

Big Data:-

Collection of datasets, which are large & complex that it becomes difficult to process using on-hand database management tools. are traditional data processing applications.

Data Science:-

is a field that uses scientific methods, algorithms, & tools to extract knowledge & insights from structure/unstructured data. involves techniques from statistics, ML, CS, to analyze data

Business Intelligence.

turning data into useful info to help business make better decisions.

Data Deluge :- overwhelming flow of large amount of data being generated constantly from various sources.

Reasons for Data Deluge.

- ① widespread use of internet & smart devices.
- ② Social media & streaming platform
- ③ Sensors & IoT Devices
- ④ Digital transactions of Businesses.
- ⑤ cloud computing & Big data technologies.
- ⑥ Healthcare products & security monitoring

Eg :- fitness tracker watch, every day x millions of users = terabytes of data. every day.

Data Analytics lifecycle

- | | |
|--------------------|-----------------------|
| ① Discovery | ④ Model Building |
| ② Data preparation | ⑤ Communicate results |
| ③ Model planning | ⑥ Operationalize |

- Iterative (to cover new informations on data)

① Discovery :-

- first phase, goal is to
- understand business problem, identify project goals & gather all necessary background info.

① Learning the Business Domain :-

understand company goals & business problems to solve

② Resources :-

identify tools, tech, team, data infrastructure for project

③ Framing the problem :-

convert the business challenges to clearly defined analytical problem

④ Identifying key stakeholders :

who's involved & who is benefited.

⑤ Interviewing the Analytics Sponsor.

discuss with clients their objectives, expectations.

⑥ Developing initial hypotheses.

assumptions about what might be happening in data.

⑦ Identify potential data sources.

external/internal sources from which data can be collected.

Activities involved. identify potential resources.

- ① understanding data needs.
- ② Identify internal data sources
- ③ Identify external data sources
- ④ check data availability & access.
- ⑤ Assess data Quality. (accurate, complete & up to date)
- ⑥ understand data formats & storage
- ⑦ Document data source. (record)

- right data found, verified & prepared.

② Data Preparation

- includes ~~pure~~ steps to explore, preprocess & condition data
- most iterative phase.

Steps in data preparation

- ①. Data collection. - gather data internal/external source
- ②. Data integration:- combine data from multiple sources
- ③. Data cleaning:- fixing issues like, missing val, duplicates
- ④. Data transformation:- Convert data to suitable form i.e. normalize, encoding, aggregation.
- ⑤. Data reduction:- Reduce data size by removing unnecessary ^{feature}
- ⑥. Data formatting:- data in standard format.
- ⑦. Data sampling:- selecting a small subset for working on ^{data}
- ⑧. Feature Engineering:- new feature from existing to increase performance

Preparing analytic Sandbox.

- ⊕ It is a separate environment or workspace for data analysis & experimentation.

purpose of sandbox.

- ①. Safe Exploration:- Team can perform testing, analysis & model building without affecting live production systems.
- ②. Expansive Data collection.
Sandbox gathers diverse & large dataset from multiple source for deeper analysis.
- ③. Supports Advanced Analytics.
complex tasks like AI, ML not just BI
- ④. Encourages Innovation:- provides risk free zone for testing new ideas.

- (5). Separation from production Data
- (6). Gaining acceptance.

Performing ETLT (Extract, Transform, load, Transform)

- Data Extracted from source.
- transformed (cleaned / structured) & then
- loaded to targeted systems.

ETLT in Analytics Sandbox

In sandbox env. Data is ~~ext~~ E & LT first & then transformed.

- early load to prevent the original form of Data.
- use case - Fraud Detection in credit card transaction.
- outliers can signal fraud.

- Hadoop is used. Since it can process large amount of data.

Data Conditioning

- includes cleaning data, normalizing dataset & performing transformations.
- viewed as data preprocessing step.
- Conditional Question like what source etc is asked.

Survey & Visualise:-

analysts use graphs, charts, visual tools to get a clear overview of data., better understanding data.

patterns, quality & readiness of data.

need. (1) Quick understanding of data.

(2) Data Quality check

(3) Distribution consistency

(4) Assess granularity

(5) population representation.

Tools for data preparation.

- ①. Hadoop. - perform parallel insight & analysis
- ②. Alpine Miner. - provides GUI for creating analytic workflow.
- ③. OpenRefine - opensource tool to work with data.
- ④. Data Wrangler - tool for data cleaning & transformation.

⑤. Model planning.

- Deciding how to approach the data analysis what techniques to use & planning the modelling process. to solve business problems.

Key activities in model planning.

- ①. Assess Data Structure. (structure / semi-structured / unstructured)
- this helps choose right tool & technique for modeling
- ②. Choose suitable Analytical Techniques.
select ~~model~~ methods that helps meet business goals.
& test hypotheses effectively
- ③. plan modeling Approach.
Decide if single model is enough, or series of them is needed.
- ④. learn from past work.
Research similar problems ^{from} past & use those insights.

need.

1. ensures right tool selection.
2. align model work with business goals.

Data Explanation & Variable Selection.

- explaining data to understand relationship b/w variables & selecting most relevant ones for modeling.

Activities involved in Explanation & variable selection

- ①. Explaining variable Relationship
- ②. Use Data Visualization tools. to find patterns, trends.
- ③. consider stake holders input to guide variable selection.
- ④. identify key predictors. :- select most imp. variables.
- ⑤. variable selection needs trial & error hence iterative.
- ⑥. plan for specific models.
- ⑦. Remove irrelevant or Redundant Data.

Model Selection.

Step where the team chooses most suitable analytical techniques. to solve business problems. based on type of data & project goals.
key points.

- ①. Goal oriented selection. - ^{choose} best analytical technique
- ②. Real world representation - models are simplified versions of real-world behavior. created using ~~the~~ rules & logic. to mimic real events.
- ③. Type of Data. Matters :- structured / un / hybrid. & select model accordingly.
- ④. Use of statistical Tools :- initial model building R, SAS, Matlab for quick testing & validation.
- ⑤. Tool ~~to~~ limitations. :- tools don't work well with large datasets.
- ⑥. Candidate Models :- select & see which performs best
- ⑦. prepare for next phase

Tools : 1] R :- complete set of modelling capabilities, 5000 packages for. D.A & graphical representation
2] SQL Analysis services :- can perform in database analysis of common data mining functions, involved aggregations,

& basic predictive models.

3] SAS / ACCESS provides integration b/w SAS & the analytics sandbox via multiple data connections.

4] Model Building.

- Team creates, test & validate the model using training & testing data, based on techniques selected in previous phase.

Key activities

①. Execute the chosen models.

- Run models planned in phase 3 - model selection.

②. Prepare Datasets.

Divide data into training, testing & production set.

③. Train the model using training dataset.

④. Test the models using testing dataset to eval. model.

Questions to consider.

①. Model Validity & Accuracy

②. Domain Expert Validation

③. Parameter Meaningfulness

④. Readiness for use. (accurate, less mistakes, desired env).

Tools.

① SAS Enterprise Miner - enterprise level computing & analytics.

② SPSS Modeler (IBM) - enterprise level

③ Matlab - high level data analytics, algo, data exploration.

④ Alpine Miner - provides GUI & backend analytics tools

Statistical & mathematical - for data mining & analytics tools.

⑤ R & PL/R - procedural lang for PostgreSQL with R

⑥ Octave - for computational modeling.

⑦ Weka - data mining software package with analytic workbench

⑧ Python - toolkit for ML & analysis

⑤ Communicate Results.

- determine if the team succeeded or failed in objectives
- Assess if results are statistically significant & valid
- Identify aspects of results that present salient features
- communicate & document the key findings & major insights of analysis
- Satisfaction of customer

⑥ Operationalize.

model is deployed in real world setting & its performance is monitored in production key points.

- ①. Communicate Project Benefits. (Business stakeholders)
- ②. Pilot Deployment :- launch model in small controlled environment to minimize risk
- ③. Risk management :- pilot run helps detect risk early & reduce risk
- ④. Database Execution for Efficiency.
for large Dataset. the algo needs to be executed inside database. instead of memory
- ⑤. Live Testing :- Run model in Real time conditions.
& test actual performance
- ⑥. performance monitoring (model accuracy, speed & op)
- ⑦. model Retraining. with new Data.
- ⑧. Ready for full deployment

Four Main Deliverables.

- ①. Presentation for project Sponsors.
 - executives & senior management.
 - focused on key takeaways.
 - high level overview.

- ②. presentation for Analysts.
 - aimed at business analysts & operational ^{users}.
 - Explains ~~of~~ business processes & reports.
- ③. Code for Technical Team.
 - complete & cleaned codebase used to build & deploy the model.
 - Replication & scaling of model.
- ④. Technical specifications!
 - Detailed doc of how code works & how to implement it
 - helps in future maintenance & upgrades.

BI

Data science.

- | | |
|--|---|
| ①. Analyzes past & present data. | ① predicts future trends & outcomes. |
| ②. Reporting & decision making based on historical data. | ② Building productive models & uncovering patterns. |
| ③. Mainly structured data used | ③ structured / unstructured data used |
| ④. Tools:- Power BI, Tableau | ④. Tools:- Python, R |
| ⑤. used ^{by} Business analysts, managers. | ⑤. Data Scientists, researchers. |
| ⑥. less complex., UI based. | ⑥. more complex., stats & ML |
| ⑦. Eg. Monthly salary Report | ⑦ predict next month customers churn |

Descriptive Analytics.

- looks at past data to understand what has already occurred.
 - helps in identifying trends, patterns & summaries.
 - answers the question "What has happened?"
- Purpose:- Summarize historical data

tools & techniques.

1. data aggregation.
2. data visualization.
3. statistical summaries.

Eg. Monthly salary Report.

Diagnostic Analytics. - why did it happen?
- investigates causes behind past outcomes.
- helps in identifying reasons or root cause of trends or problem.

Purpose.

To investigate causes behind trends or outcomes found in descriptive analytics.

Tools.

- ① drill down analysis
- ② Data mining
- ③ correlation & root cause analysis

Eg. Analyzing why sales dropped in a specific region.

Predictive Analytics. - what's likely to happen?
- uses past data, statistics, & ML to make predictions
- helps in forecasting future events.

Tools. ① Regression

② Classification.

③ forecasting models.

Eg. Predicting which customers are likely to cancel a subscription.

Key Roles for successful Analytics project.

① Business User.

- understand business domain
- provides practical insights & context for problem

② Project Sponsor.

- initiates & funds the project.
- provides high level requirements & business goal.

③ Project manager.

- manages project timelines & resources.

④ Business Intelligence Analyst.

- Bridges gap b/w business & data.
- understand data trends, reporting & KPI's

⑤ Database Administrator (DBA)

- Design & manage database infrastructure

⑥ Data Engineer.

- Handles data ETL (Extract, Transform, Load)

⑦ Data Scientists.

- Apply statistical & ML models.
- generates insights & predictions to solve problems.