

Name: Krushnakumar Patle

Email: krishnapatle128@gmail.com

Batch: Data Engineering Batch-1

*** Spark SQL:-**

- Spark SQL is component of top of the spark core.
- Spark SQL was first released in Spark 1.0
- Spark introduce a module for structure data processing.

challenges:-

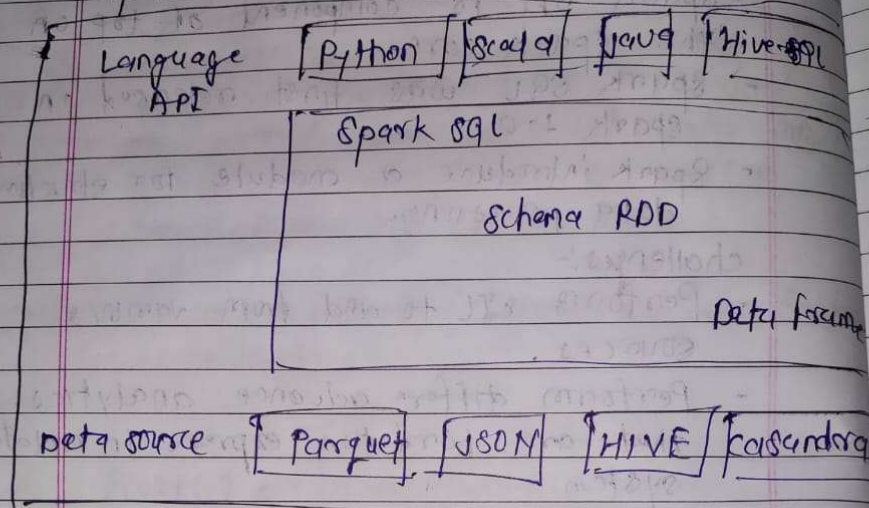
- Perform ETL to and from various sources..
- Perform differ advance analytics that are hard to express in relational system.

solution:-

- A Dataframe API that can perform relational operation on both external data source & spark built-in data source.
- A highly extensible, optimizer, catalyst that uses feature of scalar to add comparable rule, control code gen & define extension.

Spark SQL	Spark Streaming	MLlib (ML)	GraphX (graph)
Apache Spark			

Spark SQL Architecture



• Language API:-

- Spark is compatible with diff^t language & Spark SQL.

- supported by - python, scala, java, HiveSQL

• Schema RDD:-

- Spark Core is designed with special data structured called RDD.

- Spark SQL work on schema, table & records.

- Therefore, we use schema RDD.

- we call this Schema RDD & Datasource

1 Data source:-

- It is text file, Avro file etc.
- Data source for spark SQL is diff
- Those are Parquet file, JSON document, HIVE tables, & cassandra database.

1 features of spark SQL:-

1) It Integrated:-

- seamlessly mix SQL queries with spark program.
- spark SQL lets you query structured data as a RDD in spark, with integrated APIs in python, scala, & Java.
- This tight integration makes it easy to run SQL queries alongside analytical alg.

2) Unified data access:-

- Load & query data from a variety of source like Apache Hive table, Parquet.

3) Hive Compatibility:-

- Run unmodified Hive queries on existing warehouse.
- ex. `SELECT COUNT(*) FROM hiveTable WHERE hive-udf(data)`

* Hive:-

Apache hive is a distributed, fault-tolerant data warehouse system that enables analytics at a massive scale.

* User-Defined Functions (UDF):-

- UDF:- Plug in your own processing code & invoke it from Hive query.

- UDF (Plain UDF)

- Input: single row, output: single row.

- UDAF (User-defined Aggregate funⁿ)

- Input: multiple rows, output: single row.

4) Standard Connectivity:-

- connect through JDBC or ODBC

5) Scalability:-

- use same engine for both interactive & long queries.

- do not worry about using a diff engine for historical data.

fault-
m
cessive

very.

now

n)
system

after.

* Spark RDD (Resilient Distributed Dataset)

- RDD is fundamental data structure of Spark
- It is an immutable distributed collection of objects that can be stored in memory or disk across a cluster
- Each dataset in RDD is divided into logical partitions, which may be computed on diff't nodes of the cluster
- Automatically rebuilt on failure
- RDDs can contain any type of Python, Java, or Scala objects.
- Read-only, partitioned collection of records.
- created through deterministic operation on either data on stable storage or other RDDs.
- RDD is fault-tolerant collection of elements that can be operated on in parallel.
- Two ways to create RDD:-
 - ① parallelizing:-
 - ② referencing a dataset in an external storage system

* Dataset & Dataframe:-

- A distributed collection of data, which is organized into named columns.
- Dataframe can be constructed from an array of diff't sources such as Hive table, structured data files,

* Dataframe:-

Data is organised into named columns like a table in a relational database.

* DataSet: a distributed collection of data.

* Features of DataSet / Dataframe:-

- Ability to process the data in Kb to Pb
- Support diff't data format.
- Can be easily integrated with all Big data tools & frameworks
- provides API for python, Java, Scala, R, R programming.

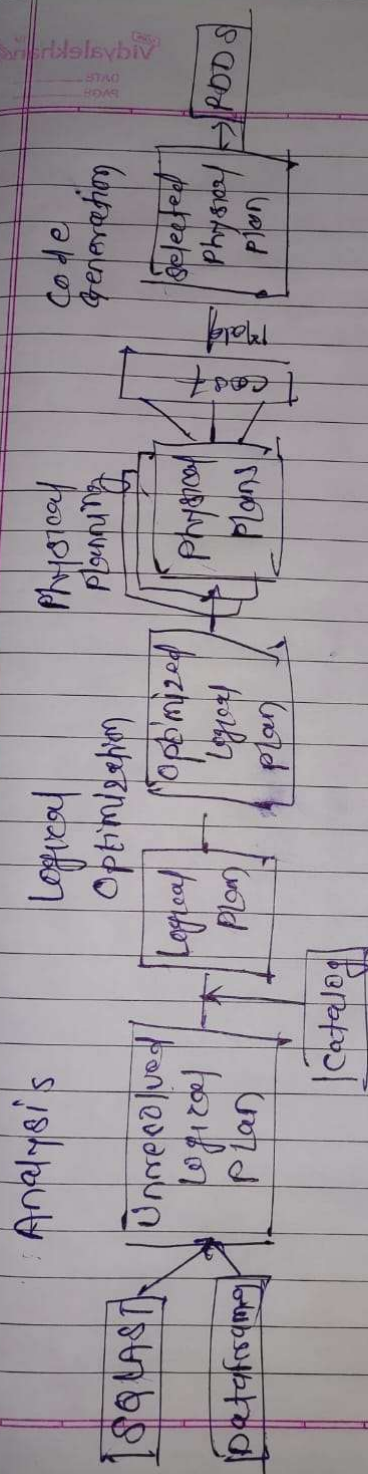
named column
database

on of data.

me:
in

mat.
with
networks
on Java,

Plan optimization & Execution



Dataframes & SQL share the same optimization/execution pipeline