# UMKC

## UNIVERSITY OF MISSOURI-KANSAS CITY

**Bigdata Hadoop Programming**

**# Lab 3 Assignment**

**Team Members:**

Pragathi Thammaneni

Sridevi Mallipudi

**Introduction:**

The main core concept for executing the Lab 3 Assignment is to implement the MapReduce Algorithm for finding Facebook common friends problem and run the MapReduce job on Apache Spark. And also, to implement the Spark Data frames and including intuitive queries like pattern recognition, topic discussion, most important terms etc.

**Objectives:**

To code for the 2 questions the below concepts are implemented.

Map reduce Algorithm for Facebook common friend's problem

Apache Spark

Spark Data frames

Datasets -queries

**Approaches /Methods:**

Apache Spark

**Workflow &Datasets/Parameters and Evaluation:**

The below each question will follow different approaches to solve. Coding is done to perform the evaluation of each individual snippet to execute the datasets which are provided as the input parameters.
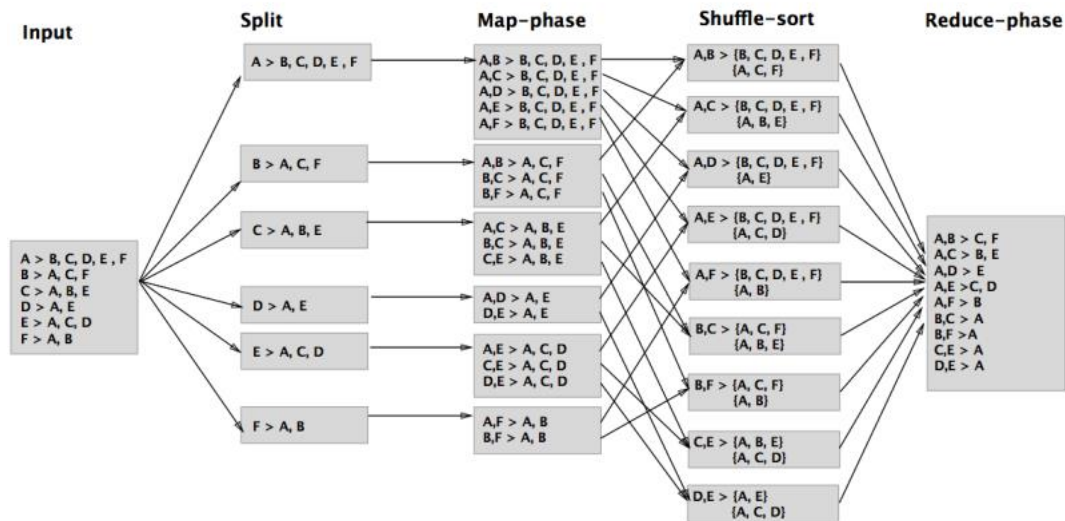
**Question 1:**

## 1. Hadoop MapReduce Algorithm

a) Implement MapReduce algorithm for finding Facebook common friends problem and run the MapReduce job on Apache Spark.
b) *Use the following dataset to run your program:*
   https://umkc.box.com/s/y6juor0fwe96f6louboy3mvbfpli6pgt
c) Write a report including your algorithm and result screenshots.

For solving the below question a map reduce algorithm is used to find out the face book common friend problem. Below is the sample use case for implementing the Map reduce.

**Use Case Diagram:**

| Input | Split | Map-phase | Shuffle-sort | Reduce-phase |
|---|---|---|---|---|

Input:
```
A > B, C, D, E , F
B > A, C, F
C > A, B, E
D > A, E
E > A, C, D
F > A, B
```

Split:
```
A > B, C, D, E , F
B > A, C, F
C > A, B, E
D > A, E
E > A, C, D
F > A, B
```

Map-phase:
```
A,B > B, C, D, E , F
A,C > B, C, D, E , F
A,D > B, C, D, E , F
A,E > B, C, D, E , F
A,F > B, C, D, E , F

A,B > A, C, F
B,C > A, C, F
B,F > A, C, F

A,C > A, B, E
B,C > A, B, E
C,E > A, B, E

A,D > A, E
D,E > A, E

A,E > A, C, D
C,E > A, C, D
D,E > A, C, D

A,F > A, B
B,F > A, B
```

Shuffle-sort:
```
A,B > {B, C, D, E , F}
      {A, C, F}

A,C > {B, C, D, E , F}
      {A, B, E}

A,D > {B, C, D, E , F}
      {A, E}

A,E > {B, C, D, E , F}
      {A, C, D}

A,F > {B, C, D, E , F}
      {A, B}

B,C > {A, C, F}
      {A, B, E}

B,F > {A, C, F}
      {A, B}

C,E > {A, B, E}
      {A, C, D}

D,E > {A, E}
      {A, C, D}
```

Reduce-phase:
```
A,B > C, F
A,C > B, E
A,D > E
A,E >C, D
A,F > B
B,C > A
B,F >A
C,E > A
D,E > A
```

**Below is the screen shot for the Map Reduce Algorithm**



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

object Friends {
  def main(args: Array[String]): Unit = {

    System.setProperty("hadoop.home.dir", "F:\\winutils")

    val sc = SparkSession
      .builder
      .appName(name = "SparkWordCount")
      .master(master = "local[*]")
      .getOrCreate().sparkContext

    val flist = sc.textFile(path = "input.txt")
    val common = flist.flatMap { x =>
      val splitlist = x.split(regex = " : ")
      val owner = splitlist(0)
      val friendslist = splitflist(1).split(regex = " ")
      friendslist.foreach(println)
      val makelist = friendslist.slice(0, friendslist.size).map(y => {
        if (owner > y) (y, owner) else (owner, y)
      })
      makelist.map(z => (z, friendslist.slice(0, friendslist.size).toSet))
    }
    val findcommon = common.reduceByKey((x, y) => x intersect y).sortByKey()
    findcommon.collect().take(10).foreach(x => {
      println(s"${x._1} -> (${x._2.mkString(" ")})")
      findcommon.saveAsTextFile(path = "output")
    })
```

**Output Screenshot:**



## Question 2:

### 2. Spark Data Frames

**Datasets:**

1. **FIFA World Cup:**
   https://www.kaggle.com/abecklas/fifa-world-cup#WorldCupMatches.csv
2. **Kickstarter Projects**
   https://www.kaggle.com/kemical/kickstarter-projects
3. **Google-Landmarks Dataset**
   https://www.kaggle.com/google/google-landmarks-dataset

   a. Create a Spark DataFrame using one of datasets, trying to use all different StructType.
   b. Perform 10 intuitive questions in Dataset (e.g.: pattern recognition, topic discussion, most important terms, etc.). Use your innovation to think out of box.
   c. Perform any 5 queries in Spark RDD's and Spark Data Frames. Compare the results

**Queries implemented:**

Below are the queries that has been implemented on the data sets FIFA World Cup.

1. Creation of Spark Data Frame using FIFA World Cup.

2. Queries regarding pattern recognition, topic discussion and most important terms has been performed.

3. Performed queries using Spark RDD's and Spark Data Frames.

Please refer the source code folder for the queries description

**Struct Type:**



**10 intuitive queries:**

**Query 1:Count of teams that reached different levels(stages) of the game**

## Query 2:Filtering the Home teams that scored goals >=3 and <=10



```
# Filtering the Home teams that scored goals >=3 and <=10
df.createOrReplaceTempView("table1")
Goals = spark.sql("SELECT Stage,Stadium,City,Home_Team_Name FROM table1 WHERE Home_Team_Goals >= 3 AND Home_Team_Goals <= 10")
Goals.show()
```

```
+----------------+-------------------+-----------+--------------+
|           Stage|            Stadium|       City|Home_Team_Name|
+----------------+-------------------+-----------+--------------+
|         Group 1|            Pocitos|Montevideo |        France|
|         Group 4|     Parque Central|Montevideo |           USA|
|         Group 3|            Pocitos|Montevideo |       Romania|
|         Group 1|     Parque Central|Montevideo |         Chile|
|         Group 2|     Parque Central|Montevideo |    Yugoslavia|
|         Group 4|     Parque Central|Montevideo |           USA|
|         Group 1|  Estadio Centenario|Montevideo |     Argentina|
|         Group 2|  Estadio Centenario|Montevideo |        Brazil|
|         Group 3|  Estadio Centenario|Montevideo |       Uruguay|
|         Group 1|  Estadio Centenario|Montevideo |     Argentina|
|     Semi-finals|  Estadio Centenario|Montevideo |     Argentina|
|     Semi-finals|  Estadio Centenario|Montevideo |       Uruguay|
|           Final|  Estadio Centenario|Montevideo |       Uruguay|
|Preliminary round|Stadio Benito Mus...|     Turin |       Austria|
|Preliminary round|  Giorgio Ascarelli|     Naples |       Hungary|
|Preliminary round|           San Siro|      Milan |   Switzerland|
|Preliminary round|           Littorale|   Bologna |        Sweden|
```

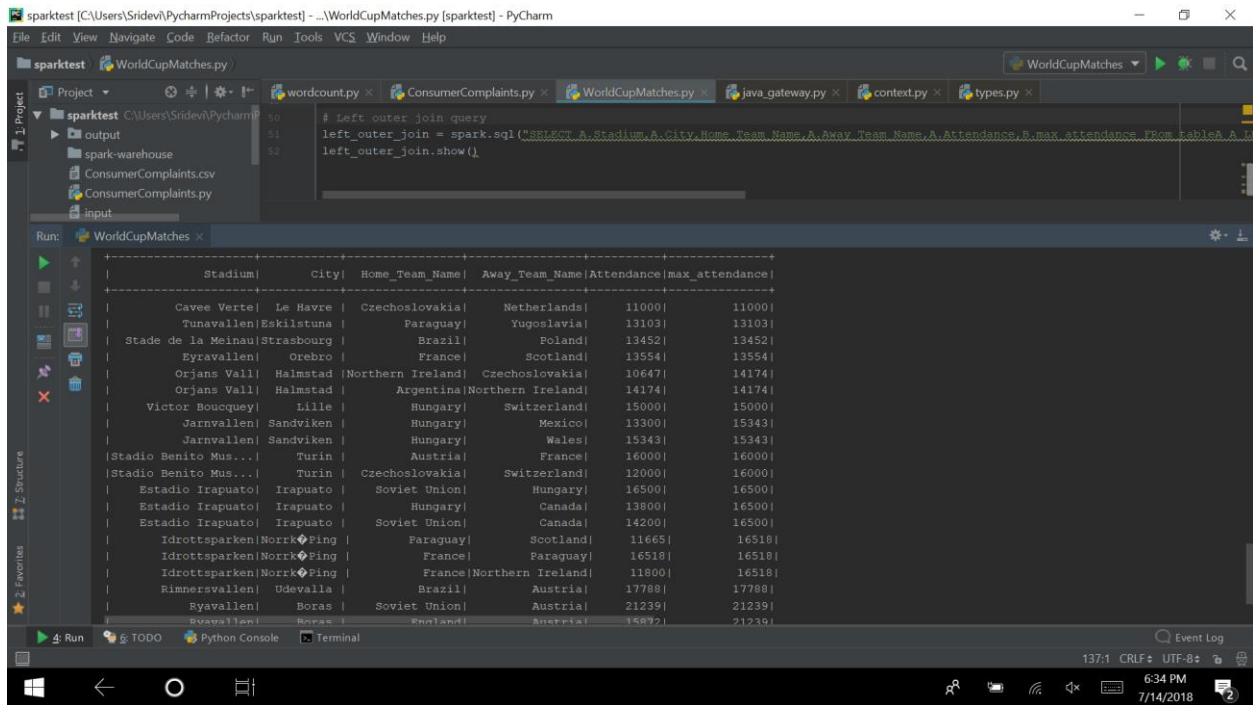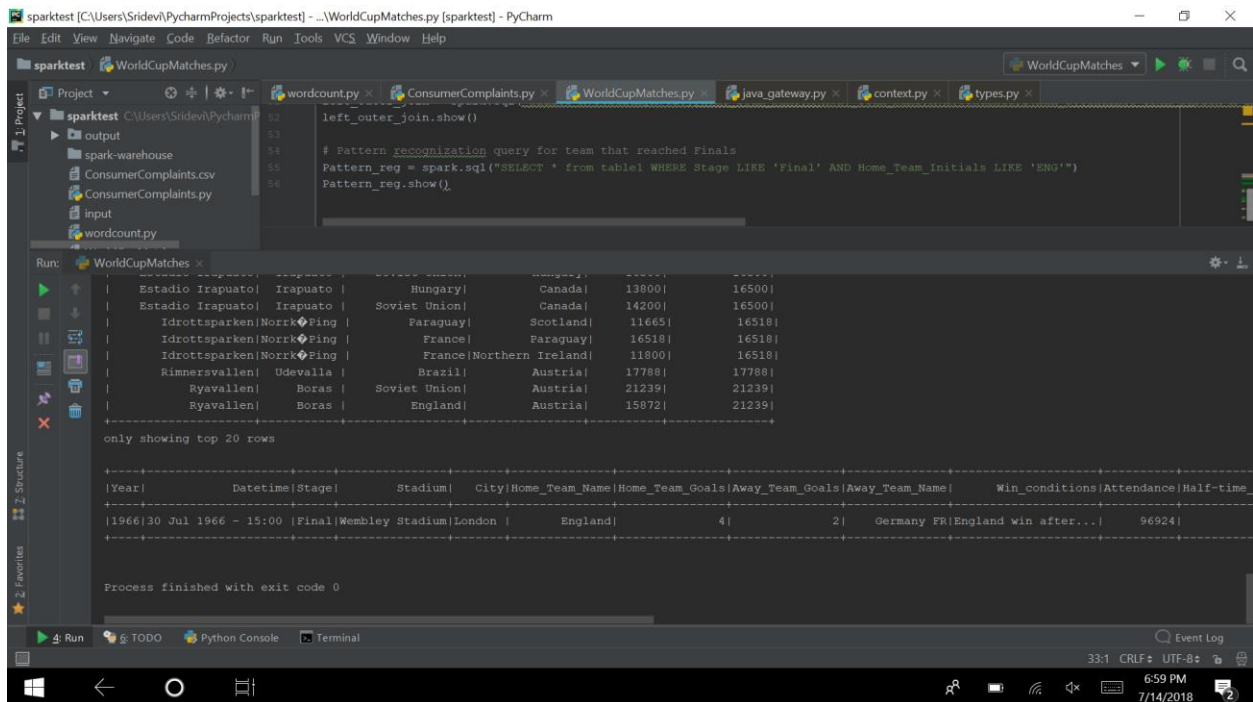## Query 3: Correlated subquery for finding attendance and max_attendance for a stadium



```
# correlated subquery for finding attendance and max_attendance for a stadium
df.createOrReplaceTempView("tableA")
df.createOrReplaceTempView("tableB")
Max_Attendance = spark.sql("SELECT A.Stage,A.Stadium,A.City,A.Home_Team_Name,A.Away_Team_Name,A.Attendance,(SELECT MAX(Attendance)
Max_Attendance.show()
```

```
+----------------+-------------------+-----------+--------------+---------------+----------+---------------+
|           Stage|            Stadium|       City| Home_Team_Name| Away_Team_Name|Attendance|max_attendance|
+----------------+-------------------+-----------+--------------+---------------+----------+---------------+
|     First round|        Cavee Verte| Le Havre | Czechoslovakia|    Netherlands|     11000|          11000|
|         Group 2|   Tunavallen|Eskilstuna |     Paraguay|     Yugoslavia|     13103|          13103|
|     First round|Stade de la Meinau|Strasbourg |       Brazil|         Poland|     13452|          13452|
|         Group 2|          Eyravallen|   Orebro |      France|       Scotland|     13554|          13554|
|         Group 1|         Orjans Vall| Halmstad |Northern Ireland| Czechoslovakia|     10647|          14174|
|         Group 1|         Orjans Vall| Halmstad |     Argentina|Northern Ireland|     14174|          14174|
|   Quarter-finals|     Victor Boucquey|    Lille |      Hungary|    Switzerland|     15000|          15000|
|         Group 3|          Jarnvallen| Sandviken |      Hungary|         Mexico|     13300|          15343|
|         Group 3|          Jarnvallen| Sandviken |      Hungary|          Wales|     15343|          15343|
|Preliminary round|Stadio Benito Mus...|     Turin |      Austria|         France|     16000|          16000|
|   Quarter-finals|Stadio Benito Mus...|     Turin | Czechoslovakia|    Switzerland|     12000|          16000|
|         Group C|     Estadio Irapuato| Irapuato |  Soviet Union|        Hungary|     16500|          16500|
|         Group C|     Estadio Irapuato| Irapuato |      Hungary|         Canada|     13800|          16500|
|         Group C|     Estadio Irapuato| Irapuato |  Soviet Union|         Canada|     14200|          16500|
|         Group 2|     Idrottsparken|Norrk◆Ping |     Paraguay|       Scotland|     11665|          16518|
|         Group 2|     Idrottsparken|Norrk◆Ping |      France|       Paraguay|     16518|          16518|
|   Quarter-finals|     Idrottsparken|Norrk◆Ping |      France|Northern Ireland|     11800|          16518|
|         Group 4|     Rimnersvallen| Udevalla |       Brazil|        Austria|     17788|          17788|
```

## Query 4: Left outer join query

## Query 5: Pattern recognization query for team that reached Finals



## Query 6: Average goals scored by a team

**Query 7: Count of rows after performing UNION ALL operation**



**Query 8: Number of times a country scores goals greater than or equal to 4**

**Query 9: Number of maximum number of goals scored by a country between years 2000 nd 2010**

## Query 10: Total number of goals scored by a country



## Spark RDD's: Apply the schema to the RDD



## Query 1: Filtering the Home teams that scored goals >=3 and <=10 with RDD

**Query 2: Correlated subquery for finding attendance and max_attendance for a stadium with RDD**



**Query 3: Pattern recognition query for team that reached Finals with RDD**

## Query 4: Average goals scored by a team with RDD



## Query 5: Number of times a country scores goals greater than or equal to 4 with RDD

**Comparison between Spark RDD's and the data frames:**

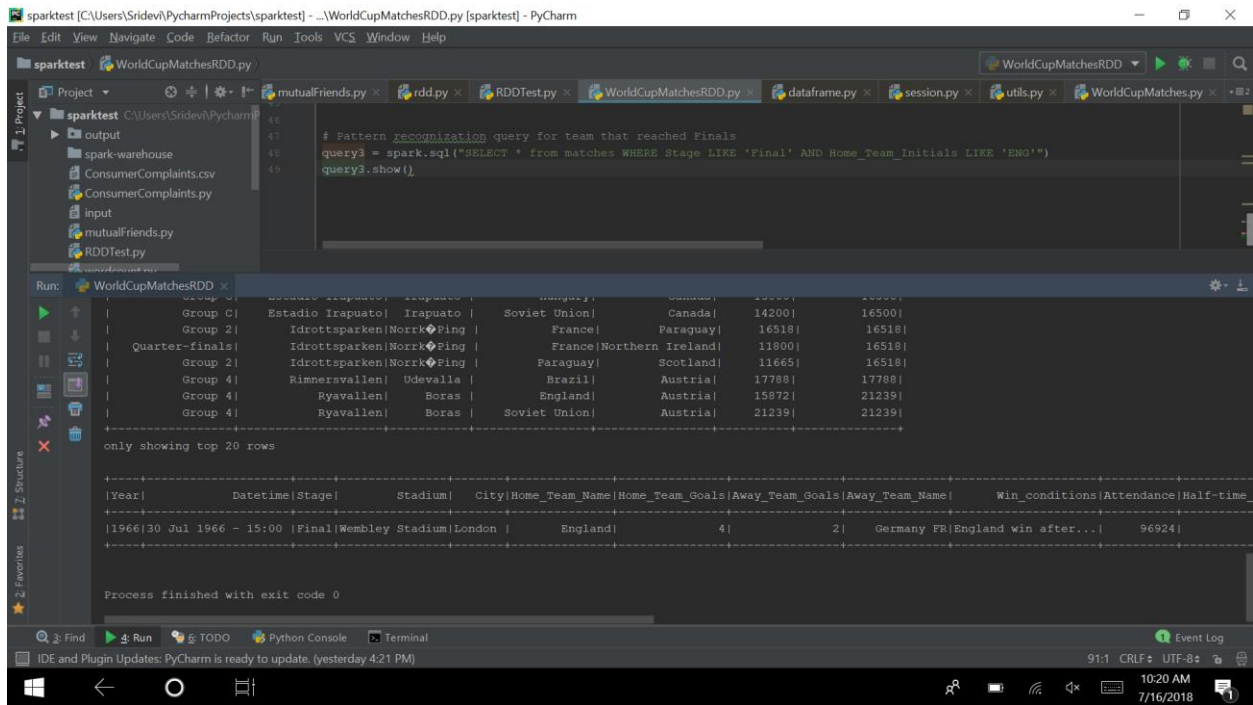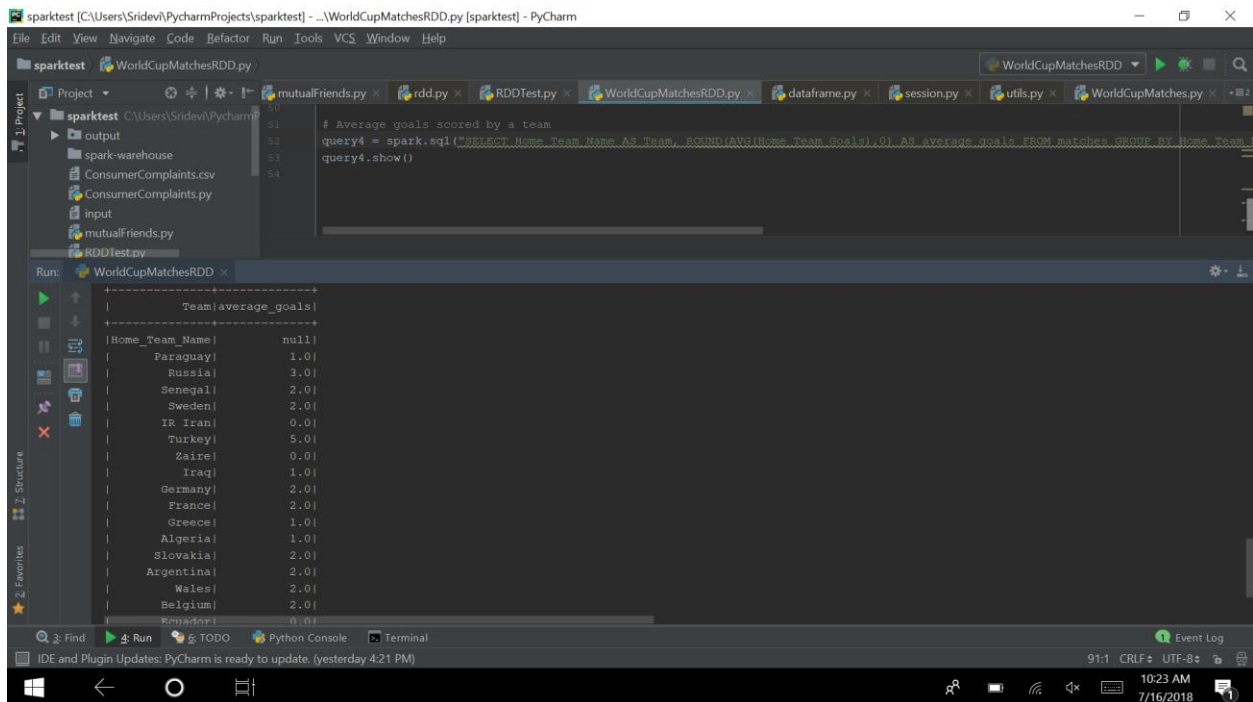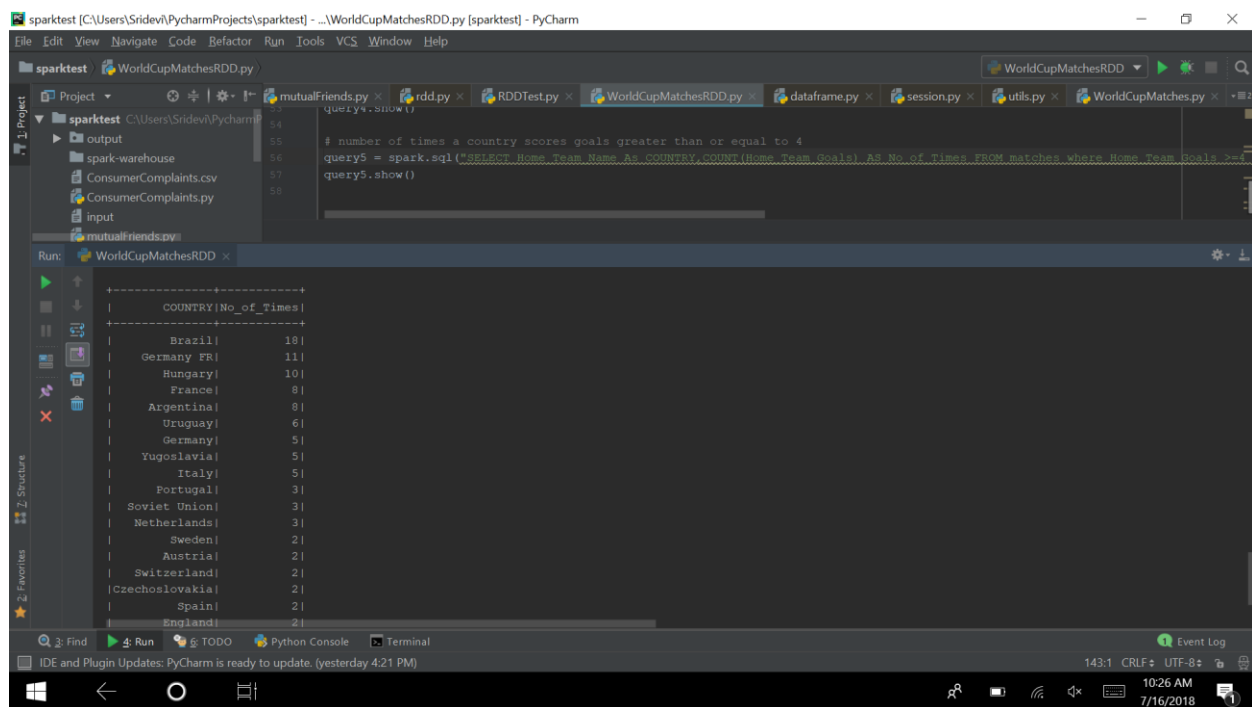| Spark RDD's | Spark Data Frames |
|---|---|
| RDD is a fault-tolerant collection of elements that can be operated on in-parallel, also we can say RDD is the fundamental data structure of Spark | It is a distributed collection of data. Basically, data is organized into named columns in data frames. |
| Basically, it is read-only partition collection of records. Moreover, it supports in-memory computations on large clusters in a fault-tolerant manner | Spark also introduced **catalyst optimizer**, along with data frame. To build an extensible query optimizer, it also leverages advanced programming features. |
| An RDD can come easily handle data with no predefined structure. | In Spark, data frame allows developers to impose a structure onto a distributed data. It also allows higher-level abstraction. |
| Compile- Time Type Safety RDD- RDD Supports *object-oriented programming* style with compile-time type safety. | If we try to access any column which is not present in the table, then an attribute error may occur at runtime. Data frame will not support compile-time type safety in such case. |
| Spark does not compute their result right away, it evaluates RDDs lazily | Computation happens only when action appears as Spark evaluates data frame lazily |

RDDs of Apache spark offers low-level functionality and control. Whereas datasets offer higher functionality. While data frame offers high-level domain-specific operations, saves space and executes at high speed. Therefore, it increases the efficiency of the system**.**

**Conclusion:**

As stated, the above workflow with certain set of parameters is followed in solving the execution by implementing the core and basic concepts Apache Spark.

**Source code link:**

https://github.com/PragathiThammaneni/Bigdata-Programming--Hadoop-Spark/tree/master/Lab%203

**Video Link: Provided in wiki** https://youtu.be/31C5DO0zQSE

**Wiki Link:**
https://github.com/PragathiThammaneni/Bigdata-Programming--Hadoop-Spark/wiki/Lab-3-Assignment