≡  BAAHU'S BLOG

# SPARK DataFrame – Java Example

👤 baahu   📅 October 2, 2016   💬 No Comments

Tweet

A `DataFrame` is a collection of data, organized into named columns. `DataFrames` are similar to tables in a traditional database `DataFrame` can be constructed from sources such as Hive tables, Structured Data files, external databases, or existing RDDs.

Under the hood, a `DataFrame` contains an RDD composed of `Row` objects with additional schema information of the types in each col.

In the Java example code below we are retrieving the details of the employee who draws the max salary(i.e get the name of the CEO 😉 )

We are going to create a `DataFrame` over a text file, every line of this file contains employee information in the below format EmployeeID,Name,Salary

1. We start with creating `SQLContext` .SQLContext is used for initializing the functionalities of Spark SQL. `SparkContext` class object is required for initializing SQLContext class object.

2. We store the schema of the table in a String(empid,name,salary) and create a ArrayList of `StructField` objects.

   I am avoiding using a for loop to create `StructField` objects since the types of columns are heterogeneous – empid,name are `String` and salary is `Integer` .

3. Create a `StructType` object using the `StructField` objects created in the above step.

4. Load the file into an RDD and create a `JavaRDD` of type `Row` for each of the entries present in the input file.

5. Create a `DataFrame` object using the `JavaRDD` in the above step and the `StructType` object created in step 3.

6. Register a temporary table "employee" , so that we can execute `sql` style queries using a `SQLContext` .

7. Execute the sql query to find details of the employee who draws the max salary and print the same using by invoking `foreach()` action.

```Java
2  import java.util.List;
3  import org.apache.spark.SparkConf;
4  import org.apache.spark.api.java.JavaRDD;
5  import org.apache.spark.api.java.JavaSparkContext;
6  import org.apache.spark.api.java.function.Function;
7  import org.apache.spark.api.java.function.VoidFunction;
8  import org.apache.spark.sql.DataFrame;
```

```java
 9  import org.apache.spark.sql.Row;
10  import org.apache.spark.sql.RowFactory;
11  import org.apache.spark.sql.SQLContext;
12  import org.apache.spark.sql.types.DataTypes;
13  import org.apache.spark.sql.types.StructType;
14  import org.apache.spark.sql.types.StructField;
15
16
17
18  public class DataFrameExample {
19
20      public static void main(String[] args) {
21
22          SparkConf conf = new SparkConf().setAppName("DataFrameExample").setMaster("local[*]");
23          JavaSparkContext jsp = new  JavaSparkContext(conf);
24
25          //Create the sql context
26          SQLContext sqlContext =  new SQLContext(jsp);
27
28          String schemaString = "name,empid,salary";
29
30          List<StructField> fieldList = new ArrayList<StructField>();
31
32          //Create 3 objects for 3 columns, please note that salary is an integer
33          fieldList.add(DataTypes.createStructField("name",DataTypes.StringType,true));
34          fieldList.add(DataTypes.createStructField("empid",DataTypes.StringType,true));
35          fieldList.add(DataTypes.createStructField("salary",DataTypes.IntegerType,true));
36
37          //Now create a StructType object from the list
38          StructType schema = DataTypes.createStructType(fieldList);
39
40          //Load the file which contains the data into an RDD
41          JavaRDD<String> rdd = jsp.textFile("file:///home/baahu/employee");
42
43          //convert the input to a list of row objects
44          JavaRDD<Row> rowsRDD = rdd.map(new Function<String,Row>(){
45
46              @Override
47              public Row call(String arg0) throws Exception {
48                  //arg0 contains each line from the rdd , i.e empid,name,salary
49                  String [] data = arg0.split(",");
50
51                  //We make sure that we convert the salary field to integer
52                  return RowFactory.create(data[0],data[1],Integer.valueOf(data[2]));
53              }});
54
55          //Create the dataframe from the rowsRDD and schema
56          DataFrame dataFrame = sqlContext.createDataFrame(rowsRDD, schema);
57          dataFrame.registerTempTable("employee");
58
59          //Run the sql query to find details of employee who draws the max salary
60          DataFrame maxSalDataFrame = sqlContext.sql("Select * from employee order by salary desc limit 1 ");
61
62          maxSalDataFrame.toJavaRDD().foreach(new VoidFunction<Row>(){
63
64              @Override
65              public void call(Row arg0) throws Exception {
66                  System.out.println(" Name is :"+arg0.get(0)+" Emp id is:"+arg0.get(1)+" Salary is:"+arg0.get(2));
67              }});
68      }
```

```
69  }
```

**Sample Input**

```
1  John,54001,10000
2  Johhny,54002,10000000
3  Janardhan,54003,5000
```

**Sample Output**

```
1  Name is :Johhny Emp id is:54002 Salary is:10000000
```

Tweet

📁 SPARK     🏷️ DataFrame, spark, SQLContext