

Data Analytics.

Date: / / Page no: _____

Unit I

Statistics :-

The branch of mathematics that transforms data into useful information for decision makers.

or
statistics is the science of conducting studies to.

↳ Collect

↳ organize

↳ Summarize

↳ analyze and

↳ draw conclusion of data.

• Why we should study statistics :-

→ Statistics is the science and also the art of learning from data. As a discipline it is concerned with the collection, analysis, and interpretation of data, as well as the effective communication and presentation of results relying on data.

→ Statistics lies at the heart of the kind of quantitative reasoning necessary for making important advances in the sciences, such as medicine and genetics, and for making important decisions in business & public policy.

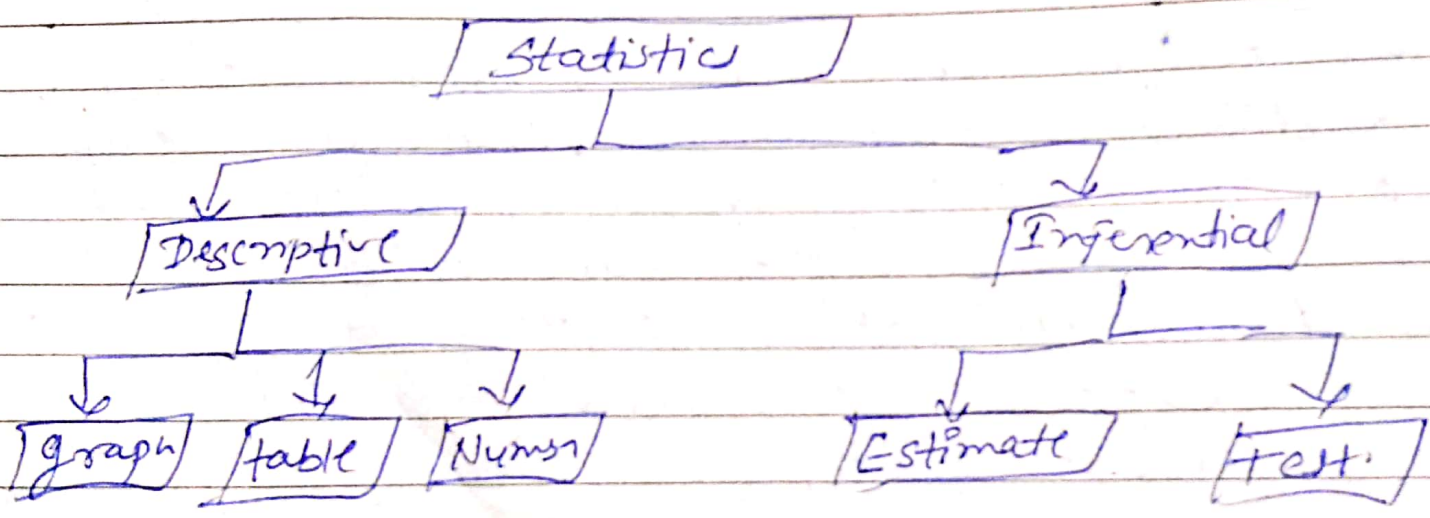
* Statistics is divided into two main areas depending on how data are used.

1) Descriptive Statistics: It describes the data. It consists of the collection, organization, summarization, and presentation of data. Table form

Ex: 30% people have A type blood 1000 people died with Cholera in 2015

2) Inferential Statistics: Consists of generalizing from sample to population, performing estimation and determining relationship among variables, and making predictions.

Ex: 30% people have A type blood group



• In inferential statistics, the answers are never 100% accurate because the calculation we use a sample taken from the population. This sample doesn't include every measurement from the population.

Difference:-

Descriptive Statistics

Inferential Statistics

- 1) It give information that describe the data in some manner.
- 2) organize, analyze and present data in a meaningful way.
- 3) To describe a situation
- 4) It explain already know data and limited to a sample or population having small size
- 5) charts, graphs and tables.

- 1) It makes inferences about populations using data drawn from the population.
- 2) Compare data, test hypothesis and make prediction (Estimation of parameters)
- 3) To explain the chance of occurrence of an event.
- 4) It attempts to reach the conclusion to learn about the population.
- 5) Probability.

Probability :- = $\frac{\text{Favorable outcome}}{\text{Total outcomes}} = \frac{12}{20} = 0.6$
 ||
 Chance

Ex:- class 12-boys
 8-girls
 20 Total Student.

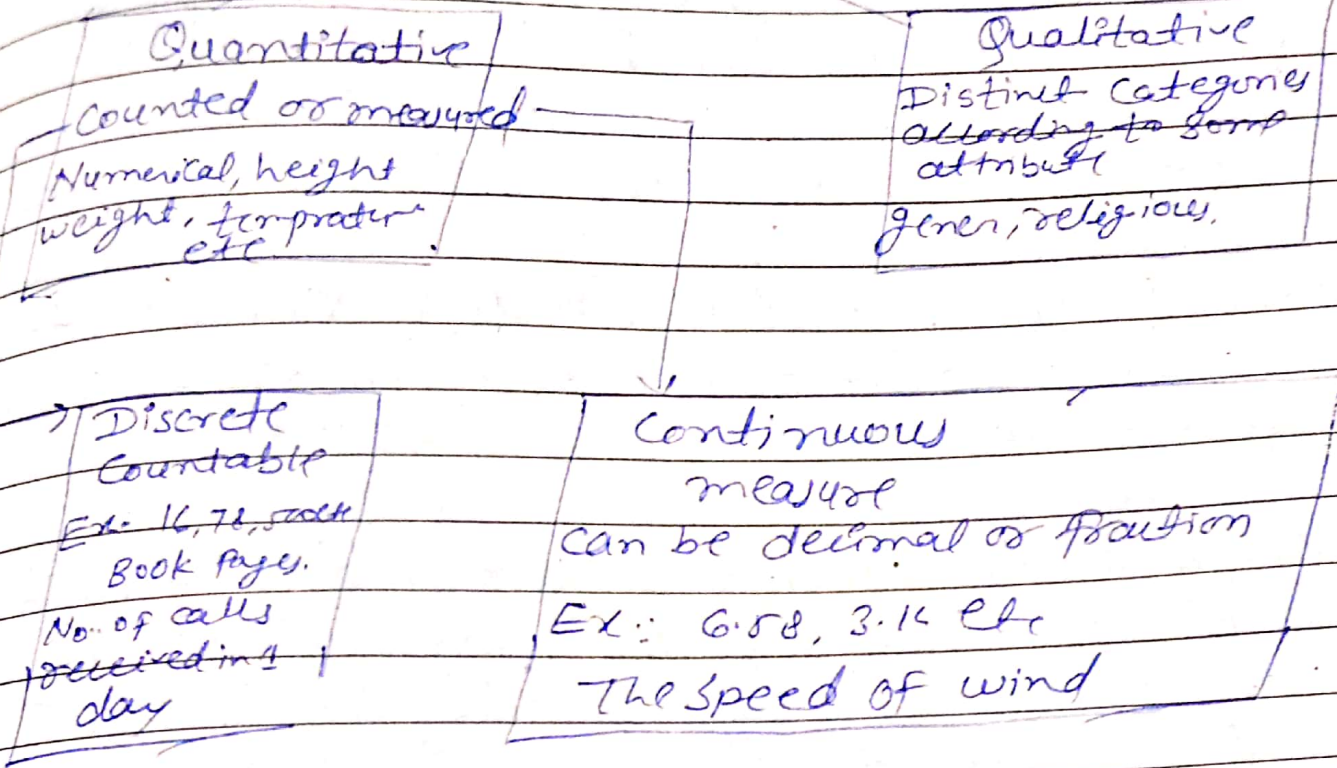
from 20 student, a student will be selected, what is the prob that he will be boys

Sum of the probability
 or
 Total Probability = 1

* Probability Distribution:-

- Variable :- A variable is a characteristic or attribute that can assume different values. x, y, z .
- Data :- The values that a variable can assume are called data.
- Population :- A population consists of all subjects (human or otherwise) that are studied, (all members of a defined group that we are studying or collecting information on for data driven decision)
Ex - Census Example.
- Sample :- A sample is a subset of the population. (A part of the population is called a sample.
Ex. - may be biased (not conform 100%.)
- Random Variable :- A variable whose values are determined by chance are called random variables.
Ex. - Automobile insurance, claim support etc. every year

Variable



* Discrete probability Distribution:-

- A discrete probability distribution consists of the values a random variable can assume and the corresponding probability of the values.

Range of probability 0-1

5th insurance Random Variable
Random variable ki kya probability hai uska discrete probability

children

Example: Construct a probability distribution for a discrete random variable of the following sample space.

3 children

BBB BBG BGB GBB GGG GGB GBS BGG

Solution:- If X is the random variable for the numbers of girls, then it assumes the value 0, 1, 2, 3.

No girls BBB	1 girls BBG, BGB, GBB	2 girls GGB, GBS, BGG	3 girls GGG
$\frac{1}{8}$	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{1}{8}$

No. of girls X	0	1	2	3
probability $P(X)$	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{1}{8}$

probability $\frac{1}{8} + \frac{3}{8} + \frac{3}{8} + \frac{1}{8}$

$$\frac{8}{8} = 1$$

probability range = 0 - 1

* What is Hypothesis :-

A hypothesis is an assumption about the relationship b/w variables. It is a tentative statement of the research problem or guess about the research outcome proposing a statement pertaining to relationship b/w the two variables is called a hypothesis. Its validity is usually unknown.

Simply
In my words, A prediction about the outcome of the research in the mind of researcher (Before starting the research work) this is called hypothesis.

* Criteria for hypothesis construction:-

- It should never formulated in form of ^{question} ~~question~~
- It can be testable.
- It ~~is~~ should be specific & precise.
- It should describe one issue only.

1972 :- If H_0 is false, even then we are accepted it
Error :- then it is said to be type 2 Error
Date: / / Page no:

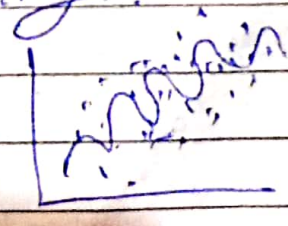
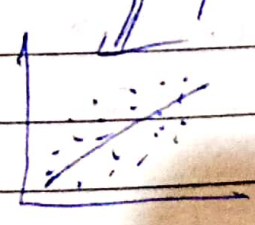
*** Regression Analysis :-** Regression analysis is one of the most used and most powerful multivariate statistical techniques for it infers the existence and form of a functional relationship in a population.

- Regression is a parametric technique used to predict continuous (dependent) variable given a set of independent variables.
- It is parametric in nature because it makes certain assumption based on the data sets. If the data set follows those assumption, regression gives incredible results.

Point
D. Regression analysis is a statistical technique used to describe relationships among variables.

- Dependent & Independent Variable
- outliers.
- multicollinearity
- Underfitting & overfitting :-

Training set pe bhi kam nahi kar pada



Training set pe kam karta but testing time me kam nahi karta

Ex:- Industry task annual sales nikalng annual sale ko affect karne wale hazar factors honge annual sale depende karta hai factors pe.

Annual Sale \rightarrow dependent Variable
 hazar factor \rightarrow Independent Variable

- Regression analysis is trying to model the relationship b/w the dependent & independent ~~mod~~ variable.

\rightarrow Outlier: Suppose mean nikalng hai
 $\frac{1+2+3+90}{4}$

90 oddone
 out 4

90-outlier, 90 ke hone

se mean cmr expected

Outlier expected result ko wrong direction dete hai

upar

\rightarrow Multicollinearity:- Independent variables when are correlated each other. non independent variable when they are sharing non linear relationship with each other (1 ko kare)

Dependent \rightarrow outcomes,
Continuous Numeric Value.

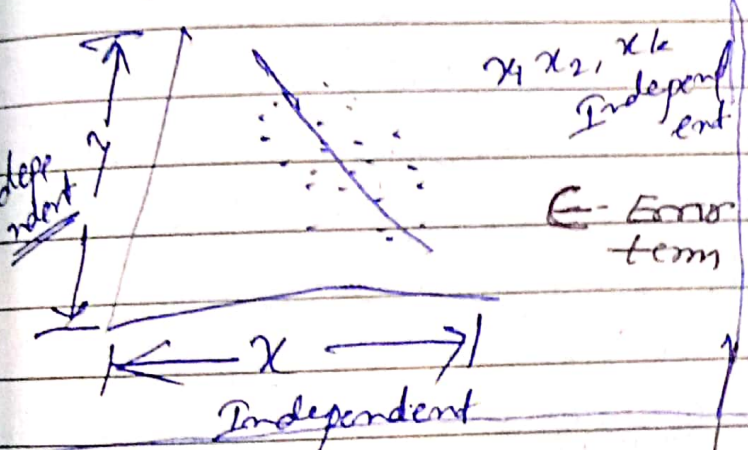
Independent - Input.
Date / / Page no

★ Linear Regression

- 1) Dependent Variable are continuous in nature (numerical value)
- 2) Linear Relationship (dependent & independent relationship)
- 3) $I \rightarrow I$ & $I \rightarrow D$
Simple Linear Regression
 $y = \alpha_0 + \alpha_1 x_i$

 α_0 Intercept (Dependent)
 α_1 Coefficient of regression (Independent)
- 4) $I - D$ & $I \rightarrow I$
Multiple Linear R.

5) $y = \beta_1 + \beta_2 x_2 + \dots + \beta_k x_k + \epsilon$



x badayenge to y bhi badhenge means linear in nature.
If x increase then y is increase.

Logistic Regression

- 1) Dependent Variable is binary
- 2) 1 (True, Success) 0 (False, Fail)
- 3) Goal is to find the best fitting model for I & D variable relationship.
- 4) Independent Variables can be continuous or binary.
- 5) Also called Logit R. expression
- 6) Used in machine learning
- 7) deal with probability to measure the relation b/w dependent & independent variable.

Probability of Support
Ya cu a bank loan use karta hai.

* Regression is a well known Statistical technique to model the predictive relationship b/w several independent variables and one dependent variable.

- The objective is to find the best-fitting curve for a dependent variable in a multidimensional space, with each independent variable being a dimension.

* The key steps for regression are simple

- List all the variables available for making the model
- Establish a dependent variable of interest.
- Examine visual relationship b/w variables of interest.
- find a way to predict DV using other variables (dependent variable)

* Regression models are simple, versatile, visual/graphical tools with high predictive ability. They include non linear as well as binary predictions. Regression model should be used in conjunction with other data mining techniques to confirm the findings.

* Advantages of regression model:-

- 1) Regression models are easy to understand as they are built upon basic statistical principles such as correlation and least square error.
- 2) Regression models provide simple algebraic equations that are easy to understand & use.
- 3) The strength of the regression model is measured in terms of the correlation coefficient and other related statistical parameters.
- 4) Regression models can include all the variables that one wants to include in the model.

★ ANOVA (Analysis of Variance) :-

In t-test we compare only two group.

ANOVA is a statistical method used to compare the mean of two or more groups.

Analysis of Variance is a collection of statistical models

- Factor (Variable) used to analyze the difference among group means and their associated procedures (Such as 'Variation' among and b/w groups)
- Level.

0mg	50mg	100mg
9	7	4
8	6	3
7	7	2
8	8	3
8	7	4
9	6	3

developed by Statistician and Evolutionary biologist "Ronald fishes"

Factor - Dosage

Levels: 0mg, 50mg, 100mg

Type of ANOVA

1) one-way ANOVA :- one factor with (general linear model (GLM-1)) at least two levels, levels are independent.

0mg	50mg	100mg
9	7	9
8	6	7
7	4	6
6	5	3

2) Repeated-measures Anova:- one factor with atleast two levels, levels are dependent.

Day 1	Day 2	Day 3
9	7	4
8	6	3
7	6	2
8	7	3

3) factorial Anova:- Two or more factors (each of which with atleast two levels), levels can be either independent, dependent or both (mixed)

gender	Three levels		
	Day 1	Day 2	Day 3
men	9	3	5
	8	7	4
	7	6	3
women	7	6	8
	8	3	3
	9	2	2

Factor - gender / day
 Level - male/female / Day 1, Day 2, Day 3.

* Assumptions in ANOVA.

1) Normality of Sampling Distribution of mean:

The distribution of sample mean is normally distributed.

2) Independence of Errors :-

Errors b/w case are independent of one another.

3) Absence of Outliers :-

Outlying scores have been removed from the data set.

4) Homogeneity of Variance :-

population Variance in different levels of each Independent Variable are equal.

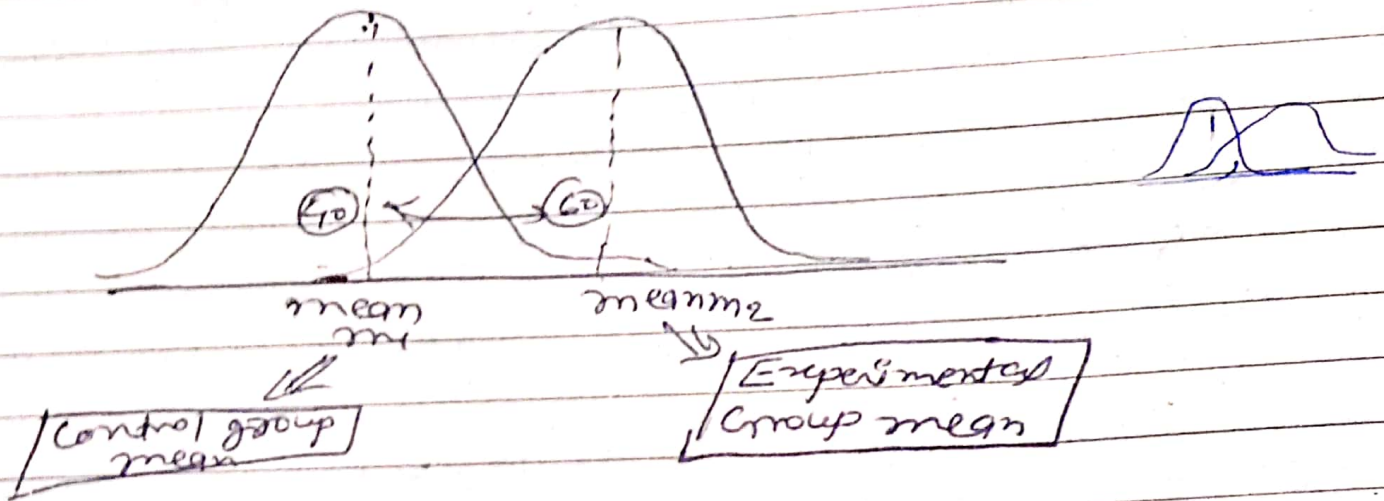
* Example of ANOVA Tests -

- A group of different therapies: Counseling, medication and biofeedback. You want to see if one therapy is better than the others.
- A manufacturer has two different processes to make light bulbs. They want to know if one process is better than the other.

* T-Test :-

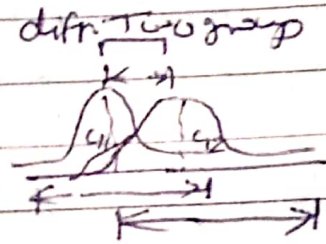
(t-जांच)

The t-test assesses whether the means of two groups are statistically different from each other.

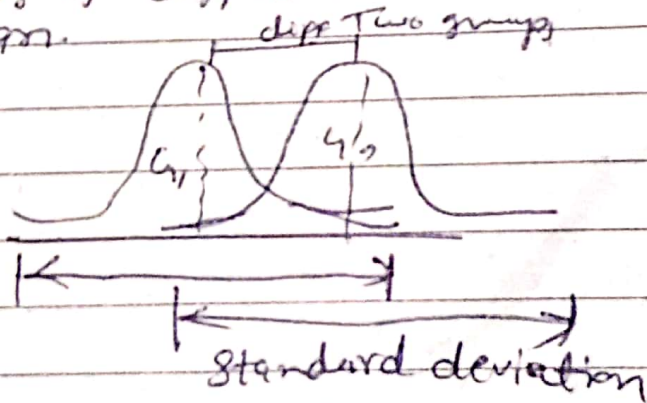


formula for t test.

$t = \frac{\text{Difference b/w mean of Two group}}{\text{Standard Error of difference b/w mean}}$



Standard Error of difference b/w mean.



$$= \frac{m_1 - m_2}{SED}$$

Unit - II

Date: / / Page no: |

★ what is Big data?

- Large amount of data. (Tera, Petta byte)
- Its a popular term used to express exponential growth of data.
- Big data is difficult to store, collect, maintain, analyze and visualize.

Store (10GB wani se, but 10 terabyte difficult) Copy hone me)
Collect (alay, alag sources se data atai hai)
maintain (store to kya hai but maintain - hude bhi Idhar udhar ^{gayaton})
analyze (meaningful information extract karna mushkil)
visualize (bar chart, graph)

- Big data is a collection of data sets so large and complex that it become difficult to process using on-hand database management tools.

The volume of data with the speed it is generated make it difficult for the current computing infrastructure to handle big data. To overcome this drawback, big data processing can be perform through a programming paradigm known as mapreduce.

- with the rapid development of science & technology there are large amounts of data generated in real life. These big data come from all areas of human life and most of them are stored in computers.

* Big data characteristics. (3Vs Concept)

- **Volume** :- ^{The volume is related to size of data. Present data is in petto bytes and in near future it will be zettabyte.} Large amount of data.
- **Velocity** :- The rate at which data is getting generated. The velocity is related to the speed of data coming from different resources.
- **Variety** :- Different types of data ^{of increasing data is not limited and is not constant.}
 - Structured data Ex: mysql
 - Semi-Structured Ex: xml
 - Unstructured data Ex: text, audio, video

* Big data Sources: There are several sources of data including some new ones. Data from outside the

- 1) Social media (organization maybe incomplete)
 - 2) Banks
 - 3) Instruments (cameras)
 - 4) websites (amazon)
 - 5) Stock market (
- and of a different quality & accuracy.

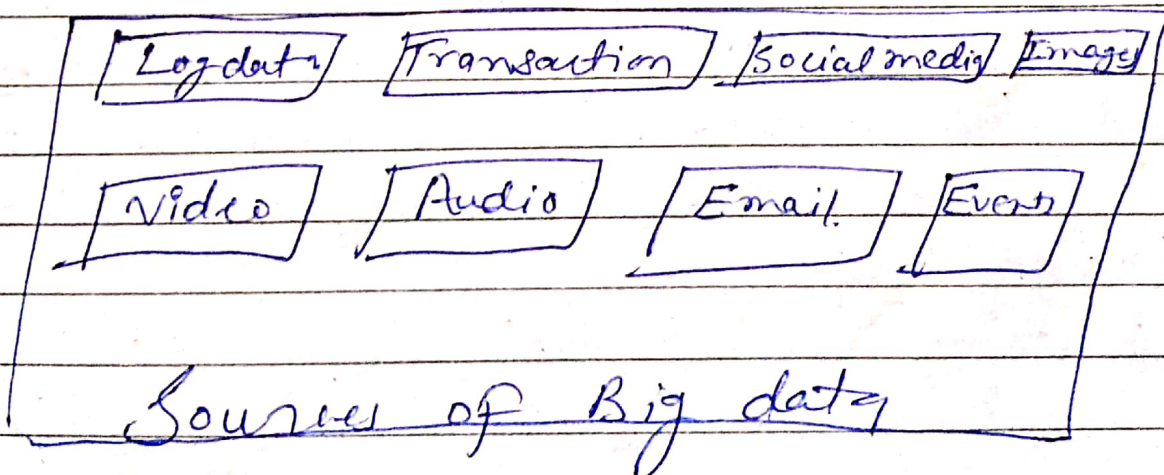
* Use cases of Big data. (comparison product scanner)

- Recommendation engine (next time ke liye recommend)
- Analysing Call Detail Record (CDR) [Call center, telecommunication]
- fraud Detection (credit card fraud, online banking)
- Market Basket Analysis (Kisi ke sath chije ke sath sale same)
- Sentimental Analysis (Text data ke analyze Ek mudda Uthra ke Social media par ke)

* There are Two types of Big data.

1) Structured Data:- are no. and words that can be easily categorized and analyzed. These data are generated by things like network sensors embedded in electronic devices, smart phones and global positioning system (GPS) devices.
→ Structured data also include things like sales figures, account balance, and transaction data.

2) Unstructured Data:- Include most complex information, such as customer reviews from commercial websites, photos and other multimedia, and comments on social networking sites. These data cannot easily be separated into categories or analyzed numerically.



* Drivers for Big data:-

* Why is Big data important?

→ The importance of big data doesn't revolve around how much data you have, but what you do with it. You can take data from any source and analyze it to find answers that enable

1) Cost reductions,

2) Time reductions.

3) New product development and optimized offerings. and

4) Smart decision making.

Sources of Big data:-

a) Social media:- All activities on the web and social media are considered stores and are accessible. Email was the first major source of new data. Google searches, facebook posts, tweets, youtube videos etc.

b) Organizations:- Business organizations & government are a major sources of data. E-Commerce systems, user-generated content, web access logs and many other sources of data generate valuable data for organization.

~~machines~~

Analyst! - $\text{Collect data} \rightarrow \text{analyzed data} = \text{Report.}$
Viewers Same valid. export Date: / / Page no: |

★ what is analytics:

Analyzing or deep understand and something using data.

→ why to do analytics

Business

Marketing

Critical Problems

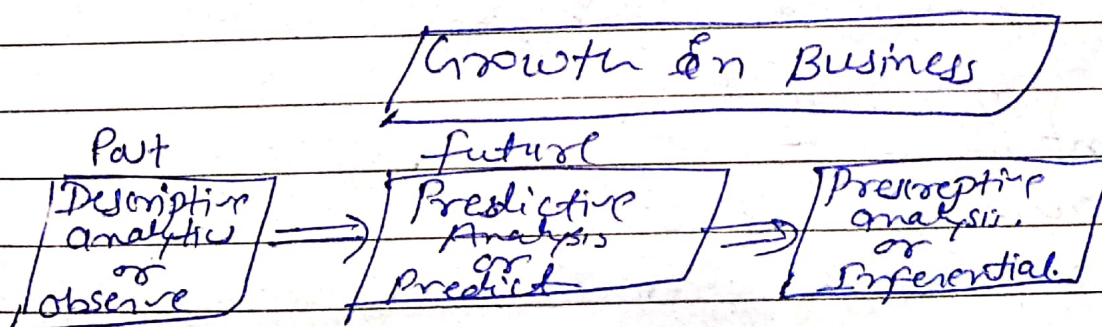
New Ideas

- Gather Hidden information insights.
- Create Reports.
- Perform market Analysis.
- Improve Business Requirements.

★ Data Analytics:

Data analytics refer to the techniques to analyse data to enhanced productivity and business gain

Business administration + Exploratory data Analysis =



Important of Bi

Date: / / Page no:

* Need for big data analytics:-

- cost reduction
- fast & Better Decision making
- New products and services.

* Big data ^{Analytics} Application :-

1) Health Care :- E-medical, purchase data, Research data, Transactional

2) Education sector :- Improving evaluation of Student Results.

- Analysing and Creating the Custom program.
- Computing the marks of Students.

3) E-Commerce :- Retailers]

- Predict Trends
- optimize pricing
- forecast demand
- personalized Stores
- Customer service.
- Sales generation.

4) Government ^{E-Government} :- ~~CP~~ - Traffic optimization

- cyber Security & Intelligence.
- Crime prediction & prevention.

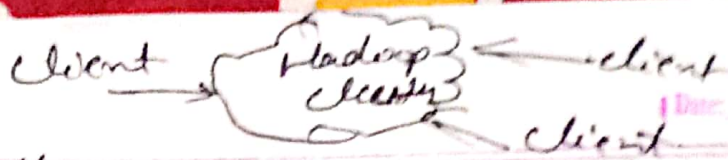
- weather forecasting
- Drug Evaluation
- Tax Compliance.

- 5) IOT :-
- Retail
 - Vehicles
 - Industry
 - Communication
 - medicare.

Smart city concept.

- 6) Media & Entertainment:-
- media scheduling
 - Audience interests.
 - Ad targetings

* What Big data Analytics:- Big data analytics is the often complex process of examining large and varied data sets or big data to uncover information including hidden patterns, unknown correlations, market trends and customer preferences that can help organizations make informed business decisions.



"Hadoop is an open source software framework which huge data storage facility"

★ Hadoop :- Hadoop is an open source, Java based programming framework that supports the processing and storage of extremely large data sets in a distributed computing environment.

- Hadoop is capable to work on clusters (thousand of interconnected machines) and handle to thousand of terabytes data.
- Hadoop provides rapid data transfer rate.
- It is developed by doug cutting and mite cafarella in 2006.
- It was inspired by google's mapreduce programming framework ^{→ GFS (google file system)}

Component present in Hadoop

- Name node (NN)
- Job tracker (JT)
- Secondary Name node (SNN)
- Data node (DN)
- Task tracker (TT)
- All this are daemons, which are present in program (necessary)

3) YARN (Yet another resource Negotiator) :- ^{JobTracker} ^{TaskTracker}
machine fail ho jaye to uska load kisi durre machine pe dena
resource manage karana

★ HDFS (Hadoop Distributed File System) :-

HDFS architecture :- It have 2 types of Nodes.

1) NameNode.

2) DataNode.

It works on master slave concept :- one master
multiple slaves

one name node multiple Data node

1) Name node :-

- Name node is also known as the master.
- Name node only stores the metadata of HDFS.
- Name node does not store the actual data or data set. The data stored in data node physically.
- Name node knows the list of the blocks and its location for any given file in HDFS. with this information Name node knows how to construct the file from blocks.
- Name Node is a single point of failure in Hadoop cluster.
- Name Node is usually configured with a lot of memory (RAM). Because the block location are kept in main memory.

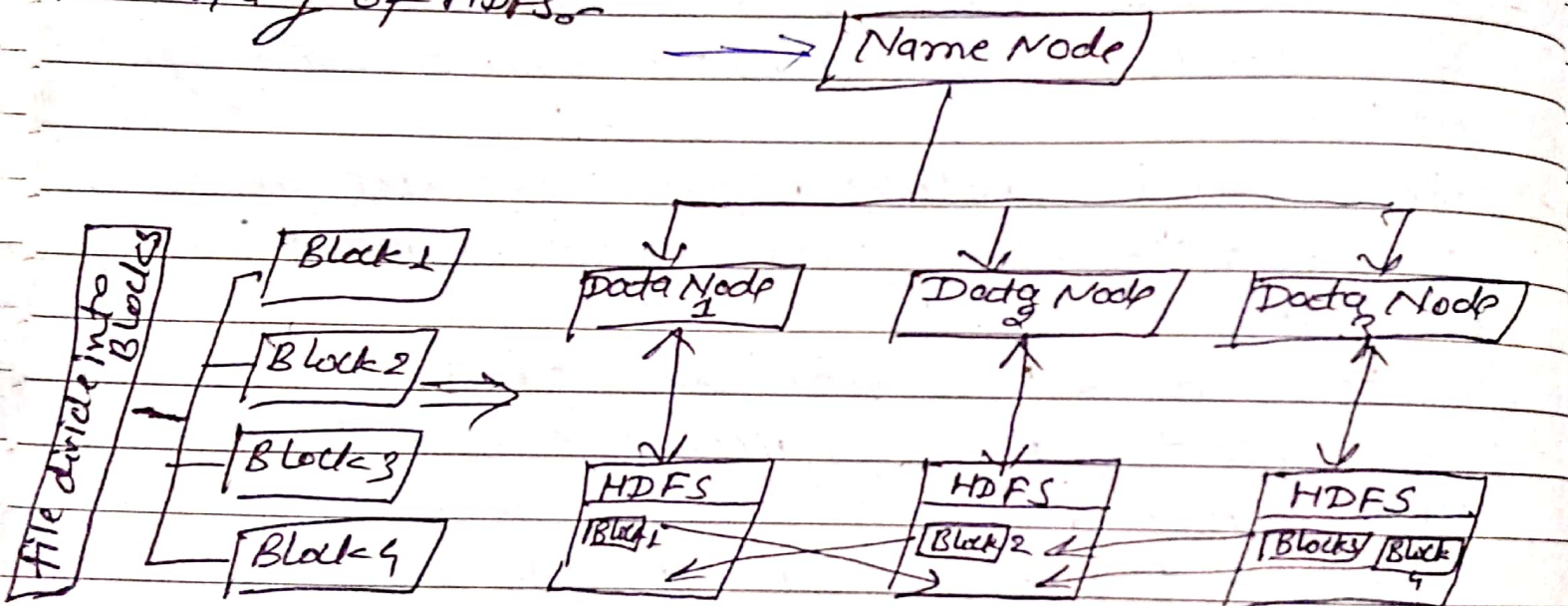
(Slave)

Q7 Data Nodes:-

- Data Node is responsible for storing the actual data in HDFS.
- DataNode is also known as the slave.
- NameNode and DataNode are in constant communication.
- When a DataNode is down, it does not affect the availability of data or the cluster.

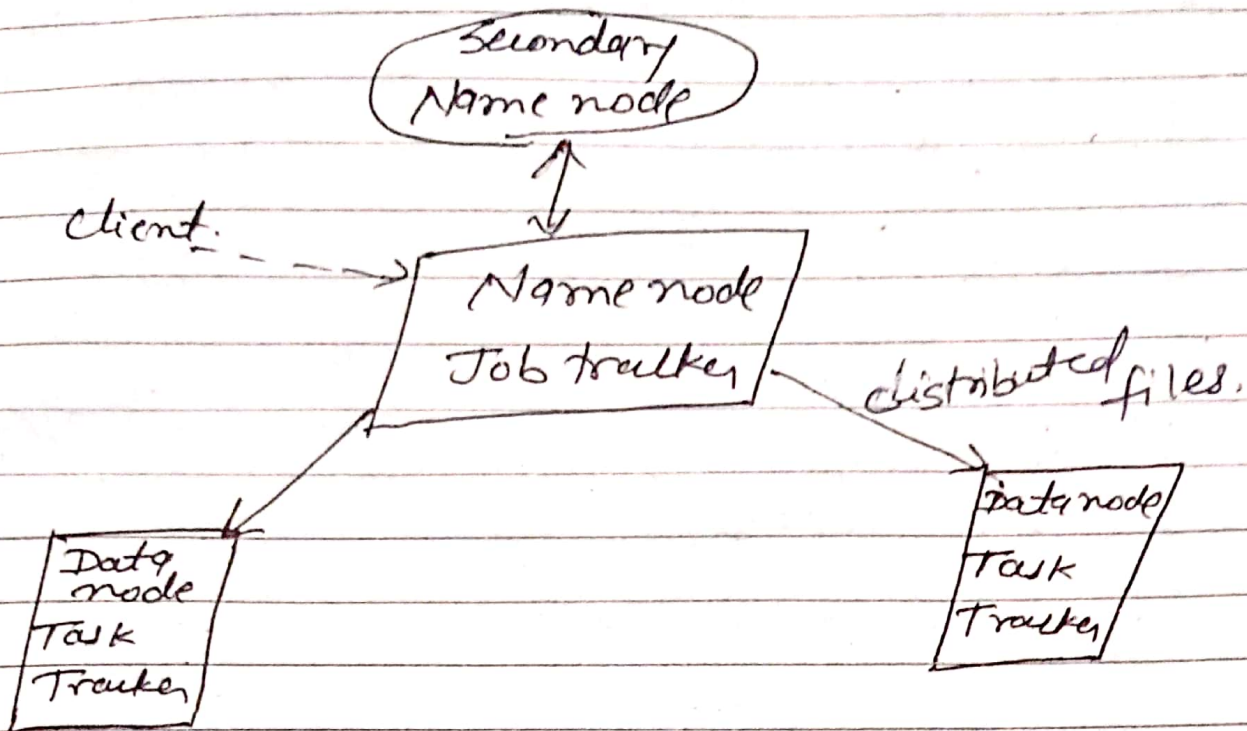
NameNode will arrange for replication for the blocks managed by the DataNode that is not available.

Working of HDFS:-



Storage & Replication of Block in HDFS.

* physical Architecture of hadoop :-



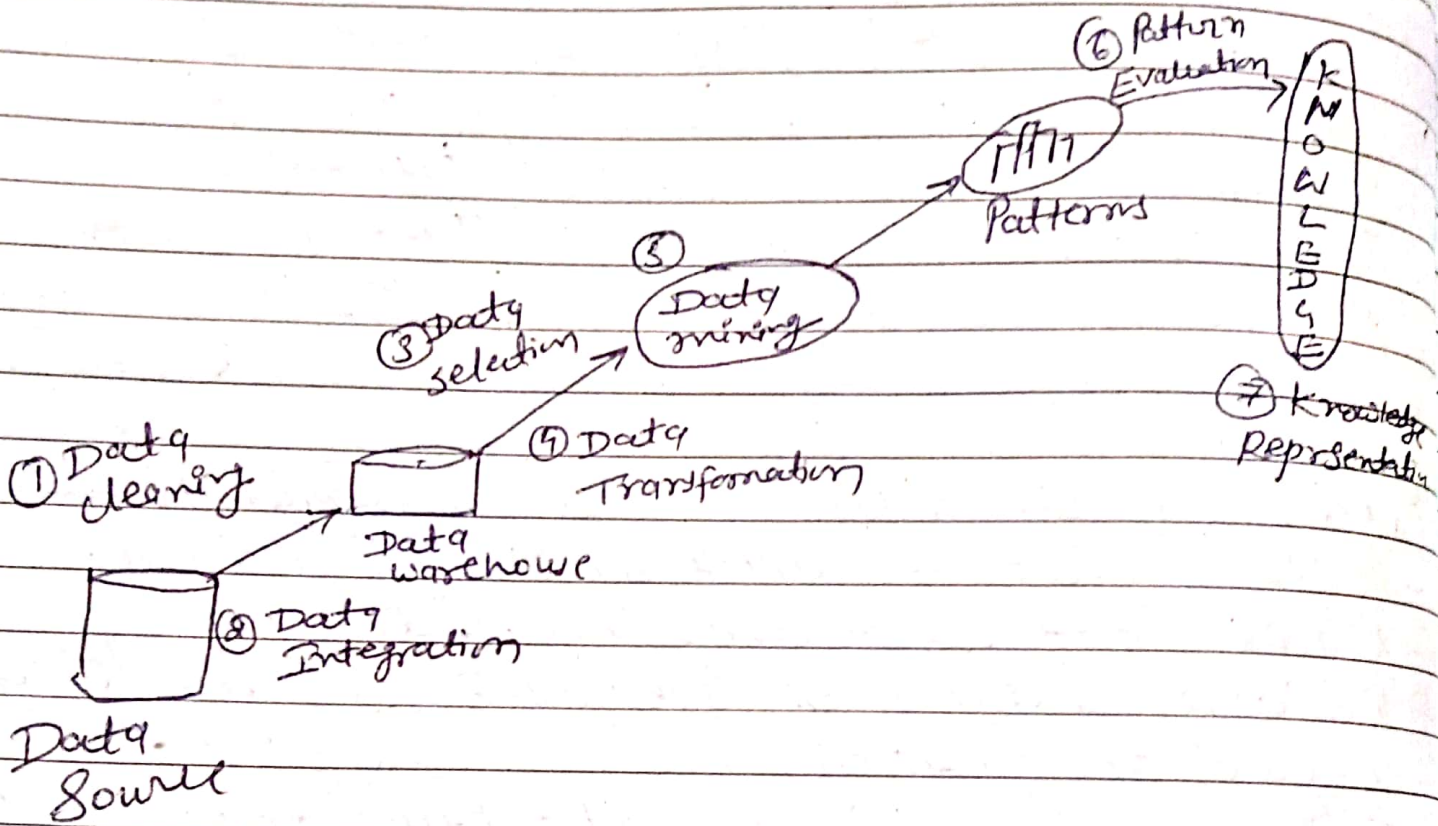
- Process:-
- client provides its job request to hadoop
 - This request is accepted by the name node.
 - Name node is master in Hadoop.
 - It also contains a Job tracker, which is again a part of master.
 - This Job is divided into tasks and tracker, provides it to the data node.
 - Now the data node is a slave and it possess task tracker which actually performs the task
 - And Job tracker continuously communicates with task tracker and if anytime it fail to reply then it assumes that the task tracker may have crashed.

Useful data ko nikalna.

* Data Discovery :-

Data discovery A combination of processes and technology that enables the detection of patterns in data.

KDD: Knowledge Discovery from Data



• Data Cleaning :-

→ To remove noise and inconsistent data example parsing the data.

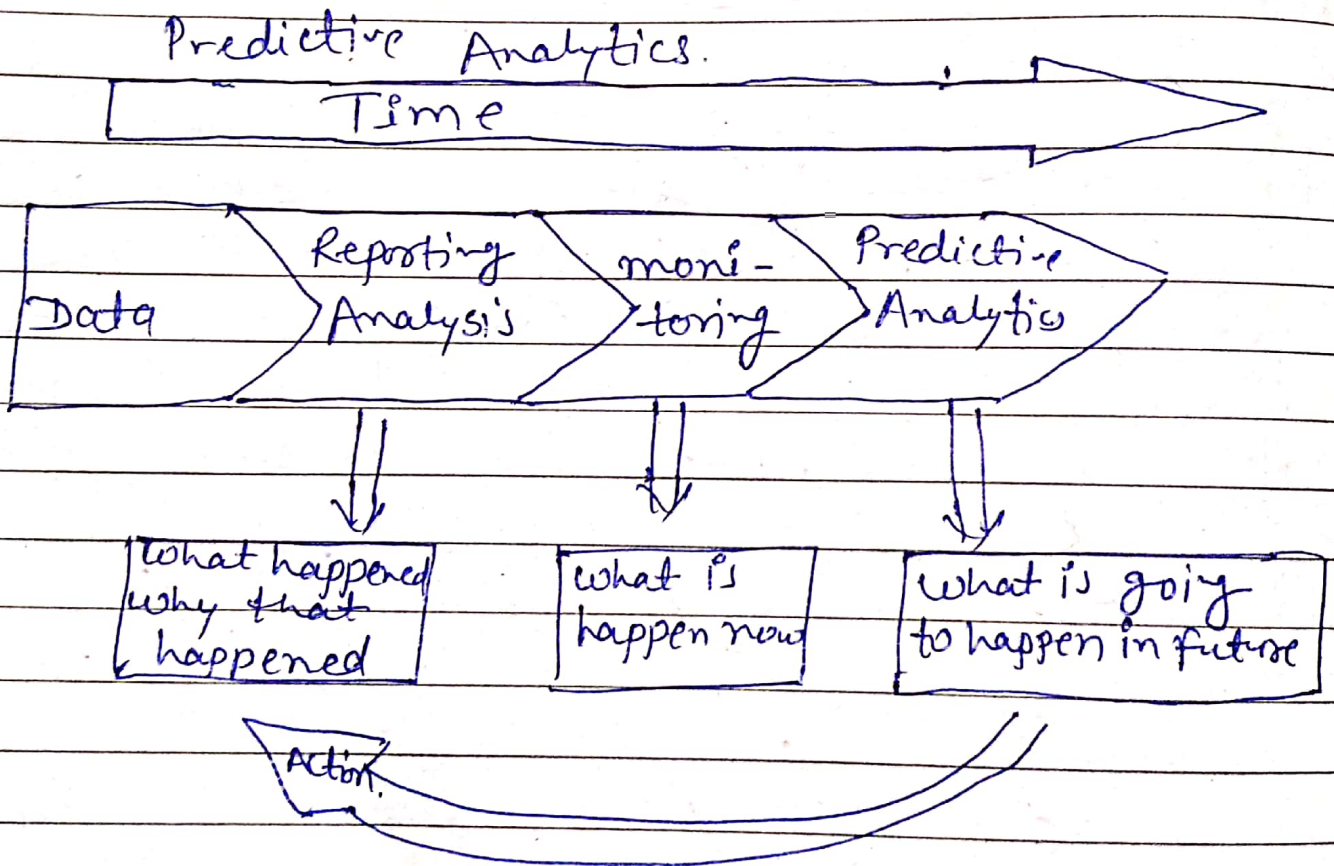
→ cleaning is performed for detection of syntax error.

- Data Integration :- Where multiple data sources are combined.
- Data selection :- Where data relevant to the analysis task are retrieved from the database.
- Data Transformation :- Where data are transformed or ~~Consolidated~~ ^(SSS) consolidated into forms appropriate for mining by performing summary or aggregation operation for instance.
- Data Mining :- An essential process where intelligent methods are applied in order to extract data patterns.
- Pattern Evaluation :- To identify the truly interesting patterns representing knowledge base on some interesting measure.
- Knowledge Representation :- Where visualization & knowledge representation techniques are used to present the mined knowledge to the user.

* Predictive analytics :-

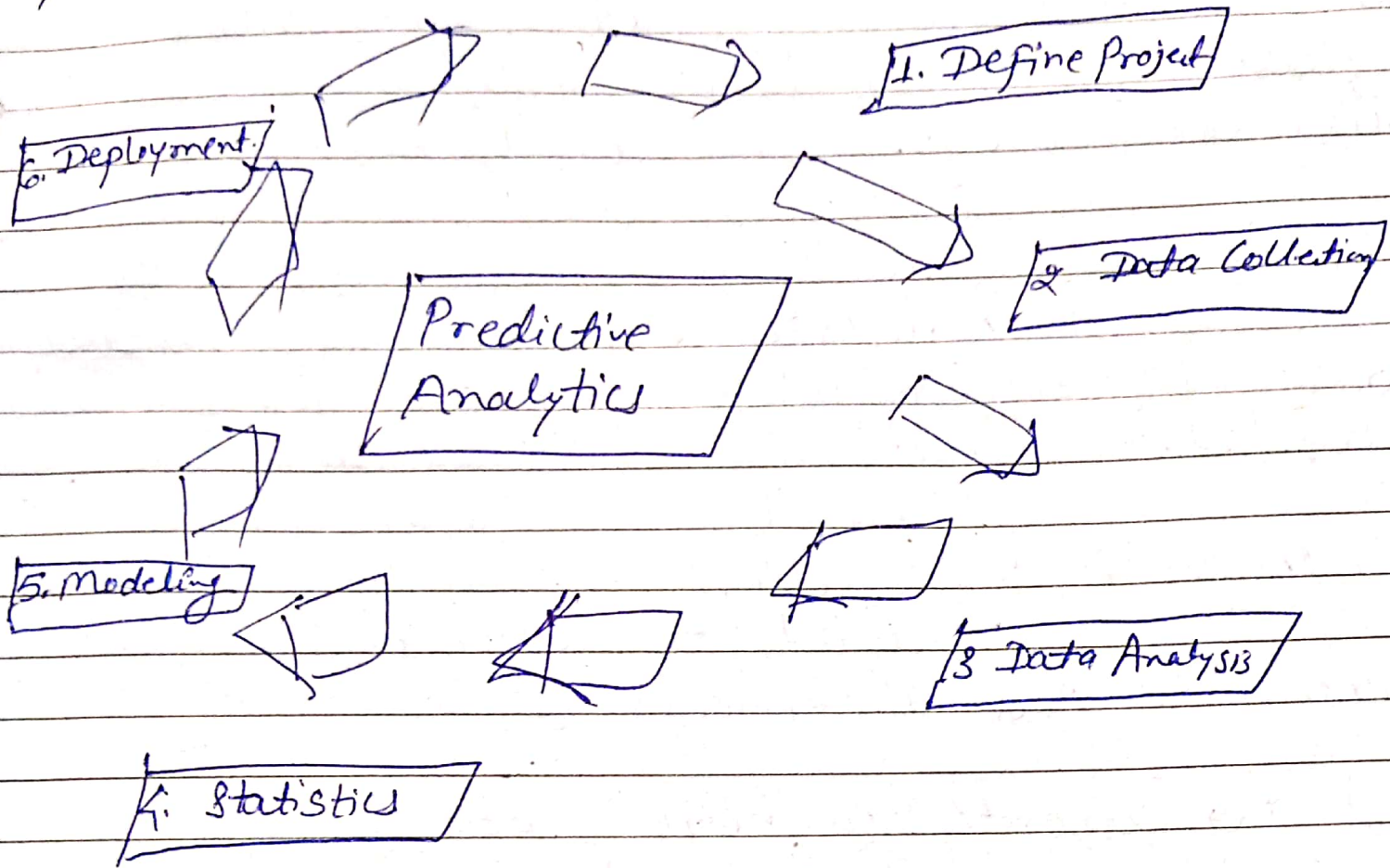
Predictive analytics is the branch of the advanced analytics which is used to make prediction about unknown future events.

- predictive analytics use many techniques from data mining, statistics, modeling, machine learning and artificial intelligence to analyze current data to make prediction about future.



[Predictive Analytics Value Chain]

⇒ Predictive analytics process:



1) Define project :- Define the project outcomes, deliverables, scoping of the effort, business objectives, identify the data sets which are going to be used.

2) Data Collection:- Data mining for predictive analytics prepares data from multiple sources for analysis. This provides a complete view of the customer interactions.

3) Data Analysis:-

Data analysis is the process of inspecting, cleaning, transforming, and modeling data with the objective of discovering useful information, arriving at conclusions.

4) Statistics:- Statistical analysis enables to validate the assumptions, hypotheses and test them with using standard statistics model.

5) Modeling:- predictive modeling provides the ability to automatically create accurate predictive model about future.

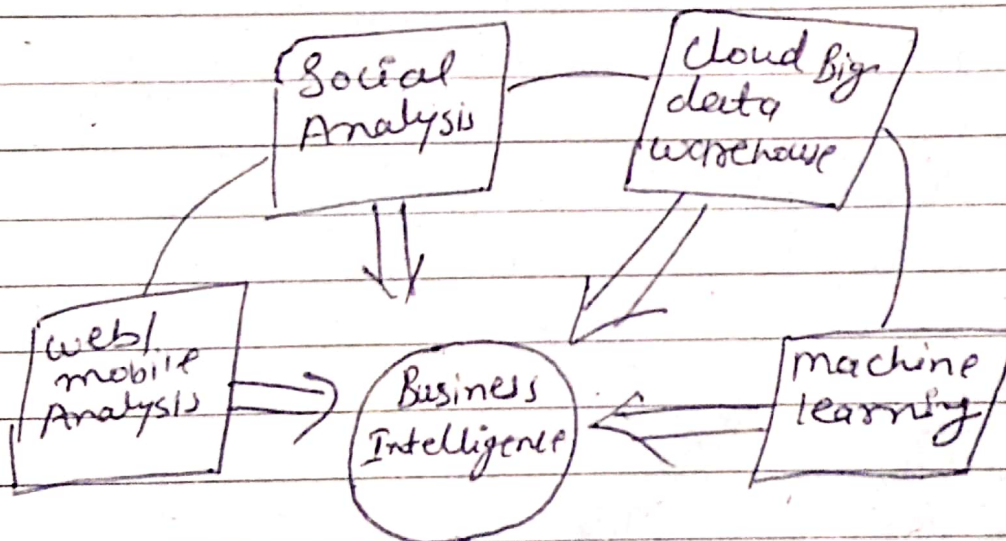
6) Deployment:- predictive model deployment provides the option to deploy the analytical results into every day decision making process to get results.

~~7) Model~~

★ Mobile Business Intelligence :-

mobile Business Intelligence (BI) is the ability to access BI-related data such as KPIs, business metrics, and dashboards on mobile devices.

- mobile BI is a system comprising both technical and organizational elements that present historical and/or real-time information to its user for analysis on mobile devices such as smartphones and tablets (not laptop), to enable effective decision making and management support, for the overall purpose of increasing firm performance.



- The definition of mobile BI refers to the access and use of information via mobile devices. With the increasing use of mobile devices for business not only in management positions, mobile BI is able to bring business intelligence and analytics closer to the user when done properly.
- whether during a train journey, ~~in the~~ airport departure or during a meeting break, information can be consumed almost anywhere and anytime with mobile BI.

* Benefits of mobile BI:-

- 1) The first major benefit is the ability for end users to access information in their mobile BI system at any time and from any location.
- 2) Improves their daily operations
- 3) React most quickly
- 4) Speeds up the decision making
- 5) Better time utilization.

★ BI	Data Science
1) Backward approach.	1) forward approach.
2) provides reports but that can't predict future	2) provides patterns and reports that state
3) preplanned and retrieved data → (no explanation)	how the data will look in future
4) helps to answer the questions you know	3) No-preplanning data is retrieved from source as per requirement.
5) Business uses.	4) help you to discover new question.
6) what happened?	5) Data Scientists.
7) Input: past data output: present solution	6) what will happen? 7) Input: Dynamic request ^{data} output: future predict.

Information technology:

main frame computing → mini computing → personal computing.

Desktop internet → mobile internet.

- Mobile BI is defined as "the capability that enables the mobile workforce to gain business insights through information analysis using application optimized for mobile devices"

★ What is Crowdsourcing?

is the process of getting work or funding, usually online, from a crowd of people.

• The word ~~it is~~ is a combination of the two words 'crowd' and 'out-sourcing'. The idea is to take work and out-source it to a crowd of workers.

~~Crowdsourcing & quality~~: The principle of

~~Crowdsourcing~~ is that more heads are better than one. By ^(कॉट मॉर्गन) canvassing a large crowd of people for ideas, skills, or participation, the quality of content and idea generation will be superior.

• The idea is generally to introduce new or more developed skill sets or a larger work force to achieve some specific goal.

★ Benefits of Crowdsourcing:

Crowdsourcing is a powerful business marketing tool as it allows an organization to leverage the creativity and resources of its own audience in promoting and growing the company for free.

- Crowd sourcing increase the productivity of a company while minimising labor expenses
- The internet is a time proven strategy for soliciting feedback from an active and passionate consumer base.

★ Inter & Trans-firewall analytics:-

Over the last 100 years, supply chain has evolved to connect multiple companies and enables them to collaborate to create ~~enormous~~ value to the end-consumer through concept like

→ CPER (Collaborative planning, forecasting and Replenishment):- A collection of business practice that leverage the internet and electronic data interchange to reduce inventories and expenses while improving customer service.

→ VMI (Vendor managed Inventory):- A

technique used by customers in which manufacturers receive sales data to forecast consumer demand more accurately.

- There are instances where a retailer and a social company can come together to share insights on consumer behaviour that will benefit both concerns.
 - Some of the more progressive companies will take this a step further and work on leveraging the large volume of data outside the firewall such as social media, location data etc.
- In other words, it will not be long before internal data and insight from within the enterprise firewall is no longer a differentiator. We call this trend the move from intra to inter and trans firewall analytics.

* Information management: Is the process of collecting, storing, managing and maintaining information in all its forms.

- Information management may also be called information asset management.
- Information can be in the form of physical data (such as papers, documents and books) or digital data assets.

(Big data technologies)

Date: / / Page no: _____

- Unstructured Information management Architecture (UIMA): This is one of the 'secret sauce' elements in the 'secret sauce' behind IBM Watson's system that reads massive amounts of data and organises for just in time processing.
- Watson beat the quiz program champion in 2011 and is now used for many business applications, like diagnosis in healthcare situations. Natural language processing is another capability that helps extend the power of big data technologies.

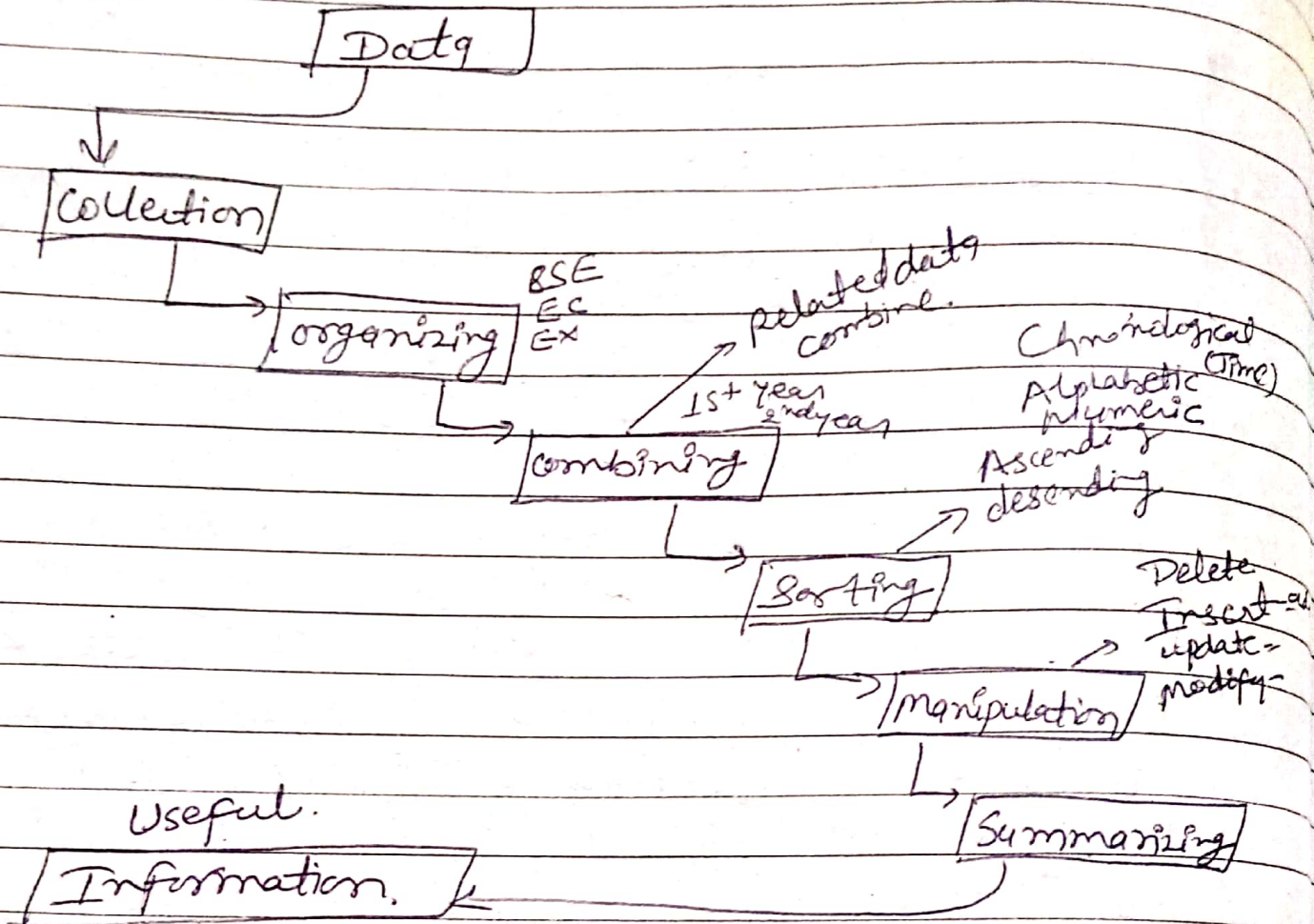
(work on concept human being dekhata hai, samajta hai thum deside best result
Virtual assistance ko bhi sawal ka sahi
Jawab dega.

Unit - III

Date: / / Page no:

★ Processing Big data:-

- Data processing:-



dissimilar
↑
different

Date: / / Page no: _____

* Integrating disparate data stores :-

• Data integration :- Data integration involves combining data from several disparate sources, which are stored using various technologies and provide a unified view of the data.

- Data integration becomes increasingly important in case of merging systems of two companies or consolidating applications within one company to provide a unified view of the company's data assets. The latter initiatives is often called a data warehouse.

* Data integration areas :- data integration is a term covering several distinct sub-area such as

- Data warehouse
- Data migration
- ~~Enterprise application~~ / Information integration
- master data management (MDM)

represents a set of tools and processes used by an enterprise to consistently manage their non-transactional data

* Data lake storage on AWS: The most Secure, durable, and scalable storage to build your data lake through Amazon S3 (Simple Storage Service).

- Amazon S3 :- is the largest and most performant storage service for structured and unstructured data and the storage service of choice to build a data lake. With Amazon S3, you can cost-effectively build and scale a data lake of any size in a secure environment where data is protected by 99.9% of durability.

- ~~with a data lake~~
- with a ^{data} lake built on Amazon S3, you can use native AWS services to run big data analytics, AI, ML, high performance computing (HPC) and media data processing applications to gain insight from your unstructured data sets.

- In the modern data marketplace, disparate data sources are largely what we refer to as unstructured in nature, making up the bulk of "big data" volumes. Databases, data warehouses, and data lakes are all

governed in unique ways. Hadoop brings different data types together in one place but does not guarantee any substantive forms of organization. (original (HSA))

* Mapping data to the programming framework :-

• The MapReduce programming framework uses two tasks common in functional programming: Map and Reduce. MapReduce is a new parallel processing framework and Hadoop is its open-source implementation on a single computing node or on clusters.

• Big data processing can be performed through a programming paradigm known as Map-Reduce. MapReduce is a programming model designed for processing large volumes of data in parallel by dividing the work into a set of independent tasks.

• It is based on "break problem up into smaller sub-problem" strategy and can be compared with SQL & SQL-Based BI tools.

- MapReduce is a programming model and an associated implementation for processing and generating large data sets.
- Users specify a Map function that processes a value pair to generate a set of intermediate value pairs, and a Reduce function that merges all intermediate values associated with the same intermediate key.

- MapReduce processing consists of two phases: map phase and Reduce phase. Each phase of MapReduce consists of key-value pairs as input and output. Hadoop divides the input to a MapReduce job into the fixed-size pieces called input splits.

- Map() function $(K_1, V_1) \rightarrow \text{List}(K_2, V_2)$
It performs filtering and sorting of tasks into queue.

- Reduce() function. It performs a summary operation of best candidate resource for task execution.

* Transforming data for processing :-
• The main purpose of data transformation Data to and feature extraction is to enhance the data in such a way that it increase the likelihood that the classification algorithm will be able to make meaningful prediction.

Data pre-processing techniques:

- Preprocessing of data involves a set of key tasks that demand extensive computational infrastructure.
- The bigger the data sets, more complex mechanism are needed to process it before analysis and visualization.
- Pre-processing prepares the data and make the analysis feasible while improving the effectiveness of the result.

• following are some crucial steps involved in data pre-processing.

- Data cleansing.
- Data Normalization
- Data Transformation
- Missing value imputation
- Noise identification.
- Minimizing the preprocessing tasks.

→ Data cleansing! -

cleansing the data is usually the first step in data processing and is done to remove the unwanted elements as well as to reduce the size of the data sets.

- This will make it easier for the algorithm to analyze it.

→ Data Normalization:- Normalization is the process of reorganizing data in a database so that it meets basic requirements:

- 1) There is no redundancy of data (all data is stored in only one place),
- 2) data dependencies are logical (all related data items are stored together).

- Normalization increase result performance.
Normalization is also known as data normalization.

NF (Normal forms)

1 NF

2 NF

3 NF

BCNF

5 NF

- Data Transformation: Data transformation is the process of converting data or information from one format to another, usually from the format of a source system into the required format of a new destination system.

- The usual reason for this data migration is the adoption of a new system that totally different from the previous one.

Data transformation ka main goal yahi hota hai ki system ko is kabil banana ki vo new system ko adopt kr sake

- Data transformation involves two key phases:

→ Data Mapping :-

The assignment of elements from the source base or system toward the destination to capture all transformation that occurs

- This is more complicated when there are complex transformations like merge or one or one-to-many rules for transformation.

→ Code generation :-

The creation of the actual transformation program. The resulting data map specification is used to create an executable program to run on computer systems.

→ Missing value imputation.

Imputation: In statistics, imputation is the process of replacing missing data with substituted values.

missing value are data that haven't been extracted or stored due to budget restrictions or other limitations in the data extraction process.

missing value is not something to be ignored as it could skew your results.

fixing the missing value issue is challenging. handling it without utmost care could easily lead to complication in data handling and wrong conclusion.

ML (machine learning technique) is used to missing value imputation.

(Intruption in data is known as noise)

➔ Noise identification: Data gathering is not always perfect, but the data mining algorithm would always assume it to be.

➔ Data with noise can seriously affect the quality of the results, ~~tackling~~ tackling this issue is crucial.

• Noise can affect the input features, output or both in most cases.

- There are two popular approaches to remove from Data

→ Data polishing methods are used to eliminate ~~noise~~ the noise.

→ The other method involves using noise filters that can identify and remove instances with noise from the data.

→ Minimizing and pre-processing data.

- preparing the data for your data analysis algorithm can involve many more processes depending on the application unique demands.

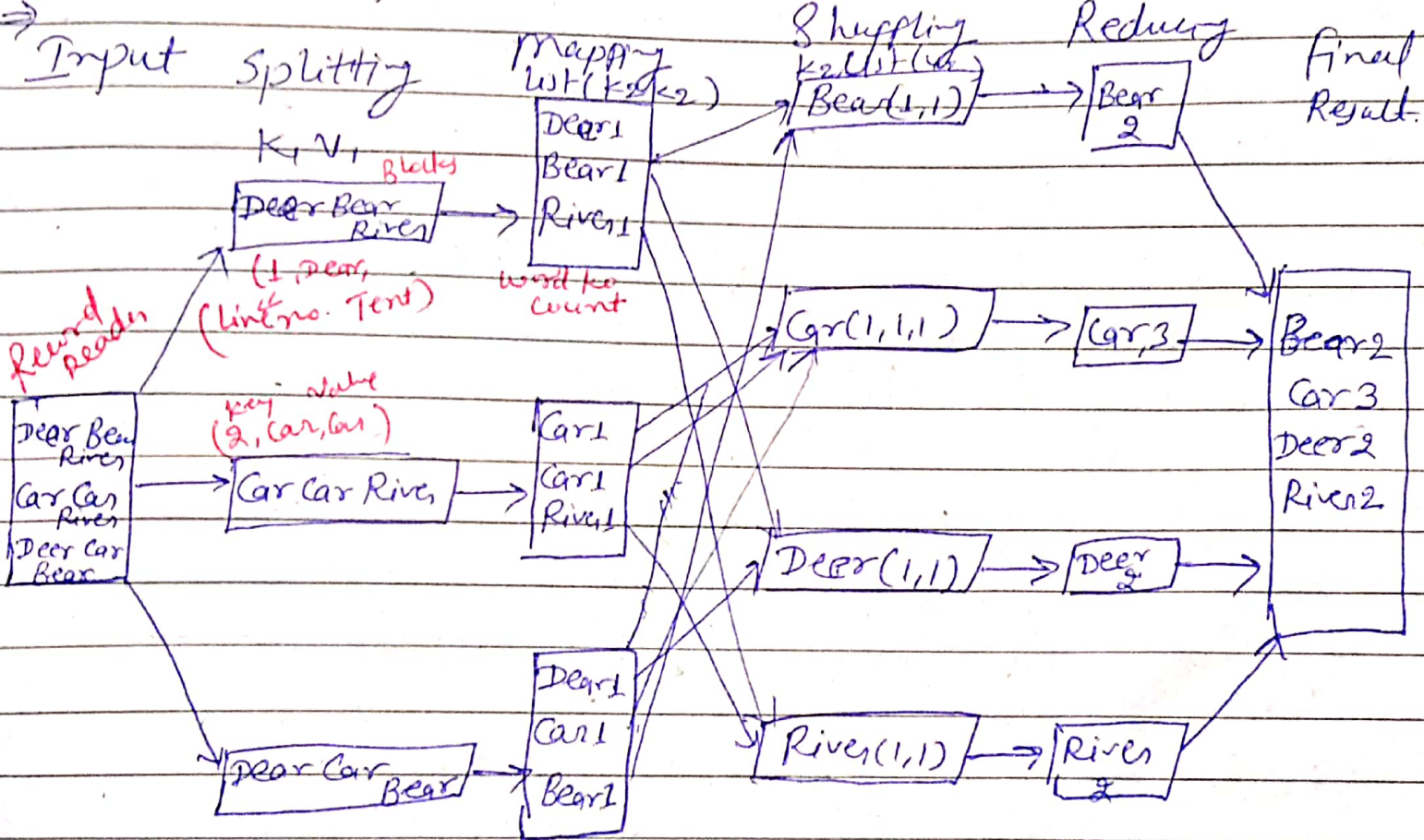
- However basic process like cleaning, deduplication and normalization can be avoided in most cases if you choose right source for data extraction.

→ But it is highly unlikely that a raw source can give you clean data.

mapReduce divides a task into small parts and assigns them to many computers. Date: / / Page no: _____

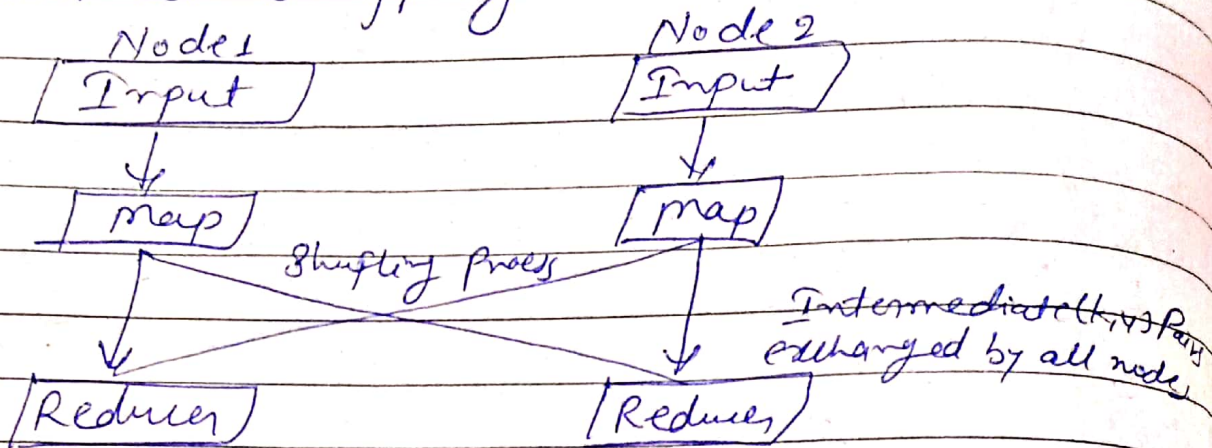
*** How MapReduce works:-**

- The Map task takes a set of data and converts it into another set of data, where individual elements are broken down into tuples (key-value pairs).
- The Reduce task takes the output from the Map as an input and combines those data tuples (key-value pairs) into a smaller set of tuples.

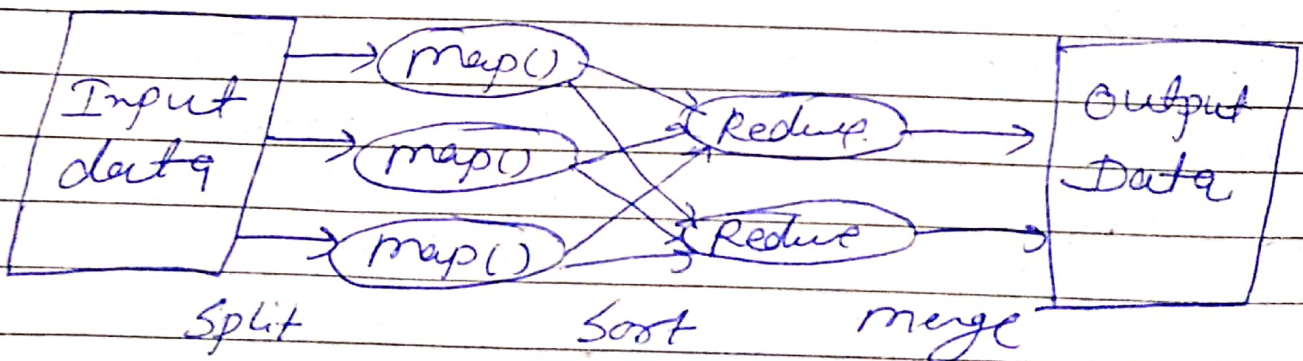


[The overall MapReduce word count process]

- Shuffling: The process of exchanging the intermediate outputs from the map tasks to where they are required by the reducers is known as "Shuffling".



- MapReduce programming model using two components: a Job tracker (master node) and many task trackers (slave nodes). The Job Tracker is responsible for accepting job requests, for splitting the data input. The task tracker executes tasks as ordered by the master node. The task can be either a map or a reduce.



★ Subdividing data in preparation for Hadoop MapReduce:-

The whole process goes through four phases of execution, namely, Splitting, mapping, Shuffling, and Reducing.

- Input Splits:- An input to a MapReduce job is divided into fixed-size pieces called input splits. Input split is a chunk of the input that is consumed by a single map.
- Mapping:- This is the very first phase in the execution of map-reduce program. In this phase data in each split is passed to a mapping function to produce output value.
- In our example, a job of mapping of mapping phase is to count number of occurrence of each word from input splits and prepare a list in the form of $\langle \text{key value} \rangle$ (word, frequency).

- **Shuffling**:- This phase consumes the output of mapping phase. Its task is to consolidate the relevant records from mapping phase output.
- our example, the same words are clustered together along with their respective frequency.

- **Reducing**:- In this phase, output values from the shuffling phase are aggregated. This phase combines values from shuffling phase and returns a single output value. In sort, this phase summarizes the complete dataset.
- In our example, this phase aggregates the values from shuffling phase i.e., calculate total occurrence of each word.

★ How MapReduce organizes work?
- Hadoop divides the job into tasks. There are two types of tasks:

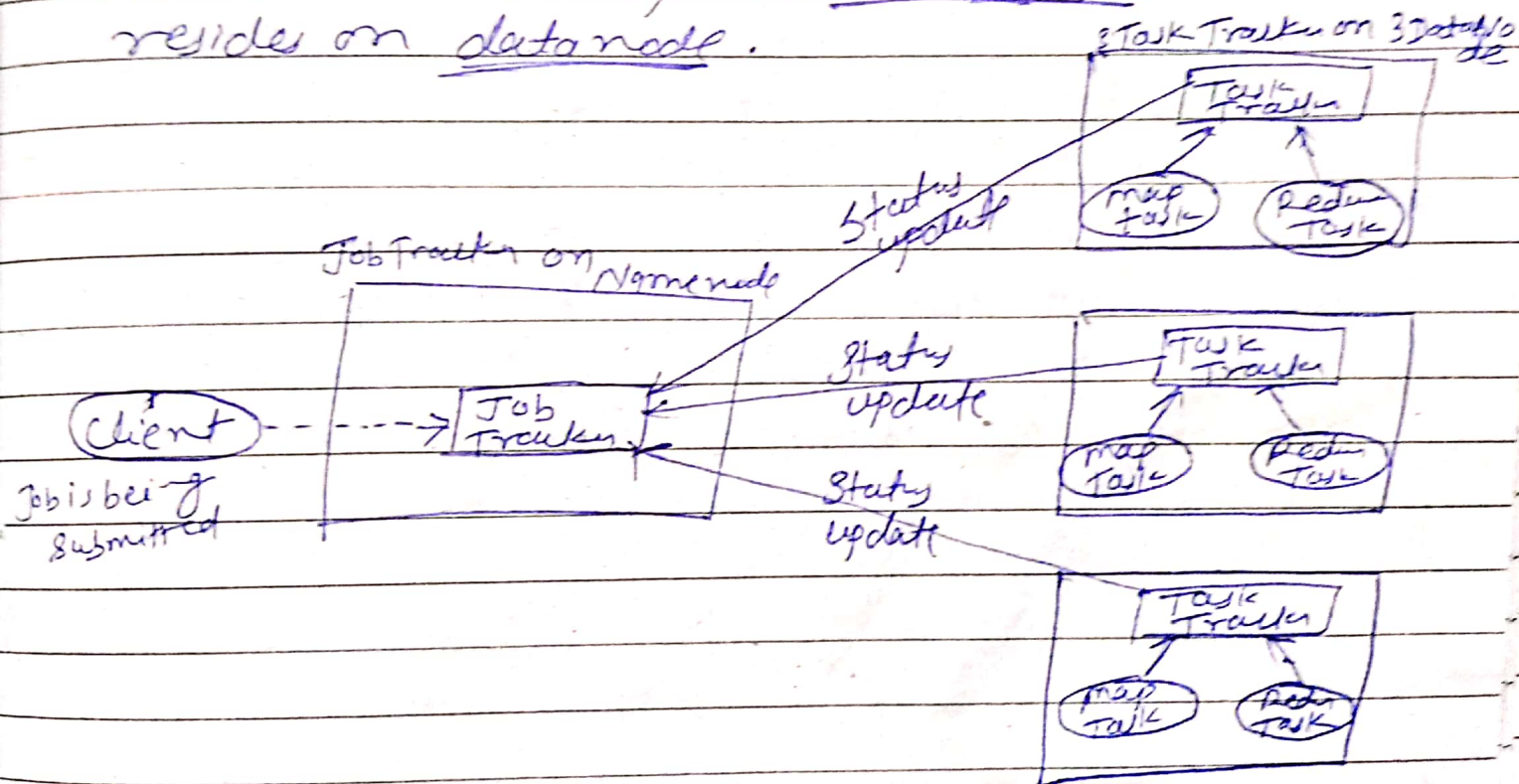
- 1) Map tasks (Splits & Mapping)
- 2) Reduce tasks (Shuffling, Reducing)

- The complete execution process (execution of map and reduce tasks, both) is controlled by two types of entities called a.

1) Job tracker: Acts like a master (responsible for complete execution of submitted job)

2) Multiple Task Tracker: Acts like slaves, each of them performing the Job

• For every job submitted for execution in the system, there is one Job tracker that resides on namenode and there are multiple task trackers which reside on datanode.



- A job is divided into multiple tasks which are then run on to multiple data nodes in a cluster.
- It is the responsibility of job tracker to coordinate the activity by scheduling tasks to run on different data nodes.
- Execution of individual task is then to look after by task tracker, which resides on every data node executing part of the job.
- Task tracker responsibility is to send the progress report to the job tracker.
- In addition, task tracker periodically sends "heartbeat" signal to the job tracker so as to notify him of the current state of the system.
- Thus job tracker keeps track of the overall progress of each job. In the event of task failure, the job tracker can reschedule it on a different task tracker.

Unit IV

Date: / / Page no: |

Hadoop MapReduce

* Employing Hadoop MapReduce:-

Configuration modes: Three Hadoop Configuration files:- settings to run hadoop cluster, mapreduce and spark.

• Hadoop cluster can be run in any of three supported models:

• Local (standalone) mode:-

By default, Hadoop is configured to run in a non-distributed mode, as a single java process where instead of HDFS, local file system is used for the purpose of input & output. Here we do not perform configurations, all services run on single JVM. This is useful for debugging.

• Pseudo-Distributed operations:-

Hadoop can also be run on a single-node in a pseudo-distributed mode, where each hadoop daemon runs in a separate java process.

• It requires configuration of three Hadoop installation files: `hdfs-site.xml`, `core-site.xml`, and `mapred-site.xml`. Here, `data node`, `namenode`, `job tracker`, `task tracker`, all supposed to be run on single node. This mode is useful for development and testing HDFS.

• Fully Distributed:-

All Hadoop services run in different JVMs belonging to one cluster. All daemons get executed in different nodes that form a multinode cluster.

Installing a Hadoop cluster typically involves unpacking the software on all the machines in the cluster or installing it via a packaging system as appropriate for the operating system.

* Creating the components of Hadoop map-reduce jobs:-

we can create a hadoop MapReduce job with Spring and apache hadoop by following these steps:

- 1) Get the required dependencies by using maven.
- 2) Create the mapper component.
- 3) Create the reducer component.
- 4) Configure the application context.
- 5) Load the application context when the application starts.

• Create the mapper component: A mapper is a component that divides the original problem into smaller problems that are easier to solve. we can create a custom mapper component by extending the class and overriding its `map()` method.

- Creating the reducer Component: A reducer is a component that removes the unwanted intermediate values and passes forward only the relevant key-value pairs. We can implement our reducer by extending the class and overriding its `reduce()` method.

- Configuring the application Context:

Because Spring does not support Java configuration, we have to configure the application context of our application by using XML.

- Loading the application Context when the application starts:

We can execute the created Hadoop job by loading the application context when our application is started. We can do this by creating a new class `ContextLoader` and providing the name of our application context configuration file as a constructor parameter.

* Distributed data processing across server farms:

Distributed data processing is a computer-networking method in which multiple computers across different locations share computer-processing capability. This is in contrast to a single, centralized server managing and providing processing capability to all connected systems.

- Computers that comprise the distributed data-processing network are located at different locations but interconnected by means of wireless or satellite links.

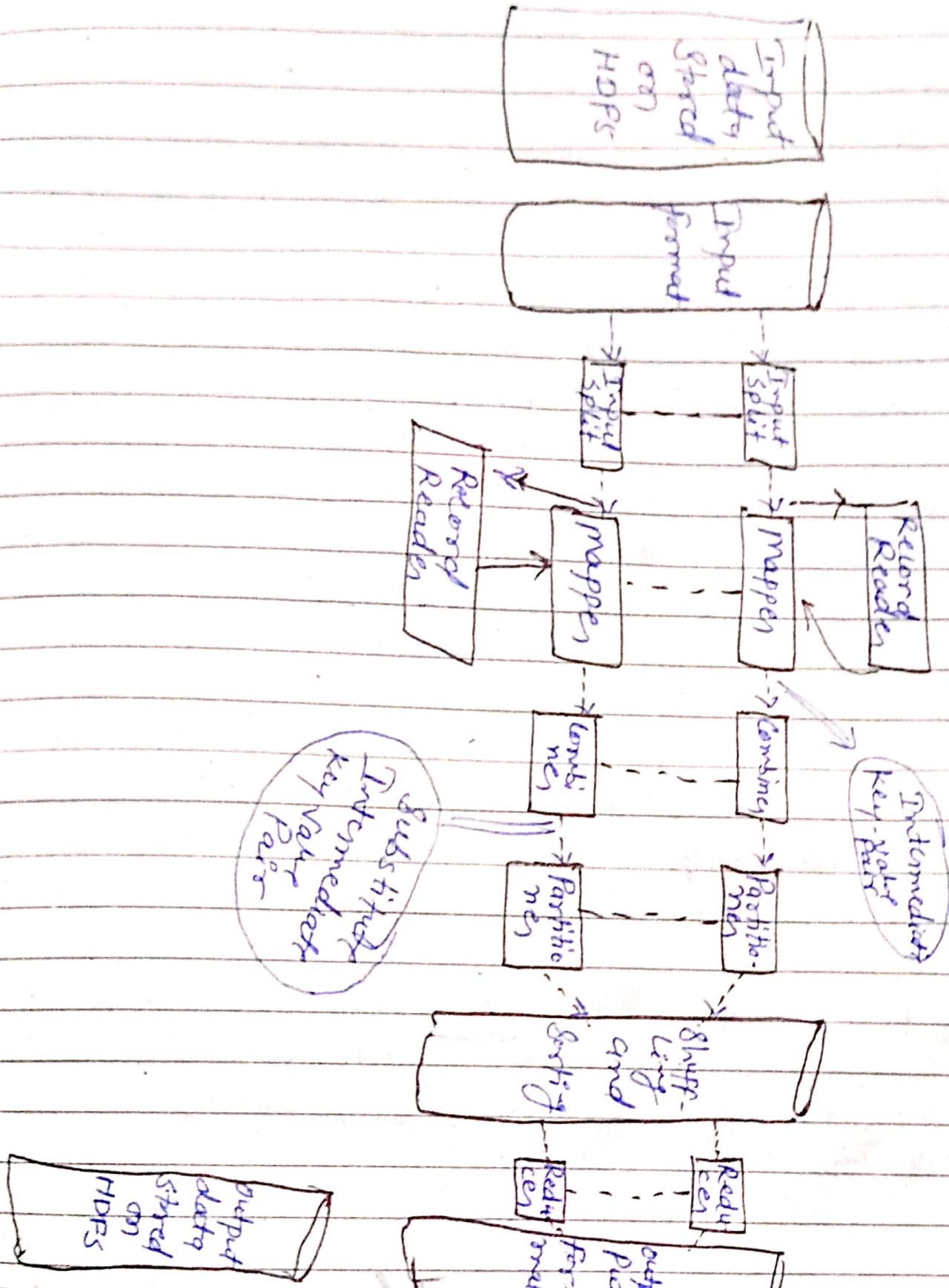
Advantages of data processing:

- Lower cost :- Larger organizations invest in expensive mainframe and supercomputers to function as centralized servers.
- Reliable :- Distributed data processing is more reliable, since multiple control centers are spread across different machines.

- Improve performance and Reduced processing times:- Single computers are limited in their performance and efficiency. An easy way to increase performance is by adding another computer to a network.

- Flexible:- Individual computers that comprise a distributed network are present at different geographical locations. For example, an organizational - distributed network comprising of three computer can have each machine in different branch.

* Executing Hadoop MapReduce Jobs:-



- Hadoop MapReduce is the data processing layer. It processes the huge amount of structured and un-structured data stored in HDFS. MapReduce processes data in parallel by dividing the job into the set of independent tasks. So, parallel processing improves speed and reliability.

- Steps of MapReduce Job execution flow:
MapReduce processes the data in various phases with the help of different components.

1) Input File :- In input file, data for MapReduce Job is stored. In HDFS, input files reside. Line based log files and binary format can also be used.

2) Input Format :- After that input format defines how to split and read these input files. It selects the files or other objects for input.

3) Input Splits :- It represents the data which will be processed by an individual mapper. For each split, one map task is created.

- 4) Record Reader:- It communicates with the InputSplit. And then converts the data into key-value pairs suitable for reading by the mapper.
- 5) Mapper:- It processes input record produced by the Record Reader and generates intermediate key-value pairs.
- 6) Combiner:- Combiner is mini-reducer which performs local aggregation on the mapper output. It minimizes the data transfer b/w mapper and reducer.
- 7) Partitioner:- Partitioner comes into the existence if we are working with more than one reducer.
- 8) Shuffling & Sorting:- After partitioning, the output is shuffled to the reducer node. The shuffling is the physical movement of the data which is done over the net.
- 9) Reducer:- Reducer then takes set of intermediate key-value pairs produced by the mapper as input.

10) Record writer :- It writes these output key-value pair from the Reducer phase to the output files.

11) Output format :- Output format defines the way how Record-writer writes these output key-value pairs in output files.

* The Building Blocks of Hadoop MapReduce
distinguishing Hadoop daemons :-

Hadoop employs a master/slave architecture for both distributed storage and distributed computation. The distributed storage system is called the Hadoop Distributed File System (HDFS).

- on a fully configured cluster, "running Hadoop" means running a set of daemons, or resident programs, on the different servers in your network.

- These daemons have specific roles, some exist only on one server, some exist across multiple servers. These daemons include:

- NameNode
- DataNode
- Secondary Name node
- Job Tracker
- Task Tracker

1) NameNode:- The NameNode is the master of HDFS that directs the slave DataNode daemons to perform the low-level I/O tasks.

It is the bookkeeper of HDFS, it keeps track of how your files are broken down into file blocks, which nodes store those blocks and the overall health of the distributed file system.

2) DataNode:- Each slave machine in your cluster will host a datanode daemon to perform the grunt work of the distributed file system - reading and writing HDFS blocks to actual files on the local file system.

3) Secondary NameNode (SNN):- The SNN is an assistant daemon for monitoring the state of the cluster HDFS. Like the NameNode, each cluster has one SNN and it typically resides on its own machine as well. No other DataNode or TaskTracker daemons run on the same server.

4) Job Tracker:- once you submit your code to your cluster, the Job Tracker determines the execution plan by determining which files to process, assigns nodes to different tasks, and monitors all the tasks as they running. There is only one Job Tracker daemon per hadoop cluster.

5) Task Tracker:- Each Task Tracker is responsible for executing the individual tasks that the Job Tracker assigns. Although there is a single Task Tracker per slave node, each Task Tracker can spawn multiple JVM to handle many map or reduce task in parallel.

* Investigating the hadoop distributed file system selecting appropriate execution models:

• Single Node (Local Mode or Standalone mode) :-

Standalone mode is the default mode in which hadoop runs. Standalone mode is mainly used for debugging where you don't really use HDFS.

- You can use input and output both as a local file system in standalone mode.

- Standalone mode is usually the fastest hadoop modes as it uses the local file system for all the input & output.

• Pseudo-distributed Mode :- The pseudo-distributed mode is also known as single-node cluster where both Name Node and Data Node will reside on the same machine.

- In pseudo-distributed mode, all the hadoop daemons will be running on a single host. Such configuration is mainly used while testing when we don't need to think about the resources and other users sharing the resource.

- In this architecture, a separate JVM is spawned for every hadoop component as they could communicate across network sockets, effectively producing a fully functioning and optimized mini-cluster on a single host.

● Fully distributed mode: As the name suggests, this mode involves the code running on an actual hadoop cluster. It is the mode in which you see the actual power of hadoop, when you run your code against a very large input on 1000s of servers.

- It is always difficult to debug a mapreduce program as you have mappers running on different machines with

different piece of input. You can never know where the mappers are going to run eventually. Also with large inputs it is likely that the data will be irregular in its format.

Unit V

Date: / / Page no: |

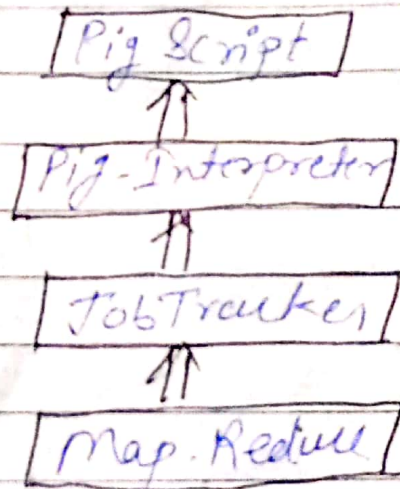
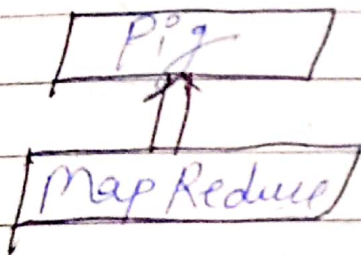
Big data tools and Techniques.

* Installing and Running Pig.

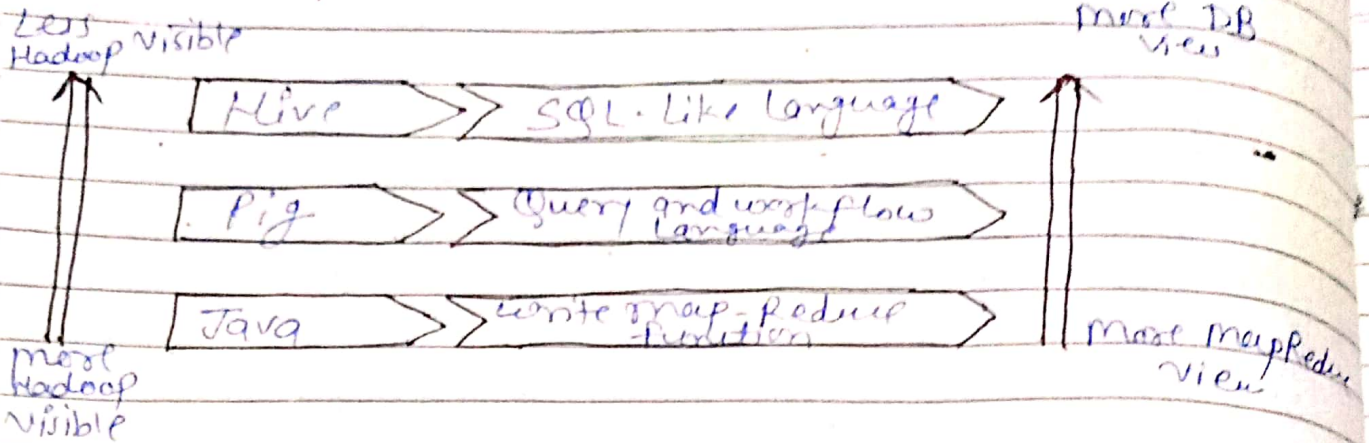
Apache pig :-

• Apache pig is an abstraction over mapReduce. It is a tool/platform which is used to analyze larger sets of data representing them as data flows.

- Pig provides a high-level language known as Pig latin.
- Compiles down to MapReduce jobs.
- Developed by Yahoo. open source language.



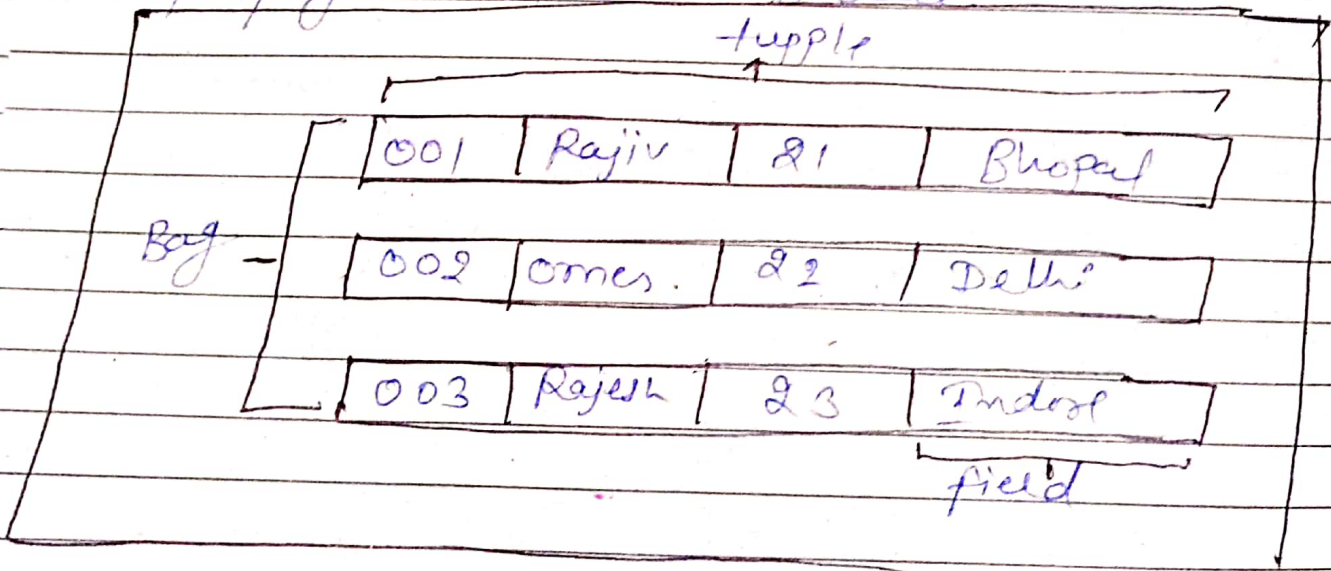
* Levels of Abstraction:-



* Pig Latin Data model:-

The data model of Pig Latin is fully nested and it allows complex non-atomic data types such as map and tuple.

- given below is the diagrammatical representation of Pig Latin data model



- Atom:- Any single value in pig Latin, irrespective of their data, type is known as an Atom. It is stored as string and can be used as string and number. int, long, float, double, chararray, and bytearray are the atomic value of pig.
- Field:- A piece of data or a simple atomic value is known as field.
- Tuple:- A record that is formed by an ordered set of fields is known as tuple, the fields can be of any type. A tuple is similar to a row in a table of RDBMS.
- Bag:- A Bag is an unordered set of tuples. In other words, a collection of tuples (non-unique) is known as a bag.

* Comparison with databases:-

- with the range of technologies in the big data world, there is often confusion to choose from them. It is required to handle huge databases efficiently with big data, and the options for managing and querying data are also needed.
- When it comes to managing databases, SQL (Structured, Query Language) is the old friend well tried and tested by everyone for data analysis. But the complicated world of Hadoop needs high-level data analysis model.
- ~~These~~
- Though old SQL still is the favorite of many and is popularly used in numerous organizations, Apache Hive and pig have become the buzz terms in the big data world today.
- These tools provide easy alternatives to carry out the complex programming of MapReduce helping data developers and analysts.

Comparison Pig, Hive, SQL.

Criteria	Hive	Pig	SQL
1) Language used	Uses hive SQL a declarative language	Uses pig latin a procedural data flow language	SQL itself is a declarative language.
2) Definition	An open source built with an analytical focus used for analytical queries.	An open source and high-level data flow language with a multi-query approach.	General purpose database language for analytical and transactional queries.
3) Suitable for.	Batch processing OLAP (online Analytical processing)	Complex & nested data structure	Business demand for fast data analysis
4) Developed by	Facebook	Yahoo	Oracle.
5) Operational for	Structured data	Structured & Semi-structured	Relational database management
6) Mainly used by	Data Analysts	Researchers & Programmers.	Data analysts Data Scientists and programmers

* Apache Pig - User defined functions :-

- Apache pig provides extensive support for User defined Functions (UDFs). Using these UDFs we can define our own functions and use them.
- The UDF support is provided in six programming languages, namely, Java, Python, JavaScript, Ruby and Groovy.
- For writing UDFs complete support is provided in java and limited support is provided in all the remaining languages.
- Using Java, you can write UDFs involving all parts of the processing like data load/store, column transformation, and aggregation.
- In apache pig, we also have a Java repository for UDFs named piggybank. Using piggybank we can access Java UDFs written by other users, and contribute our own UDFs.

- Type of UDF's in Java:

While writing UDF's in java, we can create and use the following three types of functions:

- **filter functions:** - The filter functions are used as conditions in filter statements. These functions accept a pig value as input and return a boolean value.
- **Eval functions:** - The Eval functions are used in FOREACHGENERATE statements. These functions accept a pig value as input and return a pig result.
- **Algebraic functions:** - The algebraic functions are used on inner bags in a FOREACHGENERATE statement. These functions are used to perform full MapReduce operations on an inner bag.

★ Apache Pig operators:-

The apache pig operator is a high-level procedural language for querying large data sets using hadoop and the Map Reduce platform.

- A pig latin statement is an operator that takes a relation as input and produces another relation as output.
- These operators are the main tools for pig latin provides to operate on the data. They allow you to transform it by sorting, grouping, joining, projecting and filtering.

Relational operators:-

Relational operators are the main tools pig latin provides to operate on the data. It allows you to transform the data by sorting, grouping, joining, projecting and filtering.

* Relational operations:-

loading and storing.

• LOAD :- To load the data from the file system (local/HDFS) into a relation.

• STORE :- To save a relation to the file system (local/HDFS).

Filtering.

• FILTER :- To remove unwanted rows from a relation.

• DISTINCT :- To remove duplicate rows from a relation.

• FOREACH, GENERATE :- To generate data transformation based on columns of data.

• STREAM :- To transform a relation using an external program.

Grouping & Joining.

• Join :- To join two or more relations.

- SPLIT :-
- CORROUP :- To group the data in two or more relation.
- GROUP :- To group the data in a single relation.
- CROSS :- To create the cross product of two or more relations.

Sorting :-

- ORDER :- To arrange a relation in a sorted order based on one or more field (ascending or descending).
- LIMIT :- To get a limited number of tuples from a relation.

Combining & Splitting :-

- UNION :- ~~To~~ To combine two or more relation into single relation.
- SPLIT :- To split a single relation into two or more relations.