**Spark-SQL**

- Spark SQL is Spark module for structured data processing.
- It runs on top of Spark Core.
- It offers much tighter integration between relational and procedural processing, through declarative DataFrame and Datasets API.
- These are the ways which enable users to run SQL queries over Spark.
- Apache Spark SQL integrates relational processing with Sparks functional programming.
- Spark SQL blurs the line between RDD and relational table.

**SparkSession**

The entry point into all functionality in Spark is the SparkSession class. To create a basic SparkSession, just use SparkSession.builder():

- **Using Program we have to create a spark session object**

```
import org.apache.spark.sql.SparkSession
val spark = SparkSession
 .builder()
 .appName("Spark SQL basic example")
 .config("spark.some.config.option", "some-value")
 .getOrCreate()
```

- **Using Shell it is by default available as spark**

```
23/09/30 12:36:55 WARN NativeCodeLoader: Unable to l
tform... using builtin-java classes where applicable
Spark context Web UI available at http://192.168.28.
Spark context available as 'sc' (master = local[*],
Spark session available as 'spark'.
Welcome to


      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /___/ .__/\_,_/_/ /_/\_\   version 3.5.0
      /_/
```

**Creating DataFrames**

val df = spark.read.json("examples/src/main/resources/people.json")

- **Displays the content of the DataFrame to stdout**
  df.show()

```
scala> df1.show()
+----+-------+
| age|   name|
+----+-------+
|NULL|Michael|
|  30|   Andy|
|  19| Justin|
+----+-------+
```

- **DataFrame Operations**

  // This import is needed to use the $-notation
  import spark.implicits._

  // Print the schema in a tree format
  df.printSchema()

```
scala> df1.printSchema()
root
 |-- age: long (nullable = true)
 |-- name: string (nullable = true)
```

- **Select only the "name" column**
  df.select("name").show()

```
scala> df.select("name").show()
+-------+
|   name|
+-------+
|Michael|
|   Andy|
| Justin|
+-------+
```

- **Select everybody, but increment the age by 1**
  df.select($"name", $"age" + 1).show()

```
scala> df1.select($"age"+1).show()
+---------+
|(age + 1)|
+---------+
|     NULL|
|       31|
|       20|
+---------+
```

- **Count people by age**
  df.groupBy("age").count().show()

```
scala> df1.groupBy("age").count().show()
+----+-----+
| age|count|
+----+-----+
|  19|    1|
|NULL|    1|
|  30|    1|
+----+-----+
```

Running SQL Queries Programmatically

The sql function on a SparkSession enables applications to run SQL queries programmatically and returns the result as a DataFrame.

- **Register the DataFrame as a SQL temporary view**
  df.createOrReplaceTempView("people")

  val sqlDF1 = spark.sql("SELECT * FROM people")
  sqlDF1.show()

```
scala> val sqldf1=spark.sql("SELECT * FROM people")
sqldf1: org.apache.spark.sql.DataFrame = [age: bigint, name: string]

scala> sqldf1.show()
+----+-------+
| age|   name|
+----+-------+
|NULL|Michael|
|  30|   Andy|
|  19| Justin|
+----+-------+
```

**Global Temporary View**

Temporary views in Spark SQL are session-scoped and will disappear if the session that creates it terminates. If you want to have a temporary view that is shared among all sessions and keep alive until the Spark application terminates, you can create a global temporary view. Global temporary view is tied to a system preserved database global_temp, and we must use the qualified name to refer it, e.g. SELECT * FROM global_temp.view1.


**Register the DataFrame as a global temporary view**
df.createGlobalTempView("people")

**// Global temporary view is tied to a system preserved database `global_temp`**
spark.sql("SELECT * FROM global_temp.people").show()
// +----+-------+
// | age|   name|
// +----+-------+
// |null|Michael|
// |  30|   Andy|
// |  19| Justin|
// +----+-------+

**// Global temporary view is cross-session**
spark.newSession().sql("SELECT * FROM global_temp.people").show()
// +----+-------+
// | age|   name|
// +----+-------+
// |null|Michael|
// |  30|   Andy|
// |  19| Justin|

```
// +----+-------+
```

**Creating Datasets**

Datasets are similar to RDDs, however, instead of using Java serialization or Kryo they use a specialized Encoder to serialize the objects for processing or transmitting over the network. While both encoders and standard serialization are responsible for turning an object into bytes, encoders are code generated dynamically and use a format that allows Spark to perform many operations like filtering, sorting and hashing without deserializing the bytes back into an object.

```
case class Person(name: String, age: Long)

// Encoders are created for case classes
val caseClassDS = Seq(Person("Andy", 32)).toDS()
caseClassDS.show()
// +----+---+
// |name|age|
// +----+---+
// |Andy| 32|
// +----+---+

// Encoders for most common types are automatically provided by importing spark.implicits._
val primitiveDS = Seq(1, 2, 3).toDS()
primitiveDS.map(_ + 1).collect() // Returns: Array(2, 3, 4)

// DataFrames can be converted to a Dataset by providing a class. Mapping will be done by name
val path = "examples/src/main/resources/people.json"
val peopleDS = spark.read.json(path).as[Person]
peopleDS.show()
// +----+-------+
// | age|   name|
// +----+-------+
// |null|Michael|
// |  30|   Andy|
// |  19| Justin|
// +----+-------+
```

**Interoperating with RDDs**

park SQL supports two different methods for converting existing RDDs into Datasets. The first method uses reflection to infer the schema of an RDD that contains specific types of objects. This reflection-based approach leads to more concise code and works well when you already know the schema while writing your Spark application.

The second method for creating Datasets is through a programmatic interface that allows you to construct a schema and then apply it to an existing RDD. While this method is more verbose, it allows you to construct Datasets when the columns and their types are not known until runtime.

Inferring the Schema Using Reflection
/ For implicit conversions from RDDs to DataFrames
import spark.implicits._

```
// Create an RDD of Person objects from a text file, convert it to a Dataframe
val peopleDF = spark.sparkContext
  .textFile("examples/src/main/resources/people.txt")
  .map(_.split(","))
  .map(attributes => Person(attributes(0), attributes(1).trim.toInt))
  .toDF()
// Register the DataFrame as a temporary view
peopleDF.createOrReplaceTempView("people")

// SQL statements can be run by using the sql methods provided by Spark
val teenagersDF = spark.sql("SELECT name, age FROM people WHERE age BETWEEN 13 AND 19")

// The columns of a row in the result can be accessed by field index
teenagersDF.map(teenager => "Name: " + teenager(0)).show()
// +------------+
// |       value|
// +------------+
// |Name: Justin|
// +------------+

// or by field name
teenagersDF.map(teenager => "Name: " + teenager.getAs[String]("name")).show()
// +------------+
// |       value|
// +------------+
// |Name: Justin|
```

```scala
// +------------+

// No pre-defined encoders for Dataset[Map[K,V]], define explicitly
implicit val mapEncoder = org.apache.spark.sql.Encoders.kryo[Map[String, Any]]
// Primitive types and case classes can be also defined as
// implicit val stringIntMapEncoder: Encoder[Map[String, Any]] = ExpressionEncoder()

// row.getValuesMap[T] retrieves multiple columns at once into a Map[String, T]
teenagersDF.map(teenager => teenager.getValuesMap[Any](List("name", "age"))).collect()
// Array(Map("name" -> "Justin", "age" -> 19))
```

Programmatically Specifying the Schema

```scala
import org.apache.spark.sql.Row

import org.apache.spark.sql.types._

// Create an RDD
val peopleRDD = spark.sparkContext.textFile("examples/src/main/resources/people.txt")

// The schema is encoded in a string
val schemaString = "name age"

// Generate the schema based on the string of schema
val fields = schemaString.split(" ")
  .map(fieldName => StructField(fieldName, StringType, nullable = true))
val schema = StructType(fields)

// Convert records of the RDD (people) to Rows
val rowRDD = peopleRDD
  .map(_.split(","))
  .map(attributes => Row(attributes(0), attributes(1).trim))

// Apply the schema to the RDD
val peopleDF = spark.createDataFrame(rowRDD, schema)

// Creates a temporary view using the DataFrame
peopleDF.createOrReplaceTempView("people")

// SQL can be run over a temporary view created using DataFrames
```

```
val results = spark.sql("SELECT name FROM people")

// The results of SQL queries are DataFrames and support all the normal RDD operations
// The columns of a row in the result can be accessed by field index or by field name
results.map(attributes => "Name: " + attributes(0)).show()
// +-------------+
// |        value|
// +-------------+
// |Name: Michael|
// |   Name: Andy|
// | Name: Justin|
// +-------------+
```