

Q What is data?

The Quantities, characters or symbols on which operations are performed by a computer which may be stored & transmitted in the form of electrical signals & recorded on magnetic, optical or mechanical recording media.

Q What is Big data?

Big data is a collection of data that is huge in volume yet growing exponentially with time. It is a data with so large size & complexity that none of traditional data management tools can store it or process it efficiently. Big data is also a data but with huge size.

Big data analytics.

It largely involves collecting data from different sources, manage it in a way that it becomes available to be consumed by analysts & finally deliver data products useful to the organization business.

Data Analyst

- The process of the extraction of information from a given pool of data is called data analytics. A data analyst is a person who engaged in this form of analysis.
- A data Analyst extracts the information through several methodologies like data cleaning, data conversion & data modeling.
- data analytics allows the industries to process fast queries & produce actionable results that are needed in a short duration of time.

... This restricts data analytics to a more short-term growth of the industry where quick action is required.

Data Analyse Rule

- Gathering data, using database query languages to retrieve & manipulate information.

- Perform data filtering, cleaning & early-stage transformation.
- Communicating results with the team using data visualization.
- Visualisation.
- Converting raw data into visual form.
- Data visualization is their key role.
- work with the management team to understand business requirements.

Data Engineer :-

- A data engineer is a person who specializes in preparing data for analytical usage.
- Data Engineering also involves the development of platforms & architectures for data processing.
- Data engineers have to deal with Big data where they engage in numerous operations like data cleaning, management, transformation, data duplication etc.
- A data engineer is more experienced with core programming concepts & algorithms.

Eg.

Developing a cloud infrastructure to facilitate real time analysis of data requires various development principles. Therefore, building an interface API is one of the job responsibilities of a data ~~system~~ engineer.

Tools Used By Data Engineers

- ① Hadoop ② Spark

Data Engineer Responsibilities

- development construction & maintaining of data architecture.
- Conducting testing on large scale data platform.
- Handling error logs & building robust data pipelines.
- Ability to handle raw & structured data.
- Ensure & support to the data architecture utilized by data scientist & analysts.

- Development of data process for data modeling, mining & data production.

Data Scientist

- Data Scientists are analytical experts who utilize their skills in both Technology & Social Science to find & test & manage data. They use industry knowledge, Contextual understanding, skepticism of existing assumptions - to uncover solutions to business challenges.
- A data Scientist's work typically involves making sense of messy unstructured data, from sources such as smart devices, social media feeds, & emails that don't neatly fit into a database.
- Data scientists are big data wranglers, gathering & analyzing large sets of structured & unstructured data.

Data Scientist Responsibilities :-

- Performing data preprocessing that involves data transformation as well as data cleaning.
- Using various machine learning tools to forecast & classify pattern in the data.
- Increasing the performance & accuracy of ML algorithms through fine-tuning & further performance optimization.

Applications of Big data in 10 Industry Verticals

① Banking And Securities Challenges :-

- early warning for securities fraud & Trade visibility
- Card fraud detection & credit trails.
- enterprise credit risk reporting.
- customer data transformation & analytics.

② Communication, Media & entertainment :- Challenges :-

- ① collecting, analyzing & utilizing consumer insights.

- ② Leveraging mobile & Social media content
- ③ Understanding patterns of real time media content usage.

③ Healthcare providers:

Challenges

- ① Rising medical costs.
- ② unavailability / inadequacy / unusable data.

④ Education

Challenges:

- ① Incorporating data from varied sources.
- ② Untrained staff & Institutions about Big data.
- ③ Issues of privacy & data protection.

⑤ Manufacturing & Natural Resources

Challenges:

- (i) Increase in the volume, complexity & velocity of data due to rising demands of Natural resources.
- ② large volumes of uncapped data from the manufacturing industry.
- ③ underutilization of data prevents improved quality, energy efficiency, reliability & better profit margins.

⑥ Government

Challenges

- ① Integration
- ② Interoperability of big data.

⑦ Insurance

Challenges:

- ① lack of personalized services, pricing, targeted services to new market segments.
- ② underutilization of data gathered by loss adjusters.

③ target for better insights.

Conventional System:

System which consist of 1 or more size each having either manual operated or automatic detection devices or combination of both. Big data is a huge amount of data which is beyond the processing capacity of conventional data base system to manage & analyze the data in specific time interval.

Difference b/w Conventional & Intelligent Computing

- The Conventional Computing functions logically with the sets of rules & calculation while intelligent Computing can function via images & concepts.
- Conventional Computing is often unable to manage the variability of data obtained in real world. Whereas intelligent Computing is well suited to situation that have no clear algorithmic solution and are unable to manage noisy imprecise data. This allows them to excel in those area that Conventional Computing often finds difficult.

Diff B/w Big data & Conventional data

Big data	Conventional data
① huge data set	① Data set size is in control
② Big data \Rightarrow unstructured data such as video, text, audio & images.	② normally structured data which has fixed size & specific format.
③ Hard to perform queries & analyse	③ Relatively easy to perform query & analyse.
④ need tools such as hadoop, hive, hbase, pig, sqoop.	④ Tools such as sql & excel alone may be sufficient.
⑤ use for reporting basic analysis & text mining. Advanced analytic is only in starting	⑤ used for reporting & predictive modeling

③ target for better insights.

Conventional System:

System which consist of 1 or more size each having either manual operated or automatic detection devices or combination of both. Big data is a huge amount of data which is beyond the processing capacity of conventional data base system to manage & analyze the data in specific time interval.

Difference b/w Conventional & Intelligent Computing

- The Conventional Computing functions logically with the sets of rules & calculation while intelligent Computing can function via images & concepts.
- Conventional Computing is often unable to manage the variability of data obtained in real world. Whereas intelligent Computing is well suited to situation that have no clear algorithmic solution and are unable to manage noisy imprecise data. This allows them to excel in those area that Conventional Computing often finds difficult.

Diff B/w Big data & Conventional data

Big data	Conventional data
① huge data set	① Data set size is in control
② Big data \Rightarrow unstructured data such as video, text, audio & images.	② normally structured data which has fixed size & specific format.
③ Hard to perform queries & analyse	③ Relatively easy to perform query & analyse.
④ need tools such as hadoop, hive, hbase, pig, sqoop.	④ Tools such as sql & excel alone may be sufficient.
⑤ use for reporting basic analysis & text mining. Advanced analytic is only in starting	⑤ used for reporting & predictive modeling

Life cycle of data analytics unit ①

① ~~The~~ The data analytic lifecycle is designed for Big data problems & data science projects. The cycle is iterative to represent real project.

① Discovery

- The data science team learn & investigate the problem.
- Develop context & understanding
- Come to know about data sources needed & available for the project
- the team formulates initial hypotheses that can be later tested with data.

② Data preparation

- steps to explore, preprocess & condition data prior to modeling and analysis.
- It requires the presence of an analytic sandbox. The team execute, load, & transform to get data into the sandbox

- Data preparation tasks are likely to be performed multiple times if not in predefined order.
- Several tools commonly used for this phase are - Hadoop, Himpine, Hives, Open Refine etc..

(3) Model planning

- In this phase data science team develop data sets for training & production purposes.
- Teams build & executes models based on the work done in model planning phase

(4) Model Building:-

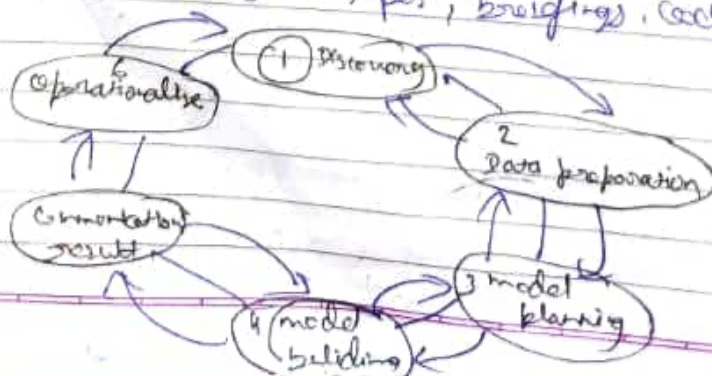
- Team develops datasets for testing, training & production for
- Team also considers whether its existing tools will suffice for running the models or if they need more robust environment for executing models.

(5) Communication Results:-

- After executing model team need to compare outcomes of modelling to criteria established for success & failure.
- Teams consider how best to articulate finding and outcomes to various team member & stakeholders, taking into account timing, assumptions.
- Teams should identify key finding, quantify business value and develop narrative to summarize & convey finding to stakeholders.

(6) Operationalize :-

- The team communicate benefits of project more broadly and sets up pilot project to deploy work in controlled way before broadcasting the work as full enterprise of user.
- This approach enables team to learn about performance and related constraints of the model in production environment on small scale & make adjustments before full deployment.
- The team deliver final report, briefings, codes.



Unit 1

Types of data analysis

Analytics is the discovery of meaningful patterns in data. Especially valuable in data rich with recorded info. analytics relies on the simultaneous application of statistics, computer programming & operations research to quality performance.

Types

① Predictive

③ Prescriptive

② Descriptive

④ Diagnostic analytics

① Predictive Analytics :-

Predictive Analytics turn the data into valuable actionable info. Predictive analytics use data to determine the probable outcome of an event or a likelihood of a situation occurring. Predictive Analytics holds a variety of statistical techniques from modeling machine learning & deep learning & game theory that analyse current & historical facts to make predictions about a future event. Techniques that are used for predictive analytics are

- Linear regression
- Time series analysis & forecasting
- Data Mining

Descriptive Analytics

Descriptive analytics looks at data & analyze past event for insight as to how to approach future events. It looks at the past performance & understands the performance by mining historical data to understand the cause of success or failure in the past. Almost all management reporting such as sales marketing, operations,

finance uses this type of analysis. The descriptive model quantifies relationships in data in a way that is often used to classify customers or projects into groups within a predictive model that focuses on predicting the behaviour of a single customer. Descriptive analytics identifies many different relationships between customer & product.

Prescriptive Analytics

Prescriptive Analytics automatically synthesizes big data, mathematical science, business rules and machine learning to make a prediction & then suggests a decision option to take advantage of the prediction. Prescriptive analytics goes beyond predicting future outcomes by also suggesting actions benefits from the predictions & showing the decision maker the implication of each

Page No. _____
Date _____

Decision option prescriptive analytics not only anticipate what will happen further prescriptive analytics can suggest decision options on how to take advantage of a future opportunity or mitigate a future risk & illustrate the implications of each decision option.

Eg - Prescriptive Analytics can benefit health care strategically planning by using analytics to leverage operational & usage data combined with data of external factors such as economic data, population demography, etc.

Diagnostic Analytics: In this analysis which generally, use historical data over other data to analyse any quo or for this resolution of any problem we try to find any dependency & pattern in the historical data of a particular problem.

Eg: Companies go for this analysis because it gives a greater insight into a problem & they also keep detailed information about their disposal. This way data collection may turn out individual for every problem & it will be more time consuming. Common technique used for diagnostic Analytics are:-

- 1) Data Discovery
- 2) Data Mining
- 3) Co-relations

5V's of Big data

Big data is a collection of data from many different sources & is often described by 5 Characteristics: volume, value, variety, velocity & veracity.

• **volume** :- the size & amounts of big data that companies manage & analyze.

• **value** :- the most imp. V from the perspective of the business & the value of big data usually comes from insight discovery & better recognition that lead to more effective operations. It improves customer relationship & offers clear & quantifiable business benefits.

• **variety** :- the diversity & range of different data types including structured data, semi-structured data & unstructured data.

• **velocity** :- the speed at which companies receive & manage data. Eg. The specific hrs of social media posts or search queries received within a day, hours or other unit of time.

• **Veracity** :- the truth or accuracy of data & info. assets which often determine executive-level confidence.

Data Analytics

Date: _____

Page No. _____

#1 Statics

Statics is a discipline dealing with collecting, analyzing and interpreting (usually) presenting that numerical data. It is used in every aspects of data science that means it is used to analyse, transformed & cleaned data.

- Population is the entire data set which is used for a particular test
- Sample is a portion of population which is used for testing.

#2 Types of Statics

① Descriptive

② Inferential

#3 Descriptive Statics

- It helps us to describe the data
- It enables us to underline the characteristics of data

- It does not predict anything & it does not make any assumption

Eg mean, mode and median.

Inferential Statistics

- It helps us to draw conclusions i.e.
- It helps to predict or make assumptions.
Eg Histogram, Bar Chart

Limitation of Descriptive Statistics.

- They are useful but they hide the important information about the data set. So we use Inferential Statistics.

Probability = $\frac{\text{no. of events}}{\text{total no. of outcomes}}$

Eg Rolling of dice

$$P = \frac{1}{6}$$

* Bias - tendency of a statistical or predictive model to over or underestimate of a parameter

- (1) Over sampling
- (2) Under sampling



Over Sampling

We use to apply ML techniques to model the data and make prediction when the sampled is biased then we perform 2 types of Sampling

- ① under sampling
- ② over sampling

① under sampling :

We select some data from the majority class as the same no. of the minority class.

Over Sampling

Multiply the minority class such that it has the count as the majority class.

What is sampling?

It is the active process of gathering observations with the intent of estimating a population.

Resampling : It is the methodology

of statistically using a data sample to improve the accuracy and quantify the uncertainty of a population parameter.

SE (Sampling error) = Calculated o/p - Observed o/p

⇒ we can perform Sampling by 2 methods

- ① Bootstrap method
- ② K-fold cross validation

Bootstrap

Samples are drawn from the data set with replacement allowing the same sample to appear more than once in the sample.

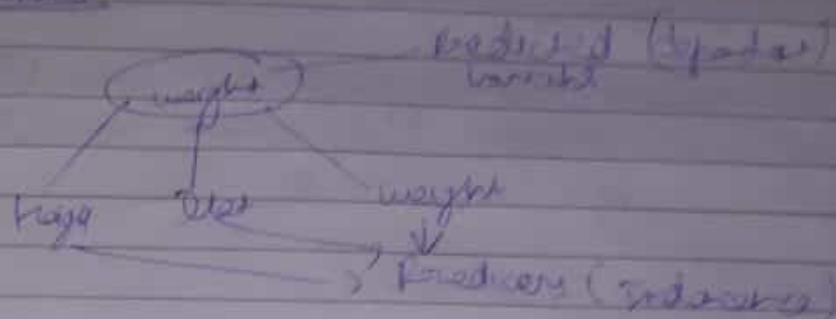
K-fold cross validation

In K-fold cross validation data set is partitioned into K-groups where each group is given the opportunity of being used as a sample.

Applications of Statistics

- ① Business Analytics
- ② Artificial Intelligence
- ③ Financial Analysis
- ④ Fraud detection
- ⑤ Machine Learning
- ⑥ Share Market
- ⑦ Pharmaceutical Sector

Regression



Regression

It is a method to mathematically formulate relationship b/w variables to search to what extent the predicted variable is affected by the predictors.

Linear Regression

$$y = b_0 + b_1x$$

Non linear regression $\rightarrow y = b_1^2x_1 + b_2^2x_2$

Logistic Regression \rightarrow It is binomial 2 values
Eg: \rightarrow happy / sad
True / false

Time series Regression \rightarrow Based on Time

Regression Type

① Linear Regression \rightarrow

It assumes that there is a linear relationship b/w predictors (x) & predicted variable (y)
 $y = b_1x + c$

multilinear regression

It assumes that a predicted variable, depends on more than 1 independent variable.
eg: $h_0, b_1, b_2, \dots, b_n$

non linear regression

It assumes that relationship b/w x & y is having a polynomial function.

logistic regression

It is useful when a target variable is binomial (1 or 0) (accept or reject).
It estimates probability at certain events.
eg: Happy or Sad, Rain or Not.
that means it has value 0 or 1.

Time Series regression

It is used to forecast behaviour of a variable based on historical data.
eg: unemployment rate

$$Y_t = x_t \beta + \epsilon_t$$

Y_t : It has value of y

Q) Calculate prediction error for linear regression model

Predicted value (\hat{y})	Actual value (y)	$(\hat{y} - y)^2$
14	12	4
15	15	0
18	20	4
19	16	9
25	20	25
18	19	1
12	16	16
12	20	64
15	16	1
22	16	36
		160

$$RMSE = \sqrt{\frac{1}{n} \sum (\hat{y} - y)^2} = \sqrt{\frac{160}{10}} = 4\sqrt{10} \approx 12.56$$

for logistic regression

Misclassification rate = $\frac{\text{Incorrect prediction}}{\text{Total prediction}}$

Predicted	Actual	Incorrect
1	0	✓
1	1	
1	1	
0	1	✓
0	0	
1	0	✓
0	0	

$$\therefore \text{Misclassification rate} = \frac{3}{7} \times 100 = 42.8\%$$

Multi Variate Analysis

It involves multiple dependent variables existing in one outcome.

uses ->

It is widely used in many industries like healthcare. Recently in (COVID) a team have more than 5 lakh COVID patients by end of July 2020. This analysis was made on multiple variables like

- Govt decision
- Public behaviour
- Occupation
- Public transport
- Healthcare services
- Overall immunity of community

Advantages

- Therefore it considers more than 1 factor of independent variables. Hence the conclusion drawn is more accurate.
- Conclusion are more realistic.

Disadvantages

It requires complex computation to arrive at a satisfactory conclusion. Very observation for a large no. of variables leads back to the problem of multicollinearity.

Bayesian model \Rightarrow It is a mathematical approach that involves the application of probability to solve statistical problems.

Prior \Rightarrow It refers to the preconceived belief we have

Posterior \Rightarrow the probability of occurrence of event on basis of prior & evidence

Bayes Rule \Rightarrow

$$P\left(\frac{A}{B}\right) = \frac{P(A \cap B)}{P(B)} \quad \text{--- (1)} \quad P(B) \neq 0$$

Conditional

Probability $P\left(\frac{B}{A}\right) = \frac{P(A \cap B)}{P(A)} \quad \text{--- (2)}$

Using (1) & (2)

$$P(A/B) = \frac{P(B/A) (P(A))}{P(B)} \quad \text{--- Bayes}$$

likelihood \rightarrow evidence

Posterior

likelihood \Rightarrow

Probability of event B being true given that event A is true.

It refers to the probability of observing what we did give that our prior are true

Ex \Rightarrow A situation where Bayesian analysis is not used is the spam filter in your mail server. Message is scrutinized for the difference of keywords which makes it likely the message is spam.

Q. of two persons in Spain
 C → Subjective line (means, the center)
 Check this out

Compute conditional probability $P(S|C)$ where
 Add. of events line span 1% of Spain results
 leave Check this out. For Subjective line, take
 0.4% of the Subjective non-span results leave
 this scenario in the subjective line

Sol: $P(S) = \frac{40}{100} = 0.4$

$P(C|S) = \frac{1}{100} = 0.001$

$P_C = \frac{0.6 \times 0.4}{100} = 0.0024$

$P(S|C) = \frac{P(C|S) \cdot P(S)}{P(C)}$

$P(S|C) = \frac{0.001 \times 0.4}{0.0024} = \frac{0.001 \times 0.4}{0.0024}$

$\frac{40 \cdot 10}{24 \cdot 10} = \frac{10}{6} = 1.67$

Q. Find out Conditional probability where
 2 → express as
 P =

$P(P) = \frac{1}{2}$

$P(Q) = \frac{1}{4}$

$$P(A) = \frac{1}{2}$$

$$P(A/B) = \frac{P(B/A) P(A)}{P(B)} = \frac{\frac{1}{2} \times \frac{1}{2}}{\frac{1}{4}} = \frac{\frac{1}{4}}{\frac{1}{4}} = 1$$

Major advantage of Bayesian statistics

It takes the prior knowledge into an consideration while calculating probability by applying Bayes's rule.