

Introduction to data science:

↳ Data Science:- is a combination of multiple disciplines like statistics, Data Analytics and Machine learning.
It is basically a data gathering, analyzing the data finding or discovering the pattern in data and predicting the future outcomes by making better decisions.

↳ Data science need:-

Data science is used in many industries like - banking
- consultancy.
- healthcare
- consumer goods.
- E-commerce
- Politics
- logistics company

In consulting - the people with prediction related to the terms.
Analyzing and which health exercises and treatment are best

Prediction of future profits and outcomes of company.

Promotion of products.
Prediction in election

Analyzing routes / predicting delays in transport.

Basically Data science is a process of using algorithm, methods and systems to extract knowledge & insight from structured & non-structured data. It uses Analytics, ML to help users make decision, prediction, enhance optimization and improve opⁿ.

Evolution of Data Science:-

Stage-I:- In 1962 J.W Turkey published article on "The future of Data Analysis". It was a relationship between statistics and Data analysis.
At that time business transaction were centralized.

Stage-2:- In 1977, J.W Turkey published another article related to hypothesis for testing & analyzing the data.

Term Data mining comes into the picture of world.

Stage-3:- Businessmen started using data analysis and using in prediction of their profit and future decision. They started getting efficiency because of data analysis.

Stage-4:- Data science started begin recognised by public.
researchers started research on D.S. Journal started commenting

on D.S. Big data comes into picture

(11)

Stage-5

The term Data Science was introduced to world. Acc. to IBM around 90% of total data available to world was produced in last two yrs.

Stage-6

D.S. started rising in various society or working culture all over the world. Everyone's firm started collecting vast data.

- Apple gave credit to Big data and data mining.
- MS, Goog, started using DS for voice recognition & speech detection.

Roles of Data Science :-

→ Data Scientist Role :

- collect & Analyze data
- Extract insight from patterns using ML model.
- explain & visualize result

→ Data Analyst

- Domain Exploration & understanding
- collect & Analyze data
- Understand data
- explain & visualize result / optimization
- Communicate with stakeholders.

→ Data Engineer

- transform raw data into usable pipeline
- Upgrade the existing version of technology with upgraded or newer version.
- Test and Deploy ML models.

→ Database Administrator

- Database backups & recoveries
- responsible for proper functioning of all database.

→ Data Architect

- creates blueprint for data management
- Also ensures the resources for engineers

→ Data & Analytics Manager

- oversees the data science operation
- Assign duties to their team.

- Business Analyst:

- Business oriented work
- identify how Big Data can be linked to actionable business insights for business growth.

- Statistician

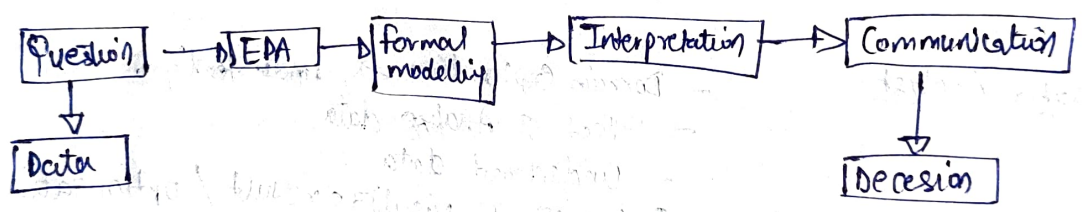
- understanding of statistical theory
- data organization
- create new methodology for engineers

- Machine Learning Engineer

- implement common ML algo.
- develop pipelines
- A/B testing.

→ Structure of Data Science Project

Q D E F I C D



5 phases

1 - Questioning phase:

- most important phase
- helps you understand your data & decide on the type of analysis.
- To extract your data from bigger dataset, one can use distributed storage like Apache Hadoop, Spark etc.

2 - EDA (Exploratory Data Analysis)

- check if the data you have is suitable to answer your question
- develop the sketch of solution
- check if your dataset carries all the data that is required.

Formal modelling:

- write down the parameter you are trying to estimate. (R)
- challenge your results through variety of approaches like sensitivity analysis.
- make sure data & algo are reproducible.

Interpretation:

- phase is to assemble all the information you've got after analysis.
- helps to filter out the result you're got.

Communication

- Project communicated to some sort of audience
- Make sure the result of your project are visualized for quick understanding.

Characteristics of data 5V's

Volume - refers to the vast amount of data generated every second.

Variety - refers to the different types of data we can use now.

Velocity - refers to speed at which new data is generated and the speed at which data moves around.

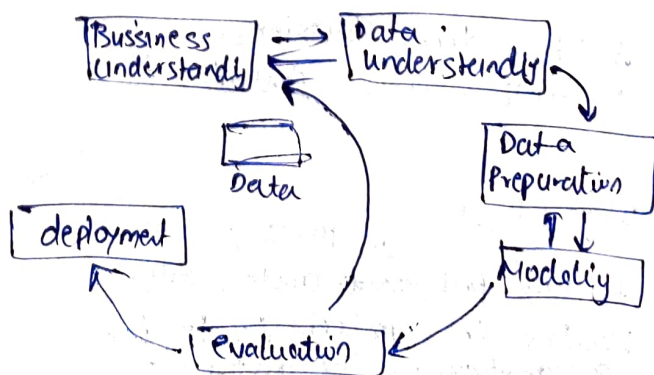
Veracity - refers to the messiness or trustworthiness of the data.

Value - refers to having access to big data is no good unless we can turn it into value.

→ Existing Data Analytic Methodologies.

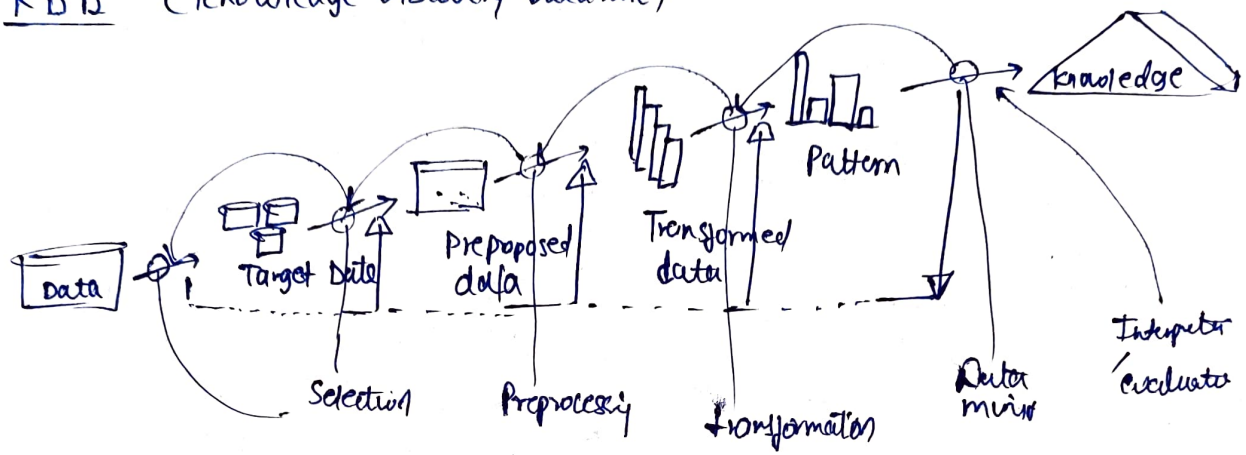
CRISP - DM

Cross-industry standard process for Data mining.



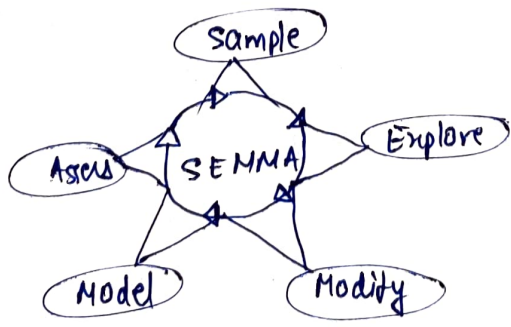
- Business understanding - What does business need?
- Data understanding - What data do we have / need?
- Data preparation - How do we organize data for modeling?
- Modeling - What modeling techniques should we apply?
- Evaluation - Which model best meet business model?
- Deployment - How do stakeholders access the result.

→ KDD (Knowledge Discovery Database)



- selection - targeted data is determined through database of compiled data.
- Pre-processing - improving the data
- transformation - converting pre-processed data to fully utilized kind.
- Data mining - focuses on shifting through transformed data to seek out pattern.
- Interpretation / Evaluation - cleaned, converted, relevant data framed into visual representation.

→ SEMMA
Sample
Explore
Modify
Model
Assess



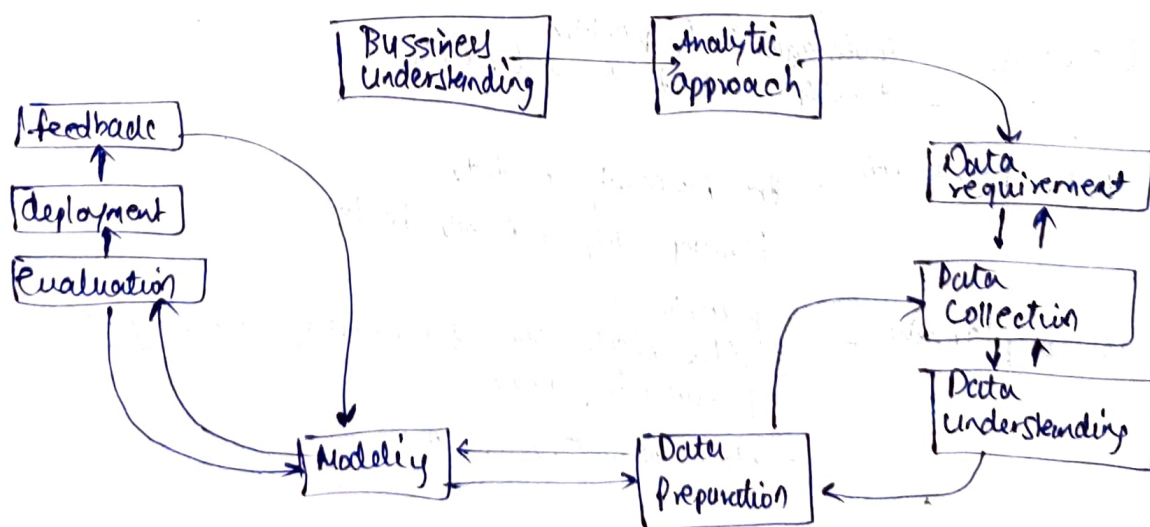
- Sample :- Generate a representative sample of data
- Explore :- visualization & basic description of data.
- Modify :- select variables, transform variable representation.
- Model :- use variety of statistics and ML model.
- Assess :- evaluate accuracy of model.

Applⁿ of Data Science:

9

- In Search engine :-
 - Data science is used to get Searches faster.
 - for ex. if we search for something on google, google will show the result of that website on the top which is visited most.
- In transport :-
 - Driverless Cars
 - training data is fed into the algorithm and with the help of DS techn., the data is analyzed like speed limit in highway n all.
- In finance :-
 - stock market
 - D.S is used to examine past behavior and to examine the future outcome.
- In E-commerce :-
 - on searching for something on E-commerce website we get suggestion similar to choices acc. to our past data and also we get recommendation.
- In healthcare :-
 - detecting tumor
 - Drug discovery
 - Medical bots.
- In Gaming :-
 - used with ML with the help of past data computer player will improve its performance.
- In delivery logistics :-
 - helps to find best route for shipment, best time suited, best mode of transport.

→ Data Science Methodologies:



Business understanding :- problem that you're trying to solve

Analytic Approach :- selecting the right analytical approach

data requirement :- identify the necessary data content.

data collection :- requirements are revised and statistical and visualization are applied.

data understanding :- activity related to construction of data

data preparation :- process of using domain knowledge of data to make ML model algo work

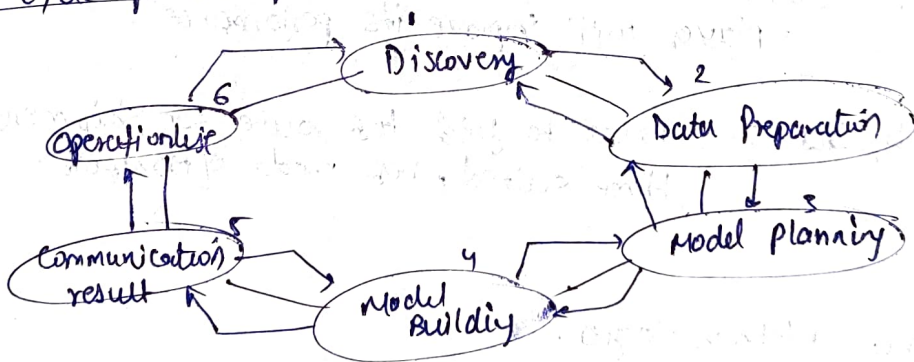
modeling :- focuses on developing models that are either descriptive or predictive.

evaluation :- to check whether the model used really answer the question of client

Deployment :- deployed to ultimate test

feedback - users will help refine the model.

→ Life Cycle phases of Data Analytics :-



Phase

1 - Discovery - investigation of problem
- develop context & understanding

2 - Data Preparation - pre-processing of data
- Hadoop, Alpine Miner tools.

3 - Model planning - planning
develop data set for training, testing & production purpose
- MATLAB, STATICA

4. Model Building - development of dataset
- Matlab, STASTICA, Octave.

5. communication result : - identify key history, summarize and convey to stakeholder

6. operationalize - deployment on small scale
- MADlib, wsc4