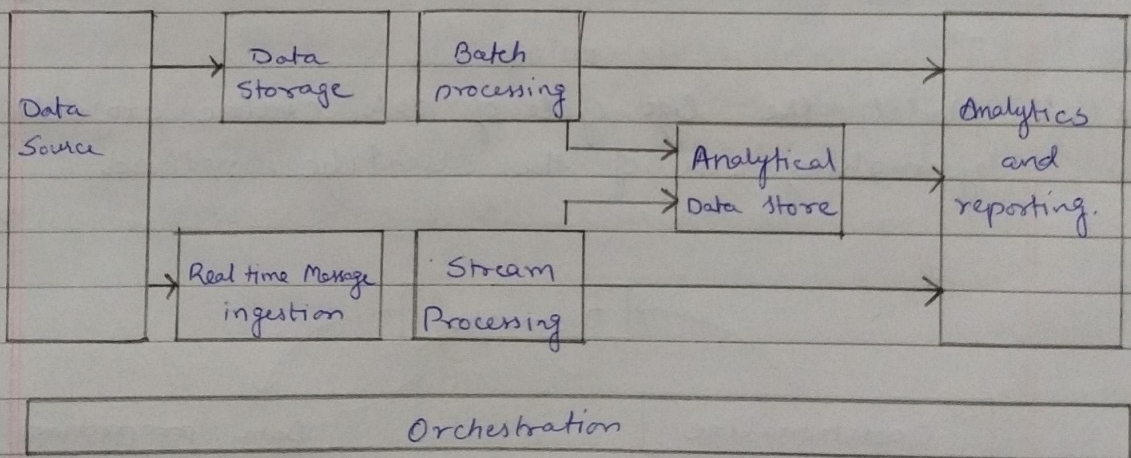Name - Vasu Kalariya
Roll - PE29 (Batch 1)
Sub - BDA

Theory Assignment - 1

**Q1 (a)** Draw and explain the Big Data architecture of the System.

Big data architecture is the foundation for big data Analytics. It is the overarching system used to manage large amounts of data so that it can be analysed for business purposes, steer data analystics and provide an in which big data analystics tools can extract vital business information from otherwise ambigueous data

Big Data Solutions typically involve one or more of following types of workhood
→ Batch processing of Big Data source at rest.
→ Real time processing of Big Data in motion
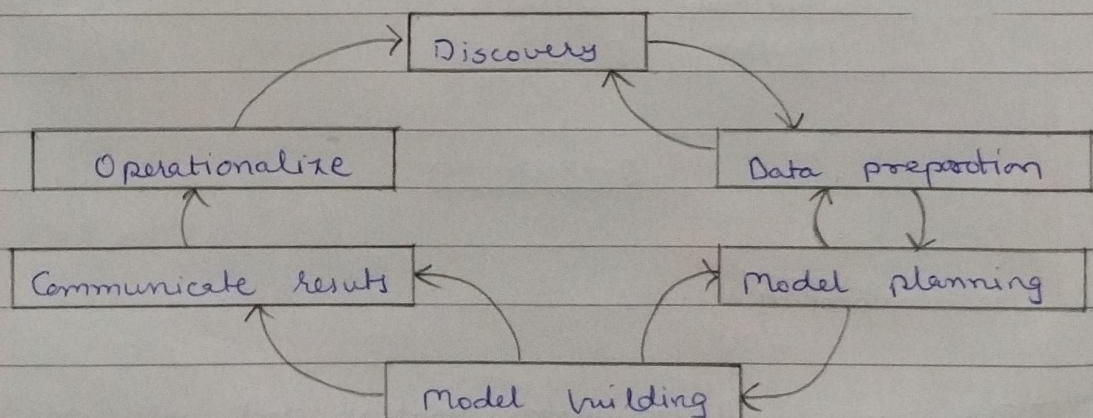→ Interaction explained of big data



The components are:

→ Data Sources: All big data solution start with one or more data source

→ **Data Storage:** Data for batch procesing operation is typically stored in distributed file store that can hold high volume of large file in various formatted. This kind of store is often called Data lake.

→ **Batch procesing:** Because data sets are so large often a big data solution must proces data files using long running batch jobs to filters aggregate and otherwise prepare the data for analysis.

→ **Real time message ingestion:** If the solution indeed real time sources, the architechue must include a way to capture and store real time messages for stream procesing.

→ **Stream preprocesing:** After capturing real time messages, the solution must process them by filtering, aggregating and otherwise ~~preprocess~~ preparing the data for analysis.

**Q1 (b)** Model the life cycle of data centric projects ~~to~~ by making use of the scientific method

Discovery → Data preparation → Model planning → Model building → Communicate results → Operationalize → Discovery
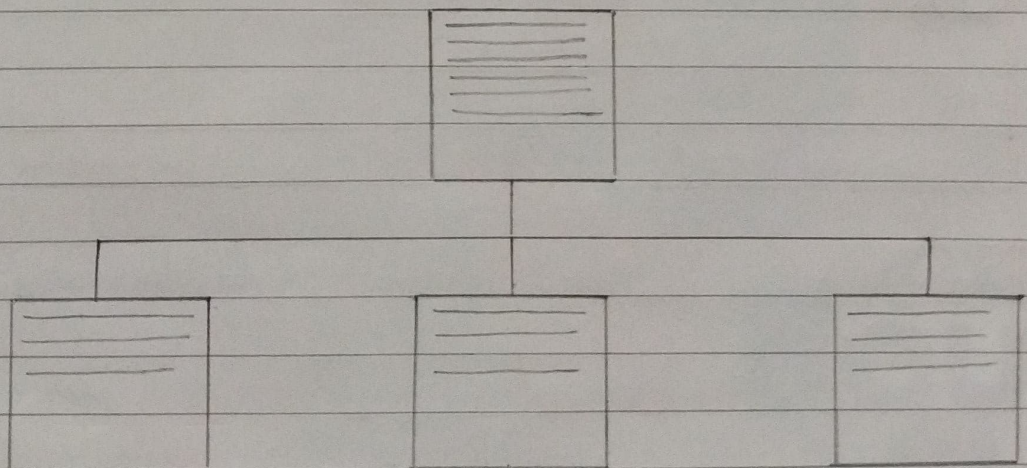
Data analysis lifecycle define analytics process best practices. spanning discovery for project completion the life cycle draws from established methods in the real of data analytics and decision since science

Phase: 1] Discovery
2] Data preparation
3] Model planning
4] Model Building
5] Communicateo results
6] Operationalize.

Q 2 @ MongoDB use horizontal scaling for handling huge amounts of data. Defend it with the help of suitable diagram

Sharading is a method for distributing data across multiple machines. MongoDB uses horizontal scalling which invobeds, dividing the system datasut and load multiple servers, adding additional servers to increase capacity as required.

Q 2 (b) Consider Employee database having.
{ Eid, Department, Fname, Address: [ street: streetname,
city: cityname, state: statename, pincode: pincode-
value }, phone: [ home contact, mobile: contact ],
age, salary }

1] Display only 3 employees residing in the
state: "Maharastra". Skip the first employee
document

→ db.Employee.find ({ "Address.City" }: "Maharastra" }). limit (3).
skip (1).

2 Eid should allow only unique value

→ db. Employee. createIndex ({"Eid": 1}, unique: true)

3 Display dis department-wise average salary of
employee having the average salary > 5000.

→ db. Employee. aggregate ( [
{ $group: { _id : "$DeptName", Avg-Sal: { $avg: "$salary"}
}},
{ $match: { Avg-sal: { $gt: 5000 }}}.
]);