

NAME : AKULA SHARATH CHANDRA

BATCH: DATA ENGINEERING

DATE : 10-02-2024

**TOPICS : CONCEPTS OF SPARK SQL AND
HANDS ON OF SPARK SQL**

CONCEPTS OF SPARK SQL:

10/09/24

1) Spark SQL

-) Introduced in Spark 1.0 (May, 2014)
-) Created by Michael Armbrust & Reynold Xin from databricks
-) introduces programming module for structured data processing called Spark SQL
-) Provides programming abstraction called DF & can distributed SQL query engine

2) Challenges

- Perform ETL to & from various data sources
- Perform advanced analytics

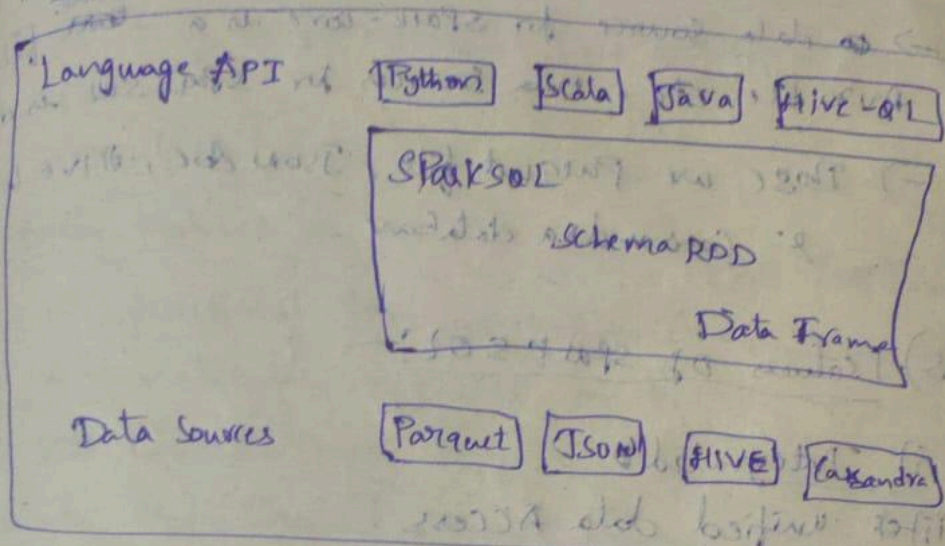
Sol'n's

- A D.F. App can perform relational operations on both external data source & Spark built-in RDD's
- A highly extensible optimizer, Catalyst, that can use features of Scala to add composable rules & define extensions

3) -) Spark SQL is Spark module for structured data processing

-) Spark SQL is component on top of Spark core that introduces.

4) Spark SQL Architecture:-



i) Language API:-

- Spark is compatible with diff lang & SparkSQL.
- It also supported by Python, Scala, Java, Hive-QL.

ii) Schema RDD:-

- Spark Core designed with special data structure called RDD.
- SparkSQL works on schemas, tables & records.
- ∴ we can use Schema RDD as temporary table.
- We can call this Schema RDD as data source.

iii) Data Sources:-

→ data source for Spark-core is a Text
 Avro, etc. Data sources for Spark SQL

-> These are Parquet file, JSON & CSV, HIVE & Cassandra database -

5) Features of SpalSOL:-

(i) Integrated unified data Access.

iii) hive compatibility:

→ run unmodified hire guides on existing warehouses.

→ Spark SQL unifies the hive frontend & metastore giving you full compatibility with existing hive data, queries & UDF's.

→ simply install it along

iv) Standard Connectivity (drivers)

v) Scalability

6) Use refined functions

• Plug in your own processing code & invoke it from a hive query

— UDF (Plain UDF) i/p - single row o/p - single row

- UDAF User defined Aggregate function

— UDIF (User defined table-generating function)
 I/P—single row O/P—multiple rows

Spark RDD

- > it is a fundamental data structure of Spark.
- > immutable distributed collection of objects.
- > each dataset is divided into logical partitions.
- > 11th functional transformations (map, filter, ...)
- > Automatically rebuilt on failure.
- > ~~It~~ contains any type of Python, Java or Scala obj, including user-defined classes.
- > RDD is read only, Partitioned collection of records.
- > ~~It~~ It can be created through deterministic operators on either data on stable storage or other RDDs.
- > It is fault-tolerant collection of elements that can be operated on in 11th.
- > 2 ways to create RDD
 - i) Parallelizing an existing collection.
 - ii) Referencing a dataset in external storage system.

8) Dataset & Dataframe:-

-) A distributed collection of data, which is organized into named columns.

→ It is equivalent to relational tables with good optimization techniques.

-) A DF can be constructed from array of diff source such as Hive tables, Datafiles, External databases etc.

-) This API was designed for modern bigdata & data science applications taking inspiration from dataframe in R Programming & Pandas in Python.

→ Dataframe:-

- data is organized into named cols like table in Relational databases.

→ Dataset:- a distributed collection of data

- Static typing & Runtime type-safety
- new interface added in Spark 1.6

10) Spark save whole story

→ Create & Run Spark Programs faster:

- i) write less code & unified interface to reading/writing data in variety of format
- ii) read less data
- iii) optimize do work

i) write less code: ip & o/p

→ unified interface to reading/writing data in variety of format

df = sqlContext.read

• format("json")

• option("samplingRatio", 0.1)

• load("Path")

df.write

• format("Parquet")

• mode("append")

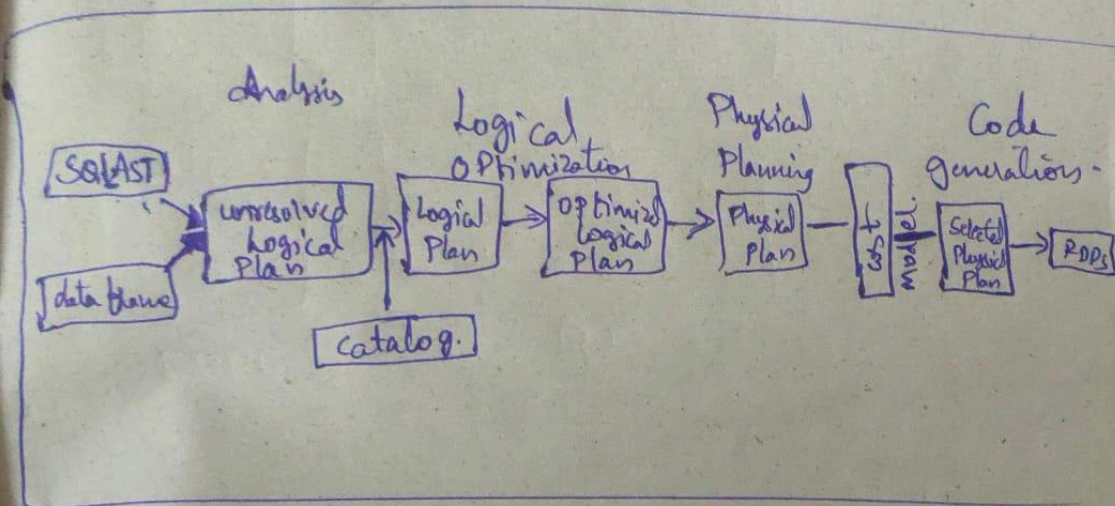
• PartitionBy("year")

• SaveAsTable("faster table")

load(...), save(...) & saveAsTable(...)

create new builders
for doing I/O.

11) Plan Optimization & Execution



HANDS ON :

Spark SQL, DataFrames and Datasets Guide

Python ☆

File Edit View Run Help Last edit was 9 minutes ago

New cell UI: OFF

Run all

Connect

Share

Publish



```
1 from pyspark.sql import SparkSession
2
3 spark = SparkSession \
4     .builder \
5     .appName("Python Spark SQL basic example") \
6     .config("spark.some.config.option", "some-value") \
7     .getOrCreate()
```

Command took 0.16 seconds -- by a user at 10/02/2024, 15:13:38 on unknown compute

Cmd 2

Python

```
1 df = spark.read.json("/FileStore/tables/bus/people.json")
2
3 df.show()
4
```

df: pyspark.sql.dataframe.DataFrame = [age: long, name: string]

```
+-----+
| age | name |
+-----+
| null | Michael |
| 30 | Andy |
| 19 | Justin |
+-----+
```

Spark SQL, DataFrames and Datasets Guide

Python ☆

File Edit View Run Help Last edit was 9 minutes ago

New cell UI: OFF

Run all

Connect

Share

Publish



```
1 # Print the schema in a tree format
2 df.printSchema()
3
4 df.select("name").show()
5
6 # Select everybody, but increment the age by 1
7 df.select(df['name'], df['age'] + 1).show()
8
9
10 # Select people older than 21
11 df.filter(df['age'] > 21).show()
12
13 # Count people by age
14 df.groupBy("age").count().show()
```

```
root
|-- age: long (nullable = true)
|-- name: string (nullable = true)
```

```
+-----+
| name |
+-----+
| Michael |
| Andy |
| Justin |
+-----+
```



```

+-----+-----+
|  name|(age + 1)|
+-----+-----+
|Michael|    null|
|  Andy|    31|
| Justin|    20|
+-----+-----+

+-----+
|age|name|
+-----+
| 30|Andy|
+-----+

+-----+
| age|count|
+-----+
|  19|    1|
| null|    1|
|  30|    1|

```

Command took 3.08 seconds -- by a user at 10/02/2024, 15:29:50 on unknown compute

Spark SQL, DataFrames and Datasets Guide

Python

File Edit View Run Help [Last edit was 10 minutes ago](#) New cell UI: OFF

Run all

Connect

Share

Publish

```
| 30| 1|
```

Command took 3.08 seconds -- by a user at 10/02/2024, 15:29:50 on unknown compute

Cmd 4

```

1 # Register the DataFrame as a SQL temporary view
2 df.createOrReplaceTempView("people")
3
4 sqlDF = spark.sql("SELECT * FROM people")
5 sqlDF.show()

```

sqlDF: pyspark.sql.dataframe.DataFrame = [age: long, name: string]

```

+-----+-----+
| age|  name|
+-----+-----+
| null|Michael|
|  30|  Andy|
|  19| Justin|
+-----+-----+

```

Spark SQL, DataFrames and Datasets Guide

Python 

File Edit View Run Help

Last edit was 11 minutes ago

New cell UI: OFF 

 Run all

 Connect 

Share

Publish





```
1 # Register the DataFrame as a global temporary view
2 df.createGlobalTempView("people")
3
4 # Global temporary view is tied to a system preserved database `global_temp`
5 spark.sql("SELECT * FROM global_temp.people").show()
6
7
8 # Global temporary view is cross-session
9 spark.newSession().sql("SELECT * FROM global_temp.people").show()
```

```
+-----+
| age |  name|
+-----+
| null|Michael|
|  30|   Andy|
|  19|  Justin|
+-----+
```

```
+-----+
| age |  name|
+-----+
| null|Michael|
|  30|   Andy|
|  19|  Justin|
+-----+
```

