

58369

Copy No.

Number of Book  
II nd used

TRIBHUVAN UNIVERSITY  
INSTITUTE OF SCIENCE & TECHNOLOGY  
**EXAMINATION BOARD**  
KIRTIPUR

Code No.:

**Book I**

Code No.:

Students are required to write their answers on both sides & a margin of about  
1 $\frac{1}{4}$  inches should be left on each page.

<b>MARKS OBTAINED</b>	
1st Q	<input type="text"/>
11th Q	<input type="text"/>
2nd Q	<input type="text"/>
12th Q	<input type="text"/>
3rd Q	<input type="text"/>
13th Q	<input type="text"/>
4th Q	<input type="text"/>
14 th Q	<input type="text"/>
5th Q	<input type="text"/>
15th Q	<input type="text"/>
6th Q	<input type="text"/>
16th Q	<input type="text"/>
7th Q	<input type="text"/>
17th Q	<input type="text"/>
8th Q	<input type="text"/>
18th Q	<input type="text"/>
9th Q	<input type="text"/>
19th Q	<input type="text"/>
10thQ	<input type="text"/>
20th Q	<input type="text"/>
TOTAL	40

**INSTRUCTIONS TO CANDIDATES**

[1] Each Candidate will write legibly on the title page his or her ROLL NO. REGISTERED NO., AND THE SCRIPT in which answers are written but not his or her name and the name of the campus from which he or she appears. This should be done in each answer - book used before beginning to write inside.

[2] No loose papers will be provided for scribbling and no other paper should be brought in for this purpose. Any candidate found with loose paper in his or her possession WILL BE EXPELLED. All work must be done in the book provided and pages MUST NOT BE TORN OUT. The book provided CANNOT BE REPLACED BY ANOTHER but, if necessary, an additional book will be given. All work intended for the examiner must be written ON BOTH SIDES of the paper. Anything cancelled will not be looked into. Should a torn leaf be discovered inside an answer book, it should not be removed but crossed out and folded after bringing it to the notice of the invigilator.

[3] Candidate is forbidden to write answers or anything else on the question - paper.

[4] No Candidate will be allowed to leave the room until one hour has passed from the time when the papers are distributed.

[5] Candidate who uses two or more answer books should see, before handing over to the invigilator, that they are properly stitched together.

**INSTRUCTIONS TO THE EXAMINER**

- Aggregate marks of each question should be placed in the given appropriate box.
- Marks for parts of a question should be totalled under the question inside the answer- book.

MDS

Level .....

Year/Part/Sem. II/II/III

SMS TV

Centre .....

Roll No. 06

Six

Arpan Sapkota

Reg. No. ....

Subject Techniques for Big Data

Paper .....

Script English

Date 2081/02/04

Signature of Scrutiny Board

Signature of Examiner

## Group 'A'

1. NoSQL stands for Not only SQL which means, SQL which can only works with relational and structured language was the limitation overcame by NoSQL. Unstructured, semi-structured data can be processed with NoSQL and it has high scalability.

In SQL, the five V's are fixed and the complex relationship can be handled by the ~~Joins~~, however those are not possible in NoSQL which are the main key difference between SQL and NoSQL. The other general differences are:

### NoSQL

i) It stands for Not Only Structural Query language.

ii) It works on non-relational databases.

iii) NoSQL are highly scalable.

iv) Horizontal partitioning are done in NoSQL.

v) It works on Unstructured and, semi-structured data.

vi) Example: MongoDB, HBase, Cassandra, etc.

### SQL

i) It stands for Structural Query language.

ii) It works only on the Relational databases.

iii) SQL are hard to scale.

iv) Vertical partitioning are done in SQL.

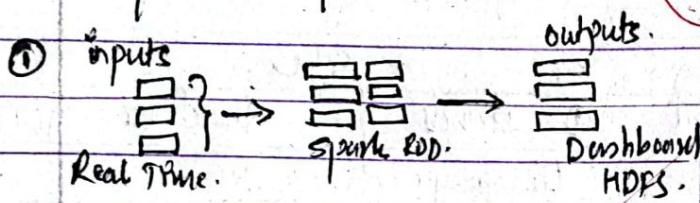
v) It works on Structured data only.

vi) Example: Oracle, MySQL, MS SQL, etc.

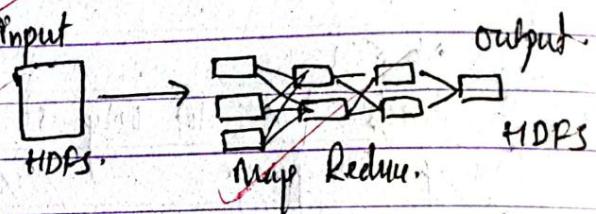
2. Apache Spark is different from Hadoop Mapreduce.

Apache Spark is the real time stream processing in big data technology where the real time data like: live video streams, sensor logs, sensor data are processed in real time. whereas Hadoop Mapreduce is the batch processing in big data where the stored data are processed in batches of jobs.

### Apache Spark



### Hadoop Mapreduce.



① It works on Real time data

It works on stored batches of data.

② Apache spark processes data quickly with the help of RDD.

Hadoop Mapreduce processes data very slow.

③ It does caching in memory.

There is no caching in Hadoop mapreduce.

④ RDD is the heart (core) component of Apache Spark.

Mapreduce framework is the core component of Hadoop.

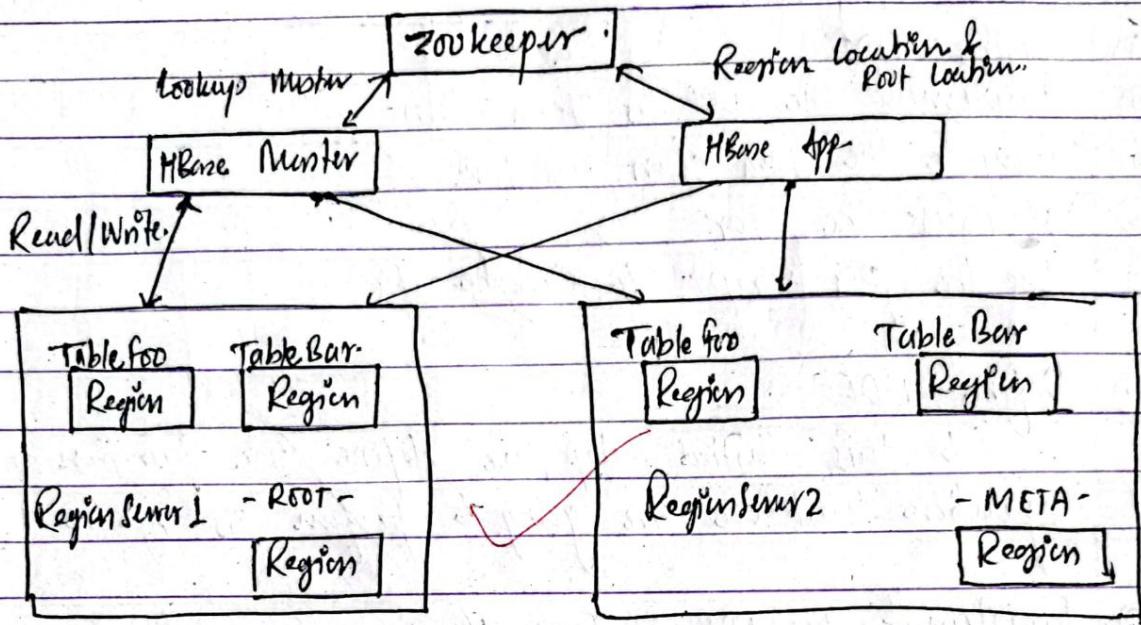
⑤ No disk is involved here.

Disk is involved here so it is slow.

⑥ Spark can be implemented in Python, Java

It is implemented only in Java.

#### 4. HBase Architecture.



Q. In HBase Architecture zookeeper is the main component which acts as the master and look after the Region servers.

Here master assign job to the Region servers. HBase is built on top of Hadoop Distributed File System (HDFS). It has following core components:

- ① Master
- ② Region
- ③ Region Server

Master assign jobs to the slaves which are Regions. In this case, Regions are the tables in the region servers which acts as slave. Master performs the read write operations on the Region servers.

5. The process to create UDFs (User Defined functions) in Apache Pig are:

- (i) Define UDF
- (ii) Implement the UDF in Java class.
- (iii) Compile the Java class into Jar.
- (iv) Register the Jar
- (v) Use the pig script to execute it

#### (i) Define UDF

In this initial step, we define the function or the operation that we are going to perform in Apache Pig.

#### (ii) Implement the UDF in Java class.

We then implement the defined function into a Java class and methods.

#### (iii) Compile the Java class into Jar

The implemented Java program is compiled into the Jar file.

#### (iv) Register the Jar

Jar file is registered for the data flow operation in pig. Here, REGISTER, DEFINE, LOAD, and DUMP operation are performed.

#### (v) Use the pig script to execute.

The registered program is executed in the pig command-line interface.

Q. The different execution modes of Apache Pig are:-

- (i) Grunt Shell Execution.
- (ii) CLI Execution.
- (iii) UI Execution.

#### (i) Shell Execution

Pig script can be executed with its shell which is known as Grunt Shell. We give commands to this shell and pig script is executed.

With  
map mode

#### (ii) CLI Execution.

Pig script can be executed with command line interface as similar to shell execution.

#### (iii) UI Execution.

Pig script can also be executed in user interface mode where we give input script in the UI and get output in the UI.

In order to execute the pig we must have to follow the following pattern.

REGISTER → Register Java class

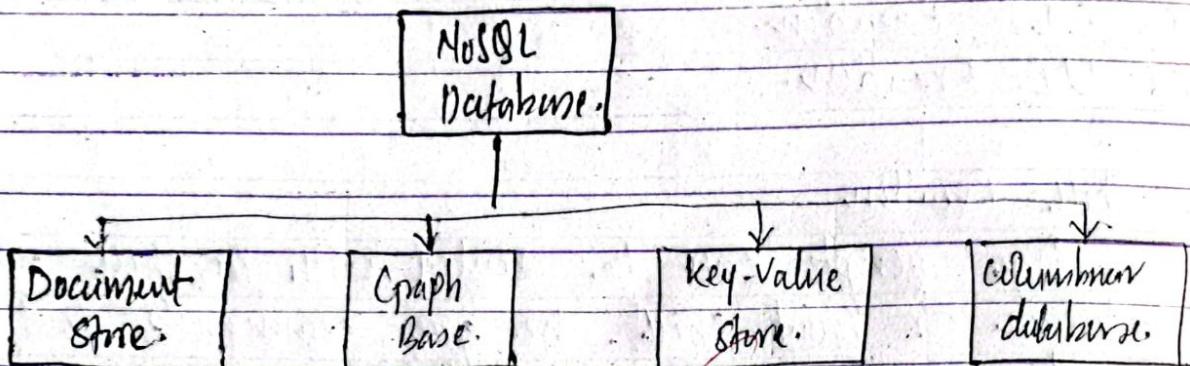
DEFINE → Define the pig class

LOAD → Load the data

DUMP → Displaying the result.

Group 'B'

6. The different types of NoSQL databases are:-



### ① Document store.

Document store are the three type of NoSQL database where the documents i.e. records are stored in a key value pair structure like in a JSON file.

MongoDB is the example of NoSQL Database.

Data stored in document store looks like:-

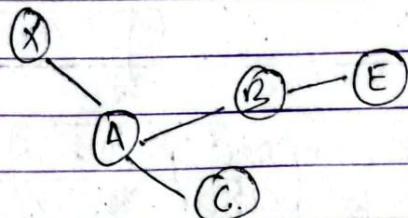
```
{ "Name" : "Arpan"  
  "Address" : "Khumaltar"  
  "website" : "www.arpansayokta.com.in"}
```

MongoDB, CouchDB, etc are the document store NoSQL databases.

### ② Graph Base.

Graph Base is the another type of NoSQL database which is used to store the data in the graphical structures. Data are stored in the vertices of the graph.

Example:-



Graph base NoSQL Database.

Hadoop is the example of Graph Base NoSQL database.

### ⑪ Key-value store.

~~key-value store are the those type of NoSQL database where the data are stored in the basis of key and values on the table. key represents the unique identifier in this database and value is the actual data pointed by the key.~~

Example:

key	Value
1	Ram
2	shyam
3	Ghanashyam

MongoDB, cassandra, etc are the key-value store databases.

### ⑫ Columnar databases.

Columnar databases are those database where the data are stored on the column families with row key as key value.

Example:

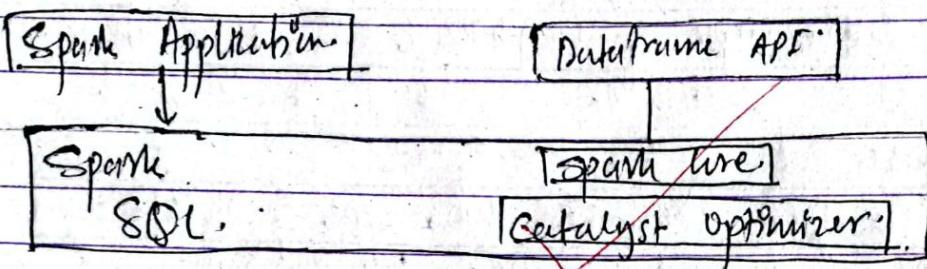
key ↓	column family - 1		column family - 2		3 cell.
	column 1	column 2	column 1	column 2	
Row 1					
Row 2					

Example: HBase is the columnar Database.

7.

### ⑨. Spark SQL.

Spark SQL is one of the components of the Spark Ecosystem which is used to perform the query operations on the Spark.



In Spark SQL, we have spark Application which is the client side of the SQL where client can perform the query operation with the spark Application which goes into the spark SQL Engine in the spark core.

Catalyst optimizer performs the optimization task on the query plan. The logical query are optimized into physical query with this optimizer.

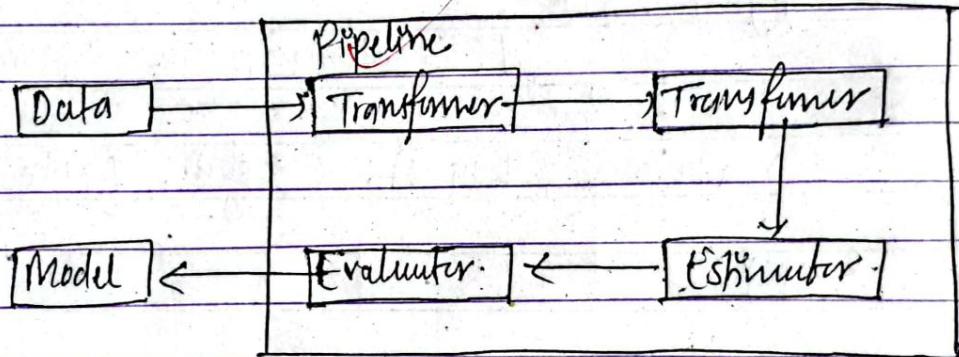
Dataframe API is an interface where the spark Query engine can be integrated into the different data sources like: HDFS, MySQL, MongoDB, etc.

Spark core is the SQL engine of the spark which performs the query operation as according to the input. In this way with the combined effort of all components spark SQL is operated.

## ⑤. Spark Mlib.

Spark Mlib is the another component of spark ecosystem which is used to perform the machine learning task in spark. It has growing set of machine learning Algorithms like PageRank, Triangle count algorithm.

spark mlib has the following pipeline



Data is provided as input to the transformer which clean and loads the data into spark and the feature extraction are done in another transformer. Once the features are extracted then, it is fitted in the estimator, These are then sends to the evaluator which evaluates the performance of the machine learning algorithm and finally it sends an output which is the model.

We have various algorithms in spark mlib and some are:-

- ① PageRank
- ② SVD++
- ③ Triangle cutting, etc

10. db.order.find({})

Use order

① Maximum price per customer

db.order.aggregate([

{

\$group : {

- id: "\$customer", price: "\$price"

}

max price : { \$max : { \$sum: "\$price" } }

}, ] )

② Total revenue per day.

db.order.aggregate([

{

\$group : {

day: "\$date", price: "\$price"

,

TotalRevenue : { \$sum : 1 }

}, ] ).

① Average price per customer per product.

db. order. aggregate ([  
])

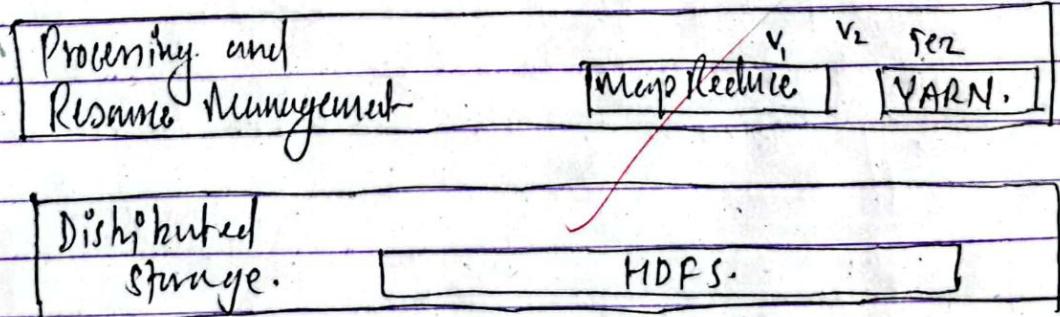
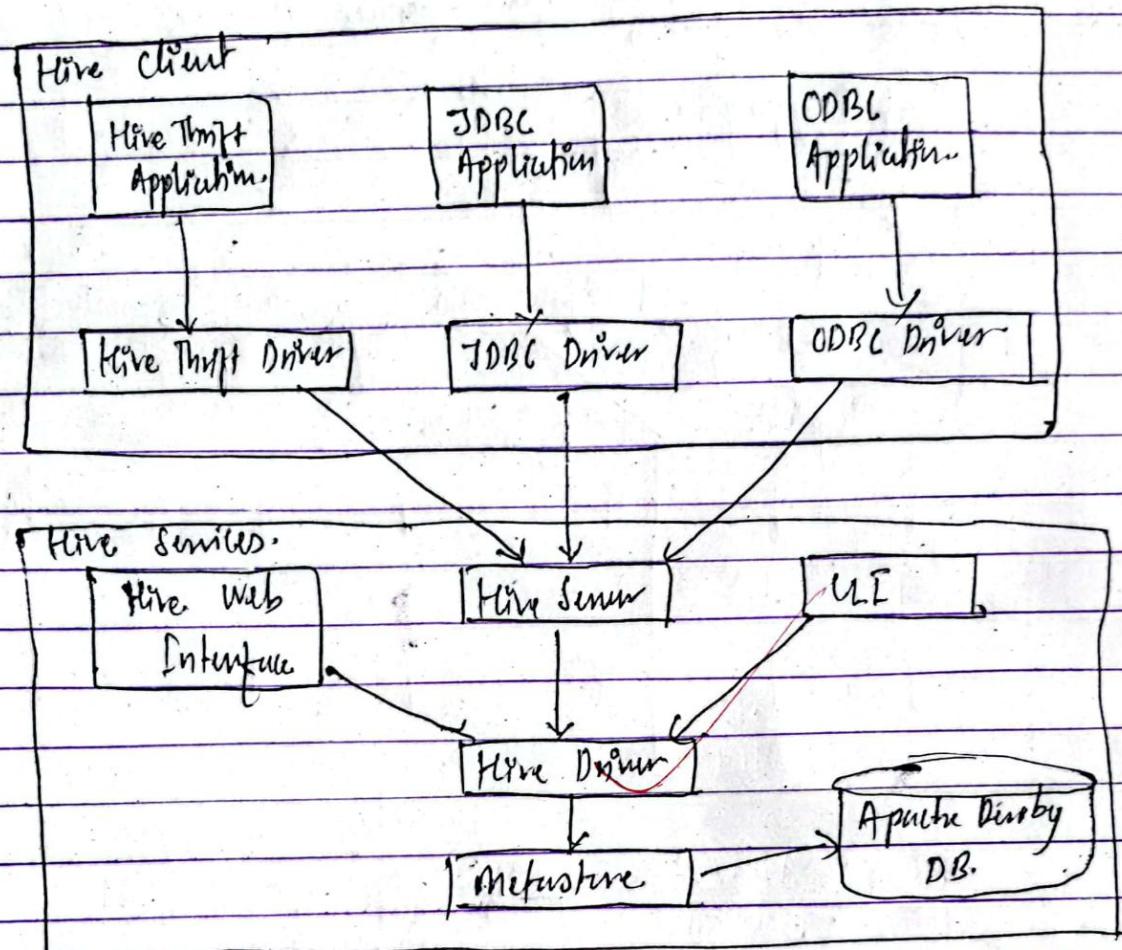
\$group : {

- id : "\$customer", product : "\$product", price : "\$price"

AvgPrice : { \$avg : "\$price" }.

}) ] )

## 9. Hive Architecture



Hive architecture has the following major components:

- (I) Hive Client
- (II) Hive Server
- (III) Processing and Resource Management.
- (IV) Distributed Storage.

#### (I) Hive Client:

It is the client side component in hive which has the applications like Thrift application, JDBC Application, ODBC Application and the respective driver. These are the database connection application to connect to the server to process the data in the server.

#### (II) Hive Server:

Hive servers are the another main component in the server side which has the following services:

- Web Interface
- Hive Server
- Command Line Interface (CLI)
- Metastore and Apache DB.

When client sends request to hive servers then the processing happens with the help of these services.

#### (III) Processing and Resource Management

Mapreduce, YARN are used for the processing and resource management task in apache hive.

#### (IV) Distributed Storage:

For storing the files, HDFS is used in hive.

In this way, Apache Hive runs with the client request on Apache Hive client side. It is normally a Data Warehouse.

Comparing Apache Hive with Apache Pig-

### Apache Pig.

① It was developed by Yahoo.

② It is used for Data Flow programming.

③ It does not have fixed schema.

④ Pig Latin is the language used in Apache Pig.

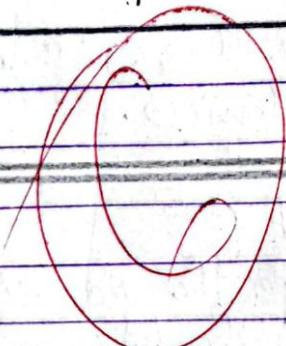
### Apache Hive.

It was developed by Facebook.

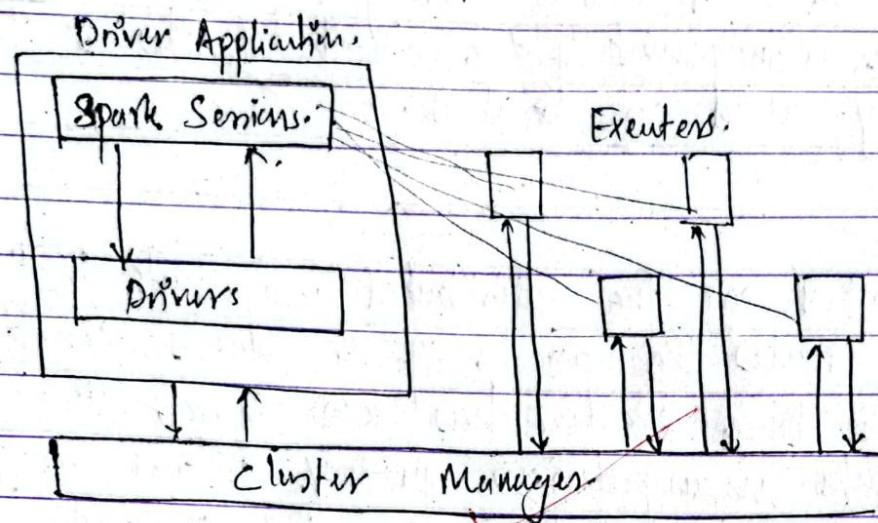
It is used for Research and Analysis.

It has schema during insertion.

HiveQL is the language used in Apache Hive.



## 8. Architecture of Apache Spark.



Architecture of Apache spark consists of following core components.

- ① Driver Application.
- ② Cluster Manager
- ③ Executors.

### ① Driver Application.

Driver application is the spark application where the spark sessions and drivers are contained. This is simply the user session. When any real time data needs to be processed, first the user submits the job from this component of the Apache spark architecture. spark sessions are handled by the drivers.

### ② Cluster Manager.

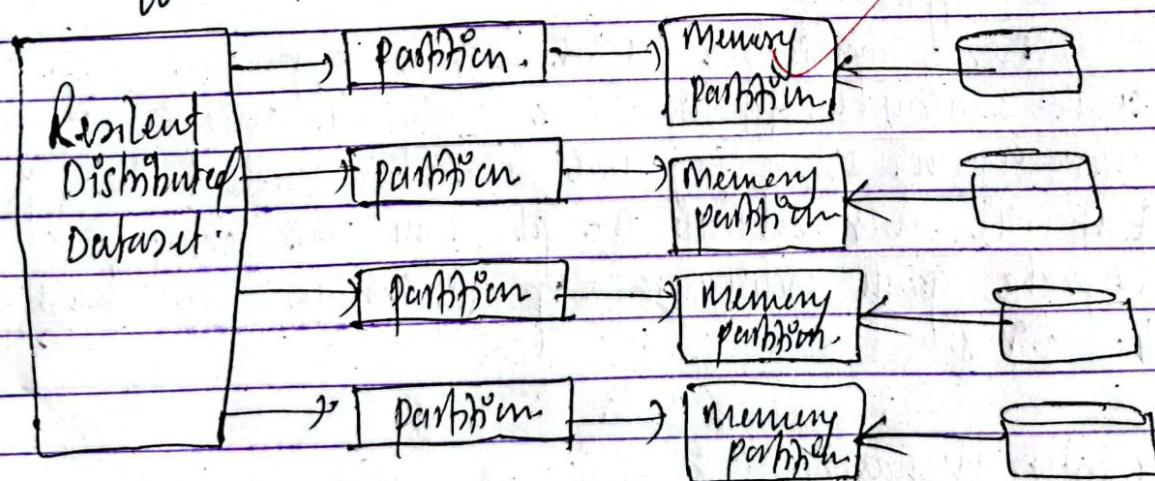
In order to process real time stream of data, there are lots of clusters where the partition are made to process the data. There are checks of RDD Persistence

Distributed Datasets) in order to manage those distributed clusters we need cluster manager. This cluster manager does the management of the distributed clusters those are generally called as executors.

### ⑦ Executors.

Executors are the distributed components in the spark which executes the apache spark actions i.e. data processing actions. These executors has RDD inside it, which performs the data manipulations quickly. It acts as the sumo since executors just performs the jobs assigned by the cluster manager.

RDD (Resilient Distributed Dataset) is the core part of the Apache Spark where the data are immutable and distributed into different partitions.



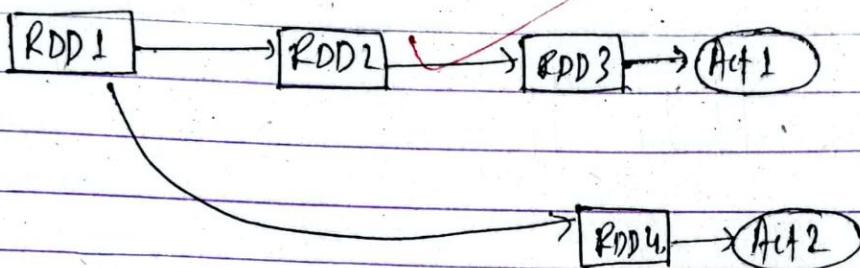
Files.

~~RDD~~ helps to form PEs chain and the Data processing

becomes very fast in Apache Spark.

Spark transformation and actions.

In Apache Spark, we have chain of RDDs  
let us consider an example.



Spark Transformation

When there needs to form another RDD in the chain then it uses previous RDD and it does not uses the actual data. It just performs the transformation only. The chain formed are just from the transformation of data which is known as spark transformation.

Example

Spark Action.

Whenever there needs to perform any task then spark will eagerly acts and perform the task which is the spark action.

In spark Evaluation is slow and task action are fast so it is also called as lazy evaluator.