

SPARK SQL

1. Inferring the Schema using Reflection

❖ create RDD → apply schema on it → convert it to Dataframe

a. using case class

1. Create case class as per the structure of data
2. Iterate on every line of rdd, split on the delimiter and apply the structure calling the case class (SchemaedRDD)
3. Convert the schemaedRDD to DF.
4. Ready to write DSL
5. register the DF to temp view.
6. Write ISO sql on top of the temp view created.

- Case classes are just regular classes that are: Immutable by default.
- Decomposable through pattern matching.
- Compared by structural equality instead of by reference.

```
case class Auction(auctionid: String, bid: Float, bidtime: Float, bidder: String, bidderrate: Integer, openbid: Float, price: Float, item: String, daystolive: Integer)
```

```
// load the data into an RDD
```

```
val auctionRDD = spark.sparkContext.textFile("file:///home/hduser/sparkdata/auctiondata")
```

```
// create an RDD of Auction objects
```

```
val ebay = auctionRDD.map(_._split("~")).map(p => Auction(p(0), p(1).toFloat, p(2).toFloat, p(3), p(4).toInt, p(5).toFloat, p(6).toFloat, p(7), p(8).toInt))
```

```
// change ebay RDD of Auction objects to a DataFrame
```

```
val auction = spark.createDataFrame(ebay)
```

and now the RDD ebay is converted into a Dataframe.

b. Using toDF() function

1. Iterate on every line of rdd, split on the delimiter and make it tuples
2. Convert the schemaedRDD to DF using toDF() function.
3. In the toDF() function pass the structure as parameter.
4. Ready to write DSL
5. register the DF to temp view.
6. Write ISO sql on top of the temp view created.

```
val auctionRDD = spark.sparkContext.textFile("file:///home/hduser/sparkdata/auctiondata")
```

//in the below command, we split the above RDD delimited by “~” and then instead of using case class, we convert the splitted RDD to data frame using toDF() function.

//While converting to dataframe using toDF() we pass the header of the data as parameter.

```
val ebay = auctionRDD.map(_._split("~")).map(p => (p(0), p(1).toFloat, p(2).toFloat, p(3),  
p(4).toInt, p(5).toFloat, p(6).toFloat, p(7), p(8).toInt)).toDF("auctionid", "bid", "bidtime",  
"bidder", "bidderrate", "openbid", "price", "item", "daystolive")
```

2. Create Dataframe programatically or using csv module

❖ Using spark.read dataframeReader function, read the file directly based on its format.

//We can use (spark.read) function to read the data directly based on format.

//But reading the data with out any information will not have much information.

//as show below, we don't have header /schema information

```
scala> val usdf = spark.read.csv("file:///home/hduser/sparkdata/usdata.csv")  
usdf: org.apache.spark.sql.DataFrame = [_c0: string, _c1: string ... 11 more fields]
```

```
scala> usdf.printSchema()  
root  
|-- _c0: string (nullable = true)  
|-- _c1: string (nullable = true)  
|-- _c2: string (nullable = true)  
|-- _c3: string (nullable = true)  
|-- _c4: string (nullable = true)  
|-- _c5: string (nullable = true)  
|-- _c6: string (nullable = true)  
|-- _c7: string (nullable = true)  
|-- _c8: string (nullable = true)  
|-- _c9: string (nullable = true)  
|-- _c10: string (nullable = true)  
|-- _c11: string (nullable = true)  
|-- _c12: string (nullable = true)
```

//instead of using just the module from read, if we use options() function we can make the data frame more informative.

```
val usschemadf = → this is using spark-csv module. We have spark modules for json,parquet etc  
spark.read.option("header","true").option("inferSchema","true").option("delimiter","").csv("file:///home/hduser/sparkdata/usdata.csv")
```

or → // above command used the direct “CSV” module.

But the below one mentions format as CSV and loads the data using load function.

```
val usmodule = → this is using load function and format as “csv”
```

```
spark.read.option("header","true").option("inferschema","true").option("delimiter","").format("csv").load("file:///home/hduser/sparkdata/usdata.csv")
```

```
scala> usschemaDF.printSchema() or scala> usmodule.printSchema()
```

```
root
|-- first_name: string (nullable = true)
|-- last_name: string (nullable = true)
|-- company_name: string (nullable = true)
|-- address: string (nullable = true)
|-- city: string (nullable = true)
|-- county: string (nullable = true)
|-- state: string (nullable = true)
|-- zip: integer (nullable = true)
|-- age: integer (nullable = true)
|-- phone1: string (nullable = true)
|-- phone2: string (nullable = true)
|-- email: string (nullable = true)
|-- web: string (nullable = true)
```

3. using STRUCT TYPE:

// StructType objects define the schema of Spark DataFrames. StructType objects contain a list of StructField objects that define the name, type, and nullable flag for each column in a DataFrame.

```
import org.apache.spark.sql.types.{StructType, StructField, StringType, IntegerType}
```

```
val custschema =
StructType(Array(StructField("first_name", StringType, true)
,StructField("last_name", StringType, true)
,StructField("company_name", StringType, true)
,StructField("address", StringType, true)
,StructField("city", StringType, true)
,StructField("country", StringType, true)
,StructField("state", StringType, true)
,StructField("zip", StringType, true)
,StructField("age", IntegerType, true)
,StructField("phone1", StringType, true)
,StructField("phone2", StringType, true)
,StructField("email", StringType, true)
,StructField("website", StringType, true))
);
```

```
val uscsvdf1 =
spark.read.option("delimiter","").schema(custschema).csv("file:///home/hduser/sparkdata/usdata.csv")
```

```
scala> uscsvdf1.printSchema
```

```

root
|-- first_name: string (nullable = true)
|-- last_name: string (nullable = true)
|-- company_name: string (nullable = true)
|-- address: string (nullable = true)
|-- city: string (nullable = true)
|-- country: string (nullable = true)
|-- state: string (nullable = true)
|-- zip: string (nullable = true)
|-- age: integer (nullable = true)
|-- phone1: string (nullable = true)
|-- phone2: string (nullable = true)
|-- email: string (nullable = true)
|-- website: string (nullable = true)

```

Now we have seen the different ways of creating data frames.. Now lets see how we can do DSL and SQL operations on it.

1. select & show

```

scala> val auctShow = auction.select("bidder").show()
+-----+
| bidder |
+-----+
| jake7870 |
| davidbresler2 |
| gladimacowgirl |
| daysrus |
| ..... |
| .... |

```

2. distinct & count

```

val count = auction.select("auctionid").distinct.count()

```

3. groupBy / agg / max / sum / alias / sort

```

import org.apache.spark.sql.functions.avg

```

```

val grpbid =
auction.groupBy("item").agg(max("bid").alias("max_bid"), sum("price").alias("sp")).sort($"sp".desc)

```

or

```

val grpbid =
auction.groupBy("item").agg(max("bid").alias("max_bid"), sum("price").alias("sp")).orderBy($"sp".desc)

```

4. Collect()

```

scala> grpbid.collect
res1: Array[org.apache.spark.sql.Row] =
Array([cartier,5400.0,1806618.5611763], [palm,290.0,1367597.7882232666], [xbox,501.77,401664.219581604])

```

5. collectAsList()

```
scala> grpbid.collectAsList
res9: java.util.List[org.apache.spark.sql.Row] =
[[cartier,5400.0,1806618.5611763], [palm,290.0,1367597.7882232666],
[xbox,501.77,401664.219581604]]
```

6. where()

```
scala> val grpbidcartier = auction.where('item ===
"cartier").groupBy("item").agg(sum("price").alias("sp"))
```

7. filter()

```
scala> val grpbidcartier = auction.filter('item ===
"cartier").groupBy("item").agg(sum("price").alias("sp"))
grpbidcartier: org.apache.spark.sql.DataFrame = [item: string, sp: double]
```

```
scala> grpbidcartier.collect
res10: Array[org.apache.spark.sql.Row] = Array([cartier,1806618.5611763])
```

8. printSchema

```
auction.printSchema
scala> auction.printSchema
root
|-- auctionid: string (nullable = true)
|-- bid: float (nullable = false)
|-- bidtime: float (nullable = false)
|-- bidder: string (nullable = true)
|-- bidderrate: integer (nullable = true)
|-- openbid: float (nullable = false)
|-- price: float (nullable = false)
|-- item: string (nullable = true)
|-- daystolive: integer (nullable = true)
```

```
grpbid.printSchema
root
|-- item: string (nullable = true)
|-- max_bid: float (nullable = true)
|-- sp: double (nullable = true)
```

9. Show()

```
scala> grpbid.show(1,false)
+-----+-----+-----+
|item   |max_bid|sp      |
+-----+-----+-----+
|cartier|5400.0 |1806618.5611763|
+-----+-----+-----+
only showing top 1 row
```

10. limit()

```
scala> auction.show(5,true)
or
scala> auction.limit(5).show()
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| auctionid| bid| bidtime| bidder|bidderrate|openbid|price|item|daystolive|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
-----+
```

```
|8213034705| 95.0|2.927373|      jake7870|          0|   95.0|117.5|xbox|
3|
|8213034705|115.0|2.943484| davidbresler2|          1|   95.0|117.5|xbox|
3|
|8213034705|100.0|2.951285|gladimacowgirl|         58|   95.0|117.5|xbox|
3|
|8213034705|117.5|2.998947|      daysrus|         10|   95.0|117.5|xbox|
3|
|8213060420|   2.0|0.065266|      donnie4814|          5|    1.0|120.0|xbox|
3|
+-----+-----+-----+-----+-----+-----+-----+-----+
-----+
only showing top 5 rows
```

11. head() / take()

```
scala> auction.head(3)
or
scala> auction.take(3)
res16: Array[org.apache.spark.sql.Row] =
Array([8213034705,95.0,2.927373,jake7870,0,95.0,117.5,xbox,3],
[8213034705,115.0,2.943484,davidbresler2,1,95.0,117.5,xbox,3],
[8213034705,100.0,2.951285,gladimacowgirl,58,95.0,117.5,xbox,3])
```

12. takeAsList()

```
scala> auction.takeAsList(3)
res21: java.util.List[org.apache.spark.sql.Row] =
[[8213034705,95.0,2.927373,jake7870,0,95.0,117.5,xbox,3],
[8213034705,115.0,2.943484,davidbresler2,1,95.0,117.5,xbox,3],
[8213034705,100.0,2.951285,gladimacowgirl,58,95.0,117.5,xbox,3]]
```

13. Join() → copy/paste the below table format output into notepad to understand better.

- Here how it works is,
- First row from left table will be joined with all rows in second table.
- If first table has 2 rows and second has 4, then output will have 8 rows.
- 1st rows from table 1 mapped with all 4 rows in second table.
- Then 2nd row from table 1 again mapped with all 4 rows in second table.

```
scala> val auct1 = auction.limit(2)
auct1: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] =
[auctionid: string, bid: float ... 7 more fields]
```

```
scala> auct1.show()
```

```
+-----+-----+-----+-----+-----+-----+-----+-----+
| auctionid|  bid| bidtime|      bidder|bidderrate|openbid|price|item|daystolive|
+-----+-----+-----+-----+-----+-----+-----+-----+
|8213034705| 95.0|2.927373|      jake7870|          0|   95.0|117.5|xbox|          3|
|8213034705|115.0|2.943484|davidbresler2|          1|   95.0|117.5|xbox|          3|
+-----+-----+-----+-----+-----+-----+-----+-----+
```

```
scala> val auct2 = auction.limit(4)
auct2: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] =
[auctionid: string, bid: float ... 7 more fields]
```

```
scala> auct2.show()
```

auctionid	bid	bidtime	bidder	bidderrate	openbid	price	item	daystolive
8213034705	95.0	2.927373	jake7870	0	95.0	117.5	xbox	3
8213034705	115.0	2.943484	davidbresler2	1	95.0	117.5	xbox	3
8213034705	100.0	2.951285	gladimacowgirl	58	95.0	117.5	xbox	3
8213034705	117.5	2.998947	daysrus	10	95.0	117.5	xbox	3

```
scala> val auctjoin = auct1.join(auct2)
auctjoin: org.apache.spark.sql.DataFrame = [auctionid: string, bid: float
... 16 more fields]
```

```
scala> auctjoin.show()
```

auctionid	bid	bidtime	bidder	bidderrate	openbid	price	item	daystolive
8213034705	95.0	2.927373	jake7870	0	95.0	117.5	xbox	3
8213034705	95.0	2.927373	jake7870	0	95.0	117.5	xbox	3
8213034705	95.0	2.927373	jake7870	0	95.0	117.5	xbox	3
8213034705	115.0	2.943484	davidbresler2	1	95.0	117.5	xbox	3
8213034705	95.0	2.927373	jake7870	0	95.0	117.5	xbox	3
8213034705	100.0	2.951285	gladimacowgirl	58	95.0	117.5	xbox	3
8213034705	95.0	2.927373	jake7870	0	95.0	117.5	xbox	3
8213034705	117.5	2.998947	daysrus	10	95.0	117.5	xbox	3
8213034705	115.0	2.943484	davidbresler2	1	95.0	117.5	xbox	3
8213034705	95.0	2.927373	jake7870	0	95.0	117.5	xbox	3
8213034705	115.0	2.943484	davidbresler2	1	95.0	117.5	xbox	3
8213034705	115.0	2.943484	davidbresler2	1	95.0	117.5	xbox	3
8213034705	115.0	2.943484	davidbresler2	1	95.0	117.5	xbox	3
8213034705	100.0	2.951285	gladimacowgirl	58	95.0	117.5	xbox	3
8213034705	115.0	2.943484	davidbresler2	1	95.0	117.5	xbox	3
8213034705	117.5	2.998947	daysrus	10	95.0	117.5	xbox	3

14. UnionAll: (just horizontal merging of data)

```
scala> val auct1 = auction.limit(2)
auct1: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] =
[auctionid: string, bid: float ... 7 more fields]
```

```
scala> val auct2 = auction.limit(4)
auct2: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] =
[auctionid: string, bid: float ... 7 more fields]
```

```
scala> val auctunion = auct1.unionAll(auct2)
warning: there was one deprecation warning; re-run with -deprecation for
details
auctunion: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] =
[auctionid: string, bid: float ... 7 more fields]
```

```
scala> auctunion.show()
```

auctionid	bid	bidtime	bidder	bidderrate	openbid	price	item	daystolive
8213034705	95.0	2.927373	jake7870	0	95.0	117.5	xbox	3
8213034705	115.0	2.943484	davidbresler2	1	95.0	117.5	xbox	3
8213034705	95.0	2.927373	jake7870	0	95.0	117.5	xbox	3
8213034705	115.0	2.943484	davidbresler2	1	95.0	117.5	xbox	3
8213034705	100.0	2.951285	gladimacowgirl	58	95.0	117.5	xbox	3
8213034705	117.5	2.998947	daysrus	10	95.0	117.5	xbox	3

15.System.out.println()

```
System.out.println(count)
```

16.//Stores the output in JSON format

```
auction.write.mode("overwrite").json("file:/home/hduser/sparkdata/auctiondata.json");  
auctunion.write.mode("overwrite").json("file:/home/hduser/sparkdata/auction.json");
```

Now lets see how to perform SQL operation on the dataframe.

- ❖ To write a SQL query on a data frame, we need to create a tempview using that dataframe.
- ❖ createOrReplaceTempView() function is used to create temp view on top of a DF.

```
val usdataDF =  
spark.read.option("header","true").option("inferSchema","true").option("delimiter",",").csv("file:///home/hduser/sparkdata/usdata.csv")
```

```
scala> usdataDF.createOrReplaceTempView("custinfo")
```

```
scala> spark.catalog.listDatabases.show(10,false);
```

name	description	locationUri
custdb	null	
hdfs://localhost:54310/user/hive/warehouse/custdb.db		
default	Default Hive database	hdfs://localhost:54310/user/hive/warehouse
kjpractice	null	
hdfs://localhost:54310/user/hive/warehouse/kjpractice.db		
retail	null	
hdfs://localhost:54310/user/hive/warehouse/retail.db		
retail_stg	null	
hdfs://localhost:54310/user/hive/warehouse/retail_stg.db		

```
scala> spark.catalog.listTables.show(10,false);
```


name	database	description	tableType	isTemporary
game	default	null	MANAGED	false
custinfo	null	null	TEMPORARY	true

- The meta data of SQL objects of dataframe are derby data base by default
- But in real time it will be stored in some remote data base.
- In our example its stored in MySQL. As show in the above highlighted .

```
spark.catalog.listTables.show(10,false);
```

```
sqlctx.sql("describe custinfo").show(10,false)
```

```
sql("select concat(first_name,' ',last_name),company_name,address,city,phone1 from custinfo
where upper(company_name) like 'CH%'").show(10,false)
```

```
sql("select count(1),city from custinfo group by city having count(1) >1 order by city").show()
```

```
sql("select count(1) as cnt,city from custinfo group by city having count(1)> 5 order by
city").show()
```

```
val agecat = sqlctx.sql("select distinct age,case when age <=10 then 'childrens' when age >10
and age < 20 then 'teen' else 'others' end as agecat from custinfo order by agecat")
```

```
scala> agecat.show()
```

age	agecat
9	childrens
10	childrens
7	childrens
8	childrens
22	others
20	others
34	others
55	others
43	others
25	others
33	others
21	others
11	teen
16	teen
19	teen
12	teen
15	teen
14	teen

//Creating and registering functions:

```
def addfunc (a:Int,b:Int):Int=
{
return a+b
```

```
}
```

```
import org.apache.spark.sql.functions.udf
```

```
spark.udf.register("addudf",addfunc _) spark.catalog.listFunctions.filter('name like  
"%addu%").show(false)
```

```
spark.sql("SELECT addudf(20,30) FROM custinfo ").show()
```

//Stores the output in Parquet and orc format, Parquet files are self-describing so the schema is preserved

```
agecat.write.mode("overwrite").parquet("file:/home/hduser/sparkdata/agecategory.parquet");
```

```
agecat.write.mode("append").orc("file:/home/hduser/sparkdata/agecategory1.orc");
```

// Read in the parquet file created above

// Parquet files are self-describing so the schema is preserved

// The result of loading a Parquet file is also a DataFrame

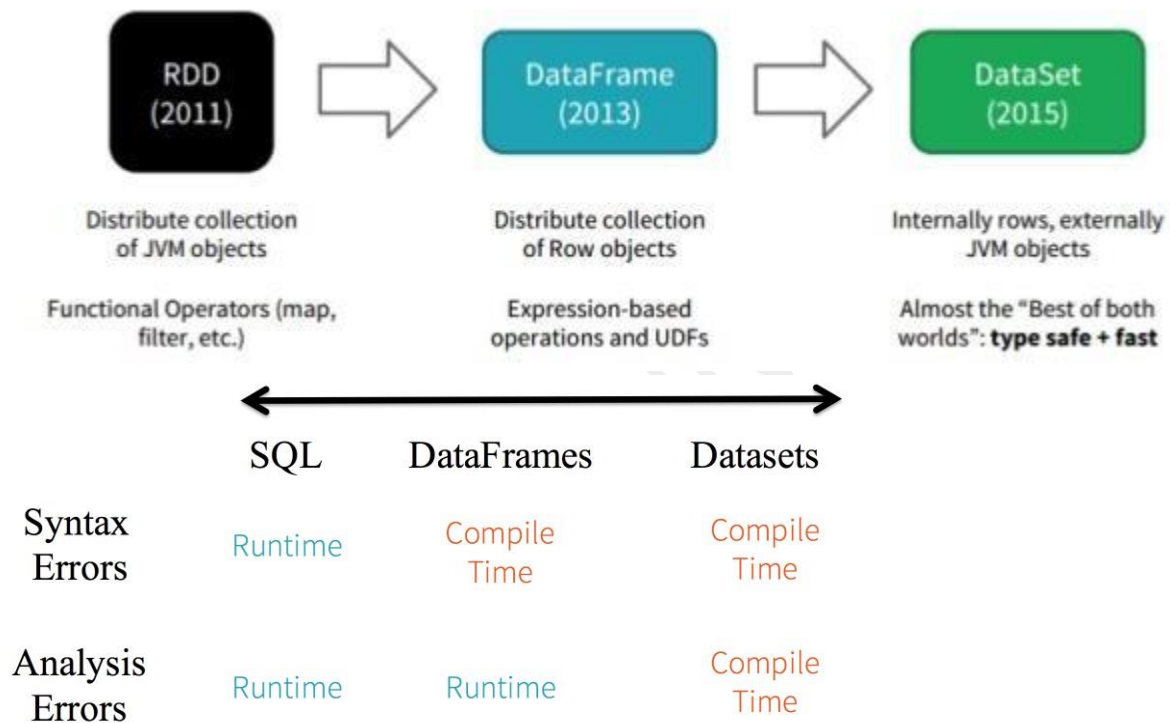
```
val parquetagecatdf =  
spark.read.parquet("file:/home/hduser/sparkdata/agecategory.parquet")
```

// Parquet files can also be used to create a temporary view and then used in SQL statements

```
parquetagecatdf.createOrReplaceTempView("parquetagecat")
```

```
spark.sql("SELECT max(age),min(age),count(distinct agecat) FROM parquetagecat ").show()
```

Dataset Vs Dataframe



A Dataset is a strongly typed collection of domain-specific objects (like object of type auction declared below) that can be transformed in parallel using functional or relational operations. Each Dataset also has an untyped view called a DataFrame, which is a Dataset of Row.

Dataframe is merged with Dataset API. So we can use any method available for dataframe in datasets.

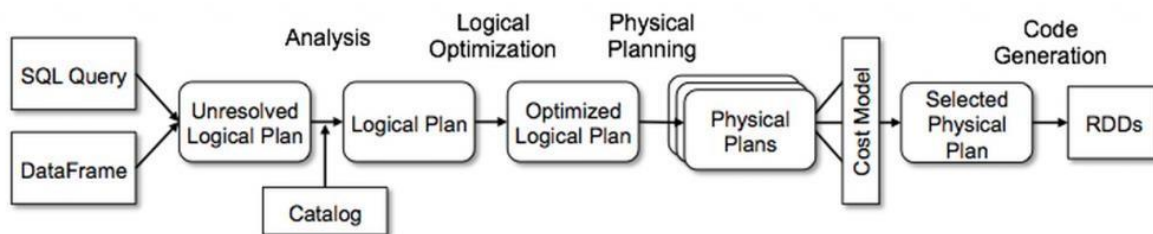
In summation, the choice of when to use RDD or DataFrame and/or Dataset seems obvious. While the former offers you low-level functionality and control, the latter allows custom view and structure, offers high-level and domain specific operations, saves space, and executes at superior speeds.

To simplify Spark for developers, how to optimize and make it performant it was decided to elevate the low-level RDD APIs to a high-level abstraction as DataFrame and Dataset and to build this unified data abstraction across libraries a top Catalyst optimizer and Tungsten.

The goal of Project Tungsten is to improve Spark execution by optimizing Spark jobs for CPU and memory efficiency (as opposed to network and disk I/O which are considered fast enough).

- Tungsten includes specialized in-memory data structures tuned for the type of operations required by Spark, improved code generation, and a specialized wire protocol.
- As Tungsten does not depend on Java objects, both on-heap and off-heap allocations are supported. Since Tungsten no longer depends on working with Java objects, you can use either on-heap (in the JVM) or off-heap storage.
- Tungsten will not do deserialization when processing data, for example of this is with sorting, a common and expensive operation using tungsten this can be done without having to deserialize the data again.
- By avoiding the memory and GC overhead of regular Java objects, Tungsten is able to process larger data sets than the same hand-written aggregations.

Catalyst Optimizer:



Catalyst supports both rule-based and cost-based optimization.

At its core, Catalyst contains a general library for representing trees and applying rules to manipulate them. On top of this framework, we have built libraries specific to relational query processing (e.g., expressions, logical query plans), and several sets of rules that handle different phases of query execution: analysis, logical optimization, physical planning, and code generation to compile parts of queries to Java bytecode. For the latter, we use another Scala feature, quasiquotes, that makes it easy to generate code at runtime from composable expressions. Finally, Catalyst offers several public extension points, including external data sources and user-defined types.

Pick one—DataFrames and/or Dataset or RDDs APIs—that meets your needs and use-case

// A JSON dataset is pointed to a path. The path can be either a single text file or a directory storing text files

//weakly typed Dataframes

```
hadoop fs -put ~/sparkdata/auctiondata.json
```

```
val auctionjson = "/user/hduser/auctiondata.json";  
val auctionjsonDF = spark.read.json(auctionjson);  
auctionjsonDF.distinct.show
```

//Strongly typed DataSets

```
case class auctionclass  
(auctionid:String,bid:Double,bidder:Int,bidderrate:Long,bidtime:Double,  
daystolive:Long,item:String,openbid:Double,price:Double)  
val auctionjson = "file:/home/hduser/sparkdata/auctiondata.json";  
val auctionjsonDS = spark.read.json(auctionjson).as[auctionclass];  
case class auctionclass  
(auctionid:String,bid:Double,bidder:String,bidderrate:Long,bidtime:Double,  
daystolive:Long,item:String,openbid:Double,price:Double)  
val auctionjsonDS = spark.read.json(auctionjson).as[auctionclass];  
// The inferred schema can be visualized using the printSchema() method  
auctionjsonDF.printSchema();  
auctionjsonDS.printSchema();
```

// Creates a temporary view using the DataSet

```
auctionjsonDS.createOrReplaceTempView("auctionjsontable");
```

```
val auctionquery = spark.sql("SELECT * FROM auctionjsontable")
auctionquery.show();
```

Hive operations:

```
import spark.implicits._
import spark.sql
```

```
//Initialize hive context wrapping spark context
```

```
val hiveContext = new org.apache.spark.sql.hive.HiveContext(sc)
```

```
//Get a hiveContext context
```

```
val sparkSession =
SparkSession.builder.enableHiveSupport.getOrCreate()
```

```
//Create a hive table
```

```
sql("create database if not exists sparkdb")
sql("use sparkdb")
```

```
sql("DROP TABLE IF EXISTS employee ")
```

```
sql("CREATE TABLE IF NOT EXISTS employee(id INT, name STRING, age INT) ROW FORMAT
DELIMITED FIELDS TERMINATED BY ',' LINES TERMINATED BY '\n'")
```

```
spark.catalog.listDatabases.show(10,false);
```

```
spark.catalog.listTables.show(10,false);
```

```
//Load data
```

```
sql("LOAD DATA LOCAL INPATH '/home/hduser/sparkdata/employee.txt' INTO TABLE employee")
```

```
//View data
```

```
val results = sql("FROM employee SELECT id, name, age")  
results.show()
```

Additional Hive Usecases (For self practices):

```
//Drop table
```

```
sql("CREATE TABLE IF NOT EXISTS src (key INT, value STRING) ROW FORMAT DELIMITED  
FIELDS TERMINATED BY ',' LINES TERMINATED BY '\n'")
```

```
sql("LOAD DATA LOCAL INPATH '/home/hduser/sparkdata/sampleddata' OVERWRITE INTO  
TABLE src")
```

```
sql("SELECT * FROM src").show()  
sql("drop table src")
```

```
//Write the table content into the given hdfs location in orc format
```

```
results.write.format("orc").save("hdfs:/user/hduser/emp_orc")
```

```
//Read the table from the given hdfs location in orc format and create a dataset out of it.
```

```
case class empclass (id :Int,name :String,age:Int);  
val empdata_orc =  
spark.read.format("orc").load("hdfs:/user/hduser/emp_orc").as[empclass];
```

```
empdata_orc.createOrReplaceTempView("orcdata")  
sql("SELECT * from orcdata").collect.foreach(println)
```

Hive Partitioning:

```
sql("""create table txnrecords(txnno INT, txndate STRING, custno INT, amount DOUBLE,
category STRING, product STRING, city STRING, state STRING, spendby STRING) row format
delimited fields terminated by ','
lines terminated by '\n'
stored as textfile""")
```

```
sql("LOAD DATA LOCAL INPATH '/home/hduser/hive/data/txns' OVERWRITE INTO TABLE
txnrecords")
```

```
sql("select * from txnrecords order by 1 limit 10").show
```

```
sql("""create external table exttxnrecsByCat(txnno INT, txndate STRING, custno INT, amount
DOUBLE, product STRING, city STRING, state STRING, spendby STRING)
partitioned by (category STRING) row format delimited
fields terminated by ','
stored as textfile
location '/user/hive/warehouse/exttxnrecsbycat'""")
```

```
sql("""Insert into table exttxnrecsbycat partition (category='Games')
select txnno,txndate,custno,amount,product,city,state,spendby
from txnrecords
where category='Games'""")
```

```
sql("select count(1) from exttxnrecsbycat").show
```

Adding UDFS in hive

```
spark-shell --jars /home/hduser/hive/replaceword.jar
sql("use sparkdb")
sql("create function repwords as 'tech.replacewords.replaceword'")
sql("select repwords(product,'Archery','Arc') from txnrecords where product='Archery' limit
10")
```

Supported Hive Features

Spark SQL supports the vast majority of Hive features, such as:

- Hive query statements, including:
 - SELECT
 - GROUP BY
 - ORDER BY
 - CLUSTER BY
 - SORT BY
- All Hive operators, including:
 - Relational operators (=, <=>, ==, <>, <, >, >=, <=, etc)
 - Arithmetic operators (+, -, *, /, %, etc)
 - Logical operators (AND, &&, OR, ||, etc)
 - Complex type constructors
 - Mathematical functions (sign, ln, cos, etc)
 - String functions (instr, length, printf, etc)
- User defined functions (UDF)
- User defined aggregation functions (UDAF)
- User defined serialization formats (SerDes)
- Window functions
- Joins
- Unions
- Sub-queries
 - SELECT col FROM (SELECT a + b AS col from t1) t2
- Sampling
- Explain
- Partitioned tables including dynamic partition insertion
- View
- All Hive DDL Functions, including:
 - CREATE TABLE
 - CREATE TABLE AS SELECT
 - ALTER TABLE
- Most Hive Data types, including:
 - TINYINT
 - SMALLINT
 - INT
 - BIGINT
 - BOOLEAN
 - FLOAT
 - DOUBLE
 - STRING
 - BINARY
 - TIMESTAMP
 - DATE
 - ARRAY<>
 - MAP<>
 - STRUCT<>

Unsupported Hive Functionality

Below is a list of Hive features that we don't support yet. Most of these features are rarely used in Hive deployments.

Major Hive Features

- Tables with buckets: bucket is the hash partitioning within a Hive table partition. Spark SQL doesn't support buckets yet.

Hive Input/Output Formats

- File format for CLI: For results showing back to the CLI, Spark SQL only supports `TextOutputFormat`.
- Hadoop archive

Hive Optimizations

A handful of Hive optimizations are not yet included in Spark. Some of these (such as indexes) are less important due to Spark SQL's in-memory computational model. Others are slotted for future releases of Spark SQL.

- Block level bitmap indexes and virtual columns (used to build indexes)
- Automatically determine the number of reducers for joins and groupbys: Currently in Spark SQL, you need to control the degree of parallelism post-shuffle using `"SET spark.sql.shuffle.partitions=[num_tasks];"`.
- Meta-data only query: For queries that can be answered by using only meta data, Spark SQL still launches tasks to compute the result.
- `STREAMTABLE` hint in join: Spark SQL does not follow the `STREAMTABLE` hint.
- Merge multiple small files for query results: if the result output contains multiple small files, Hive can optionally merge the small files into fewer large files to avoid overflowing the HDFS metadata. Spark SQL does not support that.
-