

Data Science

A field that gives insights from structured and unstructured data, using different scientific methods and algorithms. Consequently help in generating insights, making Predictions, devising data driven Solutions.

Need

We need Data Science to convert abundant Complex data into insights enabling informed decisions, innovations and solutions to challenges across various fields like healthcare, finance and environmental conservation.

Life Cycle of Data Science:

- The things go into the formation of data science.

Business Understanding \Rightarrow Data Collection \Rightarrow Data Preparation \Rightarrow

Model deployment \leftarrow Model building \leftarrow Exploratory Data Analysis and Maintenance

- Business understanding - first step - defining the business problem because a well defined business statement define a specific goal and is the key to success of the project.

- Data Collection - Second step - Data acquisition - Essential step in data science because data needs to be relevant to solve the business problem correctly. ~~and preparation~~

Data Preparation - Third step - A crucial step in DS Project as it helps clean and bring the data into shape which is required for further analysis and modeling. Also referred as data cleaning.

We treat issues like missing values, outliers and also transform the data into required format.

Exploratory data Analysis - fourth step - Data is Analysed using Summary Statistics and graphically to understand key Patterns. This is simplest step but highly effective in extracting some useful patterns.

Model Building - fifth step - Two types -

i) Descriptive Analytics

- Provides insights based on historical data

ii) Predictive modeling

- involves future Predictions.

Model deployment and Maintenance - model is ready to deploy in real world. The deployment can occur offline or the web, on the cloud, any Android or iOS app.

Database

Data

Machine learning

Big data

Scientist

Programming language

Statistics

Python

- Pandas - Data Analysis and Manipulation
 - ↳ Provide fast, flexible data structure, Numerous functions and methods
- Numpy - Scientific Computations
 - ↳ Contains multidimensional array and matrix data structure
- Seaborn / matplotlib - Data visualisation
- Power BI / Tableau - Visualization tools

Machine Learning

- ML Algorithms
 - ↳ Supervised - When we worked with labelled data means output is known
 - ↳ Unsupervised - Algorithm has to make guess means output is not known
- Two types of Supervised learning
 - ↳ Classification
 - ↳ Regression
- HyperParameter Tuning - Model optimisation
 - They Control the behaviour of Algorithm.
 - Done to improve results

- Sckit Learn - Machine learning library. It supports both supervised and unsupervised Learning.
- Provides various tools for model fitting, data Preprocessing, Model selection, Model evaluation

In R - Mir Package / Caret used

for Big data - Traditional Machine learning Algorithms are not used. we use deep learning.

Deep learning:

- Neural networks: working with different kind of layers and neural networks then deep learning is Preferred
 - There are set of Algorithms which are modeled you know loosely after human brain to recognize patterns.
- Algorithms
 - CNN → Convolution Neural network
 - RNN → Recurrent Neural network
 - LSTMs → long short term memory network
- Natural language Processing - A Subfield of AI concerned with interaction b/w Machine & humans
- Libraries - → TensorFlow
 └ PyTorch

Data Science

A field that gives insights from structured and unstructured data, using different Scientific Methods and Algorithms. Consequently help in generating insights, making Predictions, devising data driven Solutions.

Need

we need data Science to convert abundant Complex data into insights enabling informed decisions, innovations and Solutions to challenges across various fields like healthCare, finance and environmental Conservation.

Life Cycle of Data Science:

- The things go into the formation of data Science.

Business Understanding \Rightarrow Data Collection \Rightarrow Data Preparation \Rightarrow

Model deployment \leftarrow Model building \leftarrow Exploratory data Analysis and Maintenance.

- Business understanding - first step - defining the business Problem because a well defined business statement define a specific goal and is the key to success of the project.

- Data Collection - Second Step - Data acquisition - Essential step in data Science because data needs to be relevant to solve the business problem correctly. ~~and also preparation~~

- Data Preparation - Third Step - A crucial step in DS Project as it helps clean and bring the data into shape which is required for further Analysis and Modeling. Also referred as data cleaning.

We treat issues like missing values, outliers and also transform the data into required format.

Exploratory data Analysis - Fourth Step - Data is Analysed using Summary Statistics and graphically to understand key patterns. This simplest step but highly effective in catching some useful patterns.

Model Building - Fifth Step - Two types -

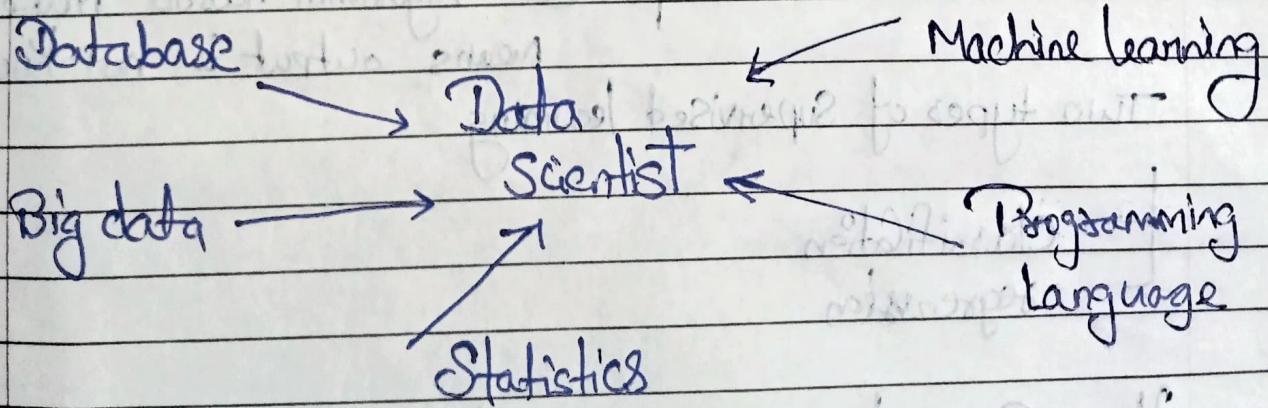
i) Descriptive Analytics

- Involves insights based on historical data.

ii) Predictive Modeling

- Involves future Predictions.

Model deployment and Maintenance - model is ready to deploy in real world. The deployment can occur offline on the web, on the cloud, any Android or iOS app.



Statistics

— Identify \Rightarrow Predict \Rightarrow Design \Rightarrow Estimate

Descriptive Statistics

\hookrightarrow It is a Summary Statistics that quantitatively summarizes features

- Variance
- Mean
- Mode
- Median
- Skewness

Inferential Statistics

\hookrightarrow Inferring something from small data assuming that this will also apply to large data.

- Hypothesis
- Probability
- Regression Analysis
- Chi-square
- Anova

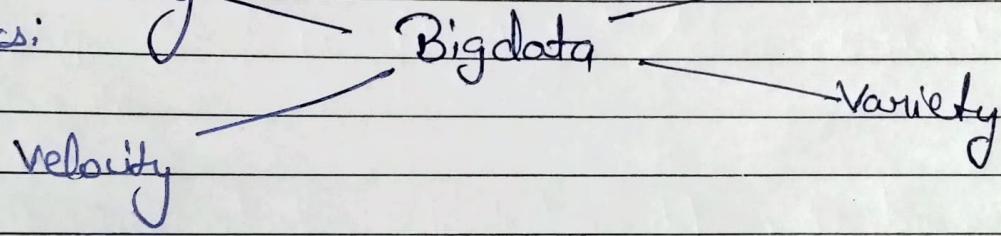
Databases

- Non-Relational DBs
 - follows no pre-set Schema
 - MongoDB, Apache Giraph
- Relational DBS
 - Data is related, in tabular format, easy to query
 - MySQL, Oracle

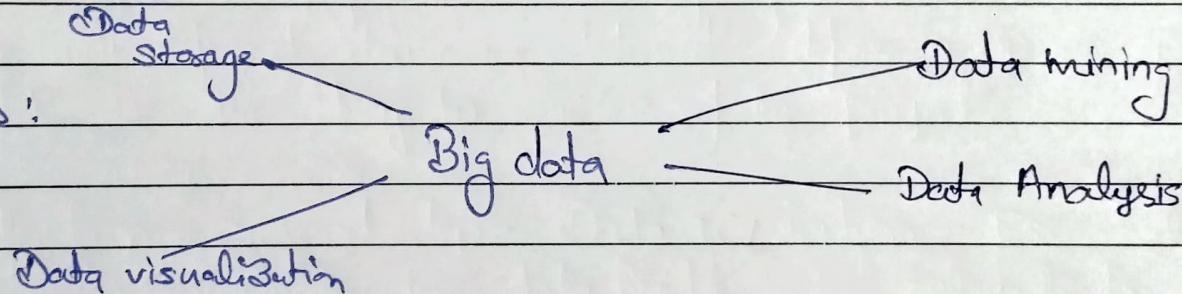
Big data

- Large and Complex Ochesta and Companies are dealing with huge volume of data especially from new Data Sources.

Characteristics:



Technologies:



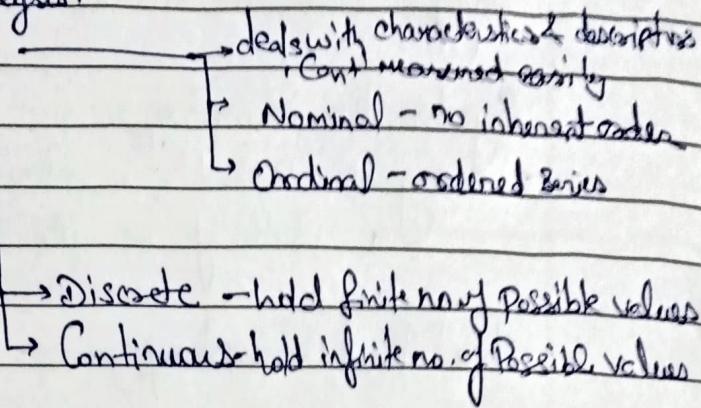
for data storage - hadoop - process data in distributed environment. it can store + Analyze data, high speed, low cost

→ Excel, Web Scapping, Linux, Git, Cloud

Data

Data refers to facts and statistics Collected together for reference and Analysis.

- ↳ Qualitative data
- ↳ Quantitative data



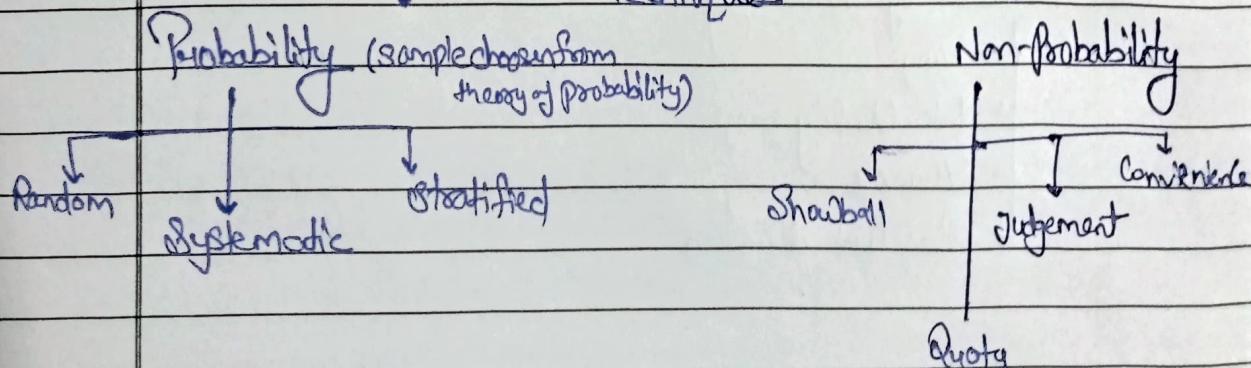
Statistics for data Science

Statistics is an area of Applied mathematics concerned with the data Collection, Analysis, interpretation, Presentation,

Terminologies:

- Population: A Collection or Set of individuals or objects or events whose properties are to be analyzed.
- Sample: A Subset of population is called 'Sample'. A well Chosen Sample Contain most of information about a Particular population Parameter.

Sampling Techniques



Random Sampling

- Each member of Population has equal chance of being Selected in the Sample.

Systematic Sampling

- Every n^{th} record is chosen from Population to be a Part of Sample.

Stratified Sampling

A Stratum is a Subset of Population that shares at least one Common characteristics.

Random Sampling is used to Select a Sufficient no. of Subjects from each stratum.

Types of Statistics

Descriptive Statistics

- uses the data to provide description of the Population either through Numerical GIC, graphs.

→ Measure of Central tendency

Measure of Variability

→ range

→ Inter-Quartile range

→ Variance

→ Standard deviation

Inferential Statistics

- Makes inferences and Predictions about a Population based on Sample of data taken from Population.

Mean :

Measure of Average of all the values in a sample is called Mean -

Median :

Center of value, arranged in Asc order -

Mode :

The most recurrent in Sample set.

Range :

- The given measure of how spread apart the values in dataset are -

$$\text{Range} = \text{Max}(x_i) - \text{Min}(x_i)$$

Quartiles :

- Tell us about the spread of dataset by breaking dataset into quarters.

Variance :

- how much a random variable differs from its expected value.

$$S^2 = \frac{\sum_{i=1}^{n-1} (x_i - \bar{x})^2}{n}$$

Standard deviation :

- Difference b/w each element from the mean

$$\text{Deviation} = (x_i - \bar{x})$$

Example.

Data:- 9, 2, 5, 4, 12, 7, 8, 11, 9, 3, 7, 4, 12, 5, 4, 10, 9, 6, 9, 4

- Mean:- $\frac{9+2+5+4+12+7+8+11+9+3+7+4+12+5+4+10+9+6+9+4}{20}$

$$\boxed{M = 7}$$

$$(x_i - M)^2 \Rightarrow (9-7)^2, (2-7)^2, (5-7)^2, \dots \text{ So on}$$

We get,
 $\sigma_m = \sqrt{9, 25, 4, 9, 25, 0, 1, 16, 4, 16, 0, 9, 25, 4, 9, 9, 4, 1, 4, 9}$

Standard deviation,

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - M)^2}$$

$$= \sqrt{\frac{1}{20} (80)}$$

$$\sigma^2 = 8.0$$

$$\boxed{\sigma = 2.983}$$

Bayes's theorem:

- Shows the relation between one Conditional probability and its inverse. \rightarrow Likelihood ratio

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Prior

Posterior

- Conditional Probability, $P(A|x) = \frac{P_r(A \cap x)}{P_r(x)}$

- Interval Estimate:

An interval or range of values used to estimate a population parameter is called interval Estimate.

- Margin of Errors:

The greatest possible distance between the Point Estimate and value of Parameter it is estimating.

$$E = Z_c \frac{\sigma}{\sqrt{n}}$$

Hypothesis testing

- used to formally check whether hypothesis is accepted or rejected.

State hypotheses: Formulating a null or alternative hypotheses formulate Analysis plan! Construction of Analysis plan.

Analyses of Sample data: Calculation and interpretation of Test Statistics as described in Analysis plan.

Interpret result: - Application of decision rule described in Analysis plan.