

# Difference betn data science & big data

- |  |  |
|--|--|
| <ul style="list-style-type: none"><li>① It is field of scientific analysis of data in order to solve the complex problem &amp; necessary cleaning &amp; preparing the data.</li><li>② It is used in bio-tech, energy, gaming &amp; insurance</li><li>③ Goals:- data classification, anomaly, prediction, scoring &amp; ranking</li></ul> | <ul style="list-style-type: none"><li>① Big data is storing &amp; processing large volume of structured &amp; unstructured data that can not be possible with traditional application</li><li>② Used in retail, education, healthcare &amp; social-media</li><li>③ Goal:- To provide better customer service, revenue identification &amp; marketing</li></ul> |
|--|--|

Applications of data science:-

- ① healthcare
- ② Gaming
- ③ Image recognition
- ④ Logistics
- ⑤ Predict future market trends
- ⑥ Recommendation sys
- ⑦ Streamline manufacturing

ngilpms  
hpqinls

data - explosion:-

- Business model transformation
- Globalization
- personalization of services
- New sources of data

BGNA

## \* 5 v's of Big data :-

- ① Volume
- ② Velocity
- ③ Variety
- ④ Value
- ⑤ Veracity

Volume:- Machine data, App logs, clickstreams log, Emails, Contracts, geographical info sys & geo-spatial data

Velocity:- Speed of creation data

Variety:- nature of data structured & unstructured data.

- ① Sensor data
- ② Mob networks
- ③ Amazon, facebook, yahoo
- ④ Social Media



## # 5 V's of Data :-

- ① Volume :-
  - i) Larger than conventional relational database
  - ii) consist of terabyte or petabyte of data.
- ② Velocity :-
  - i) refers to speed of generation of data.
  - ii) How fast data is generated & processed to meet the demands.
  - iii) Determines real potential in data.
- ③ Variety :-
  - i) Heterogeneous source
  - ii) Nature of data - structure & unstructured
- ④ Value :-
  - i) Business value derived from data.
  - ii) Objective :- Generate some sort of value for the company doing an analysis.
- ⑤ Veracity :- True Data relates to assurance of data quality, integrity, credibility & accuracy.



## Business Intelligent

- ① Business Intelligence tends to provide reports, dashboard, & queries on business question for the current period or in the past
- ② BI take it easy to answer the question like quarter-to-date rev, progress towards quarterly targets, & how product performed
- ③ Understand history
- ④ power BI

CC

- ① it provide resources on demand
- ② It refers to internet services from SaaS, PaaS to IaaS
- ③ Cloud is used to store data & info on remote servers
- ④ low maintenance
- ⑤ on demand soln
- ⑥ eg:- aep & Salesforce

## Data Science

- ① ds tends to use disaggregated data & looking forward to extract info in such a way to get further predictions
- ② data science tends to be more exploratory in nature & may use scenario optimization to deal with more open-ended question.
- ③ definition.
- ④ application - HP4IRLS

## Big data

- ① handle huge volume of data
- ② It can be stored, used & retrieved
- ③ It deals huge volume of data & info
- ④ big data is highly scalable, robust & cost-effective
- ⑤ eg:- cloudera, apache.



# data science life cycle - BDE DVDC PDHFE

BDE DVDC PDHFE		Parameter	Structured data	Unstructured data
Business understanding	Representation		Stored in row & column format	Unstructured data doesn't follow specified format
data Exploration	Metadata		Syntax	Schemas
Data Visualization	Storage		database Management system	Unstructured file structure
Data cleaning	Standard		SQL, ADO.net, ODBC	open XML, STIX, STS
Predictive Modeling	Tool for integration		ETL	Manual data integration
data Mining	Used by ORG		Low vol op <sup>r</sup>	High vol op <sup>r</sup>
Perks etc.				



## Information

- ① It is processed data
- ② Info is specific
- ③ Info depends on data
- ④ Info is output of computer
- ⑤ refers to knowledge of interpretation of data

## data

- ① Raw data
- ② data is not specific
- ③ doesn't depend on info
- ④ data is input to computer
- ⑤ data refers to facts, measurements etc.

## Qualitative data

data which cannot be measured

Type:- Nominal data & Ordinal data

Makes use of appearance, color, texture & other qualities

These are descriptive form

ex: Jean is bad

## Quantitative data

Can be measured

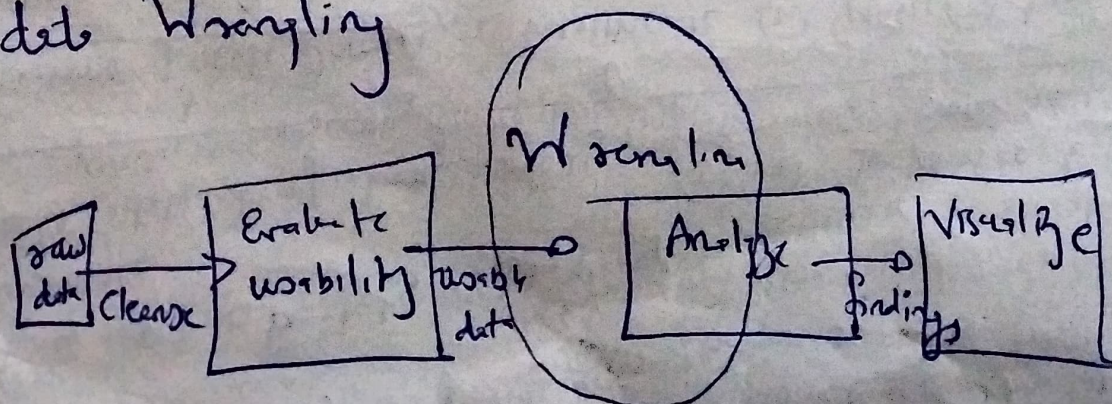
Type:- Interval data & Ratio data

Measures quantities such as length, size, amount, price & even distance

These are numerical form

ex: there is 7 people

## data Wrangling





## # Data Wrangling :-

### ① Definition :-

Data wrangling is the process of cleaning, organizing and transforming raw data into desired format for better decision making in less time.

### ② 5 STEP Data Wrangling Process :-

1] Discovering

2] Structuring

3] Cleaning

4] Enriching

5] Validating

6] Publishing

### ③ Benefits :-

- 1] Enables <sup>business</sup> users to make timely & accurate decisions.
- 2] Processes large ~~amount~~ volume of data easily.
- 3] Merges several datasets into one for analysis.