

## \* Big Data &amp; sources -

large volumes of data available at various sources in varying degree of complexity, generated at different speed i.e. velocities & technologies, processing modules.

- term → describe collection of data i.e. huge in size & yet growing exponentially with time.
- data → large & complex  
→ no traditional method/tools work or to process/store it.
- consists of extensive datasets that require a scalable architecture for efficient storage, manipulation & analysis.

Sources of Big Data -

- 1) Social Media - Fb → 500+ TB of data everyday  
→ status msg, photos & videos.
- 2) Stock exchange → TB/sec → trade data of users & companies
- 3) aviation industry → Jet engine → 10 TB of data → 30 min flight
- 4) Survey Data → online or offline survey → 100 & 1000's response  
→ processed for analysis & visualization
- 5) Compliance Data → Healthcare, hospitals, life sciences & finance, have file compliance reports, etc.

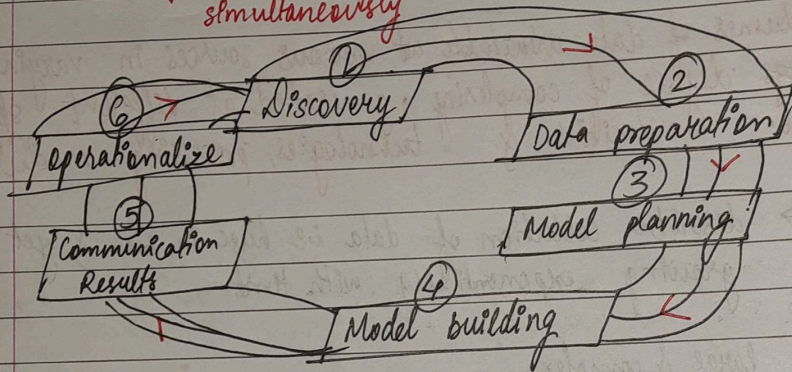
appl<sup>n</sup> of Big data →

- govt
- medical
- energy
- finance
- manufacturers & transportation
- media & entertainment
- weather
- e-commerce.



## \* Data Analytics Lifecycle -

→ data problems & data science projects.  
→ six phases the project work can occur in several phases simultaneously.



- 1) Discovery
  - defining data's process.
  - critical objective of business → mapping out data.
  - term learns → business domain & previous similar project.
  - team evaluates - technology, people, data & time.
- 2) Data prep
  - business req → data req
  - collecting, processing & clean data
  - analytics → workspace.
  - sandbox → collects all kind of data.
  - ETLT → extract, transform, load & transform.
  - data → sandbox.
  - data acquisition → existing data from outside.
  - data entry → creating new data values from data.
  - data conditioning → cleaning data performing transformation.
  - common tools for data prep.

Tools → • Hadoop • Open Refine  
• Alphere Miner



## 3) Model Planning -

- assess structure of data.
- model selection
- common tools → SQL/SAS/ACCESS.

## 4) Model Building -

- design model
- execute model
- SPSS modular (IBM) → enterprise level computing.
- STATISTICA & MATHEMATICA → data mining & analytics tools.

## 5) Communicate Results -

- project report results → success or failure
- vital findings → analysis
- communicate/document → key findings & major insights.

## 6) Operationalize -

- sandbox → live environment.
- team → communication → benefit of project
- risk management → efficiently.
- pilot project → execute algorithms more efficiently.
- test model → produc<sup>n</sup> env<sup>i</sup> → monitor accuracy is  
retrain model if  
necessary.



## \* 5V's of Big Data -

- **Volume** → • Big Data characterised → enormous volume → petabytes  
• large data needs specialized tools & technologies to store process & analyze.  
• observes & tracks data from various sources.
- **Velocity** → • data streams → high speed & dealt timing  
• processing of data → streamed data → real time results → fast  
• insights generated → relevant & actionable.  
• speed of generation of data.
- **Variety** → • heterogeneous sources  
• nature of data → structured & unstructured.  
• social media data → sensors, audio, video, etc.  
• data → insights → manage & analyze.
- **Value** → • business value → Big Data.  
• generate some sort of value → company doing analysis  
• insights → formal decisions  
→ optimize operations  
→ gain competitive edge in marketplace.
- **Veracity** → • inconsistency / uncertainty in data  
led to • challenges in data quality & reliability.  
• insights → accurate & trustworthy.



663 2484523

MRP ₹ 150/-

Inclusive of all taxes

Official Use



83

## Exploratory Data Analysis (EDA) -

- crucial step in Data Analysis → understand data before modelling
- visually & statistically exploring data → gain insights  
 identify patterns  
 detect anomalies  
 formulate hypotheses
- analysis to develop a intuition about dataset & informed decisions

EDA's → 5 M framework

### Mindset

- curious & open about dataset
- allow data to reveal characteristics & patterns

### Methods

- to explore & summarize data
- summary statistics
- Data Visualization
- Dimensional Redu<sup>n</sup>
- Data Profiling
- Correlation Analysis

### Missing Data

- missingness
- pattern (underlying)
- strategies to handle

### Multivariate Analysis

- Rel<sup>n</sup> b/w multiple variables
- diff variables interaction
- identify clusters

### Models

- potential predictors
- variable importance
- formulate hypothesis
- ML models

- uncover pattern & anomalies
- generate hypothesis
- guide analysis
- modelling decisions