterns correlation & deviations based on visual & descriptive techniques to give a deep understanding of the data.

5. Model building

Select the variables to build the model & also a modeling technique. Building models with a goal of making predictions classifying objects.

6. Presentation & automation

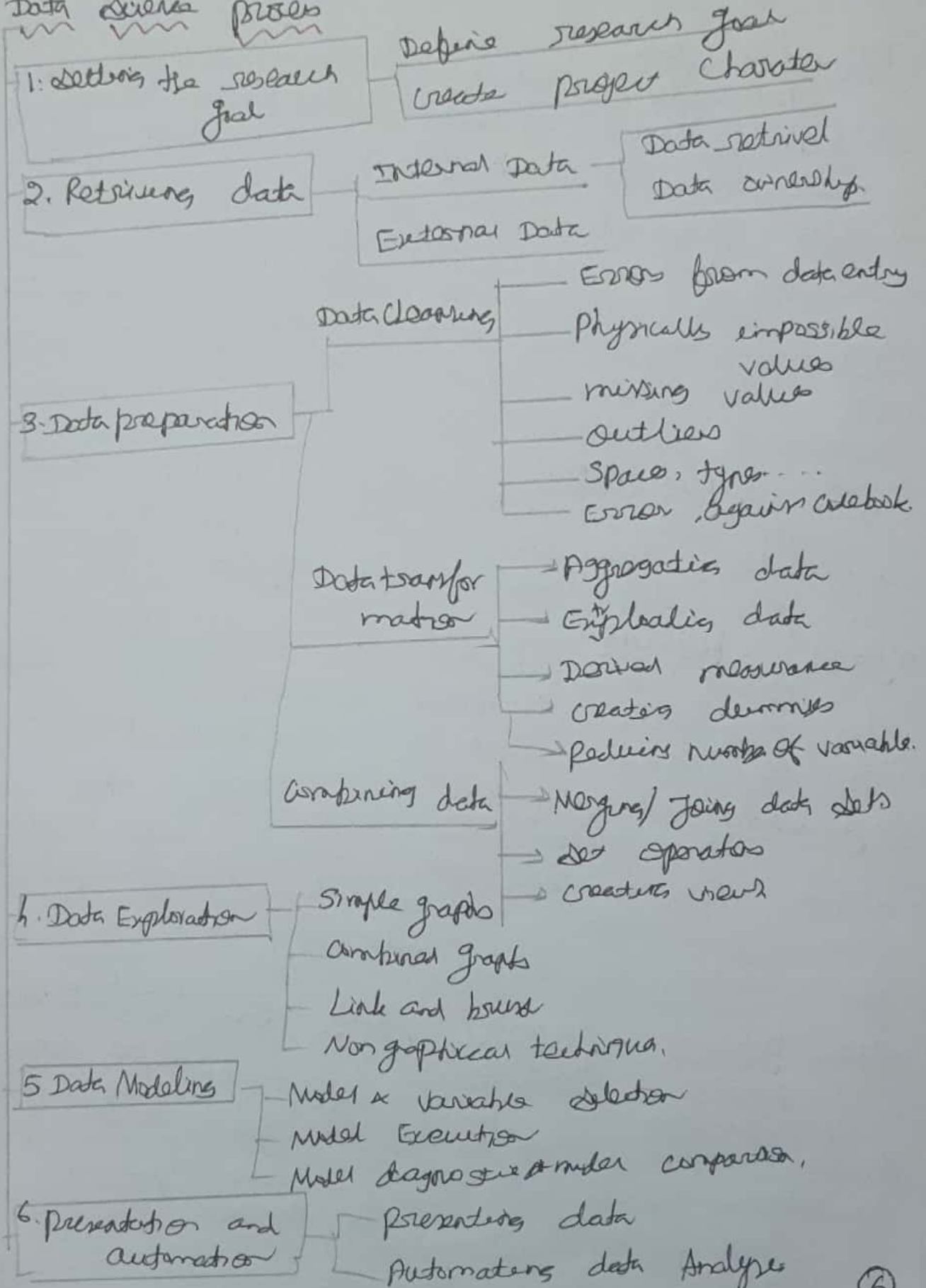
Finally present the result to the business results can take many forms ranging from presentation to research reports. Automate the execution of the process if needed.

Step 1: Defining research goal and creating a project character.

understanding the what, why & how of the project answering questions is the goal of first phase. outcome should be a clear research goal just understanding of the content well defined deliverables plan of action with a time table.

Overview of the data science process
 The following figure summarizes the data science process.

Data Science Process



Creating a project character:

- A clear research goal
- the project mission and context
- how do you perform the analysis
- what resources are expected to use.
- Proof of concepts
- Deliverables and a measure of success A timetable

A budget can use this information to estimate the project cost, data and people required for completion of project.

Step 2: Retrieving data

- second step of data science process is retrieve the required data.
- many companies have already collected & stored the data for use internally
- many organizations are making high quality data freely available for public & commercial use
- Data can be stored in official data repositories such as databases, data marts, data warehouses & data lakes.
- Getting access to data is difficult task. So organizations have policies which controls the access of data everyone in that organization [ie Access rights - who can access what data]

A list of open data provides:

- 1) Data.gov - the home of the US governments open data
- 2) Open-data.europa.eu - the home at European commission's open data
- 3) Data.worldbank.org - open data initiative from the world Bank.

Step 3: Cleansing, Integrating and Transforming

The data received from the data retrieval phase is like a diamond in the rough" so it must be sanitized & prepared for use in modeling & reporting phase.

Cleansing data

Data cleansing is a sub-process of data science process that focus on removing errors in the data to make the data become true & consistent base.

An overview of Error Description	Common errors: possible solution:
<ul style="list-style-type: none"> → mistakes during data entry → Redundant white spaces → Impossible values → Missing values → outliers → Deviates from a case book → Different units of measurement → Different levels of aggregation 	<ul style="list-style-type: none"> → manual overrules Use string functions manual overrules Remove observation (or) value Validate, if error treat as missing. match on key (or) else use manual overrules Recalculate Bring to same level of measurement.

possible errors with examples.

- ① Data entry errors - Entries as "Grade" instead of "Good"
- ② Redundant white space - white space at the end of a string causes errors, very hard to detect.
- ③ Impossible values - people taller than 3 meters (or) people within age of 299 year.
- ④ outliers - is a data that lies abnormally far away from other values in a data set
15, 25, 30, 20, 35 ⇒ 98 outlier.
- ⑤ Deviates from a case book - detecting errors in large data base
eg Salary / week and Salary / work week
A good practice is to correct errors as early as possible week.

combining data from different sources.
different way to combine data.

- joining
- Appending table.

① Joining:

↳ Considering an observation from one table with information from another table

eg/

Client	Item	month
John	pen	January
Jane	pen	January

Client	Region
John	Tamilnadu
Jane	Karnataka

Result of joining two tables

Client	Item	month	Region
John	pen	January	Tamilnadu
Jane	pen	January	Karnataka

② Appending

table to those of another table. Adds the observation of one table to those of another table

Client	Item	month
John	cup	January
Jane	paper	January

Client	Item	month
Jane	pen	February
Jack	fruit	February
Ben	fruit	February

Client	Item	month
Jack	fruit	February
Ben	fruit	March

↳ To avoid duplication of data, one can virtually combine data with views.
 ↳ Also data redundancy can be done by adding calculated information

Step 3: Transformers

↳ It makes the data suitable for data modeling.

eg: Transformers x to $log x$

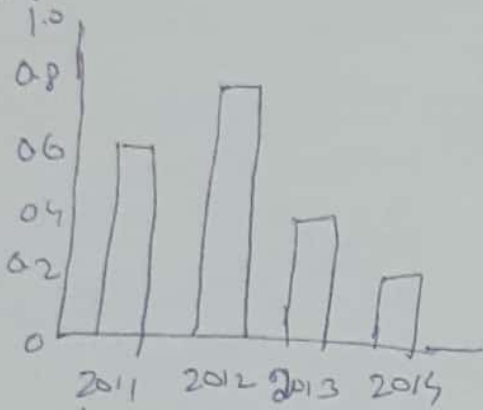
↳ Reducing the no of variables (Principle of Dimensional Analysis)

↳ Threshold variables with dummies. Dummies variables can only take two values true (1) or false (0)

Step: Exploratory Data Analysis.

→ Helps to gain deep understanding of data & the interactions between variables.

→ A bar chart, a line chart, & histograms used in exploratory analysis.



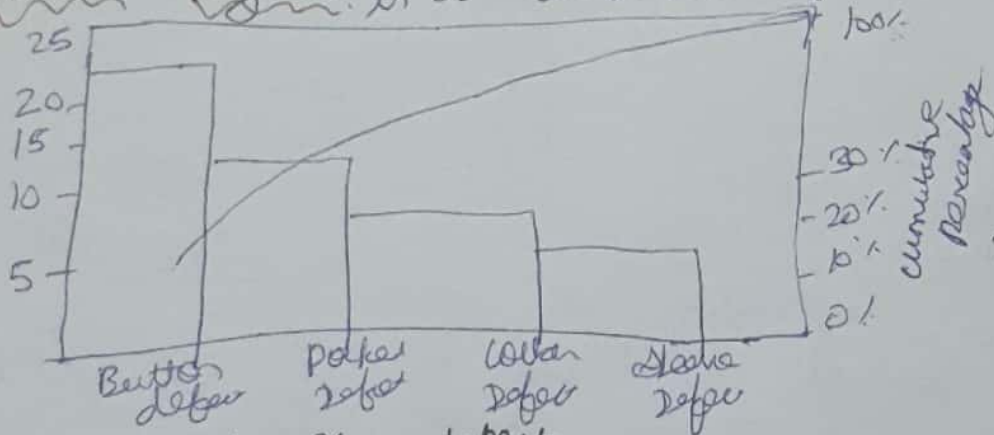
bar chart



Line Chart

Combined graph - Pareto diagram

Pareto diagram: is a combination of a bar graph & line graph

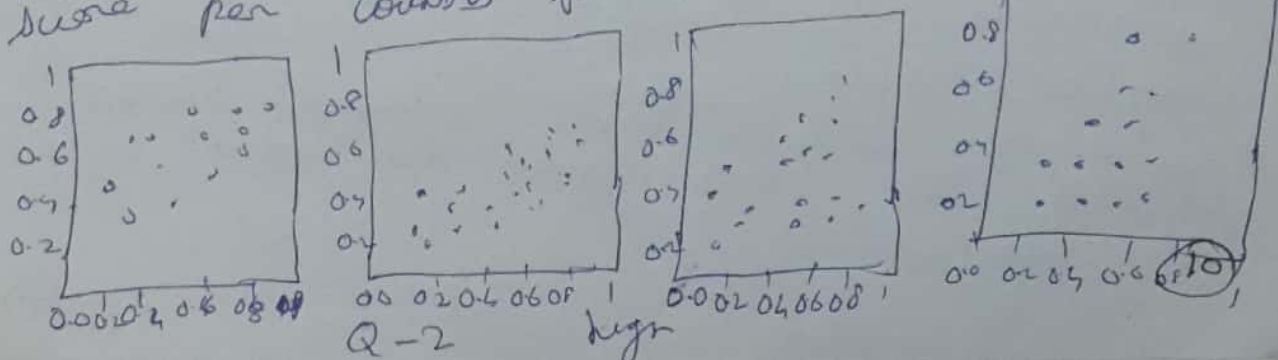


Pareto charts depicts in short

Link and

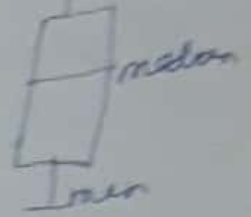
→ helps to brush combine & link different graphs & are automatically transferred to the other graphs.

→ the following transfer shows the average score per country for questions



It shows the correlation between answers & also to see the bipolar points in the graphs

Boxplot: Shows the maximum, minimum, median & other characterizing measures
 → Other are tabulation, clustering & other modeling techniques.



Steps build the model:

- building models with the goal of making better predictions. Classifying objects.
 - builds a model in an iterative process
 - It depends on statistics & machine learning techniques
 - The components of data modeling
- a) model x variable subsets:
- many modelling techniques are available.
 - Choosing the right model based on factors such as model performance, projects requirements, easy to implement, difficulty in maintenance, & easy to explain.

b) model execution

- Implementation of chosen model in code
- Python has libraries such as StatsModels(e)
- scikit-learn to implement models.

Implementation of Linear Regression Model:

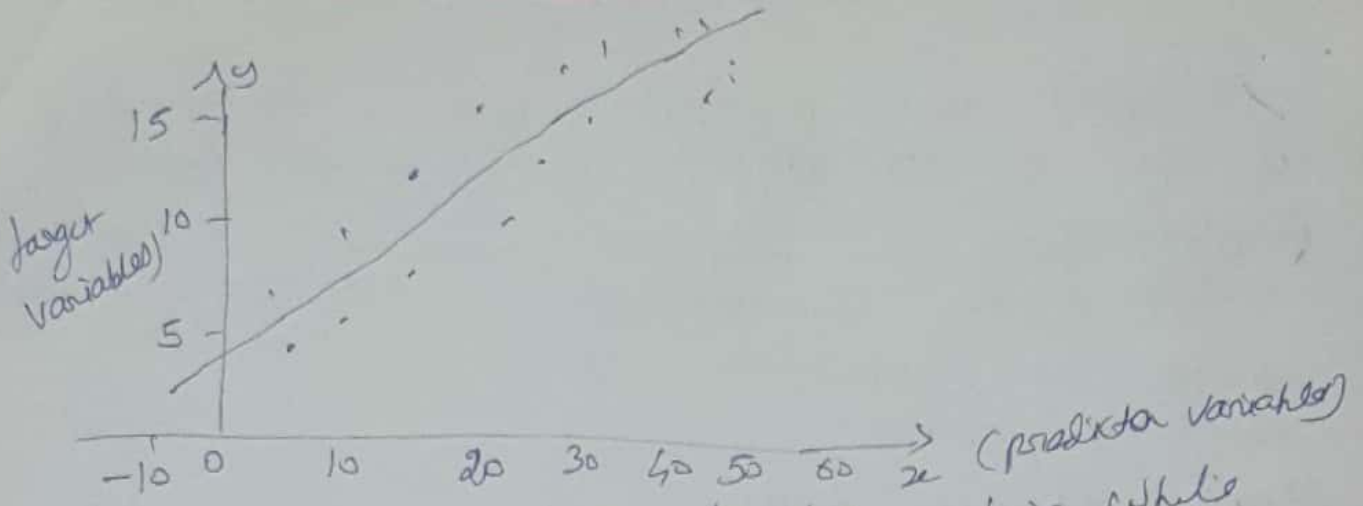
```
import statsmodels.api as sm
import numpy as np
import pandas as pd
import sklearn.linear_model as lm
```

```
predictors = np.random.random(1000).reshape(500)
target = predictors.dot(np.array([0.4, 0.6])) + np.random.random(500)
```

```
smf = sm.OLS(target, predictors) # fits Linear Regression on data
```

```
res = smf.fit()
```

res.summary() → Shows model statistics



Linear regression tries to fit a line while maintaining the distance to each point. Linear equation is $y = 0.7858x + 1.1252$

Some important output parts of linear regression model

- model fit: R-squared or adjusted R-squared is used \rightarrow It is an indication of the amount of variation in the data that gets captured by the model.
- predictor variables have a coefficient
- predictor significance: how significant the predictor is.

Model diagnosis and model comparison

- Choose the best model based on multiple criteria
- multiple error measures are available one is mean square error (MSE)

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

- It checks for each predictor how far it is from the actual.
- Choose the model with the lowest error
- Verification of meet out of model assumptions is called model diagnostics.

Step 6: Presenting findings and building applications

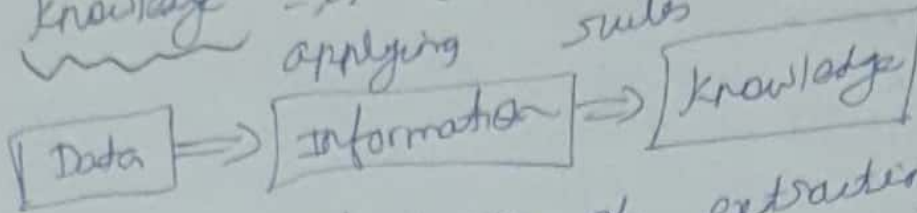
presentations of presenting data

powerpoint presentations. Automate the models in order to do the tasks repeatedly.

Automating data analysis

Data Mining:

Data is raw facts (or) disconnected fact
Information is the processed data
Knowledge - is derived from information by applying rules to it



Data mining: is the process of extracting hidden, valid and potentially useful patterns in huge datasets.

Steps in data mining:

1. Data cleaning: Remove noise & irrelevant multiples data sources may be combined
2. Data Integration: where data relevant to the analysis.
3. Data selection: where task are retrieved from the databases.
4. Data transformation: where data are transformed consolidated into forms appropriate for mining by performing summaries or aggregation operations.
5. Data Mining: An essential process where int methods are applied to extract data patterns
6. Pattern evaluation: to identify the truly interesting patterns representing knowledge.
7. Knowledge presentation: where visualization & knowledge mixed knowledge to user.

Types of data for mining:

1. Flat Files: Transactions measurements.
2. Data base data: Relational databases.
3. Data Warehouse data: is a repository of information collected from multiple sources, stored under a unified schema.
4. Transactional data: A transaction includes a unique ID and a list of items making up the transactions.

Data Warehousing:
It is the process of compiling and transforming data and making it available to users in a timely manner.

Data Warehouse:
A single complex and consistent store of data obtained from a variety of different sources made available to end users in a way understandable and use in a business context.
eg. Walmart - 2.4 TB (Tera Bytes)
National medical records - 10¹¹ bytes in order of Exabyte

Key Characteristics of Data Warehouse
A data warehouse is a object oriented - provides top level information
→ Integrated - Data from varied sources in a consistent format

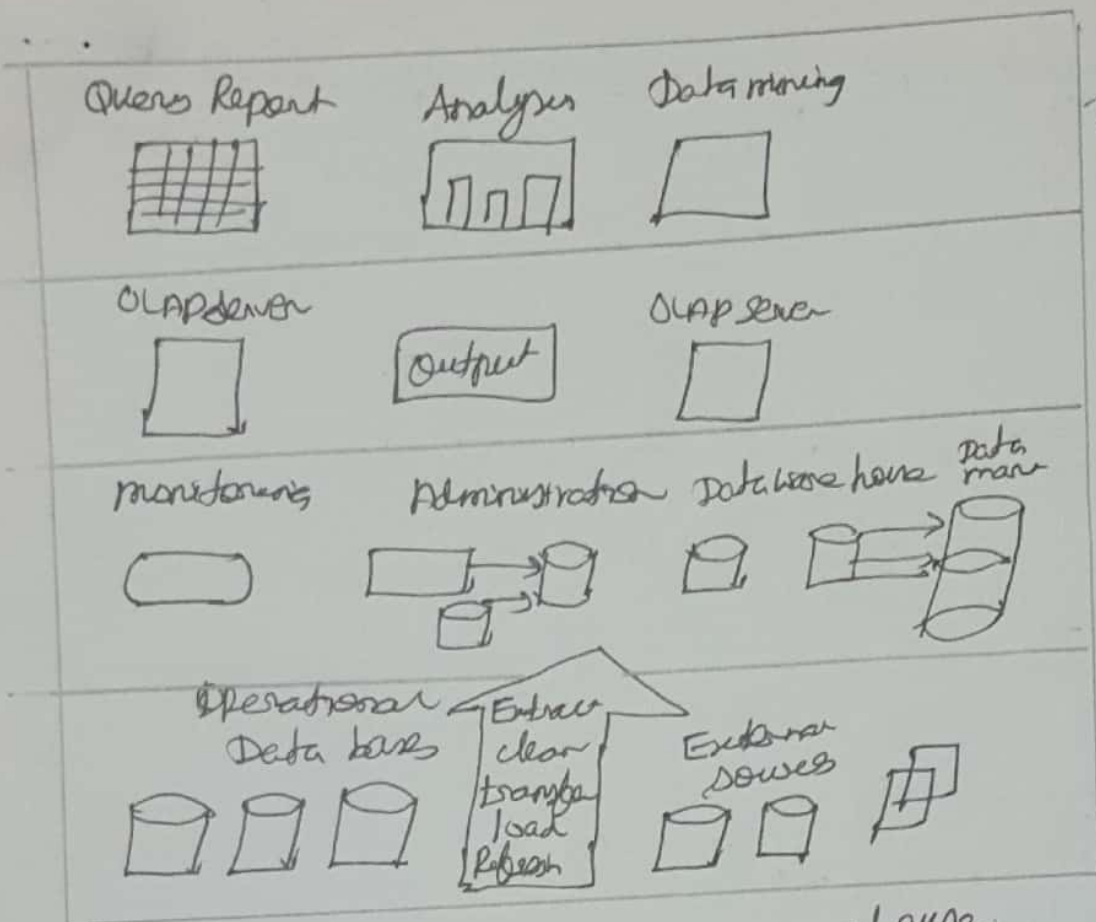
→ Time Varying - Data with time
→ Non-volatile - Data remain unchanged collection of data that is used primarily in organizational decision making.

Data Warehousing Architecture
Generally a data warehouse adopts a three-tier architecture.

Bottom Tier: It is the database server. It is the relational database system. back end tools and utilities are used to feed data into the bottom tier. They perform extract, clean load & suppress functions

Middle Tier: here OLAP server can be implemented in either Relational OLAP or multidimensional OLAP

Top Tier: This is front end client layer. This layer holds the query tools and reporting tools, analysis tools and data mining tools.



Process flow in Data Warehouse

- There are 4 main process
1. Extract and load the data
 2. Cleaning and transforming the data
 3. Backup & archive the data
 4. Managing queries and directing them to the appropriate data sources.

Extract & Load
 This process extract data from operational data bases & External sources & load it into the data warehouse in consistent form.

Clean & Transform
 There are 3 steps in cleaning & transforming process
 clean & transform data into a structured partition the data - partition each table into multiple partitions.
 Aggregator - It is required to speed up common queries.

Backup & Archive the Data - [warehouse manager]
 In order to recover the data in the event of data loss, software failure or hardware failure it is necessary to keep regular backup.

Query Management Process (overs managers)

- This process performs the following functions
- managers the queries.
 - speed up the execution time of queries.
 - directs the query to their most effective data sources.
 - monitors the actual query profile.

Online Analytical Processing (OLAP)

- OLAP allows managers and analysts to get an insight of the information through fast consistent and interactive access to information.
- It enables end users to perform better analysis of data in multiple dimensions.

OLAP Applications

a) Finance and accounting

- budgeting
- Financial performance analysis
- Financial modeling.

b) Sales and marketing

- sales analysis and forecasting
- market research analysis
- customer analysis.

c) Production

- Production planning
- Defect analysis

→ OLAP pre-calculates most of the queries
names aggregation / joins and grouping

Thank you

(16)

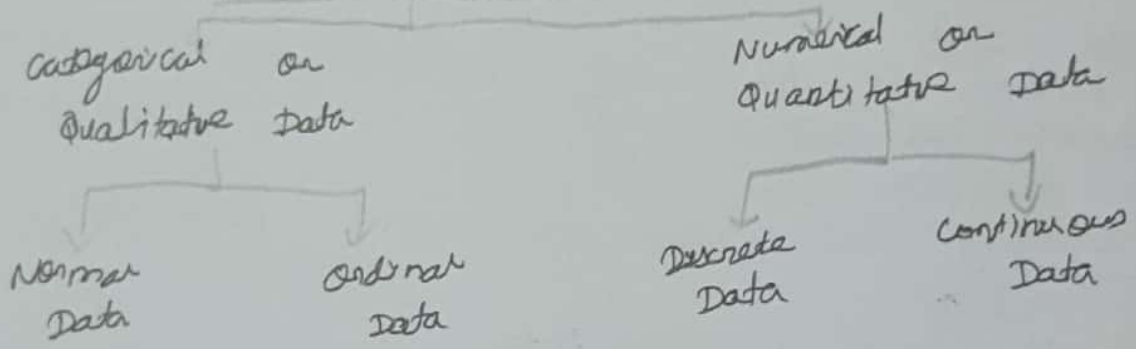
Unit - II

Types of Data - Types of Variables - Describing Data with Tables and Graphs - Describing Data with Average - Describing Variability - Normal Distributions and Standard (z) Scores.

Types of Data:

Data is the collection of actual observations or scores in a survey or experiment which can be used for statistical analysis. It is structured and stored in a data base, etc. in table format - containing Numerical or text values.

Types of Data



Qualitative (or) Categorical Data:

- Qualitative data also known as Categorical Data.
- Describes the data that fits into the categories.
- Qualitative data are not Numerical.

→ Categorical variables describe features such as persons - gender, home, town, etc.

→ Categorical data are birthdate, favourite sport, school, post code. Qualitative data consist of word (yes or no), letters (Y or N), Numerical code (0 or 1).

eg. Do you like cartoons. Replies from few students

Students	Result
1	Y
2	N
3	Y
4	N

Nominal Data which helps to label the variable without providing the numerical value. but the data can be qualitative and quantitative.

eg: Nominal data are letters, symbols, words, gender.

Ordinal Data:

Ordinal data is a qualitative data that variable is type of data that follows a natural order. (group the variable into ordinal categories). Ex: this variable is found in surveys, finance, economics, questionnaires. The ordinal data is commonly represented using a bar chart. These data are investigated and interpreted.

Ranked Data:

Ranking is the data transformation in which numerical (or) ordinal values are replaced by their rank when the data are sorted.

→ arranging in ascending order and making from low to high.

2. Quantitative (or) Numerical Data

Quantitative data represents the numerical value like how, much, how, often, how many.

eg: Numerical data are height, size, weight.

eg: Height of 5 students.

Students	Height (in cm)
1	160
2	165
3	170

Discrete Data

Discrete information contains only a finite number of possible values (these values cannot be subdivided meaningfully). It can be counted in whole numbers.

Discrete data Key Characteristics.

You can count the data if usually units or whole numbers. The values cannot be divided into smaller pieces and are additional markings.

Discrete Data can be represented by

→ Bar Graph

→ Frequency table

→ Line Plot (Number Line)

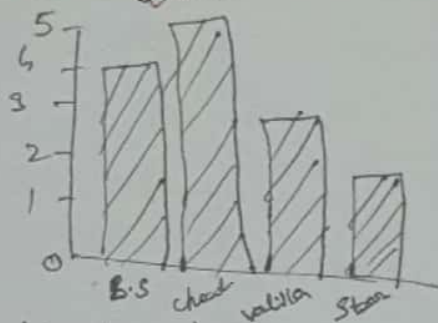
Bar Graph is the most suitable way to represent discrete data, as finite values can be presented clearly through vertical bars.

Example.

In a survey with 14 children on their favourite ice cream flavour, it was found that 4 children like butterscotch flavour, 5 children like chocolate flavour, 3 children like vanilla flavour & 2 children like strawberries flavour.

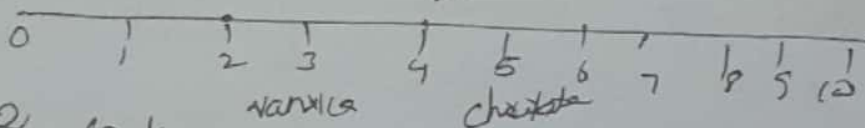
→ It is an example of discrete data. We can count the number of children who like a particular flavour of ice cream.

Bar Graph.



Number Line:

We mark the value of each variable on the number line. Strawberries, butterscotch.



2.2 Continuous data:

→ Continuous data is data that can be calculated. It has an infinite number of probable values that can be selected within a given specific range.

eg: Temperature Range.

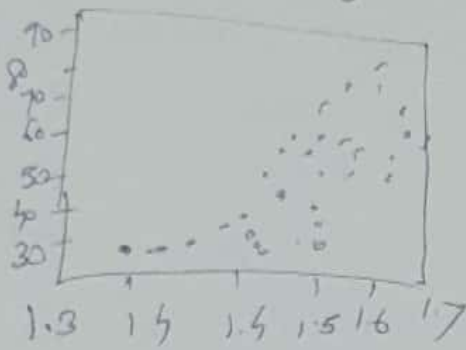
→ The amount of time required to complete a project
The height of children.

Frequency table.

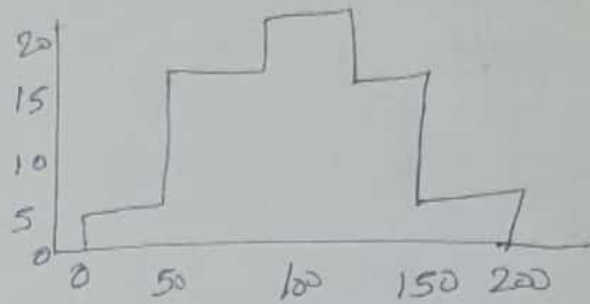
Flavour	Tally marks	Number of children
butterscotch		4
Chocolate		5
Vanilla		3
Strawberries		2

Values are represented by tally marks.

Histogram & scatter plot to give continuous data these graphs are designed to handle values



height & weight are continuous variables



Annual Revenue

Difference between Discrete & Continuous Data.

Discrete Data	Continuous Data
<ul style="list-style-type: none"> ① Values can take only specific values ② We can count the values in a discrete data series ③ Whole numbers are used to write discrete data values 	<ul style="list-style-type: none"> ① Variable can take any value within a range. ② We cannot count the values in a continuous data series ③ Real numbers are used to write continuous data values.

Types of Variable...

these are four types of variables.

- i) Numerical
- ii) Categorical
- iii) Date time
- iv) Mixed

Numerical Variable...

This category of variables deals with the numbers only. It is subdivided into

- i) Discrete
- ii) Continuous.

Categorical Variables.

This category deals with the categories.

- i) Ordinal
- ii) Nominal

Date & Time Variable

This category of variables deals with the date & time aspects. This category can obtain types of values.

- ⇒ only having date
- ⇒ only having Time
- ⇒ Having both date & time.

Mixed Variables.
It deals with the collection of multiple values for the multiple observation of a specific variable.

- i) Number or label / strings in different observation
e.g. performance of a student (CPI, percentage, marks)
- ii) Number or label / strings in the same observation
e.g. car number, value representation number.

Independent Variable.
is a singular characteristic that the other variables in your experiment cannot change. e.g. Age.

Dependent Variable.
It starts on and can be changed by other components.
e.g. grade or exam, time for a specific student (dependent)

2.3. Describing data with table and graphs.

There are 4 ways to represent the data that are:

① Tables ② Pictograph ③ Line graph ④ Bar graph

Tables are the simplest way to represent data

A pictograph uses images to represent a certain number of items.

A line graph plots individual data points as dots and connects them with lines.

Bar graph use bars of different heights to represent data

Table (Frequency Distribution)

Organization of Data.

→ Statistics refers to the collection, organization, distribution, & interpretation of data on a set of observations

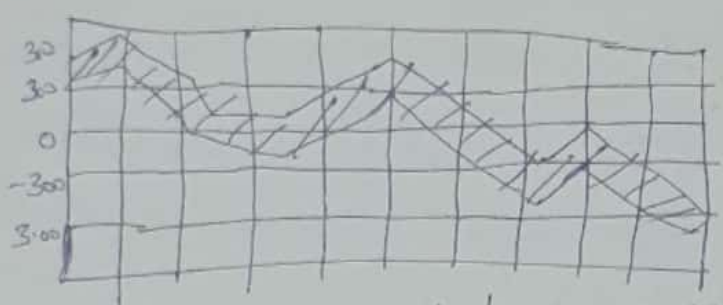
→ A frequency distribution is a collection of observations produced by sorting observations into classes and showing their frequencies (f) of occurrence in each class

→ there are two frequency distributions they are.

Types of frequency distribution.

- i) ungrouped frequency distribution
- ii) grouped frequency distribution
- iii) Cumulative frequency distribution
- iv) Relative frequency distribution
- v) Relative cumulative frequency distribution

Ungrouped frequency distribution



→ frequency distribution provided with raw data or random values. In this case, we may have to make either for individual observations or class intervals.

→ Let the test scores of all 20 students be as follows: 23, 26, 11, 18, 09, 21, 23, 30, 30, 11, 20, 11, 13, 23, 11, 29, 25, 21, 26.

Marks obtained in the test	Tally marks	No. of students (Frequency)
09		1
11		4
13		1
18		1
20		1
21		2
22		1
23		3
25		1
26		3
29		1
30		1
total		20

Grouped Frequency Distribution

It shows the scores by grouping the observations into intervals and then lists these intervals in the frequency distribution table. The intervals in grouped frequency distribution are called class intervals or limits.

→ Grouped of data in the form of class intervals to tally the frequency for the data that belongs to that particular class interval. The lower number in a class interval is called lower limit. The higher number is upper limit.

Marks obtained in the test (class-interval)	No. of students (frequency)
0-5	3
5-10	11
10-15	38
15-20	34
20-25	9
25-30	5
Total	100

Cumulative Frequency Distribution

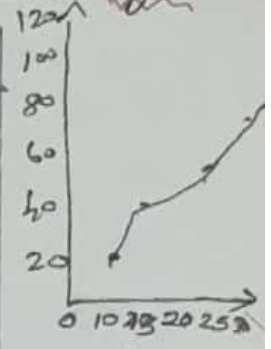
The cumulative frequency is the total of frequencies in which the frequency of first class interval is added to the frequency of the second class interval and then the sum is added to the frequency.

- Steps to construct less than C.F.
- ① Mark the upper limits on the horizontal axis or x axis.
 - ② Mark the C.F. on the vertical axis or y axis.
 - ③ Plot the point (x, y) in the coordinate plane where x represents the upper limit value.

- ④ Finally, join the points and draw the smooth curve.
- ⑤ The curve so obtained gives a cumulative frequency distribution graph of less than type.
- ⑥ To draw the frequency distribution of more than type.

Table cumulative frequency distribution table of more than type.

Level of Error	Age group class intervals	Age Class	Number of Participants	Cumulative frequency
Level 1	10-15	Less than 15	30	20
Level 2	15-20	Less than 20	32	52
Level 3	20-25	Less than 25	16	70
Level 4	25-30	Less than 30	30	100



Sum is added to the frequency of the third class interval and so on.

- Types
- i) Less than
 - ii) Greater than

Cumulative frequency begins from lowest to highest. Cumulative frequency total frequencies starting from lowest class.

Relative frequency distribution shows the proportion of the total number of observations associated with each value of values and is related to a probability distribution.

Height	frequency	Rel. frequency.
57 or less	1	0.025
57.1 to 58.6	1	0.025
58.6 to 60.1	3	0.075
60.1 to 61.7	6	0.15
61.7 to 63.3	8	0.2
63.3 to 64.8	11	0.275
64.8 to 66.4	3	0.075
66.4 to 68.0	7	0.175
Total	40	1

Cumulative Relative Frequency.

It can be found by dividing the cumulative frequencies of each interval by the total number of observations.

$$C.R.F.D = \frac{\text{Cumulative frequency of each interval}}{\text{Total number of observations}}$$

eg) following datasets.

{ 1, 1, 1, 1, 1, 1, 3, 3, 3, 3, 3, 5, 5, 5, 5, 5, 5, 11, 11, 11, 11, 11, 11, 11, 11, 11 }

first one \rightarrow find frequencies table then we find its cumulative relative frequency

Count	frequency	Cumulative frequency	Cumulative relative frequency.
1	5	5	$5/25 = 0.2$ 20%
3	5	10	$10/25 = 0.4$ 40%
5	6	16	$16/25 = 0.64$ 64%
7	1	17	$17/25 = 0.68$ 68%
11	8	25	$25/25 = 1.0$ 100%

Qualitative data from frequency table are presented by histogram.

Histogram:

It is used to visualize the distribution of data among different intervals as a series of vertical bars

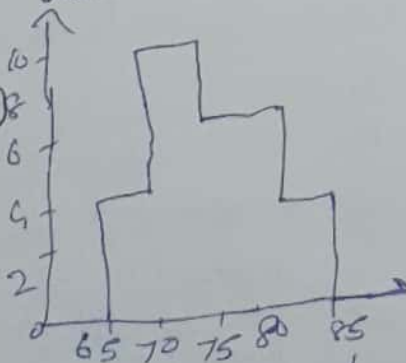
eg) Construct a histogram for the following frequency distribution table distribute the frequency of weight of

25 Student in a class.

Weight frequency

(No. of student)

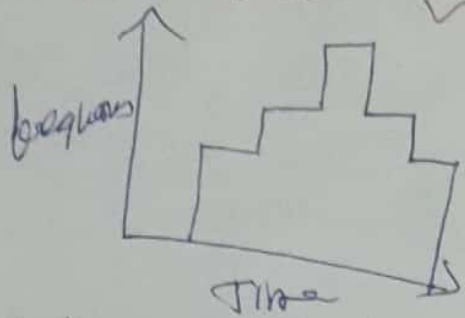
65-70	4
70-75	10
75-80	8
80-85	4



horizontal axis \rightarrow Start from 65 not from 0
unit = 15

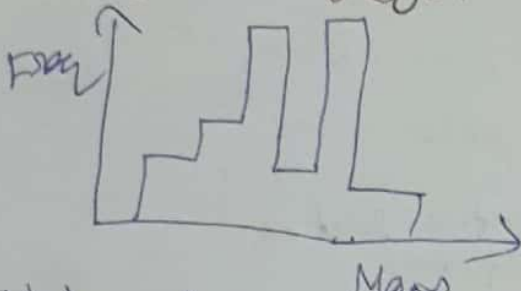
Vertical axis = frequency varies from 4 to 10 unit = 2

1) Bees shaped histogram



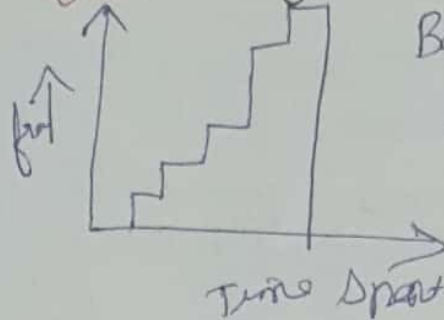
→ It has a single peak
→ has only one peak at the time.

2) Bimodal histogram



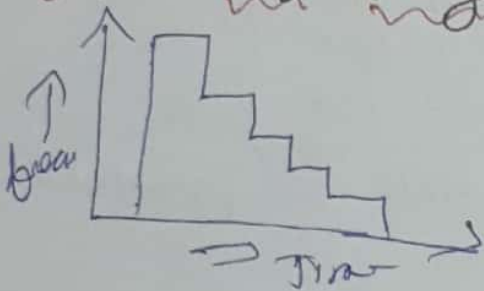
It has two peaks

3) Skewed Right Histogram

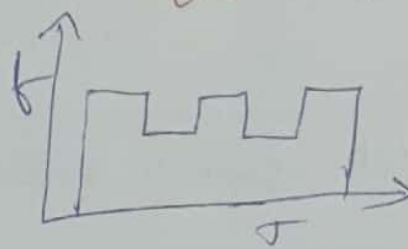


Bars of the histogram are skewed to Right.

4) Skewed Left histogram



5) Uniform histogram



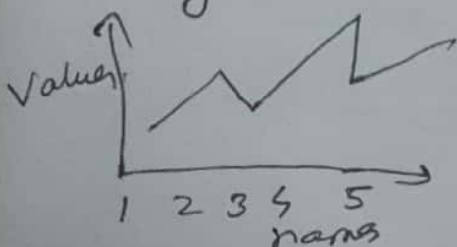
All the bars are same length

Picto graph:

Interval representation of data using images
keys or symbols.

Flavour	Number of children	Key: ○ Represents 1 child
Chocolate	○○△	$2 \times 4 + \frac{1}{4} \times \frac{1}{4} = 9$
butter scotch	○○○○△	$8 + 1 = 9$
Vanilla	○○○	
Strawberry	○○	

Line Graph

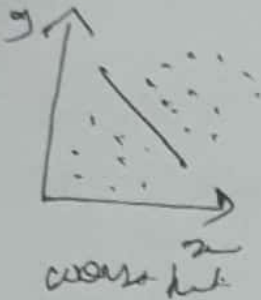


Months	1	2	3	4	5
Runs	10	30	50	20	40

x axis → name
y axis → values

Scatter Graph

A scatter xy plot has points that show the relationship between two sets of data

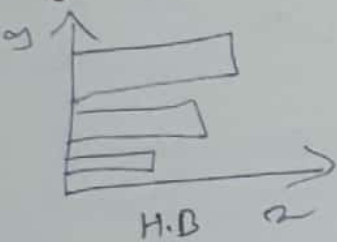


The line is drawn which is nearest to almost all the points.

Bar Graph

It is the pictorial representation of data in the form of vertical & horizontal bars

Types of Bar Graph:



H.B



V.B

- i) horizontal bar chart
- ii) vertical bar chart

2. Describing Data with Averages.

Different measures can be computed to describe how data are distributed. There are 3 categories.

- i) Measures of location
- ii) Measures of variation
- iii) Measures of shape.

→ location measure is understood as the "center" of the data eg) Average price of a new product middle of a new product, or the most frequent price of a new product.

→ Most widely used measures of the data "center" are the

- ⇒ Mean (average)
- ⇒ Median
- ⇒ Mode.

Average:

Average represents the whole value and it lies between the minimum & maximum value of the data. Hence it is the representative figure of the entire data.

Definition:

Clark and Sekhade defines as "average is an attempt to find one single figure to describe whole of figures."

Objective of an average.

- The representative value explains the character of the whole data. Its one value may represent even thousand and millions of values.
- policy decision may be taken with the help of one representative figure.
- It helps to compare with other data.

Characteristics of an average

- Everyone can easily understand the single representative value.
- It is very simple to calculate.
- It should represent the entire data.

Types of Average.

- 1) Arithmetic Mean
 - Simple Arithmetic Mean
 - Weighted Arithmetic Mean
- 2) Median
- 3) Mode
- 4) Geometric Mean
- 5) Harmonic Mean

Arithmetic Mean
The Mean of data indicates an average of the given collection of data.

It is equal to the sum of all the values in the graph data divided by the total number of values $\bar{x} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n}$

$$= \frac{\sum_{i=1}^n x_i}{n}$$

Calculating the Mean when the frequencies of the observations is given.

$$\bar{x} = \frac{f_1 x_1 + f_2 x_2 + f_3 x_3 + \dots + f_n x_n}{f_1 + f_2 + f_3 + \dots + f_n} = \frac{\sum f_i x_i}{\sum f_i}$$

In a class 30 students, marks obtained by students in mathematics out of 50 is tabulated below. Calculate the Mean of the data.

Marks obtained	Number of student	Class mark	$f_i x_i$
10-20	5	15	75
20-30	5	25	125
30-40	8	35	280
40-50	12	45	540
Total	$\sum f_i = 30$		$\sum f_i x_i = 1020$

$$\bar{x} = \frac{\sum f_i x_i}{\sum f_i} = \frac{1020}{30} = 34$$

The Mean of the given data is 34.

Median: Median is the middle most value in a given dataset after it is arranged in ascending order.

If the number of items in the dataset is odd $\left[\frac{n+1}{2} \right]^{\text{th}}$ term

$$\text{even} = \left[\frac{[n/2]^{\text{th}} \text{ term} + [n/2 + 1]^{\text{th}} \text{ term}}{2} \right]$$

$n \rightarrow$ total number of observations

Find the Median of the following data 48, 20, 50, 67, 73
 Arranged in ascending order no of observations = 5
 \Rightarrow Median of odd data = $\left[\frac{n+1}{2} \right] = \left[\frac{5+1}{2} \right] = \frac{6}{2} = 3$

Median is 3 observation

$$\boxed{\text{Median} = 50}$$

Find the median of the above grouped data, to find median, we need cumulative frequency.

Class intervals	No of girls	Cumulative frequency
120-130	2	2
130-140	8	$2+8=10$
140-150	12	$10+12=22$
150-160	20	$20+22=42$
160-170	8	$42+8=50$

$$n = \text{sum of Cf} = 50 \quad n/2 = 25$$

N : median class = 150-160

$$\text{low } l = 150, \quad c = 22, \quad b = 20, \quad h = 10$$

$$\text{Median} = l + \left[\frac{(n/2 - c)}{f} \right] \times h$$

$$= 150 + \left[\frac{(50/2 - 22)}{20} \right] \times 10 = 150 + 1.5$$

$$\text{Median} = 151.5$$

Ques:

Weekly Expenditure	0-1000	1000-2000	2000-3000	3000-4000	4000-5000	Total
No of families	34	12	43	60	51	200

Mode:

It is the value which has the maximum frequency
 It is possible to have more than one value
 which has the same maximum frequency.

Find the mode of ungrouped data.

Wickets taken by a bowler in 10 cricket matches are
 2, 6, 4, 5, 0, 9, 1, 3, 2, 3

to find the mode of ungrouped data.

Number of wickets	0	1	2	3	4	5	6
Number of matches	1	1	3	2	1	1	1

maximum no of matches = 3

Mode of the given data = 2

$$\text{Mode} = f + \left(\frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right) \times h$$

Where,

- f is the frequency of the modal class
- f_0 is the frequency of the class preceding the modal class
- f_2 is the frequency of the class succeeding the modal class
- h is the size of the class interval
- f_1 is the lower limit of the modal class

It is a ungrouped data.

mode of grouped data.

no of families	7	8	2	2	1
class interval	1-3	3-5	5-7	7-9	9-11
	f_0	f_1	f_2		

maximum class frequency = 8

Class interval = 3-5
 lower limit of the modal class = 3
 class size (or) size of the interval $h = 2$
 frequency of the modal class $f_1 = 8$

→ frequency of the class preceding to modal class $f_0 = 7$
 → frequency of the class successor of modal class $f_2 = 2$
 Know that the formula to find the mode of the grouped data.

$$\text{Mode} = l + \left(\frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right) \times h$$

$$\text{Mode} = 3 + \left(\frac{8 - 7}{2(8) - 7 - 2} \right) \times 2 \Rightarrow 3 + \left(\frac{1}{16 - 7 - 2} \right) \times 2$$

$$= 3 + \left(\frac{1}{7} \times 2 \right) \Rightarrow 3 + \frac{2}{7} = \frac{21 + 2}{7}$$

$$= 23\frac{2}{7}$$

3.5 Describing variability:

→ Variability (also called spread or dispersion) refers to how spread out a set of data is.
 → Variability gives you a way to describe how much data sets vary and allows you to use statistics to compare your data to other sets of data. The four main ways to describe variability.

- ① Range
- ② Interquartile Range
- ③ Variance
- ④ Standard deviation

Range: The Range is the simplest method of judging dispersion. It is defined as the difference between the value of the largest item and the value of the smallest item included in the distribution.

$$\text{Range} = L - S$$

$L \rightarrow$ Largest Item
 $S \rightarrow$ Smaller item

The relative measure of corresponding to the range is called its coefficient of range.

$$\text{Coefficient of range} = \frac{L-S}{L+S}$$

eg) the following are the prices of shares of ABC company Ltd, from Monday to Saturday, Monday, Tuesday, Wednesday, Thursday, Friday, Saturday.

Price: 200, 210, 208, 160, 220, 250.

Calculate Range and its coefficient

$$\text{Range} = L - S$$

$$\text{Range} = 250 - 160 = 90$$

$$\text{Coefficient of Range} = \frac{L-S}{L+S} = \frac{250-160}{250+160} = \frac{90}{410} = 0.22$$

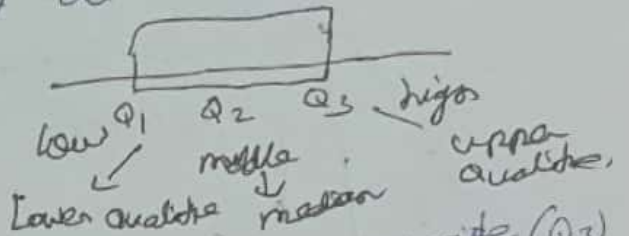
Interquartile Range (IQR)

→ It is the range for the middle 50 percent of the scores. They are values that divide the data set into 4 equal parts.

$$IQR = Q_3 - Q_1$$

→ Interquartile range gives you the spread of the middle of your distribution.

→ The interquartile range is the third quartile (Q_3) minus the first quartile (Q_1). It will provide the range of the middle half of a dataset.



$$\text{Quartile deviation } QD = \frac{Q_3 - Q_1}{2}$$

$$\text{Coefficient } Q = \frac{Q_3 - Q_1}{Q_3 + Q_1}$$

- ⇒ Lowest value
- ⇒ Q_1 : 25th percentile
- ⇒ Q_2 : the median
- ⇒ Q_3 : 75th percentile
- ⇒ Higher value (Q_4)

eg) find out the value of QD and its coefficient from the following data.

Row No: 1, 2, 3, 4, 5, 6, 7

Marks: 20, 28, 40, 10, 30, 15, 50

Calculating Quartile Deviation

$n = 7$

Marks arranged in ascending data: 10, 15, 20, 28, 30, 40, 50

$$Q_1 = \frac{1}{4}(n+1)\text{th term} = \frac{1}{4}(7+1) = \frac{3}{4}$$

odd = $n+1$ term
even = n term

$Q_1 = 2^{\text{th}}$ term is (15)

$Q_3 = 3/4 (N+1)^{\text{th}} = 3/4 \times 8 = 6^{\text{th}}$ term is (40)

$QD = (40 - 15) / 2 = 25/2 = 12.5$

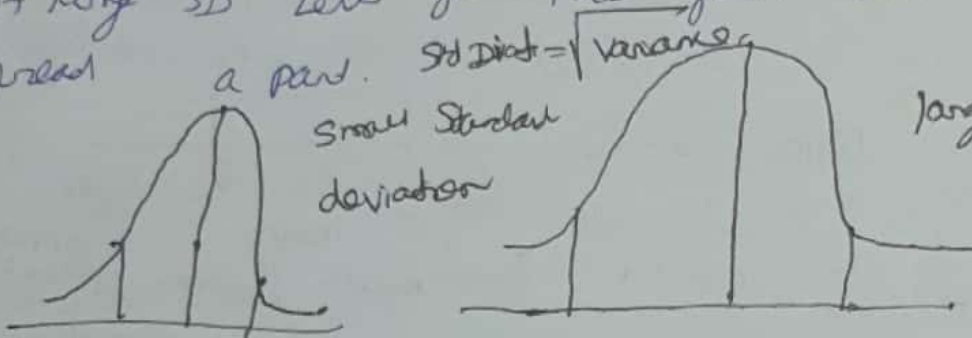
Coefficient of QD = $\frac{Q_3 - Q_1}{Q_3 + Q_1} = \frac{40 - 15}{40 + 15} = \frac{25}{55} = 0.45$

Standard deviation:

The std deviation tells you how tightly your data is clustered around the mean (the average)

A small SD indicates that your data is tightly clustered

A large SD tells you that your data is more spread



Standard Deviation formula for Population & sample.

Formula Explanation

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

- σ - Population Std deviation
- \sum - sum of
- x - each value
- μ - population mean
- N - Number of values in the population

Standard deviation

formula for sample.

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

- s - sample standard deviation
- \sum - sum of
- x - each value
- \bar{x} - sample mean
- n - number of values in the sample.

Variance:

It is the square root of standard deviation
 variance of a data set is how spread out your data

A small no. of variance means:
 dataset is tightly clustered together

A large no. for the variance means:
 the values in the dataset is more spread a part

The mean of all squared deviations squares

$$\sigma^2 = \frac{S^2}{n} = \frac{\text{sum of squares}}{n}, \text{ Standard Deviation} = \sqrt{\text{variance}}$$

Variance formula for populations → entire or complete set

$$\sigma^2 = \frac{\sum (x - \mu)^2}{n}$$

σ^2 - population variance
 \sum - sum of each value
 μ - population mean
 n - no. of values in the population

Variance formula for samples → subset of population

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

s^2 - sample variance
 \sum → sum of each value
 \bar{x} = sample mean
 n - number of values in the sample

data values are
 7, 11, 11, 15,

arranged in ascending order
 7, 11, 11, 15, 20, 20, 28

x	\bar{x}	$x - \bar{x}$	$(x - \bar{x})^2$
7	16	-9	81
11	16	-5	25
11	16	-5	25
15	16	-1	1
20	16	4	16
20	16	4	16
28	16	12	144

$$\text{Mean} = \bar{x} = \frac{\sum x}{n}$$

$$\bar{x} = \frac{7+11+11+15+20+20+28}{7}$$

$$\bar{x} = \frac{112}{7} = 16$$

Variance formula for sample

$$s^2 = \frac{\sum (x - \bar{x})^2}{n-1} = \frac{308}{7-1} = \frac{308}{6}$$

$$s^2 = 51.3333$$

$$\sum (x - \bar{x})^2 = 308$$

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

$$\sigma = \sqrt{\frac{308}{6}} = \sqrt{51.333}$$

$$\sigma = 7.165$$

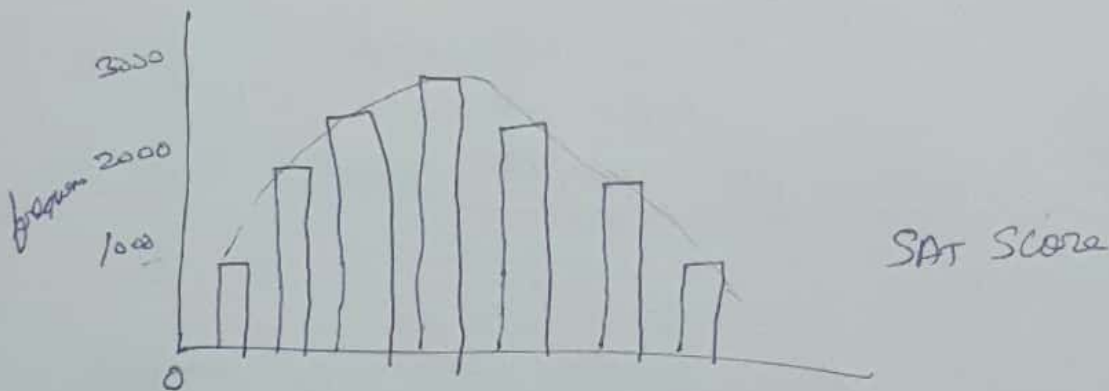
(A)

2.6. Normal Distributions and Standard (z) scores.

→ Normal distribution also known as Gaussian distribution, is a probability distribution that is symmetric about the mean, showing the data near the mean are more frequent in occurrence than data far from the mean.

→ In graphical form the normal distribution appears as a "bell curve".

→ In a normal distribution, data is symmetrically distributed with no skew. When plotted on a graph the data follows a bell shape, with most values clustering around a central region.

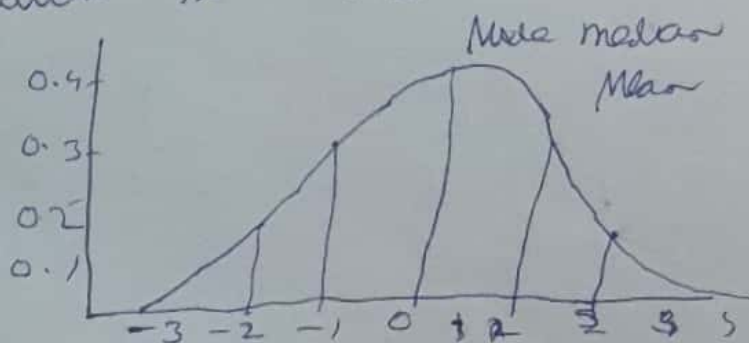


Properties of Normal Distribution

→ The mean, median, mode are exactly the same

→ The distribution is symmetric about the mean, half the values fall below the mean and half above the mean

→ The distribution can be described by two values: the mean and the standard deviation.



examples.

① Movie ratings reflect ordinal measurement because they can be ordered from most to least restrictive NC17, R, PG13 and G, the ratings of some films shown as follows.

PG	PG	PG	PG13	G
G	PG13	R	PG	PG
R	PG	R	PG	R
NC17	NC17	PG	G	PG13

- a) construct freq distribution
- b) convert to relative frequencies, exp as %
- c) construct a cumulative dist

d) find the approx percentiles rank for these films with a PG rating g.

Rating	(a) frequencies	(b) Rel. frequencies	(c) CF	d. CF %
NC17	2	$\frac{2}{20} \times 100 = 10$	$18 + 2 = 20$	$\frac{20}{20} \times 100 = 100$
R	4	$\frac{4}{20} \times 100 = 20$	$14 + 4 = 18$	$\frac{18}{20} \times 100 = 90$
PG13	3	$\frac{3}{20} \times 100 = 15$	$11 + 3 = 14$	$\frac{14}{20} \times 100 = 70$
PG	8	$\frac{8}{20} \times 100 = 40$	$3 + 8 = 11$	$\frac{11}{20} \times 100 = 55$
G	3	$\frac{3}{20} \times 100 = 15$	3	$\frac{3}{20} \times 100 = 15$
	$\Sigma f = 20$			

② 2, 17, 5, 5, 3, 28, 7, 5, 8, 5, 6, 2, 12, 10, 4, 3

1) find Mean, median mode.

2) without constructing a freq distribution or graph would you characterize the shape of the dist as balanced, truly skewed or -modskew.

Sol Sort the list -

2, 2, 3, 3, 4, 5, 5, 5, 6, 7, 8, 10, 12, 17, 28

$$\text{Mean} = \frac{\Sigma x}{n} = \frac{2+2+3+3+4+5+5+5+6+7+8+10+12+17+28}{15}$$

Mean = 7.8

① Mode = 5

Median \rightarrow odd $\Rightarrow \frac{(n+1)}{2}$

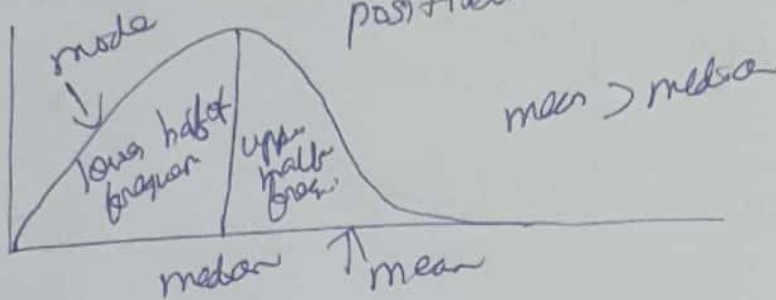
$= \frac{(15+1)}{2} = 16/2 = 8^{th} \quad (5)$

Median = 5

② Mean > median

$7.8 > 5$ \rightarrow tvels skewed.

positivells skewed distributo.



③ construct stems & leaf displas for lts group of h-wr following 10 scores obtained from a group of h-wr

old childrens.

Remains digit / unit digit

Stem Leaf

120	98	118	117	79	111
126	85	88	124	109	113
106	141	123	137	78	96
102	132	109	106	143	

seprate unit digram.

Soln

$120 \rightarrow 12-0$

$98 \rightarrow 9-8$

$118 \rightarrow 11-8$

$117 \rightarrow 11-7$

$79 \rightarrow 9-9$

$126 \rightarrow 12-6$

$85 \rightarrow 8-5$

$88 \rightarrow 8-8$

$124 \rightarrow 12-4$

$109 \rightarrow 10-9$

$113 \rightarrow 11-3$

$108 \rightarrow 10-8$

$141 \rightarrow 14-1$

$123 \rightarrow 12-3$

$137 \rightarrow 13-7$

$78 \rightarrow 7-8$

$96 \rightarrow 9-6$

$102 \rightarrow 10-2$

$132 \rightarrow 13-2$

$109 \rightarrow 10-9$

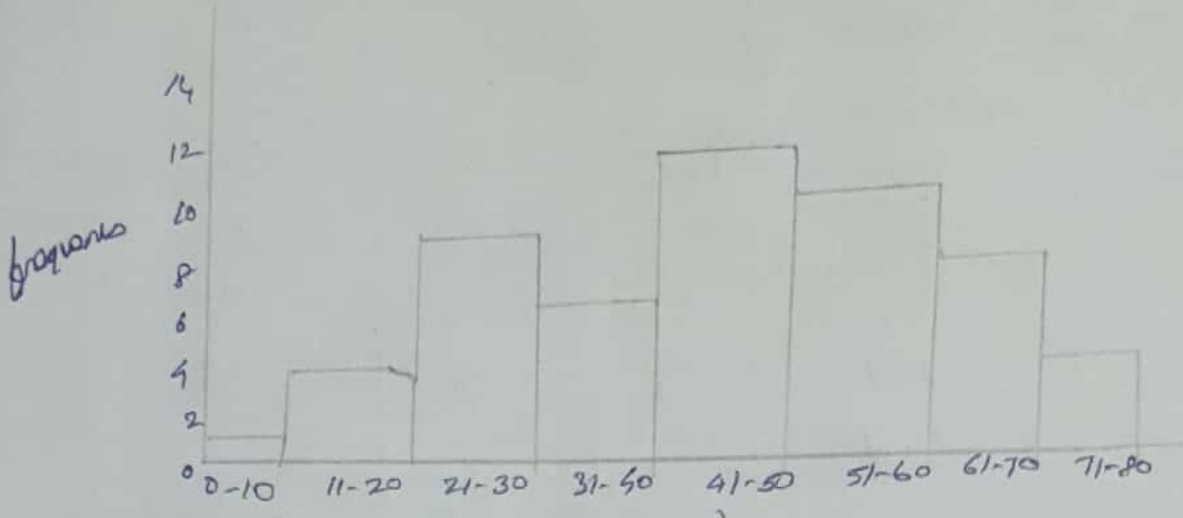
$106 \rightarrow 10-6$

$143 \rightarrow 14-3$

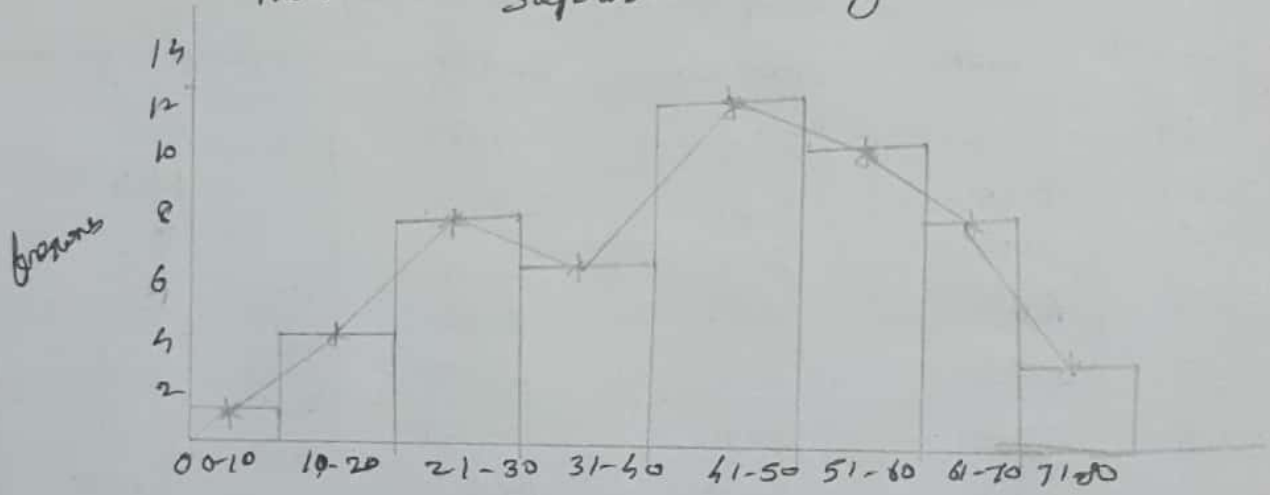
Stem	Leaf
7	8
8	5 8
9	8 9 6
10	8 2 9 6 5
11	8 7 1 3
12	0 6 3 4
13	2 7
14	1 3

④ Construct a histogram for the following data.

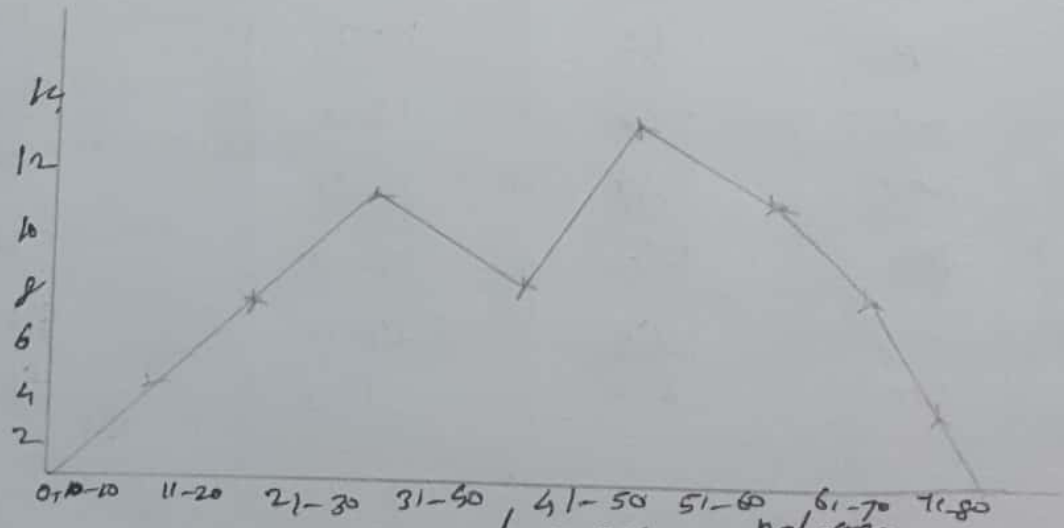
Age	0-10	11-20	21-30	31-40	41-50	51-60	61-70	71-80
frequency	1	4	8	6	12	10	7	3



histogram



connecting age midpoints.



constructing frequency polygon.

5) or Computation of population variance.

no of items (N)	observation (x)	μ	$x - \mu$	$(x - \mu)^2$
1	105	101	4	16
2	100		-1	1
3	102		1	1
4	95		-6	36
5	100		-1	1
6	98		-3	9
7	107		6	36
Σ	1429			

6) Computation of Mean, sum of squares, population variance and sample variance.

Year	Month	New No	$x - \text{mean}$	Standard Difference
2019	Jan	242	9.0	81
2019	Feb	221	-12.0	144
2019	Mar	225	-8.0	64
2019	Apr	233	0.0	0
2019	May	211	-22.0	484
2019	Jun	229	-9.0	81
2019	Jul	243	10.0	100
2019	Aug	213	-20.0	400
2019	Sep	244	11.0	121
2019	Oct	248	15.0	225
2019	Nov	249	16.0	256
2019	Dec	247	14.0	196
2020	Jan	246	13.0	169
2020	Feb	216	-17.0	289

5

Mean = 333.000

sum of squares : 2610.00

Variance of population : 186.429 $\sigma^2 = \sum_{i=1}^n \frac{(x_i - \mu)^2}{n}$

Variance of sample : 200.769 $s^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$

⑦ Computation of Standard Deviation.

No	Returns	Return - Mean	(Return - Mean) ²
1	72	12	144
2	45	-15	225
3	58	-2	4
4	84	24	576
5	60	0	0
6	10	-50	2500
7	91	31	961
8	65	5	25
9	55	-5	25
10	60	0	0
Total	600		4460
Mean	60		
Standard Deviation	22.26		

Standard deviation = $\sqrt{\text{variance}}$

The standard deviation is a rough measure of the average (or standard) amount of by which scores deviate on either side of the mean.

⑧ In a test 12% got D grade, 50% got C grade, 30% got B grade and 8% got A grade.

Grade	Frequency	Cumulative Frequency	Frequency percentage
A	16	200	8
B	60	184	30
C	100	124	50
D	24	24	12
Total	200		

Percentile Rank: Percentile Rank B: 77%

D	C	B	A
12%	50%	30%	8%

12% got D

50% got C

Half of the 30% got B

for the total percentile Rank of

$$12\% + 50\% + 15\% = 77\%$$

using the formula:

$$RR = \frac{CF - (0.5 \times F)}{N} \times 100$$

the frequency F of B = 60, the cumulative frequency of B is 184. the N value is 200. substitute these values in the above formula

$$RR = \frac{184 - (0.5 \times 60)}{200} \times 100$$

$$RR = 77\%$$

unit - III

Describing Relationships
 correlation - scatter plots - correlation coefficient
 for quantitative - computational formal for
 correlation coefficient - regression - regression line
 least square regression line - standard error
 of estimate - interpretation of r^2
 multiple regression equations - regression toward
 the mean

Correlation Definition:

→ According to Smith, correlation is the relationship that exists between two variables, when one variable value changes, the other variable changes as well (decreases or increases)
 → Mathematically we can say a function has a purpose to predict a value, by converting input (x) to output (f(x))

Types of Correlation:

- i) positive (+) and negative (-) correlation
- ii) simple correlation and multiple correlation
- iii) partial and Total correlation
- iv) Linear and non Linear correlation.

x, y - relation
 x - increase
 y - decrease
 x - decrease
 y - increase
 correlation
 x - increase, y - decrease
 x - decrease, y - increase
 no correlation (0)

Positive and negative correlation:

→ In positive correlation, increasing or decreasing the values of one variable necessarily leads to increasing (or decreasing) the values of another variable.
 ⇒ In two cases of negative correlation

(i) increasing the values of one variable certainly leads to decreasing values of another variable.
 (ii) decreasing the values of one variable leads to increasing values of another variable.

ex: Parrots - Rupee
 20 - 20
 1 - 80
 neg. walkis depend
 ↓ 50 15 katon
 120 12
 1200 1200

Positive Correlation: two variables change in the same direction

Negative Correlation: two variables change in opposite direction

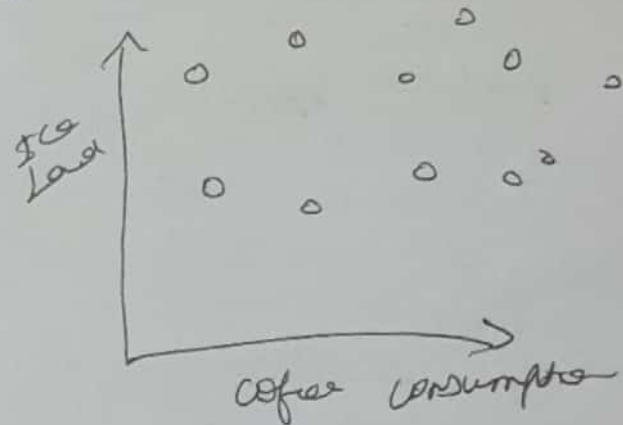
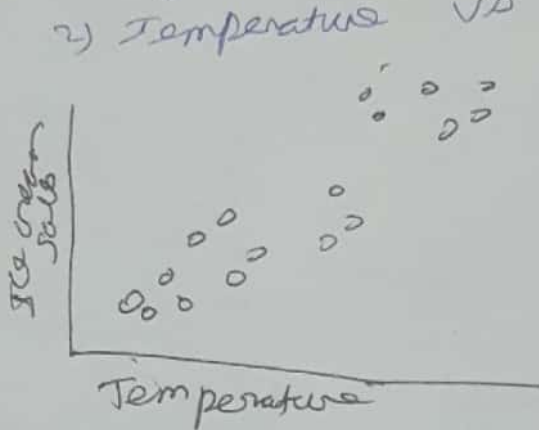
No Correlation: There is no association or relevant relationship between the two variables.

Positive Correlation . eg.

1) height vs weight

2) Temperature vs Ice cream sales

Negative Correlation



Amount of coffee and ice level has a correlation of zero

Simple and multiple correlation

Simple: When we study the relationship between only two variables, it is called simple correlation

eg: price & Demand

multiple: When we study the relationship between more than two variables, it is called multiple correlation eg: price / Demand & supply of a product

Partial and total correlation
When we study the relationship between only two variables by eliminating other variables is called partial correlation.

eg: Studying Price & Demand of a product eliminating supply

total correlation: all the facts are taken into account.

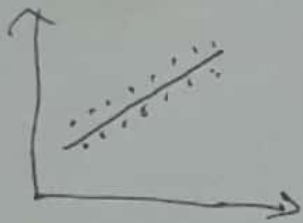
Linear and non Linear Correlation

Linear - When the ratio of changes between two variables is uniform it is called Linear

non Linear - If the amount of change in one variable is not the same like the other, it is called Non Linear correlation.

Scatter plots

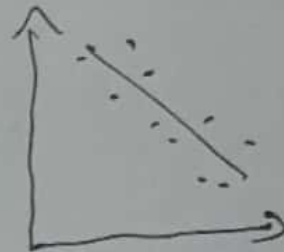
A scatter diagram is a diagram that shows the values of the variables x and y , along with the way in which these two variables relate to each other. The values of variable x are given along the horizontal axis and with the values of variable y given on the vertical axis.



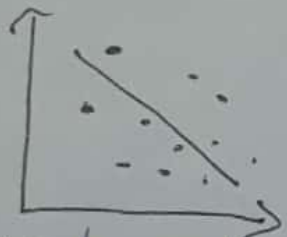
Strong positive correlation



Weak positive correlation



Strong negative correlation.



Weak negative correlation

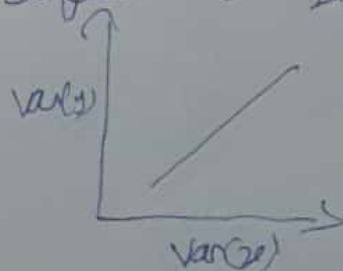


No correlation

Linear Relationship

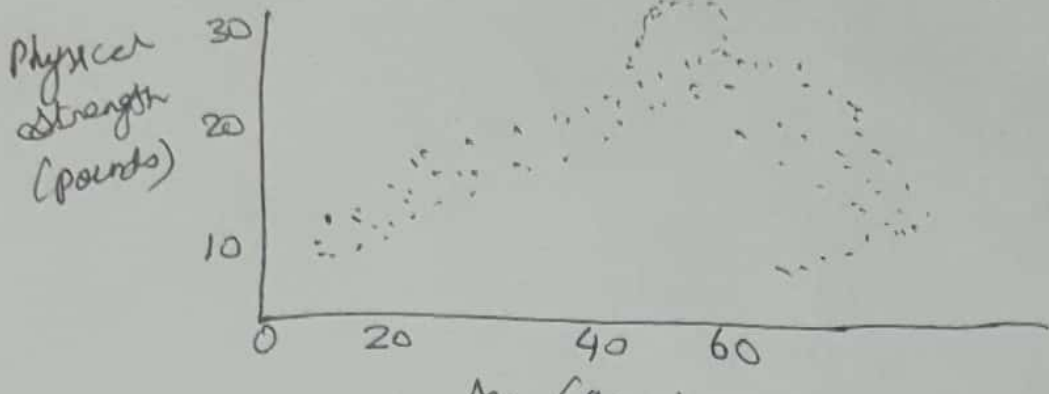
If the dot cluster approximates a straight line and it reflects a linear relationship between x & y

x	y
1	10
2	20
3	30



Curvilinear Relationship

If the dots cluster approximately a bent or curved line and therefore reflects a curvilinear relationship



Correlation Coefficient for Quantitative Data

→ Quantitative or numerical data can be discrete or continuous.

→ In statistics, correlation coefficients are a quantitative assessment that measures both the direction and strength of the tendency to vary together.

→ Correlation coefficient can be calculated using "Pearson correlation coefficient" formula, both variables are quantitative and normally distributed with no outliers.

→ Pearson's correlation coefficient is represented by the Greek letter rho (ρ) → population parameter

→ Correlation coefficient is a single number that measures both the strength and direction of the linear relationship between two continuous variables. Values can range from -1 to $+1$.

Strength:

Greater the absolute value of the Pearson correlation coefficient → stronger relationship

Direction:

The sign of the Pearson correlation coefficient represents the direction of the relationship.

Positive Coefficients.

→ It indicates that when the value of one variable increases, the value of the other variable also tends to increase.

→ positive relationships produce an upward slope on a scatterplot.

Negative Coefficients.

It represents cases when the value of one variable increases, the value of the other variable tends to decrease.

negative relationships produce a downward slope.

Pearson Correlation Coefficient Formula:

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}$$

When n = Quantity of information

$\sum x$ = total of the first variable value

$\sum y$ = total of the second variable value

$\sum xy$ = sum of the product of first & second value

$\sum x^2$ = sum of the squares of the first value

$\sum y^2$ = sum of the squares of the second value

Correlation Coefficient measures the relationship between two variables.

1 → there is a perfect linear relationship between variables

0 → there is a NO linear relationship between variables.

-1 → there is a perfect negative linear relationship between variables.

Find the least Pearson's correlation coefficient

x	x ²	y	y ²	xy
6	36	9	81	54
2	4	11	121	22
10	100	5	25	50
4	16	8	64	32
8	64	7	49	56
30	220	40	340	214

x = 6, 2, 10, 4, 8
y = 9, 11, 5, 8, 7

$$r = \frac{N \sum xy - (\sum x)(\sum y)}{\sqrt{N \sum x^2 - (\sum x)^2} \sqrt{N \sum y^2 - (\sum y)^2}}$$

$$= \frac{5 \times 214 - 30 \times 40}{\sqrt{5 \times 220 - (30)^2} \sqrt{5 \times 340 - (40)^2}} = \frac{1070 - 1200}{\sqrt{1100 - 900} \sqrt{1700 - 1600}}$$

$$= \frac{-130}{\sqrt{200} \sqrt{100}} = \frac{-130}{14.14(10)} = \frac{-130}{141.4}$$

$$= -0.9194$$

3) Computational formula for correlation coefficient are given by the four types of correlation coefficient are

- i) Pearson correlation coefficient
- ii) Linear correlation coefficient
- iii) Sample correlation coefficient
- iv) Population correlation coefficient

⇒ Correlation shows its relation between two variables. The measure of correlation coefficient shows two datasets, we can use the correlation formulas to compare them.

Pearson Correlation Coefficient Formula.

The common formula is the Pearson correlation coefficient used for linear dependence between that dataset. The value of the coefficient lies between -1 to +1. When the coefficient was down to zero, the data is considered as not related. If we get the value of +1 \rightarrow then the data are positively correlated. -1 \rightarrow Negative correlated.

$$r = \frac{n \sum(xy) - (\sum x)(\sum y)}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}$$

Where n = quantity of information

$\sum x$ \rightarrow total of the first variable value

$\sum y$ \rightarrow total of the second variable value

$\sum xy$ \rightarrow sum of the product of first and second value

$\sum x^2$ = sum of the squares of the first value

$\sum y^2$ \rightarrow sum of the squares of the second value.

Linear Correlation Coefficient Formula.

The formula for the linear correlation coefficient is given by,

$$r_{xy} = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{\sqrt{[n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2][n \sum_{i=1}^n y_i^2 - (\sum_{i=1}^n y_i)^2]}}$$

Sample Correlation Coefficient Formula.

The formula is given by

$$r_{xy} = \frac{S_{xy}}{S_x S_y} \quad \text{where } S_x \text{ \& } S_y \text{ are sample standard deviations}$$

S_{xy} - is the sample covariance. ⑦

Population Correlation Coefficient Formula
 The population correlation coefficient is given by

$$r_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$
 where σ_x and σ_y are the population standard deviations and σ_{xy} is the population covariance.

Rank Correlation: Charles Spearman developed this method of Rank Correlation coefficient where the population is not known. This method can be applied.

$$R = 1 - \frac{6 \sum D^2}{N^3 - N}$$

eg: 3 two judges in the beauty contest rank the 12 entries as follows.
 $x: 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12$
 $y: 12, 9, 6, 10, 3, 5, 4, 7, 8, 2, 11, 1$

between the two coefficient.

Soln: Calculation of Rank correlation

Rank x	Rank y	$D = R(x) - R(y)$	D^2
1	12	-11	121
2	9	-7	49
3	6	-3	9
4	10	-6	36
5	3	2	4
6	5	-1	1
7	4	3	9
8	7	1	1
9	8	1	1
10	2	8	64
11	11	0	0
12	1	11	121
$N = 12$			$\sum D^2 = 416$

$$\begin{aligned}
 R &= 1 - \frac{6 \sum D^2}{N^3 - N} \\
 &= 1 - \frac{6 \times 416}{12^3 - 12} \\
 &= 1 - \frac{2496}{1728 - 12} \\
 &= 1 - \frac{2496}{1716} \\
 &= 1 - 1.4535 \\
 &= -0.4535
 \end{aligned}$$

$$R = -0.4535$$

Regression.

→ Regression is defined as a statistical method that helps us to analyze and understand the relationship between two or more variables of interest.

→ In regression, we normally have one dependent variable and more or more independent variables.

Regression Analysis.

Regression Analysis is a branch of statistical theory that is widely used in almost all the scientific disciplines.

<u>Correlation</u>	<u>Regression</u>
Relationship between two or more variables	It is a mathematical measure showing the average relationship between variables.
x and y are random variables.	x is a Random variable y is a fixed variable
It gives limited information verifying the relationship between the variables.	It is used for the prediction of one value, in relationship to other given value.
The strength of relationship between the variables.	Its Regression value is absolute figure.
It studies linear relationship between variables.	It studies linear and non linear relationship between variables
If the coefficient of correlation is positive then the two variables are positively correlated and vice versa.	The regression line also explains that the decrease in one variable is associated with its increase in another variable.

Dependent variable:
we are trying

to understand or forecast this is the variable that

Independent Variable:

→ These are factors that influence the analysis or target variable and provide us with information regarding the relationship of the variable with the target variable.

Applications of Regression Analysis:

→ The regression analysis method of forecasting generally involves five basic applications.

- 1) Predictive Analysis
- 2) Operational Efficiency
- 3) Supporting decision
- 4) Correcting errors
- 5) New Insights.

Regression Line:

→ The Regression Line is a straight line rather than a curved line because of the linear relationship between the two variables.

→ The regression line is often referred to as the least squares regression line.

→ If we take two variables x & y we have 2 regression lines.

Regression equation:

→ The algebraic equation of the two regression lines are called regression equation.

→ Regression of x on y $\boxed{2x = a + by}$

→ The values of a and b can be calculated by solving the two normal equations.

$$\sum x = Na + b \sum y$$

$$\sum xy = a \sum y + b \sum y^2$$

Regression of y on x

$$y = a + bx$$

The values of a and b can be calculated by solving the two normal equations. (1)

$$\sum y = Na + b \sum x$$

$$\sum xy = a \sum x + b \sum x^2$$

find regression line:

x: 3, 5, 6, 8, 9, 11, y: 2, 3, 4, 6, 5, 10

x	$x - \bar{x}$	x^2	y	$y - \bar{y}$	y^2	xy
3	-4	16	2	-3	9	12
5	-2	4	3	-2	4	4
6	-1	1	4	-1	1	1
8	1	1	6	1	1	6
9	2	4	5	0	0	0
11	4	16	10	5	25	20
$\sum x = 42$	$\sum x = 0$	$\sum x^2 = 62$	$\sum y = 30$	$\sum y = 0$	$\sum y^2 = 60$	$\sum xy = 38$

$$\bar{x} = \frac{\sum x}{n} = \frac{42}{6} = 7 \quad \boxed{\bar{x} = 7}$$

$$\bar{y} = \frac{\sum y}{n} = \frac{30}{6} = 5 \quad \boxed{\bar{y} = 5}$$

Regression equation of x on y $(x - \bar{x}) = r \frac{\sigma_x}{\sigma_y} (y - \bar{y})$

$$r \times \frac{\sigma_x}{\sigma_y} = \frac{\sum xy}{\sum y^2} = 0.95$$

$$(x - 7) = 0.95 (y - 5)$$

$$x - 7 = 0.95y - 4.75$$

$$x = 0.95y - 4.75 + 7$$

$$\boxed{x = 0.95y + 2.25}$$

Regression equation of y on x $y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$

$$r \frac{\sigma_y}{\sigma_x} = \frac{\sum xy}{\sum x^2} = 0.904$$

$$y - 5 = 0.904 (x - 7)$$

$$y = 0.904x - 6.328 + 5$$

$$\boxed{y = 0.904x - 1.228} \quad (11)$$

3.7 Least Squares Regression Equation:

→ An equation pinpoints the exact least squares regression line for any scatter plot.

where $\Rightarrow y$ represents the predicted value

$\Rightarrow x$ represents the known value

$\Rightarrow b + a$ represents the calculated from the original correlation analysis.

→ using linear regression we can find the line the best "fits" our data. This line is known as the least squares regression line and it can be used to help us understand the relationship between weight and height.

The formula for the line of best fit is written as $y = b + ax$

$y \Rightarrow$ predicted value of the response variable

$b \Rightarrow$ intercept

$a \Rightarrow$ Regression coefficient

$x \Rightarrow$ value of the predictor variable.

Notice how our data points are scattered closely around the line. That's because the least squares regression line is the best fitting line for our data.

Fitting Least Squares Regression Lines.

Eg. predictor value: 140, 155, 179, 192, 200, 212

Response value: 60, 62, 67, 70, 71, 72, 75

Find a regression line using the calculated also answer

- for a person who weights 170 pounds, how tall would we expect them to be?
- for a person who weights 150 pounds, how tall would we expect that to be?

Linear Regression equation:

$\hat{y} = 32.7830 + (0.2001) * X$
the calculator automatically finds the least squares

Regression Line:

$y = 32.7830 + 0.2001x$
how to use the least squares regression line
using the least squares regression line, we can answer the questions like

a) to answer this, we can simply substitute $x = 170$ into our regression line for x and solve for y

$$y = 32.7830 + 0.2001(170) = 66.8 \text{ inches}$$

b) to answer this we can substitute for $x = 150$ into our regression line for x and solve for y .

Standard Error of estimate.

→ it is the measure of variation of an observation made around the computed regression line. Simply it is used to check the accuracy of predictions made with the regression line.

→ the regression equation helps us to predict the values of y for values of x or

→ the value of x for values of y

→ the deviation of each dot from the regression line is symbolized by $y - \hat{y}$. Thus the square root of the mean of the squared deviation.

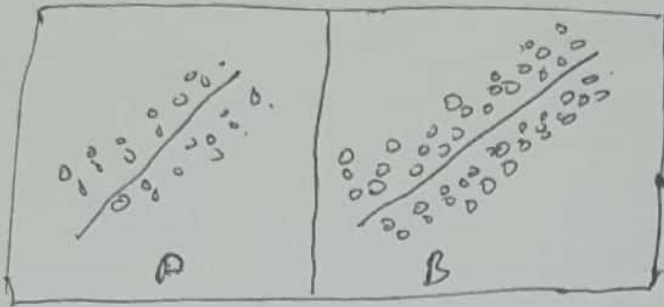
$$S_y = \sqrt{\frac{\sum (y - \hat{y}_c)^2}{N}} \quad S_x = \sqrt{\frac{\sum (x - x_c)^2}{N}} \quad \text{Similar}$$

→ In fact $(x - x_c)$ and $(y - \hat{y}_c)$ represent the unexplained variation in x and y respectively.

$$S_y = \sqrt{\frac{\text{unexplained variation in } y}{N}} \quad S_x = \sqrt{\frac{\text{unexplained variation in } x}{N}}$$

→ The computation of standard error of estimate by the above formula is quite tedious. More convenient formula.

$$\sigma_{est} = \sqrt{\frac{\sum x^2 - a \sum x - b \sum xy}{N}} \quad \sigma_y = \sqrt{\frac{\sum y^2 - a \sum y - b \sum xy}{N}}$$



Regression differing in accuracy of prediction
 → You can see that group A. The points are closer to the line than they are in Graph B, are more accurate.
 → The prediction in Graph A than in Graph B.
 → The standard error of the estimate is a measure of the accuracy of predictions.
 → The regression line is the line that minimizes the sum of squared deviations of predictions (also called the sum of squares error)

$$\sigma_{est} = \sqrt{\frac{\sum (y - y')^2}{N}}$$

σ_{est} is the standard error of the estimate
 y is an actual score
 y' is a predicted score
 N is the number pair of scores.
 $(y - y')$ is the error of prediction

$$\sigma_{est} = \sqrt{\frac{\sum (y - y')^2}{N}} \text{ for population}$$

similar formula are used when the standard error of the estimate is computed from a sample rather than a population

$$s_{est} = \sqrt{\frac{\sum (y - y')^2}{N - 2}}$$

Interpretation of r^2 multiple regression

R-squared:

R-squared (R^2 or the coefficient of determination) is the statistical measure in a regression model that determines the proportion of variance in the dependent variable that can be explained by the independent variable. In other words, R-squared shows how well the data fit the regression model.

Regression statistics

Multiple R	0.939450536
R Square	0.882548529
Adjusted R Square	0.87981709
Standard Error	207.2051602

0-100%
(0-100%)
0% = line to represent
100% = all show
(no) larger line fit
hard person.

R-squared can take any value between 0 and 1. In addition, it does not indicate the correctness of the regression model. Therefore the user should always draw conclusions about the model by analyzing R-squared together with the other variable in statistical model.

How to Calculate R-squared.

The formula for calculating R-squared is

$$R - \text{Squared} = \frac{SS_{\text{regression}}}{SS_{\text{total}}}$$

SS regression is the sum of squares due to regression (Explained sum of squares)

SS total is the total sum of squares.

Multiple Regression:

It is the statistical techniques that can be used to analyze the relationship between a single dependent variable and several independent variables.

Independent variable → Its value is independent of other variables cause.

Dependent variable → It is the effect, its value depends on change in the independent variable.

Multiple Regression analyze in case whenever we wish to model the relationship between one response variable and more than one explanatory variables.

Multiple Linear Regression formula is multiple Linear Regression

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + \dots + b_k * x_k$$

The variables in this equation are y is the predicted or expected value of the dependent variable.

x_1, x_2, x_p are three independent or predictor variables.
 b_0 is the value of y when all the independent variables are zero.
 b_1, b_2, \dots, b_k are the estimated regression coefficients.
Each regression coefficient represents the change in y relative to a one-unit change in the respective independent variable.

Public health:

If we want to predict the future spread of these illness based upon current known infectious multiple independent variables can affect number of future infections including population size, population density, air temperature, asymptomatic carriers.

Difference between Simple Linear and Multiple Linear Regression.

Simple Linear regression has only one x and one y variable. Multiple Linear regression has one y and two more x variables.

Advantages of multiple regression.

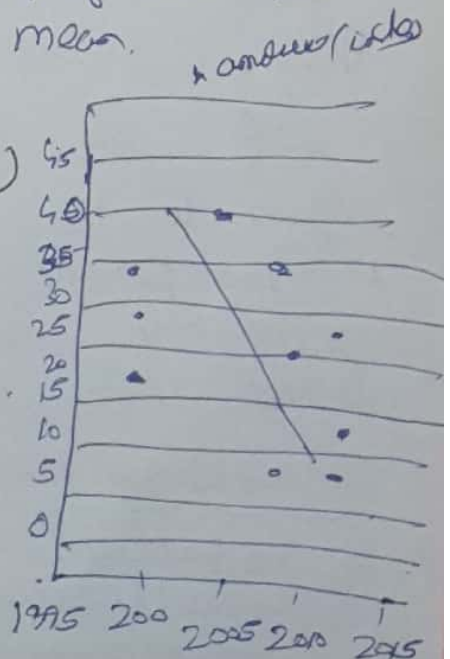
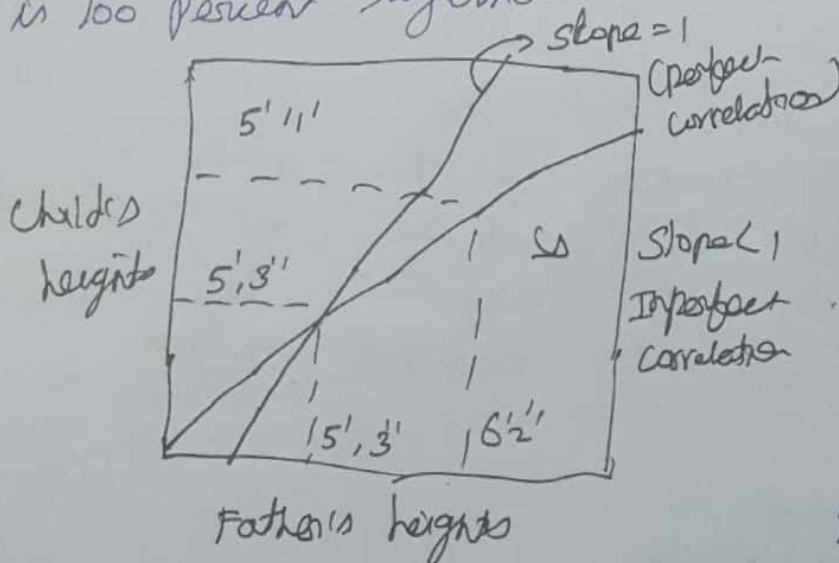
- more accurate
- Better understanding of the association of each individual factor with the outcome.

Regression to the Mean.

RTM is a statistical phenomenon that can make natural variation in repeated data look like real change. It happens when unusually large or small measurements tend to be followed by measurements that are closer to the mean.

→ mathematically, the strength of the "Regression" effect is dependent on whether or not all of the random variables are drawn from the same distribution.

→ If data has perfect correlation, it will have regression to the mean. With an r of zero, there is 100 percent regression to the mean.



Regression analysis
Linear analysis with

The fallacy:

The general statistical rule is that whenever the correlation between two variables is imperfect there will be regression to the mean.

Importance of Regression to the mean:

It is important to minimize instances of bad judgements and address the weak spot in our reasoning.

Regression to the mean fallacy:

The regression or regressive fallacy is an informal fallacy. It assumes that something has returned to normal because of corrective actions. This fails to account for natural fluctuations.

Why Regression to the mean happens:

Regression to the mean usually happens because of sampling error. A good sampling technique is to randomly sample from the population.

Percent of Regression to the mean = $100(1-r)$

⇒ Correlation coefficient & regression to Mean:

⇒ $r=1$ - correlation coefficient. If $r=1$ → perfect correlation.

Then $1-1=0$ and the regression to the mean is zero.

⇒ If your data has perfect correlation, it will regress to the mean with an r of zero. There is 100 percent regression to the mean.

⇒ In other words, data with an r of zero will always regress to the mean.

Terminologies related to regression analysis:

→ outliers

→ multicollinearity

→ heteroscedasticity

→ underfitting and

overfitting.

Normal distributions important.

All kinds of variables in natural and social sciences are normally or approximately normally distributed height, birth weight, Reading ability Job satisfaction or Entrance scores.



Normal distributions with different means.

① Outliers:

An unusual high or low value of data is called an outlier.

Start your data.

To identify outliers in your data, when you allow you to see as unusual data points within your information.

for say 3, 6, 7, 10 and 54 arranged in order you can see 54 is the largest than the rest of the data point.

2. Check your data.

To find outliers, graph your data in scatter plots or histograms.

useful for visualizing outliers, because one dot is far away from the other dots.

Histogram:

It displays data in groups called "bins".
Most of your data points, are on the Right side of the graph, and one bin of data is on the left side of the graph.
Left bin \rightarrow outlier.

3. Calculate the z-scores:

To calculate from the raw by the standard

z-scores, subtract the mean measurement and divide it

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

x - raw measurement
 μ - mean
 σ - standard deviation.

That's for you

UNIT IV PYTHON Libraries for Data Wrangling
Basic of Numpy arrays - aggregations - computations
on array - comparisons, masks, boolean logic -
 fancy indexing - structured arrays - Data manipulation
with pandas - data indexing and deletion - operating
on data - missing data - Hierarchical indexing -
 combining datasets - aggregation and grouping -
 Pivot tables.

Basic of Numpy Arrays:

Numpy - stands for Numerical python.

- It is a python library used for working with an array. Numpy is a powerful N-dimensional array object & has functions for working in domain of linear algebra, Fourier transform & matrices.
- In python lists serve the purpose of array but lists are slow to process.
- Numpy provide an array object called ndarray which is 10x faster than traditional python lists.

Why is numpy faster than lists?

- Numpy arrays are stored at one continuous place in memory. Unlike lists, do process can access & manipulate them very efficiently.
- Numpy library is partially written in C/C++ & partially written in python.
- Data manipulation in python is done with Numpy array manipulation.
- data manipulation includes access data & subarray split, reshape and join the array.

Few categories of basic array manipulations.

- ① Attributes of arrays: Determining the array shape, memory consumption & data type of array.
- ② Indexing of arrays: Creating & deleting the values of individual array elements

③ Slicing of arrays: Creating & getting smaller subarray within a larger array

④ Reshaping of arrays: Changing the shape of a given array.

⑤ Joining & splitting of arrays: Combining multiple arrays into one, and splitting one array into many.

1. Numpy Array Attributes:

Determining the size, shape, dimensions, memory consumption & data types of arrays.

Each array has the following attributes.

1. ndim - number of dimensions eg. one dimension (2D, 3D)
2. shape - size of each dimension eg (3,3) or (2,3) or (3,3,5)
3. size - ^{total} size of the array or no of elements in the array
4. dtype - data type of the array
5. itemsize - size of each item in bytes.
6. nbytes - total size of the array in bytes.

eg. sample python program.

```
import numpy as np
x1 = np.random.randint(10, size=6) # one dimensional array
x2 = np.random.randint(10, size=(3,4)) # create 3x4 2D array with random values with in the range 10
x3 = np.random.randint(10, size=(3,4,5)) # create 3D array

print ("x3 ndim : " x3. ndim)
print ("x3 shape : " x3. shape)
print ("x3 size : " x3. size)
print ("x3 dtype : " x3. dtype)
print ("x3 itemsize : " x3. itemsize)
print ("x3 total size : " x3. nbytes)

Output x3 ndim : 3
x3 shape : (3,4,5)
x3 size : 60
x3 dtype : int 64
x3 itemsize : 8
x3 total size : 480
in bytes
```


II. Array indexing.

- Accessing single element.
- The indexing specifies the desired index in square bracket.
- The indexes in NumPy array starts with '0'
- negative indexes are used to index from the end of the array.
- Creating & setting the value of array element indexing is used.

ex Creating a values.

```
import numpy as np
a1 = np.array([1, 2, 3, 4]) # Create 1D array [1, 2, 3, 4]
a11 = np.array([1, 2, 3, 6], [6, 7, 8, 9, 10])
print(a1[0]) # print the first item in a array 1
print(a11[0, 1]) # print the item on 1st row 2nd column 2
print(a1[-1]) # print the last item in a1 array 4
print(a11[1, -1]) # print the last item in a11 array 10
```

ex setting a value.

```
a1[0] = 5 # modify (or) set the first item in a1 array as '5'
a11[0, 1] = 10 # set the item in 1st row 2nd column as '10'
a1[-1] = 20 # set the last item in a1 array as '20'
```

output

```
print(a1[0]) 5
print(a11[0, 1]) 10
print(a1[-1]) 20
```

III Array slicing.

Array slicing

Syntax: `arrayname [start : stop : step]`

Start - is assumed as '0' if it is unspecified

Stop - is assumed as index of dimensions if it is unspecified

Step - is assumed as '1' if it is unspecified.

Note: if the stop value is '-ve', default values of start and stop are swapped.

eg: `print(a[::-1])` # print the values of array a in reverse order

Example

import numpy as np

a = np.array([1, 2, 3, 4, 5]) # create ndarray [1, 2, 3, 4, 5] # case 2 arrays

b = np.array([1, 3, 4, 5], [5, 6, 7, 8, 9], [9, 10, 11, 12]) [1 2 3 4] [5 6 7 8] [9 10 11 12]

print(a[:3]) # print the first 3 items of array a
Starts with index 0 and ends at 2

O/P => [1, 2, 3]

print(a[2:])

O/P => [3, 4, 5]

print(a[1:3])

O/P => [2, 3]

print(b[1:2, :3]) # print the items in first 2 rows & first 3 columns of array b

O/P => [[1, 2, 3], [5, 6, 7]]

Fancy Indexing

Access and modify the portions of arrays by using simple indices, slices, Boolean mask these are array indices

Fancy indexing is like a simple indexing, but we pass the arrays of indices in place of simple scalars. This allows us to very quickly access and modify complicated subsets of an array values.

Exploring Fancy Indexing

It means passing an array of indices to access multiple array element at once.

In [1]: import numpy as np
rand = np.random.RandomState(42)

x = rand.randn(100, size=10)

print(x)

[51, 92, 14, 71, 60, 20, 82, 86, 76, 74]
Suppose we want to access three different elements

In [2]: x[3], x[7], x[2]

Out [2]: [71, 86, 14]

Alternatively we pass a single list or array of indices to obtain the same result.

In [3] = ind = [3, 7, 4]

x [ind]

out [3] : array([71, 86, 60])

with fancy indexing.
The shape of the result reflects the shape of the index array.

In [4] = ind = np. array([3, 7], [1, 5])

x [ind]

out [4] = array([[71, 86], [60, 20]])

Structured Array.
with some compound data types.
here the structure provides efficient storage for compound heterogeneous data.

Name (string)	Age int 32	height float 64
Alice	23	173.4
Bob	35	165.5
Anto	18	175.8
Duke	28	179.4

dt = np. data type (['Name', 'U10'], ('Age', 'i4'), ('height', 'f8'])

data = np. array(['Alice', 23, 173.4], ('Bob', 35, 165.5)

('Anto', 18, 175.8), ('Duke', 28, 179.4)

print (data.dtype) → 'name' U10 'Age' i4 'height' f8

introspect table ⇒ print data

print (data['Name'])

print (np.sort (data, order = 'Age'))

data = np. zeros (4, dtype = {'name': ('Name', 'U10'), ('Age', 'i4'), ('height', 'f8')})

create a structured Array with fields such as

employee, qualification, designation and store them printed detail in the structured array normal order as well as sorted order. (5)

Name (strings)	Qualification	Salary	Designation
Jeno	BE	20000	Employee
Anu	BE	30000	Assistant
Sabin	ME	50000	Manager
Sabi	MSE	25000	Employee 2
Radha	ME	30000	Employee 3

```
dt = np. data type ([ ('Name', 'UID'), ('Salary', 'f'),
('Qualification', 'UID', 'Designation', 'UID')
data = np. array ([ ('Jeno', 20000, BE, Employee), ('Anu', 30000, BE,
Assistant), ('Sabin', 50000, Manager), ('Sabi', 25000, MSE,
Employee 2), ('Radha', 30000, ME, Employee 3)
```

print (data)

The Pandas Index Object:
 we have seen both the series and dataframe object
 contains explicit index → used to reference & modify data.
 Indexed Object → contain interesting data structure itself

It may be either
 ordered → Technically multiset, Index object may
 contains repeated values.
 Some operations available on Index object.
 eg) Let's construct an index from a list of integers

```
in [30]: ind = pd. Index ([2, 3, 5, 7, 11])
out [30]: Int 64 Index [2, 3, 5, 7, 11], dtype = 'int 64'
```

Index as Immutable Array:
 Index object in many ways operates like an
 array. eg) we can use std Python indexing notation to
 retrieve the values on slices.

```
in [31]: ind [1]
out [31]: 3
in [32]: ind [::2]
out [32]: Int 64 Index [2, 5, 11], dtype = 'int 64'
```

Index object also have many of the attributes
 familiar from NumPy Arrays.

In [33]: print (ind.size, ind.shape, ind.ndim, ind.dtype)

5 (5,) 1 int 64

Difference between Index Object and Numpy arrays

- Index are immutable
- do it cannot be modified in a normal way

In [34]: ind [1] = 0

Type Error: Index doesnot support mutable operators.

Index as operand set:

- Pandas object are designed to facilitate operators such as join across dataset. It depends on many aspect of set arithmetic.
- Index object are used in python built-in set data structure - such as unions, intersections, differences and other combinations.

In [35]: ind A = pd.Index([1, 3, 5, 7, 9])

ind B = pd.Index([2, 3, 5, 7, 11])

In [36]: ind A & ind B # intersection

out [36]: Int 64 Index([3, 5, 7], dtype='int64')

In [37]: ind A | ind B # union

out [37]: Int 64 Index([1, 2, 3, 5, 7, 9, 11], dtype='int64')

In [38]: ind A ^ ind B # symmetric difference

out [38]: Int 64 Index([1, 2, 9, 11], dtype='int64')

These operations may also be accessed via object method eg... ind A.intersection(ind B)

Data Indexing and selection.

Data selection in series:

- i) Series object act like \rightarrow one dimensional Numpy Arrays
- ii) Series object provides Dictionary style data selection mechanism. \rightarrow Standard python Dictionary.

1) Series as Dictionary:
 Like a dictionary series object provides a mapping from a collection of keys to a collection of values.

```
IN[1]: import pandas as pd
data = pd.Series [0.25, 0.5, 0.75, 1.0],
index = ['a', 'b', 'c', 'd']
```

```
data
out[1]: a 0.25      Var [start index : end index]
        b 0.5       Value
        c 0.75      Included      Value
        d 1.0       Excluded
dtype: float64      :: 2 -> Add index + 2
```

```
IN[2]: data['b']
```

```
out[2]: 0.5
```

→ we can also use dictionary like python expressions and methods to examine the key / indices and values.

```
IN[3]: 'a' in data
```

```
out[3]: True
```

```
IN[4]: data.keys()
```

```
out[4]: Index(['a', 'b', 'c', 'd'], dtype='object')
```

```
IN[5]: list(data.items())
```

```
out[5]: [('a', 0.25), ('b', 0.5), ('c', 0.75), ('d', 1.0)]
```

Series Objects:

It can be modified with a dictionary like syntax. Dictionary is extended by new key in the same way series can be extended by assigning a new index value.

```
IN[6]: data['e'] = 1.25
```

```
data
out[6]: a 0.25
        b 0.50
        c 0.75
        d 1.00
        e 1.25
```

```
data: float64
```

1) Series as one-dimensional Arrays
 provides Arrays style item selection mechanisms
 ce) & slices & masking
 & fancy indexing

In [7]: # slicing by Explicit Index

data ['a': 'c']

Out [7]: a 0.25

b 0.50

c 0.75

dtype: float64

In [8]: # slicing by implicit integer index

data [0:2]

Out [8]: a 0.25

b 0.50

dtype: float64

In [9]: # masking

data [(data > 0.3) & (data < 0.8)]

Out [9]: b 0.50

c 0.75

dtype: float64

In [10]: # fancy indexing

data [['a', 'e']]

Out [10]: a 0.25

e 1.25

dtype: float64

When you are slicing with an explicit index (data ['a': 'e']), the final index is included in the slice.

When you are slicing with an implicit index data [0:2], the final index is excluded from the slice.

Indexes, loc, iloc and ix

loc = location & row attribute to return one or more specified rows.

While a indexing operation data [2] will be use explicit index.

While a slicing operation like data [1:3] will use implicit index.

In Python code → 'Explicit is better than implicit'
 The explicit nature like `loc` and `iloc` make very useful in maintaining clean readable code.

Data Selection in DataFrame
 Data frame acts like a two dimensional or structured arrays or in other way like a dictionary.

1) data frame as a dictionary:
 It related series acts [let's consider the eg]

```
In[18]: area = pd.Series({'California': 4329, 'India': 628653, 'New York': 1497, 'Florida': 17032})
```

```
pop = pd.Series({'California': 3833251, 'India': 2644893, 'New York': 1965127, 'Florida': 19552868})
```

```
data = pd.DataFrame({'area': area, 'pop': pop})
```

data	area	pop
California	4329	3833251
Florida	17032	19552868
India	628653	2644893
New York	1497	1965127

The individual series that makes up the columns of the data frame can be accessed via dictionary-style indexing of the column name.

```
Name: area, dtype: int64
```

Equivalently we can use attribute-style access with column names.

```
In[20]: data.area
```

California	4329
Florida	17032
India	628653
New York	1497

```
Name: area, dtype: int64
```


→ Attribute style column access have the same object as the dictionary style access.

In [21]: data.area is data["area"]

out [21]: True

→ Attribute style access is not possible, when the column names are not strings, column name conflicts with the method of the data frame.

In [22]: data.pop is data["pop"]

out [22]: False

Like series object:

→ Dictionary style syntax add a new column

In [23]: data["density"] = data["pop"] / data["area"]

out [23]:

	area	pop	density
California	433987	38332521	90413921
Florida	170312	19552850	11480421
India	62652	12882135	8588276
New York	141297	19651127	13901876

ii) Dataframe

As two-dimensional Arrays

In [24]: data.values

out [24]: Array([[4.23, 3.83, 9.05], [1.7, 1.955, 1.14], [1.49, 1.28, 8.58], [1.41, 1.96, 1.39]])

Arrays like Observations on Data frame:
 → transpose for data frame.
 ↳ swap rows and columns

In [25]: data.T

	California	Florida	India	New York
area	4.23	1.70	1.49	1.41
pop	3.83	1.95	1.28	1.96
density	9.05	1.14	8.58	1.39

Indexes of columns to access to a row in a array.

In [26]: data.values [0]
 out [26]: array ([423, 3.83, 9.04])
 Passes a single 'index' to a dataframe
 across a column.

In [27]: data ['area']
 out [27]: California 423
 Florida 170
 India 149
 Newyork 191
 Name: area, dtype: int64

Array style indexing
 uses loc, iloc, ix
 iloc indexes uses implicit Python, like style index.

In [28]: data.iloc [2, :2]
 out [28]: California 423 383
 Florida 170 195
 India 149 128

In [29]: data.loc [! 'india', : 'pop']
 out [29]: California 423 383
 Florida 170 195
 India 149 128

ix index allows a hybrid of two approaches.

In [30]: data.ix [2, : 'pop']
 out [30]: California 423 383
 Florida 170 195
 India 149 128

Numpy style data access patterns also used
 up indexes like index combine masking and
 fancy indexing.

In [31]: data.loc [data.density > 1, ['pop', 'density']]
 out [31]: Florida 195 1.16
 Newyork 196 1.39

Additional Indexing Conventions:

In [33]: data ['Florida', 'India']

	area	pop	density
Florida	170	195	1.14
India	159	128	0.8

Rows while indexing refers to columns.
 Columns refers to rows.
 → Slice refers to row by number rather than by index.

Out [35]: In [34]: data [1:3]

	area	pop	density
Florida	170	195	1.14
India	159	128	0.8

Operators on Data in pandas:

Data manipulation with pandas:
 → Pandas is a high level data manipulation tool developed by Wes McKinney. It is built on the NumPy package and its key data structure is called DataFrames.

→ Data frame allows you to store and manipulate tabular data in rows of observations and columns are variables.
 → Pandas is built on top of the NumPy package, meaning a lot of structure of NumPy is used for replication in pandas. Data in pandas operations.

Operations:
 → NumPy has ability to perform basic arithmetic operations (addition, subtraction, multiple color), and with more sophisticated operators trigonometric functions, exponential, logarithmic functions.
 → Pandas inherits the functionalities from NumPy and use

"Computation on NumPy Arrays: Universal Functions"

Pandas includes:

Binary operations: Series or addition, multiplication
Pandas will automatically align indices when
performing the operations. like negation and trigonometric functions

like negation and trigonometric functions
Pandas preserve indices and column labels.

1) ufuncs: index preservation
→ Pandas is designed to work with
Numpy, Numpy ufunc with array or
Pandas Series and Dataframe objects.
→ Lets us start by defining a simple
Series and Dataframe on which to
demonstrate the

```
In[1]: import pandas as pd
→ import numpy as np
```

```
In[2]: rng = np.random.RandomState(42)
ser = pd.Series(rng.rand(5), index=[0, 1, 2, 3, 4])
ser
```

```
Out[2]:
0    6
1    8
2    7
3    5
dtype: int64
```

Pandas is the library for data
manipulation and analysis. Usually it is the
starting point for your data science tasks. It
allows you to read/write data from/to
multiple sources, process the missing data
align your data reshape it, merge and
join it with other data sources data
groups it, etc.

Missing Data:

→ Data can have missing value for a number of reasons such as observations that were not recorded and data corruption. Handling missing data is important as many machine learning algorithms do not support data with missing values.

→ You can load the dataset as a pandas DataFrame and print summary statistics on each attribute.

```
# load and summarize the dataset
dataset = read_csv('read_csv_file_name.csv', header=None)
# summarize the dataset
print(dataset.describe())
```

In Python, we mark missing values as NaN. We learn with a NaN value are ignored from operators like sum, count, etc.

Missingness:

The `fillna()` function is used to fill NaN values using the specified method. `DataFrame.fillna(value=None, method=None, axis=None, inplace=False, limit=None, downcast=None, **kwargs)`

- 1. value: It is a value that is used to fill its null value.
- 2. method: Used to fill its null value.
- 3. axis: It takes int or string value for row/column.
- 4. inplace: If it is true, it fill value at an empty place.
- 5. limit: It is an integer value that specifies what to do as like float64 to int64.

Hierarchical indexing: group is a method to create structured data. A multi index on DataFrame has more than two dimensions. A multi index adds at least one more dimension to the data. A hierarchical index has more than two dimensions. To create DataFrame with playing ratings of a few players from the Fifa 19 dataset.

```
In [1]: import pandas as pd
In [2]: data = {'position': ['GK', 'GK', 'GK', 'DF', 'DF', 'DF',
                             'MF', 'MF', 'MF', 'CF', 'CF', 'CF'],
                'Name': ['DeGroot', 'Center', 'Allis', 'Vandri', 'Ram', 'Credni',
                          'Korad', 'De', 'Ronald', 'Menni', 'Neyron'],
                'Overall': [91, 88, 90, 80, 96, 97, 92, 93, 95],
                'Rank': [1st, 3rd, 2nd, 3rd, 1st, 2nd, 3rd, 1st, 2nd, 3rd]}
In [3]: fifa19 = pd.DataFrame(data, columns = ['position', 'Name', 'Overall', 'Rank'])
```

Out [5]

	position	Name	Overall	Rank
0	GK	DeGroot	91	1st
1	GK	Center	88	3rd
2	GK	Allis	90	2nd
3	DF	Vandri	80	3rd
4	DF	Ram	96	1st
5	DF	Credni	97	2nd
6	DF	Korad	92	3rd
7	MF	keun	93	1st
8	MF	De	95	2nd
9	CF	Ronald	99	1st
10	CF	Menni	97	2nd
11	CF	Neyron	96	3rd

Combining Dataset. concat and append.
 → Data comes from different Data sources.
 → Operation involve combination of two different data set, and more complicated.
 → Pandas includes functions and methods that make some of data wrangling fast and straightforward. Simple concatenation of Series and Data frames with pd.concat function. We begin with the Standard import.

```
In[1]: import pandas as pd
import numpy as np
In[2]: def make_df(cols, ind):
    """quickly make a Dataframe"""
    data = {c: [str(i) + str(i) for i in ind]
            for c in cols}
    return pd.DataFrame(data, ind)
```

```
#example Dataframe
make_df('ABC', range(3))
```

Out[3]:

	A	B	C
0	A0	B0	C0
1	A1	B1	C1
2	A2	B2	C2

Recall: Concatenation of Numpy arrays & Dataframe
 Concatenation of objects similar to the concatenation of Numpy Arrays, done by np.concatenate of arrays into a single Array

```
In [4] x = [1, 2, 3]
y = [4, 5, 6]
z = [7, 8, 9]
np.concatenate([x, y, z])
Out [5]: array([1, 2, 3, 4, 5, 6, 7, 8, 9])
```

In [5]: $x = [1, 2], [3, 4]$

np. concatenate ([x, x], axis=1)

out [5]: array([[1, 2], [3, 4]])

Simple Pandas has a Concatenation function with pd. concat(). Similar to np. concatenate but Pandas contains several options with this.

signature in Pandas v0.18

pd.concat(objs, axis=0, join='outer', join_axes=None, ignore_index=False, keys=None, levels=None, name=None, verify_integrity=False, copy=True)

pd.concat() → It can be used for a simple concatenation of series or DataFrame objects.

np.concatenate() → It can be used for simple concatenation of arrays

In [6]: ser1 = pd.Series(['A', 'B', 'C'], index=[1, 2, 3])
ser2 = pd.Series(['D', 'E', 'F'], index=[4, 5, 6])

pd.concat([ser1, ser2])

Out [6]

1	A
2	B
3	C
4	D
5	E
6	F

dtype = object

Append () method:

Simply call df1.append(df2):
Append() method in pandas does not modify the original object. It creates a new object with the combined data.

Combining Pandas provides Merge and Join members. Join as a high performance merge operation.

pd.merge() function
Join() function

Categorize → one to one
→ many to one
→ many to many

Aggregation and Grouping

In large data set, the data are summarized or aggregated for efficiency.

aggregation \rightarrow it will return single number for a large dataset, some aggregations \rightarrow sum(), mean(), median(), max(), min(),

Simple Aggregation in pandas

For a pandas series the aggregate returns a single value

```
In [4]: mg = np.random.RandomState(42)
```

```
ser = pd.Series(rng.rand(5)) ser
```

```
Out [4]: 0    0.3
         1    0.9
         2    0.7
         3    0.5
         4    0.1
```

dtype: float64

```
In [5]: ser.sum()
```

```
Out [5]: 2.5
```

```
Out [6]: ser.mean()
```

```
Out [6]: 0.5
```

for a DataFrame, by default the aggregation can be taken in each column.

```
In [7]: df = pd.DataFrame({'A': rng.rand(5), 'B': rng.rand(5)})
```

Out [7]:

	A	B
0	0.1	0.6
1	0.5	0.9
2	0.08	0.8
3	0.6	0.2
4	0.1	0.1

Out [8]:

0	0.08
1	0.5
2	0.8
3	0.9
4	0.4

```
In [8]: df.mean()
```

Out [8]:

A	0.45
B	0.4

dtype: float64

```
In [9]: df.mean(axis='columns')
```

Some other built in pandas Aggregations.

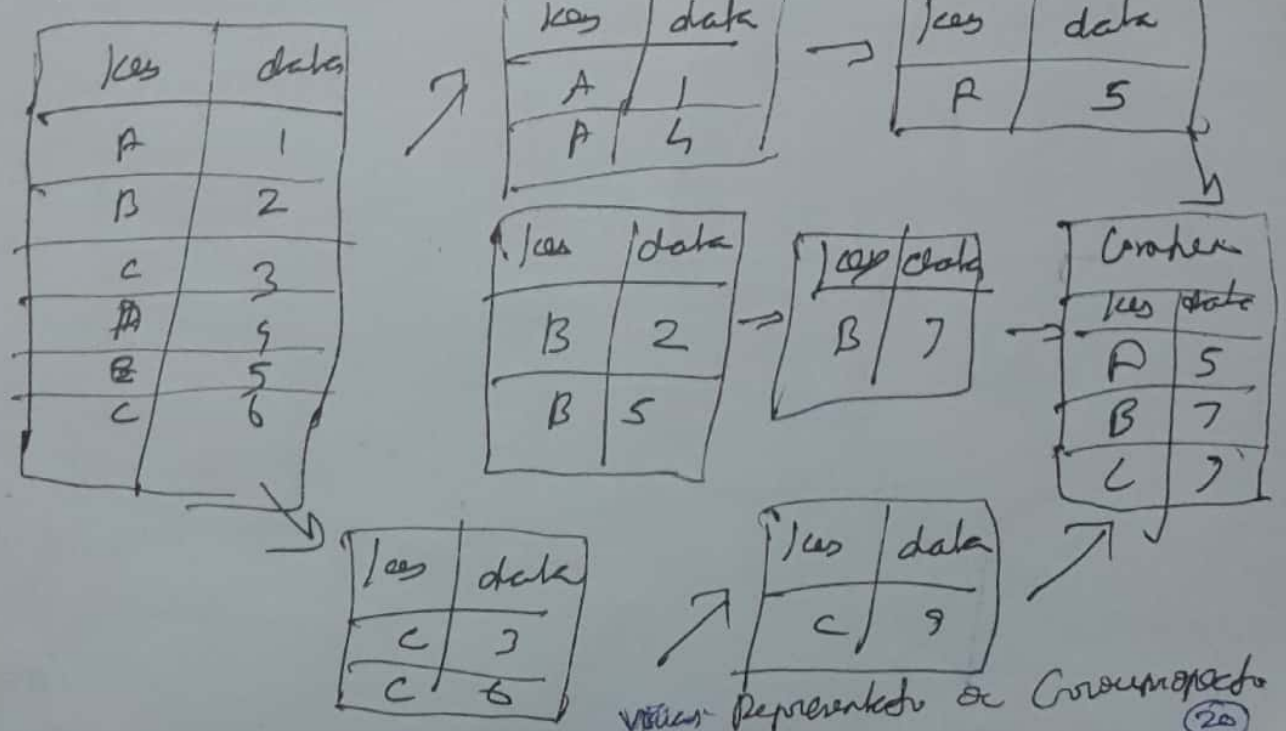
Aggregation	Description.
count	Total number of items
first(), last()	First and last items
mean(), median()	Mean and Median
min(), max()	Minimum & maximum
std(), var()	Standard deviation and variance
mad()	Mean absolute deviation
prod()	Product of all items
sum()	Sum of all items

Group by, Split, Apply, Combine.

Group by is accomplished by the step.
 a split step involve breaking up and grouping of the data frame depends on the value specified in the step.

Apply step involves computers some function usually as:
 1) Aggregate 2) transformation 3) filtering within the individual groups.

The combine step merges the result of this operation into split an order apply order.



values represented as Group objects (20)

Data Wrangling:

→ Data wrangling is the process of transforming data from its original "raw" form into a more digestible format and organizing sets from various sources into a singular coherent whole for further processing.

→ Data wrangling is also called as data munging. The primary purpose of data wrangling can be described as getting data in coherent shape. In other words, it is making raw data usable.

Data Wrangling processes:

- i) Getting the data from the various source into one place.
 - ii) Piling the data together according to the determined setting.
 - iii) Cleaning the data from the noise or erroneous, missing elements.
- "Data Wrangling is the process of cleaning & structuring and a desired format for better decision making in less time."
- Six iterative steps that make wrangling process

① DISCOVERING:

Before you can dive deep, you must better understand what is your data. What will inform how you want to analyze it, how you wrangle customer data, as may be informed by where they are located, what they bought, or what promotions they received.

② Structuring:

→ This means organizing the data, which necessitates because raw data comes in many different shapes and sizes.

→ A single column may turn into several rows for easier analysis. One column may become two. Movement of data is made for easier computation and analysis.

③ Cleaning:

→ What happens when errors and outliers skew your data? you clean the data.

→ What happens when state data is entered as AP or Andhra Pradesh or Assam instead of states? you clean the data. Null values are changed to standard formatting implemented, ultimately increases data quality.

④ Enriching: here you take stock in your data strategies about how other additional data might be. Wrangling step might be: what new types of data can be derived from what I already have? What other information.

⑤ Validating: validation rules are repetitive programming sequences that verify data consistency, quality & security. end of validation includes ensuring uniform distribution of attributes that should be distributed normally (eg. height data) or confirming accuracy of bills through a check across data.

⑥ Publishing: Analyze prepare the wrangled data for use downstream, whether a particular use or deployment requires any particular steps taken or logs to user wrangled data. undertaken for implementation of insights relies upon the ease with which it can be accessed and utilized by others.

Pivot Tables:

→ A Pivot table is an interactive way to quickly summarize huge amounts of data.
→ It can be used to analyze numerical data in detail and answer unanticipated questions about the data.
→ A pivot table is especially designed for querying large amount of data in many user friendly ways. In pandas pivot table is used to calculate, aggregate and summarize your data. It is defined as a powerful tool that aggregates data with calculations such as sum, count, average, max & min. It also allows to sort and filter data when the pivot table has been created.

Parameters:

data: A dataframe

values: It is an optional parameter

index: It refers to the column to aggregate.

Syntax:

```
pandas.pivot_table(data, value=None, index=None, columns=None, aggfunc='mean', fill_value=None, margins=False, dropna=True, margins_name='All')
```

Data: Dataframe

values: Column to aggregate, optional

index: Column, Group, axis, or list of the previous

columns: Column, Group, axis or list of the previous

aggfunc: function, list of functions, dict, default numpy.mean

fill_value [scalar, default None]: Value replace missing values with

margins [boolean, default, False]: Add all row/columns

dropna [boolean, default, True]: Do not include columns and all NaN

margins_name [string, default 'All']: Name of the row/column that will continue the total when margin is true

Consider a data frame:
 # create a simple data frame
 # importing pandas as pd
 import pandas as pd
 import numpy as np
 # creating a dataframe.

```
df = pd.DataFrame({'A': ['Alex', 'Bays', 'Sunil', 'Taylor', 'Nicks'],
                   'B': ['Master', 'Graduate', 'Graduate', 'Master', 'Graduate'],
                   'C': [27, 23, 21, 23, 29]})
```

df
OUT:

	A	B	C
0	Alex	Master	27
1	Bays	Graduate	23
2	Sunil	Graduate	21
3	Taylor	Master	23
4	Nicks	Graduate	29

Simplest Pivot table must have a dataframe
 # and an index / axis of index

```
table = pd.pivot_table(df, index = ['A', 'B'])
```

table
OUT:

	A	B	C
Alex	Master	27	
Bays	Graduate	23	
Nicks	Graduate	29	
Sunil	Graduate	21	
Taylor	Master	23	

Applications of Pivot table:
 summarizing data
 average sales for each region for each product from a product data table.
 Lists unique values in any column of table.
 creates a pivot report with sub-totals and custom formats
 the reports: filtering, sorting, drilling, down data in

Unit - V
Visualization

<p><u>Data</u> Importing Visualizing Histograms Annotation Geographic data base.</p>	<p>Matplotlib - Line plots - scatter plots - error - density and contour plot - Legends - colors - subplot - text and - customization - three dimensional plotting data with Basemap - visualization with</p>
--	---

Introduction : Matplotlib:

→ Matplotlib tool for visualization on Python. Matplotlib is a multiplatform data visualization library built on Numpy arrays and designed to work with its broader scientific stack. Main important features is ability to play well with most operating systems and graphics backends. This cross platform everything to everyone approach has been great strength of Matplotlib.

5:1 Importing Matplotlib:

→ We use np shorthand for Numpy and the plt shorthand for plots. → we will use some standard shorthand for Matplotlib imports.

```
In [1]: import matplotlib as mpl
import matplotlib.pyplot as plt
```

Setting styles
we will use the plt style directive to choose appropriate aesthetic style for our figures. we will use plt and use the class matplotlib style.

```
In [2]: plt.style.use('classic')
```

We will adjust the style as needed from the matplotlib version 1.5 \Rightarrow customizing matplotlib configurations and style sheets.

Plotting from a script if you are using Matplotlib from within a script, the function `plt.show()` opens one or more interactive windows and displays your figure as figures.

```
# ---- file: myplot.py
import matplotlib.pyplot as plt
x = np.linspace(0, 10, 100)
plt.plot(x, np.sin(x))
plt.plot(x, np.cos(x))
plt.show()
```

When you are run the script, which will result in a window opening with figure displayed.

Plotting from an IPython shell

It can be very convenient to use Matplotlib interactively within an IPython shell. IPython is built to work with Matplotlib if you speak matplotlib. \rightarrow to enable this mode you can use `%matplotlib magic` command after starting ipython.

```
In [1]: %matplotlib
```

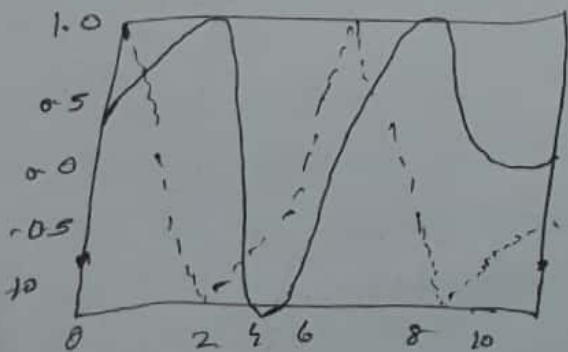
```
using matplotlib backend: TkAgg
```

```
In [2]: import matplotlib.pyplot as plt
```


In the `print` and `plt.plot` command will cause a figure window to open and further commands can be run to update the plot. Some changes (such as modifying proper ties of lines that are already drawn) will not draw automatically. To force an update, use `plt.draw()` using `plt.show()` in `matplotlib` lib mode is not required.

Plotting from an IPython Notebook
 IPython Notebook is a browser-based interactive data analysis tool that can combine narrative text, code, graphs, and executable documents. Plotting interactively within an IPython Notebook can be done with the `%matplotlib` command.

- In IPython Notebook, two possible options of embedding graphs will lead to interactive plots embedded within the notebook.
- `%matplotlib inline` will lead to static images of your plot embedded in the notebook.
- `In[3]: %matplotlib inline`
- After you run this command, the results of the plot will embed a PNG image of the figure.



```
In [5]: import numpy as np
x = np.linspace(0, 10, 10)
fig = plt.figure()
plt.plot(x, np.sin(x), '-')
plt.plot(x, np.cos(x), '--')
```

Saving figures to files

one nice feature of matplotlib is its ability to save figures in a wide variety of formats. you can save the figure using the `save_fig()` commands.
for eg) to save the previous figure as a PNG file you can run this:

```
In [5]: fig.savefig('my-figure.png')
```

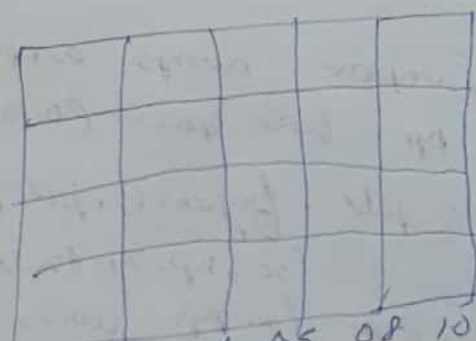
we have now have a file called `my-figure.png` in the current working directory.

```
In [6]: !ls -lh my-figure.png
```

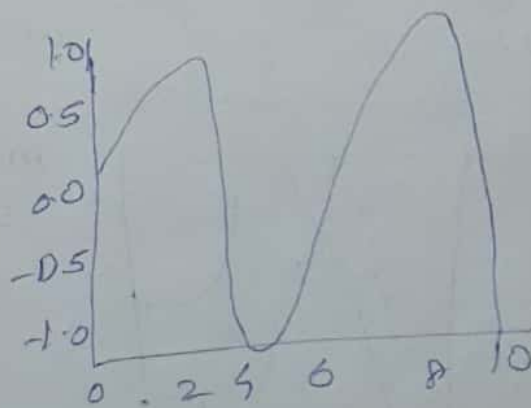
```
ms-r--r-- 1 jakevdp staff 16K Aug 11 10:59 my-figure.png
```

Let's use Python Image Object to display the contents of this file.
Simple Line plots:
The visualization of a single plot of all plots is the creation of a simple plot of function $y = \sin(x)$ for plotting and importing functions.
In [1]: `import matplotlib.pyplot as plt`
(Seaborn - Whitegrid)

we start by creating a figure and axes.
In [2]: `fig = plt.figure() .ax = plt.axes()`



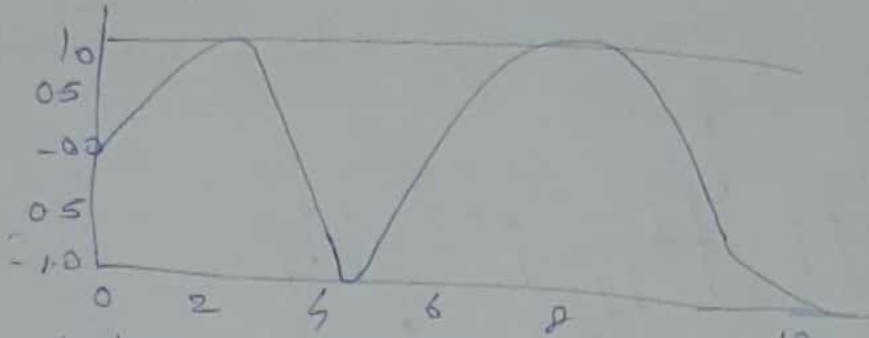
An empty gridded axes.



```

In [3]: fig = plt.figure() ax = plt.axes()
x = np.linspace(0, 10, 100) ax.plot(x, np.sin(x));
In [4]: plt.plot(x, np.sin(x))
plt.plot(x, np.cos(x));

```

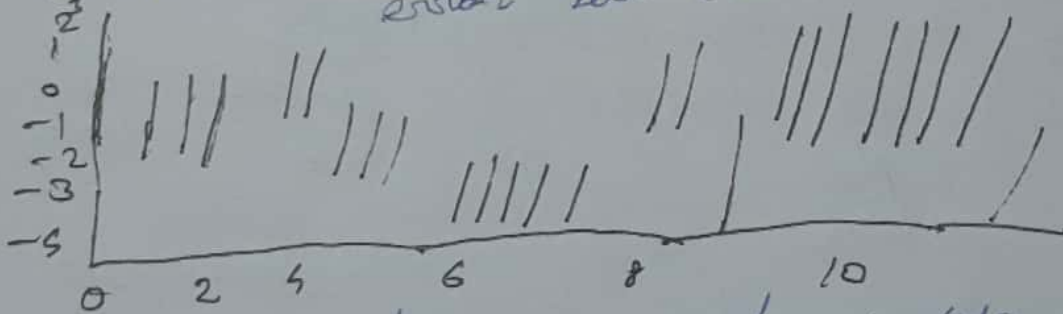


A simple standard 10^10 object-oriented interface
 Visualizing Errors:
 Showing the visualization of data and results
 comes much more effectively can make a plot
 Base Error bars. Complete information

```

In [1]: %matplotlib inline
import matplotlib.pyplot as plt plt.style.use('dark_background')
import numpy as np
In [2]: x = np.linspace(0, 10, 50)
dy = 0.8
y = np.sin(x) + dy * np.random.randn(50).plt.
error bar (x, y, yerr = dy, fmt = 'k');

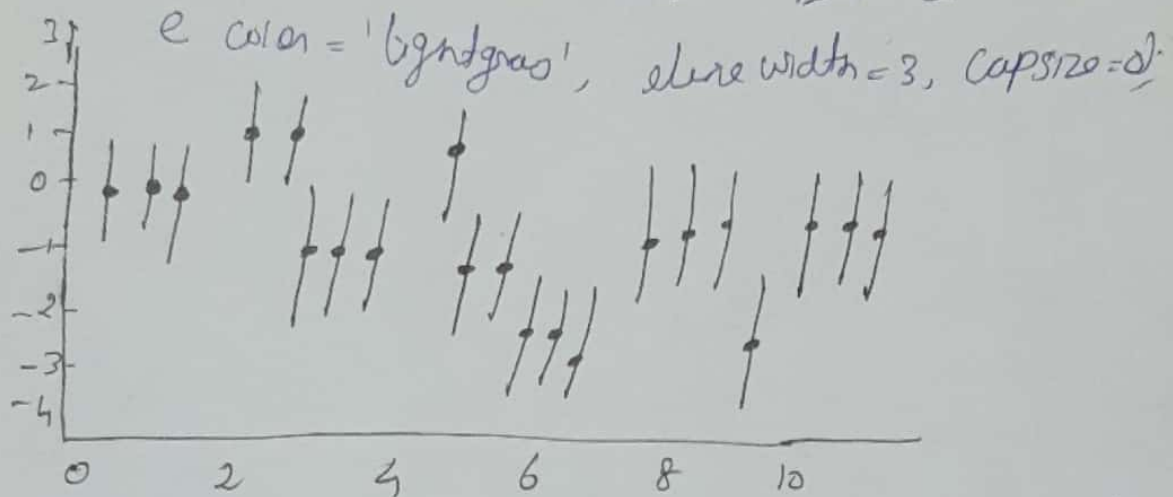
```



here the fmt is a format code controlling
 the appearance of lines and points, and
 has the same syntax as the shorthand
 used in plt.plot outlined in "Simple
 Line plots".

The errorbar function has many options to fine-tune its outputs.

IN[35]: plt.errorbar(x, y, yerr=dy, fmt='o', color='black')



In addition to these options you can specify horizontal errorbars (xerr), one-sided errorbars and many other variants for more information on the options available, refer to the documentation of plt.errorbar

Density and contour plots:

→ Sometimes it is useful to display three-dimensional data in two dimensions using contours or color-coded regions.

→ there are Matplotlib's functions that can be helpful for this task.

plt.contour for contour plots.

plt.contourf for filled contour plots.

plt.imshow for showing images

Notebook for plotting and importing functions we will use.

In [1]: %matplotlib inline

import matplotlib.pyplot as plt
plt.style.use('dark_background-white')

visualizing a three dimensional function.

→ we will start by demonstrating a contour plot using a function $z = f(x, y)$ using the following particular choice.

In [2]: def f(x, y)

return np.sin(x) * np.cos(y)

→ return np.sin(x) * np.cos(y)

A contour plot can be created with the plt.contour function. It takes three arguments

A grid of x values

A grid of y values

A grid of z values

→ the x and y values represent the positions on the plot.

z values represented by the contour levels.

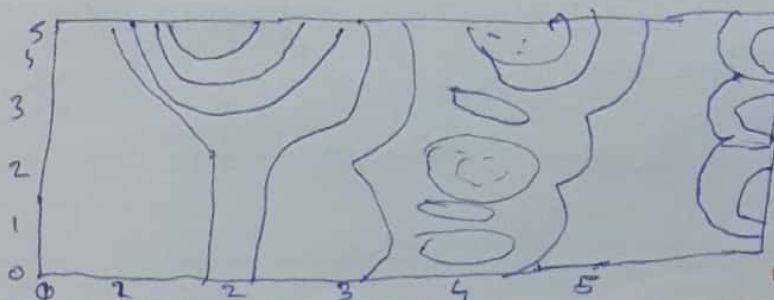
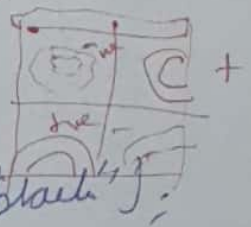
np.meshgrid function → which builds two-dimensional grids from one-dimensional arrays.

In [3]: x = np.linspace(0, 5, 5) x, y = np.meshgrid(x, y)

y = np.linspace(0, 5, 4) z = f(x, y)

x, y = np.meshgrid(x, y) z = f(x, y)

In [4]: plt.contour(x, y, z, colors = 'black')



Visualizing three dimensional data with contours

→ arrays 4x4 or map.

Histograms, Binning and Density:

A simple Histogram can be great first step in understanding data. Earlier we saw a preview of Matplotlib's histogram function which creates a basic histogram in one line once the normal boiler plate imports are done.

In [1]: %matplotlib inline

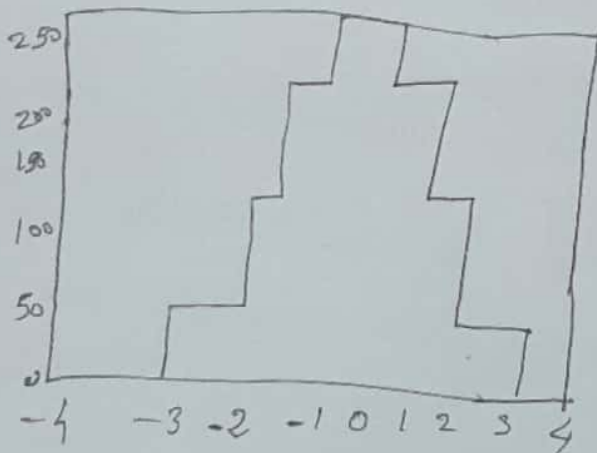
```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
plt.style.use('seaborn-white')
```

```
data = np.random.randn(1000)
```

In [2]: plt.hist(data);



```
x1 = np.random.normal(0, 0.8, 1000)
```

```
plt.hist(x1);
```

```
# contour can used to fill
```

```
plt.contour(x, y, z, 20,
```

```
    cmap = 'RdGy');
```

```
plt.colorbar();
```

```
contoursf (fill )
```

```
plt.imshow(z, extent =  
    [0, 5, 0, 5], origin = 'Lower',  
    cmap = 'RdGy');
```

```
plt.colorbar();
```

```
plt.axes(aspect = 'image');
```

upper left corner input data not

The hist() function has many options to tune both the calculation and the display

In [3]: plt.hist(data, bins = 30, normed = True,

```
alpha = 0.5, histtype = 'stepfilled', color =
```

```
'steel blue', edgecolor = 'none');
```

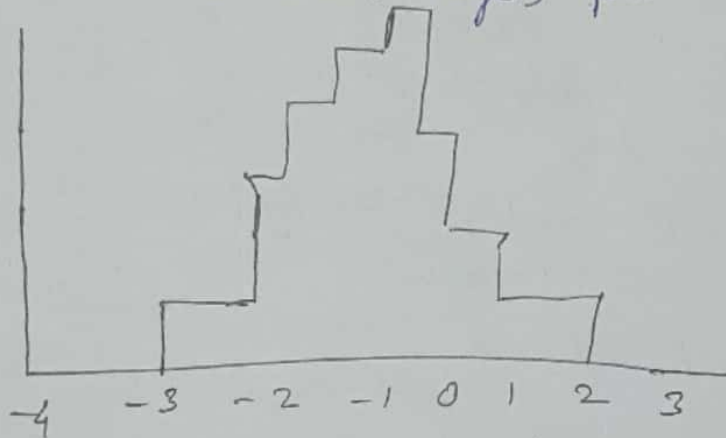
To find the combination of histtype = 'stepfilled' along with some transparency

In[4]: $x_1 = np.random.normal(0, 0.8, 1000)$

$x_2 = np.random.normal(-2, 1, 1000)$

$x_3 = np.random.normal(3, 2, 1000)$

`kwargs = dict(histtype='stepped', alpha=0.3,
normed=True, bins=40)`
`plt.hist(x1, **kwargs)`
`plt.hist(x2, **kwargs)`
`plt.hist(x3, **kwargs)`



A customized histogram.

visualization with seaborn:

→ Matplotlib has proven to be incredibly useful and avid users will admit it often leaves much to be desired. There are several solid complaints about Matplotlib that often come up

→ Prior to version 2.0 Matplotlib's default are not exactly the best choice. It was based off of matlab circa 1999.

→ Matplotlib's API is relatively low level doing sophisticated statistical visualization is possible but often requires a lot of boilerplate code.

→ Matplotlib predates pandas by more than a decade, and this is not

designed for use with pandas Dataframe, you must extract each series and often concatenate them together into a single format. It would be nice to have a plotting library that can intelligently use the Dataframe labels in a plot.

Seaborn versus Matplotlib: here is an example of a simple-random walk plot in Matplotlib using its class plot formatting and colors

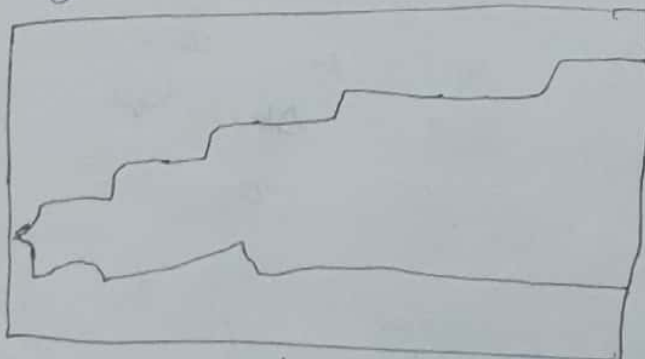
```
IN[1]: import matplotlib.pyplot as plt
plt.style.use('classic')
```

Now we create some random walk data -

```
IN[2]: # create some data
mg = np.random.RandomState(0) * np.linspace(0, 10, 500)
y = np.cumsum(mg.randn(500, 1), 0)
```

And do a simple plot.

```
IN[3]: plot the data with Matplotlib defaults
plt.plot(x, y)
plt.legend('ABCDEF', ncol=2, loc='upper left');
```



Data in Matplotlib's default style.

Two Interfaces of Matplotlib.

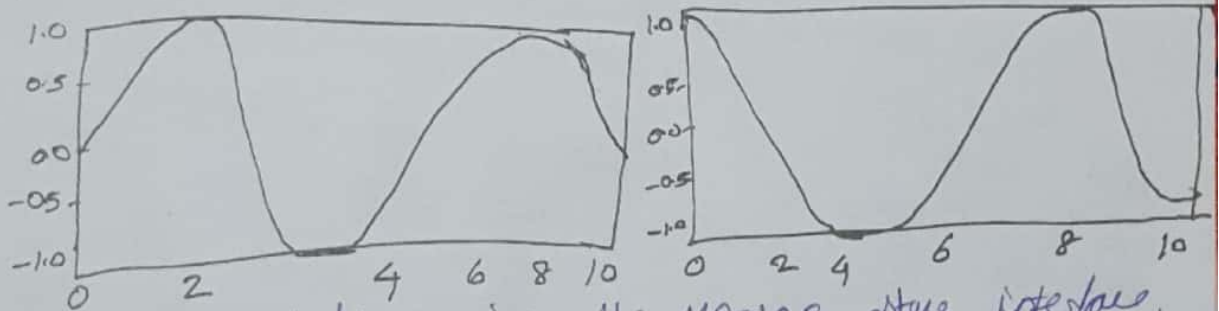
- 1) Matlab - style state-based interface
- 2) Object-oriented interface.

MATLAB - style interface.

Matplotlib was originally written as a Python alternative for MATLAB users.

The MATLAB - style tools are contained in the pyplot (plt) interface.

```
plt.figure() # create a plot figure
# create the first of two panels and set current axis
plt.subplot(2, 1, 1) # (rows, columns, panel number)
plt.plot(x, np.sin(x))
# create the second panel and set current axis.
plt.subplot(2, 1, 2)
plt.plot(x, np.cos(x))!
```



→ Subplots using the MATLAB - style interface.

This interface is stateful

Keeps track of the "current" figure and axes.

Reference to this can be obtained using

the `plt.gcf()` (get current figure) and

`plt.gca()` (get current axes) routines.

Stateful interface is fast and convenient for simple plots.

Adding something to the first panel after creating the second panel cannot be done.

Object-oriented interface.

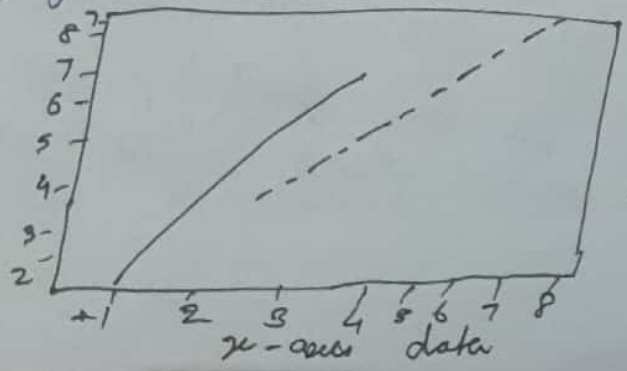
The object-oriented interface is available for more complicated situations and provides more control over the figure, plotting and functions axes objects and object-oriented interface are more readable and explicit. can structure the code and other tasks with split plotting, labeling and other tasks into code blocks.

Example.

first create a grid of plots.
 # as will be an array of two Axes objects
 fig, ax = plt.subplots(2)
 # call plot() method on the appropriate object
 ax[0].plot(x, np.sin(x))
 ax[1].plot(x, np.cos(x))

Line plots:

- The pyplot a sublibrary of matplotlib, is a collection of functions that helps in creating a variety of charts
- Line plots are used to represent two data x and y on a 2D plane
- Import matplotlib.pyplot library for plotting functions.
- Import requirements. The numpy library as per dependencies
- Data values multiple plots x and y.

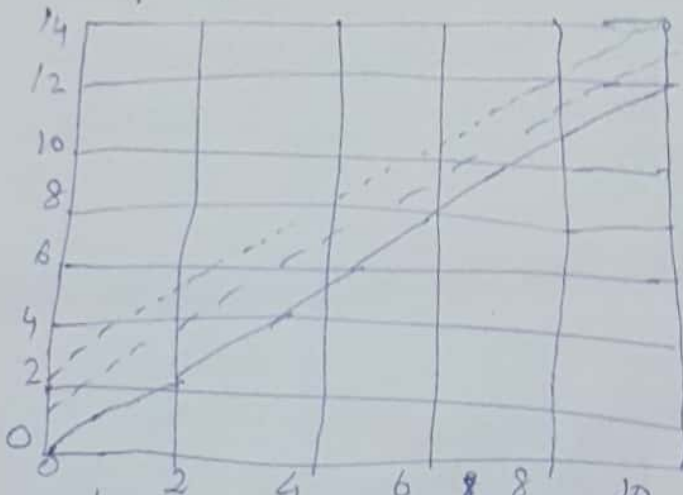


Multiple plots.

Setting Line colours and style
 The `plt.plot()` function takes additional arguments to specify line colors and styles like color, color keyword is used which accepts a color. The color can be specified in a variable of ways

```

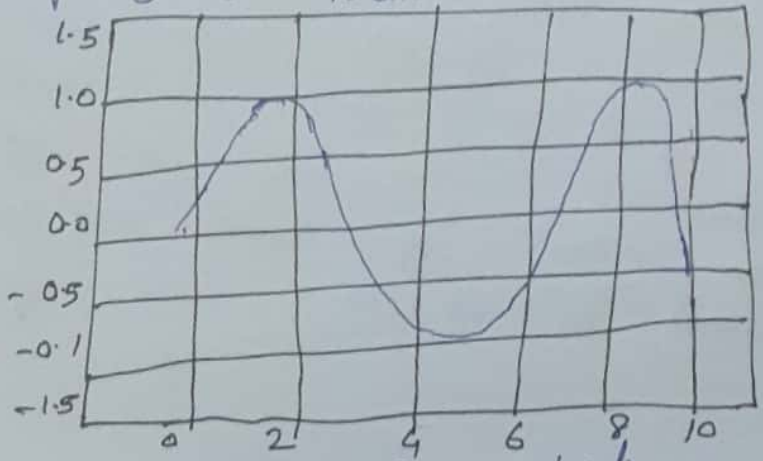
plt.plot(x, x+0, '-g') # solid green
plt.plot(x, x+1, '--c') # dashed cyan
plt.plot(x, x+2, '-k') # dotted black
plt.plot(x, x+3, ':r') # dotted red
  
```



controlling colors and styles with the Matplotlib syntax.

Setting the axes limits of plots

Default axes limits are chosen for the plot automatically. It can also be controlled by adjusting the axes limits using `plt.xlim()` and `plt.ylim()` methods.



set the axis limits

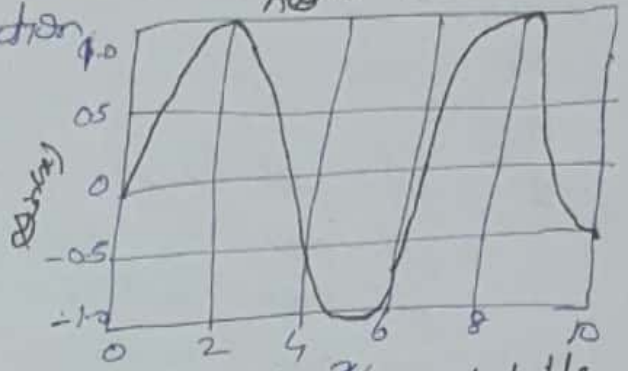
Labeling plots:

Labelling plots with the title, axis labels and simple legends can be done with the Pyplot as below. The position, size, and style of these labels can be altered using optional arguments to the function.

```

plt.plot(x, np.sin(x))
plt.title("A sine curve")
plt.xlabel('x')
plt.ylabel('sin(x)')

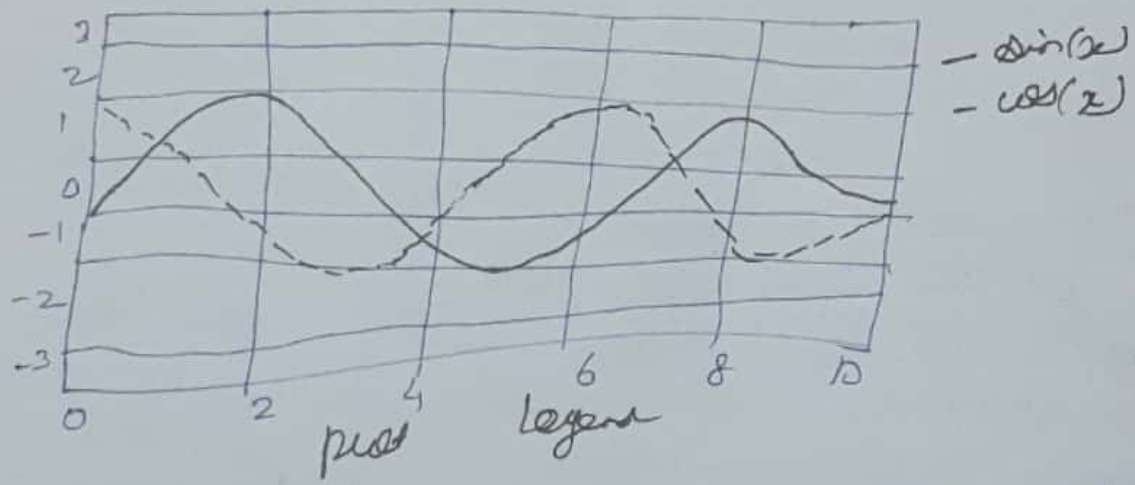
```



```

In [5]: plt.plot(x, np.sin(x), '-g', label = 'sin(x)')
plt.plot(x, np.cos(x), ':b', label = 'cos(x)')
plt.axis('equal')
plt.legend()

```



scatter plots:

scatter plots are used to observe relationships between variables. scatter plots are a type of plot in which the points are represented by individual shapes instead of point being joined by line segments as in the case of a line plot.

The `scatter()` method in the `matplotlib` library is used to draw a scatter plot. Scatter plots are used to visualise the relation among variables and how change in one affects the other variable.

Syntax: The syntax for `scatter()` method is given below:
`matplotlib.pyplot.scatter(x-axis-data, y-axis-data, S=None, C=None, marker=None, cmap=None, vmin=None, vmax=None, alpha=None, linewidths=None, edgecolors=None)`

where x-axis-data - An array containing x axis data
y-axis-data - An array containing y axis data.
S-marker size (can be scalar or array of size equal to size of x or y)

C-color of sequence of colors for markers
marker - marker style

Cmap - cmap name

linewidths - width of marker border

edgecolor - marker border color

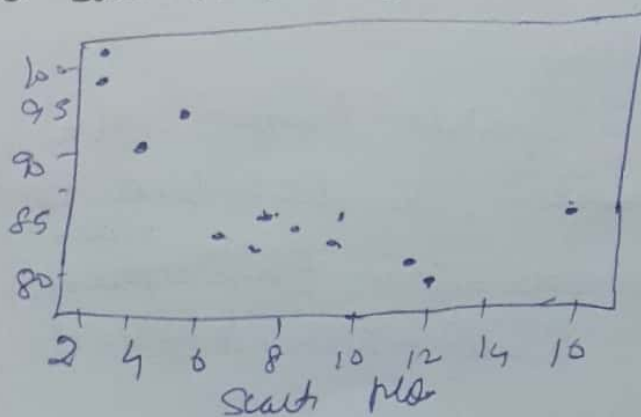
alpha - blending value between 0 (transparent) and 1 (opaque)

Simple scatter plot

`import matplotlib.pyplot as plt`

`x = [5, 7, 8, 7, 2, 17, 2, 9, 5, 11, 12, 9, 6]` `y = [99, 86, 87, 88, 100, 81, 83, 81, 99, 78, 77, 85, 86]`

`plt.scatter(x, y, c="blue")` # to show the plot plot show



Legends:

plot legends give meaning to a visualization, assigning labels to the various plot elements, the simplest legend can be created with the `plt.legend()` command. When automatically creates a legend for any labeled plot elements.

```
In [1]: import matplotlib.pyplot as plt
plt.style.use('classic')
```

```
In [2]: %matplotlib inline
import numpy as np
```

```
In [3]: x = np.linspace(0, 10, 1000)
```

```
fig, ax = plt.subplots()
ax.plot(x, np.sin(x), '-b', label='(sin x)')
ax.plot(x, np.cos(x), '--r', label='(cos x)')
ax.axis('equal')
leg = ax.legend():
```



Default plot legend.

```
In [4]: ax.legend(loc='upper left', frameon=False)
```

```
In [5]: ax.legend(frameon=False, loc='lower center',
ncol=2)
```

```
In [6]: ax.legend(frameon=True, framealpha=1,
Shadow=True, borderpad=1)
```

Choosing

Elements for the Legend

The legend includes all labeled elements by default. The labels appear in the legend by using the objects returned by plot commands.

The `plt.plot()` commands is able to create multiple lines at once, and returns a list of created line instances, passing one of them to `plt.legend()` will end it which to identify, along with the labels.

```
In [1]: y = np.linspace(0, np.pi, 100) + np.pi * np.arange(0, 2, 2.05)
```

```
lines = plt.plot(x, y)
```

#lines is a list of plt.Line2D instances

```
plt.legend(lines[:2], ['bird', 'dend']);
```

COLORS:

plt legends identify discrete labels of discrete points for continuous labels based on the color of points, lines, or regions. a labeled colorbar can be used. In matplotlib a colorbar is a separate axis that can provide a key for the meaning of colors in a plot.

for plotting and importing the function:

```
In [1]: import matplotlib.pyplot as plt
```

```
plt.style.use('classic')
```

```
In [2]: %matplotlib inline
```

```
import numpy as np
```

The simplest colorbar can be created with the `plt.colorbar` function.

```
In [3]: x = np.linspace(0, 10, 1000)
```

```
I = np.sin(x) * np.cos(x) * np.new_axis()
```

```
plt.imshow(I)
```

```
plt.colorbar();
```

Choosing the colormap:

The three categories of colormaps:

- I) Sequential colormaps: these consist of one continuous sequence of colors
eg (bans or viridis).
- II) Divergent colormaps: these usually contain two distinct colors, which show positive & negative deviations from a mean (eg: RdBu or PuOr)
- III) Qualitative colormaps: these mix colors with no particular sequence (eg. rainbow or Jet)

The Jet colormap, which was the default in matplotlib prior to version 2.0 is an example of a qualitative colormap. Its status as the default was quite unfortunate, because qualitative maps are often poor choice for representing quantitative data.

Color limits and extensions of Colorbar in an instance of plt.Axes are applicable. A large range of customization like colorbar are applicable.

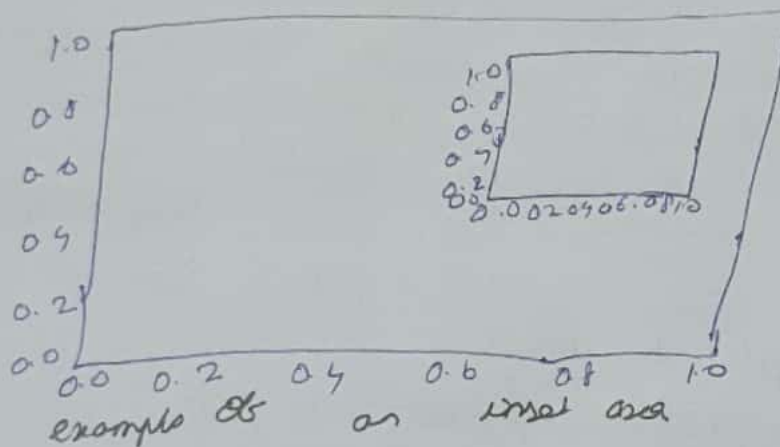
Discrete Color Bars:
but can be represented by default continuous values. It is easier was to do this and pass the name of a suitable colormap along with the number of desired bins:
In []: `plt.imshow(I, cmap = plt.cm.get_cmap('Blues', 6))`
`plt.colorbar()`
`plt.clim(-1, 1);`

Subplots:

It is necessary to compare different views of data side by side. Matplotlib has the concept of subplots: group of smaller areas that can exist together within a single figure.

```
In [1]: %matplotlib inline
import matplotlib.pyplot as plt
plt.style.use('seaborn-white')
```

```
import numpy as np
In [2]: ax1 = plt.axes() # standard axes
ax2 = plt.axes([0.65, 0.65, 0.9, 0.9])
```

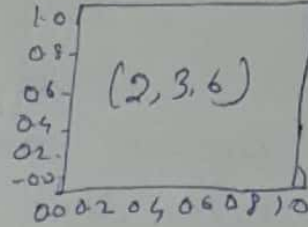
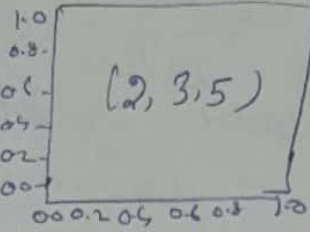
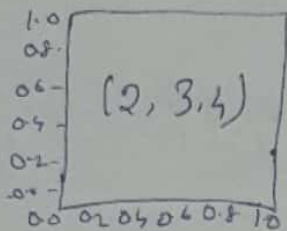
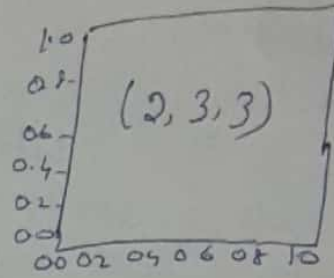
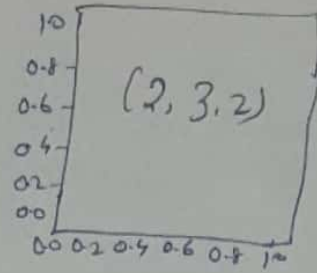
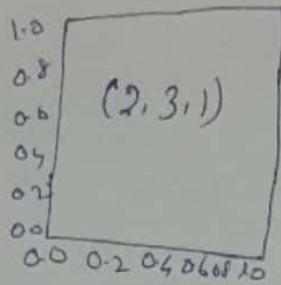


plt.subplot: Simple Grids of subplots.

\Rightarrow `plt.subplot()` creates a single subplot within a grid. This command takes three integer arguments: the number of rows, the number of columns, and the index of the plot to be created. In this scheme, which starts from the upper left to the bottom right.

```
In [3]: for i in range(1,7):
plt.subplot(2,3,i);
plt.text(0.5, 0.5, str(i), (2,3,i),
fontsize=18, ha='center')
```

```
fig, ax = plt.subplots(3,4, sharex='col', sharey='row')
```



plt.subplot() example

Text and Annotation

Good visualization guides the user so that the figure conveys a story. In some cases, this story can be expressed in an entirely visual manner. without the need for added text, but in others, small textual cues and labels are necessary. The basic types of annotations are axis labels and titles.

In [1]: %matplotlib inline.

```
import matplotlib.pyplot as plt
import matplotlib as mpl
plt.style.use('seaborn-whitegrid')
import numpy as np
import pandas as pd.
```

Transforms and Text Position:

- ax.transData: Transform data coordinates associated with
- ax.transAxes: Transform axes (in units of axis dimensions) associated with the
- fig.transFigure: Transform figure (in units of figure dimensions) associated with the

Arrows and Annotation

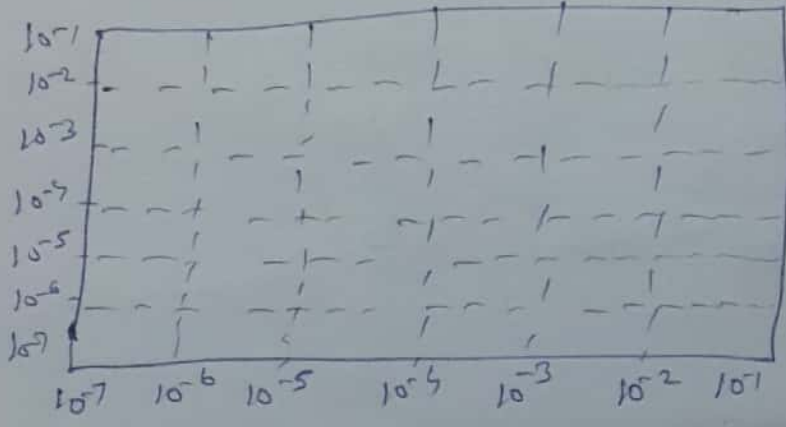
Drawing arrows in Matplotlib is much harder. `plt.arrow()` function: the objects that will be subject to the varying aspect ratio of the plot. `plt.annotate()` function: creates some text and an arrow, and the arrows are black. The arrow style is controlled through the `arrowprops` dictionary.

CUSTOMIZATION:

Customizing Ticks: Matplotlib's locators and formatters are generally different in many common situations, but are in no way optimal for every plot. Major and minor ticks:

within each axis, there is the concept of a major tick mark and a minor tick mark. As the names would imply, major ticks are usually bigger or more pronounced while minor ticks are usually smaller. By default, Matplotlib's `yticks` make use of minor ticks, but one place you can see them is within logarithmic plots.

```
In [1]: %matplotlib inline
import matplotlib.pyplot as plt
plt.style.use('dark_background')
import numpy as np
In [2]: ax = plt.axes(yscale='log', yscale='log')
```



Geographic Data with Basemap

A common type of visualization in data science is that of the geographic data. Matplotlib's main for this type of visualization is its Basemap toolkit.

There are several Matplotlib toolkits that are under the `mpl-` toolkits namespace, Basemap is a useful tool for python users to have in their virtual toolbelts.

Installation of Basemap `pip install basemap`

```
$ conda install basemap  
In [1]: %matplotlib inline.
```

```
import numpy as np  
import matplotlib.pyplot as plt  
from mpl_toolkits.basemap import Basemap
```

```
In [2]: plt.figure(figsize=(8,8))
```

```
m = Basemap(projection='ortho', resolution=None,
```

```
            lat_0=50, lon_0=-100) 27.77  
m.blueshade(scale=0.5)
```

Map projections:

It is used to project a spherical map such that of the earth, onto a flat surface without distortions or breaking its continuity.

Depends on the intended use of the map projection, there are certain map features

```
In [3]:
```

```
from matplotlib import axes:
```

```
def draw_map(m, scale=0.2):
```

```
# from a shaded-relief image
```

```
m.shadedrelief(scale=scale)
```

lats and longs are returned as a dictionary
lats = m.drawparallels (np.linspace (-90, 90, 13))
longs = m.drawmeridians (np.linspace (-180, 180, 13))

keys contain the plt.Line 2D instances.

lat_lines = Chain (* (tup [0][0] for tup in lat_items))
lon_lines = Chain (* (tup [1][0] for tup in lon_items))

all_lines = Chain (lat_lines, lon_lines)

cycle through these lines and set the desired style,
for line in all_lines;

line.set (line_style = '-', alpha = 0.3, color = 'w')

Cylindrical projections:

The simplest of map projections are cylindrical projections in which lines of constant latitude and longitude are mapped to horizontal and vertical lines respectively.

fig = plt.figure (figsize = (8, 6), edgecolor = 'w')

m = Basemap (projection = 'cyl', resolution = None,
central_lat = -90, central_lon = 90,
min_lat = -180, min_lon = 180)

draw_map(m)

Pseudo-cylindrical projections:

It gives better properties near the poles of the projection.

The extra arguments to Basemap refers to the central latitude (lat = 0) and longitude (lon = 0) for the desired map

Orthographic projection

In [5]: fig = plt.figure (figsize = (8, 8))

m = Basemap (projection = 'ortho', resolution = None,
lat = 0 = 50, lon = 0 = 0)

draw_map(m);

Drawing a Map Background

The Basemap package contains functions for drawing borders of physical features like continents, oceans, lakes and rivers, as well as political boundaries.

The following ~~are~~ ^{are} of its available drawing functions using IPython's help features:

Physical boundaries and bodies of water

draw_coastlines(): Draw continental coast lines.

draw_landmask(): Draw a mask between the land and sea, for use with projecting arrays on a map.

draw_rivers(): Draw rivers on the map.

Political boundaries

draw_countries(): Draw country boundaries.

draw_states(): Draw US state boundaries.

draw_counties(): Draw US county boundaries.

Map features draw_great_circle(): Draw a great circle between two points.

draw parallels(): Draw lines of constant latitude.

draw meridians(): Draw lines of constant longitude.

draw_mapscale(): Draw a linear scale on the map.

Whole-globe images:

blue_marble(): Project NASA's blue marble image onto the map.

shade_relief(): Project a shaded relief image onto the map.

etop(): Draw an elevation image onto the map.

map_image(): Project a user-provided image onto a map.

Thank you.