Name : Parth Nandedkar
Date : 10 Feb 2024
Topics : PySpark
Batch : Data Engineering Batch-1

**Handwritten Notes :**

- Pyspark
- Setting up Azure Lab.

## Spark SQL

- It is a module for structured data processing
- Spark SQL is a component Spark Core that introduces a new data abstraction called SchemaRDD

Challenges ⟶ Solution

Perform ETL to and from ⟷ A DataFrame API that can perform relational
Various (Semi/unstructured) data      operations on both external data source and Spark's
Source      built-in RDD.

## Spark SQL Architecture

| Python | Scala | Java | Hive-QL |
|---|---|---|---|

Language API

| Spark SQL |
|---|
| Schema RDD |
| Data Frame |

| Data Source | Parquet | JSON | HIVE | Cassandra |
|---|---|---|---|---|

(Hive tables are here)

- Language API -
- Schema RDD -
- Data Sources -

Features of Spark SQL

1] Integrated : (mix with SQL & has APIs Python, Scala and Java) can -
2] Unified Data Access ( Load and query data from variety of sources)
3] Hive Compatibility :

⌐ Run unmodified Hive query on existing warehouse
5] Scalability : - Spark SQL the reuse the Hive front end
Use the same engine for and Meta Store, giving you full compatibility
both interactive & long quaries with existing Hive data, queries & UDFs.
Simply

UDF ( User defined functions)

Tableau ⟶ Qlik

Features

Spark RDD :⟶

- It is fundamental data structure of Spark
- It is immutable distributed collection of objects that can be stored in memory or disk across cluster
- Each dataset in RDD is divided into logical partitions, which may be computed on diff node of cluster
- Parallel functional transformation.

- Automatically rebuild failures
- RDD contains any type Python, Java, Scala objects including user defined classes.

- Formally, an RDD is read-only, partitioned collection of records.

- RDDs can be created through deterministic operations on either data on stable storage or other RDD.

- RDD is fault-tolerant collection of elements that can be operated on in parallel.

Dataset & Dataframe:

- A distributed collection of data, which organized into named columns.
- Conceptually, it is equivalent to relational tables with good optimization techniques.
- A Dataframe can be constructed from an array of different sources such as Hive tables, Structured data files, external databases.
- This API was designed for modern Big data.

Data frame →
         Data is organized into named columns, like table in a relational database.

DataFrames

Features of DataFrame:

- Ability to process the data in the size of kilobytes to Petabytes on a single node s cluster to large cluster.

- Supports diff data formates (Avro, csv, elastic search, Cassandra) and storage systems (HDFS, HIVE, mysql etc)

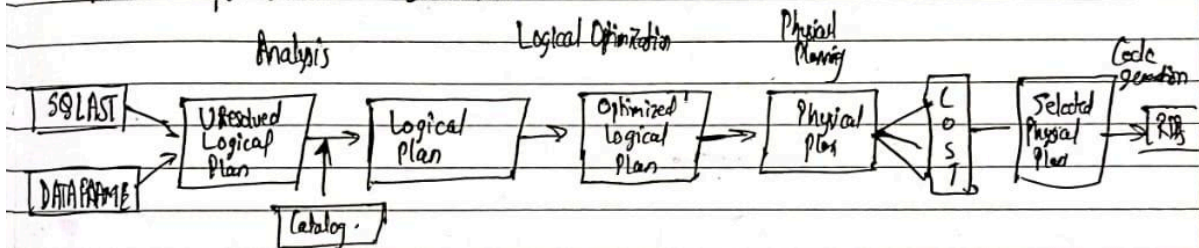* State of art optimization & code generation through the Spark SQL Catalyt optimizer (tree transformation framework)

- Can be easily integrated with all Big data tools and frameworks via Spark core.

- Provides API for Python, Java, Scala & R programming.

# Spark SQL

1. Write less code
2. Read less data
3.

## Plan Optimization & Execution:



Analysis     Logical Optimization     Physical Planning     Code generation

SQL AST → Unresolved Logical Plan → Logical Plan → Optimized Logical Plan → Physical Plan → COST → Selected Physical Plan → RDD

DATA FRAME

Catalog

# Creating Database using Spark SQL :



# Creating Table using Spark SQL :

## Adding data to table :

# Fetching Data from table :



# Filtering Data from table :