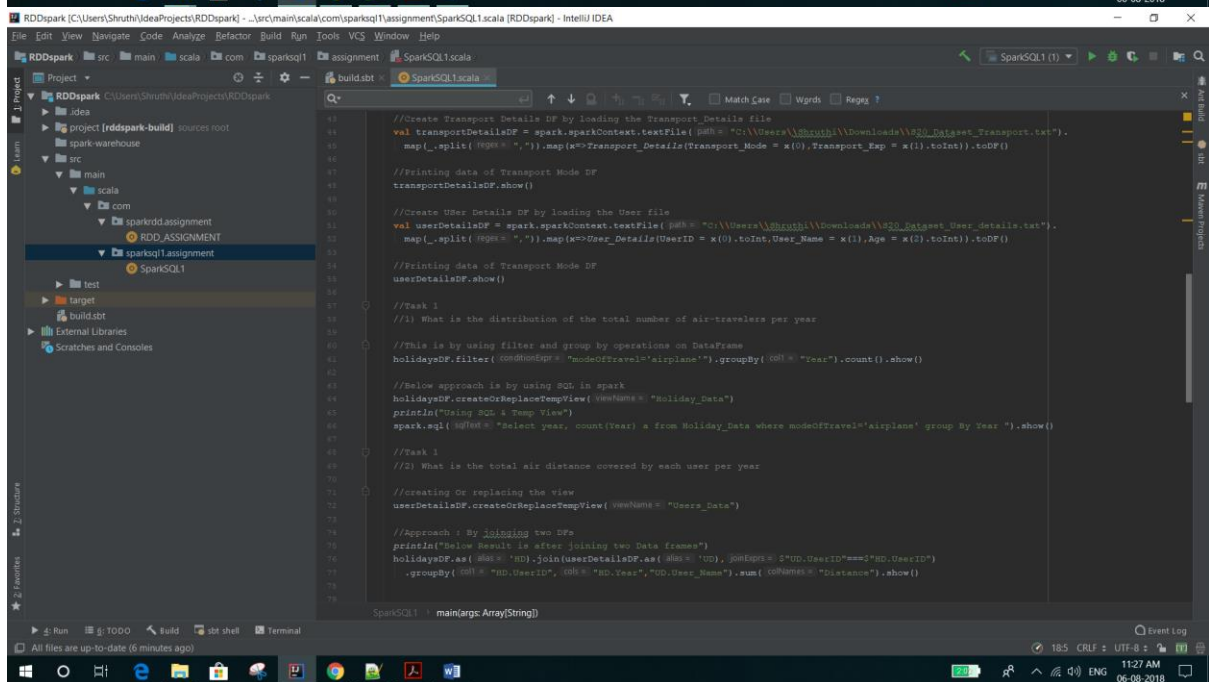
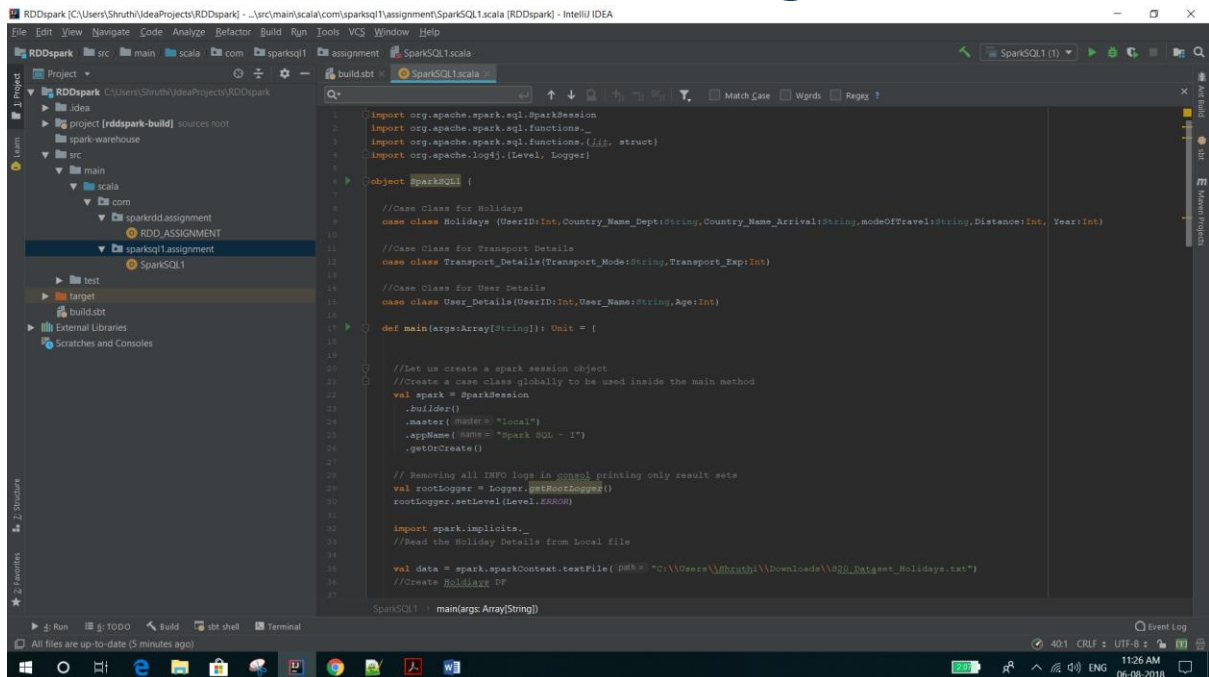


SPARK SQL 1



SPARK SQL 1

```
Run: SparkSQL (1) x
18/08/06 11:21:25 INFO Utils: Successfully started service 'SparkUI' on port 4040.
18/08/06 11:21:25 INFO SparkUI: Bound SparkUI to 0.0.0.0, and started at http://192.168.56.1:4040
18/08/06 11:21:25 INFO Executor: Starting executor ID driver on host localhost
18/08/06 11:21:25 INFO Utils: Successfully started service 'org.apache.spark.network.netty.NettyBlockTransferService' on port 62772.
18/08/06 11:21:25 INFO NettyBlockTransferService: Server created on 192.168.56.1:62772
18/08/06 11:21:25 INFO BlockManager: Using org.apache.spark.storage.RandomBlockReplicationPolicy for block replication policy
18/08/06 11:21:25 INFO BlockManagerMaster: Registering BlockManager BlockManagerId(driver, 192.168.56.1, 62772, None)
18/08/06 11:21:25 INFO BlockManagerMasterEndpoint: Registering block manager 192.168.56.1:62772 with 859.7 MB RAM, BlockManagerId(driver, 192.168.56.1, 62772, None)
18/08/06 11:21:25 INFO BlockManager: Registered BlockManager BlockManagerId(driver, 192.168.56.1, 62772, None)
18/08/06 11:21:25 INFO BlockManager: Initialized BlockManager: BlockManagerId(driver, 192.168.56.1, 62772, None)

-----+-----+
|UserID|Country_Name_Dept|Country_Name_Arrival|modeOfTravel|Distance|Year|
-----+-----+
| 1|CHN|IND|airplane|200|1990|
| 2|IND|CHN|airplane|200|1991|
| 3|CHN|airplane|200|1992|
| 4|RUS|IND|airplane|200|1990|
| 5|CHN|RUS|airplane|200|1992|
| 6|AUS|PAR|airplane|200|1991|
| 7|RUS|AUS|airplane|200|1990|
| 8|IND|RUS|airplane|200|1991|
| 9|RUS|airplane|200|1992|
| 10|AUS|CHN|airplane|200|1993|
| 1|AUS|CHN|airplane|200|1993|
| 2|CHN|IND|airplane|200|1993|
| 3|CHN|IND|airplane|200|1993|
| 4|IND|AUS|airplane|200|1991|
| 5|AUS|IND|airplane|200|1992|
| 6|RUS|CHN|airplane|200|1993|
| 7|CHN|RUS|airplane|200|1990|
| 8|AUS|CHN|airplane|200|1990|
| 9|IND|AUS|airplane|200|1991|
| 10|RUS|CHN|airplane|200|1992|
-----+-----+
only showing top 20 rows

-----+-----+
|Transport_Mode|Transport_Exp|
-----+-----+
|airplane|170|
|car|140|
|train|120|
|ship|200|
-----+-----+

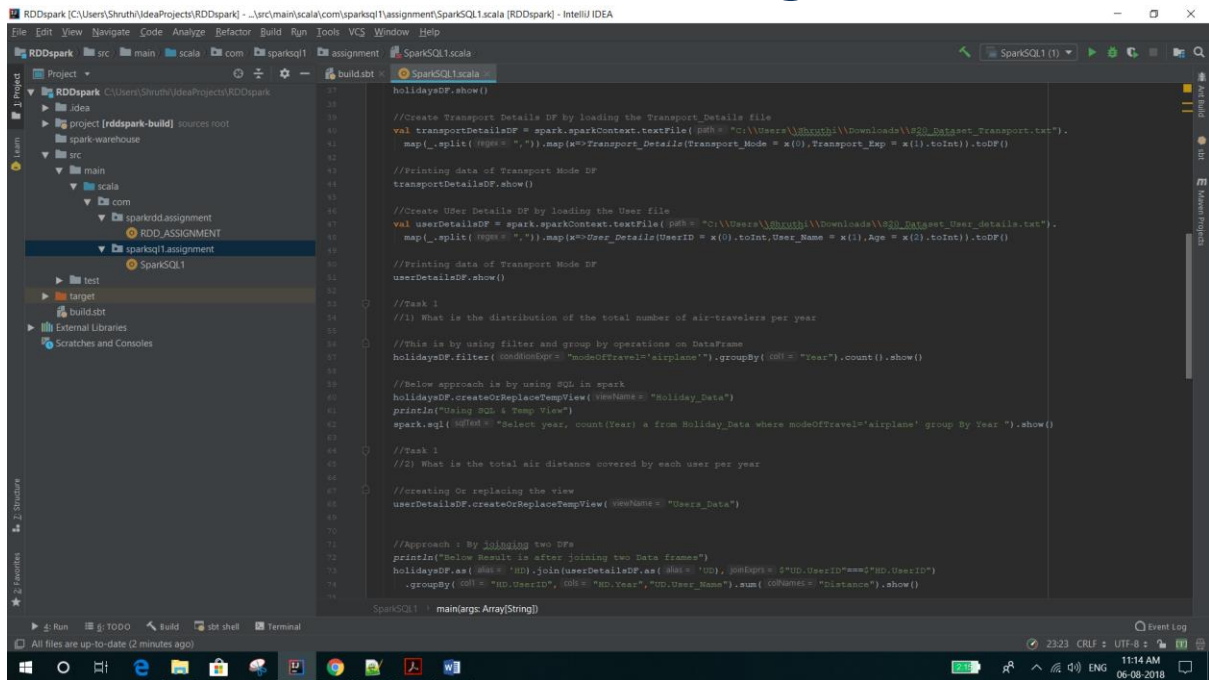
-----+-----+
|UserID|User_Name|Age|
-----+-----+
| 1|mark|15|
| 2|john|16|
| 3|luke|17|
| 4|lisa|27|
| 5|mark|25|
| 6|peter|22|
| 7|james|21|
| 8|andrew|50|
| 9|thomas|46|
| 10|annie|44|
-----+-----+
```

```
Run: SparkSQL (1) x
| 9|CHN|RUS|airplane|200|1992|
| 10|AUS|CHN|airplane|200|1993|
| 1|AUS|CHN|airplane|200|1993|
| 2|CHN|IND|airplane|200|1993|
| 3|CHN|IND|airplane|200|1993|
| 4|IND|AUS|airplane|200|1991|
| 5|AUS|IND|airplane|200|1992|
| 6|RUS|CHN|airplane|200|1993|
| 7|CHN|RUS|airplane|200|1990|
| 8|AUS|CHN|airplane|200|1990|
| 9|IND|AUS|airplane|200|1991|
| 10|RUS|CHN|airplane|200|1992|
-----+-----+
only showing top 20 rows

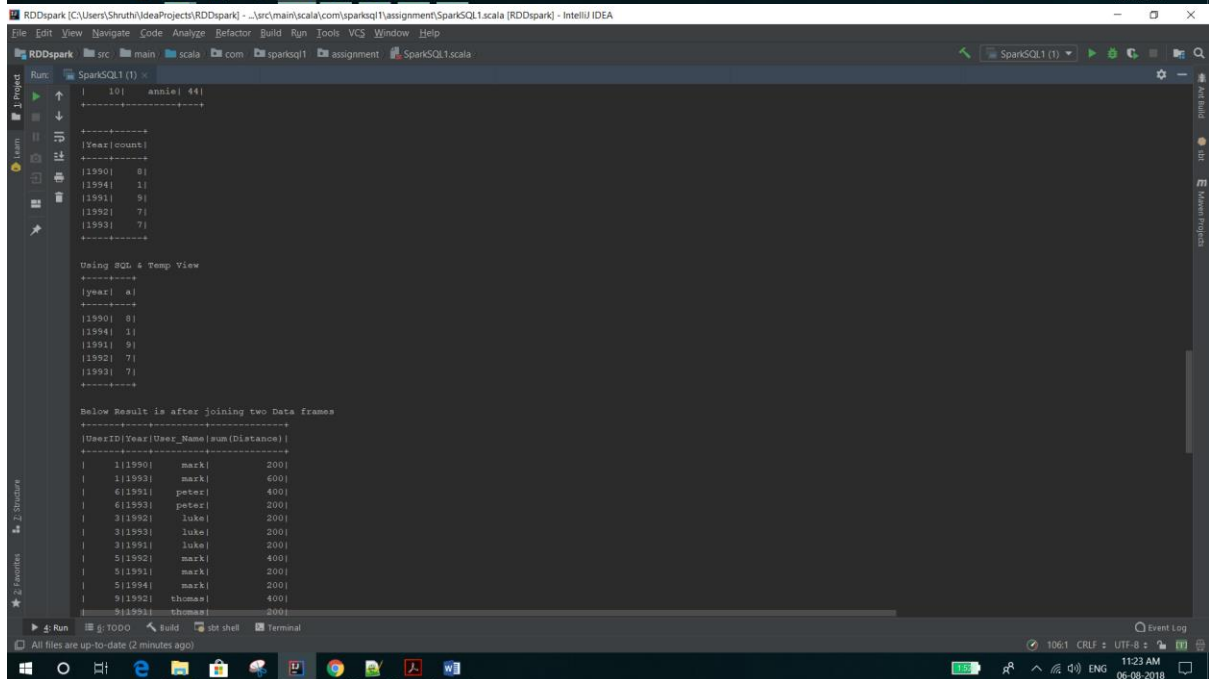
-----+-----+
|Transport_Mode|Transport_Exp|
-----+-----+
|airplane|170|
|car|140|
|train|120|
|ship|200|
-----+-----+

-----+-----+
|UserID|User_Name|Age|
-----+-----+
| 1|mark|15|
| 2|john|16|
| 3|luke|17|
| 4|lisa|27|
| 5|mark|25|
| 6|peter|22|
| 7|james|21|
| 8|andrew|50|
| 9|thomas|46|
| 10|annie|44|
-----+-----+
```

SPARK SQL 1



```
37 holidaysDF.show()
38
39 //Create Transport Details DF by loading the Transport_Details file
40 val transportDetailsDF = spark.sparkContext.textFile(URI("C:\\Users\\Shruthi\\Downloads\\SP2_Dataset_Transport.txt")).
41   map(_.split(",")).map(x=>Transport_Details(Transport_Mode = x(0),Transport_Exp = x(1).toInt)).toDF()
42
43 //Printing data of Transport Mode DF
44 transportDetailsDF.show()
45
46 //Create User Details DF by loading the User file
47 val userDetailsDF = spark.sparkContext.textFile(URI("C:\\Users\\Shruthi\\Downloads\\SP2_Dataset_User_details.txt")).
48   map(_.split(",")).map(x=>User_Details(userID = x(0).toInt,User_Name = x(1),Age = x(2).toInt)).toDF()
49
50 //Printing data of Transport Mode DF
51 userDetailsDF.show()
52
53 //Task 1
54 //1) What is the distribution of the total number of air-travelers per year
55
56 //This is by using filter and group by operations on DataFrame
57 holidaysDF.filter( (colName: String) => "modeOfTravel"=="airplane").groupBy( (col) => "Year").count().show()
58
59 //Below approach is by using SQL in spark
60 holidaysDF.createOrReplaceTempView( (viewName: String) => "Holiday_Data")
61 println("Using SQL & Temp View")
62 spark.sql(URI("SELECT year, count(Year) a from Holiday_Data where modeOfTravel='airplane' group By Year")).show()
63
64 //Task 1
65 //2) What is the total air distance covered by each user per year
66
67 //creating Or replacing the view
68 userDetailsDF.createOrReplaceTempView( (viewName: String) => "Users_Data")
69
70 //Approach : By joining two DFE
71 println("Below Result is after joining two Data frames")
72 holidaysDF.as( (viewName: String) => "HD").join(userDetailsDF.as( (viewName: String) => "UD"), (colName: String) => s"$UD.UserID"===s"$HD.UserID")
73   .groupBy( (colName: String) => "HD_Year", (colName: String) => "UD_User_Name").sum( (colName: String) => "Distance").show()
74
75 SparkSQL1 main(args: Array[String])
```

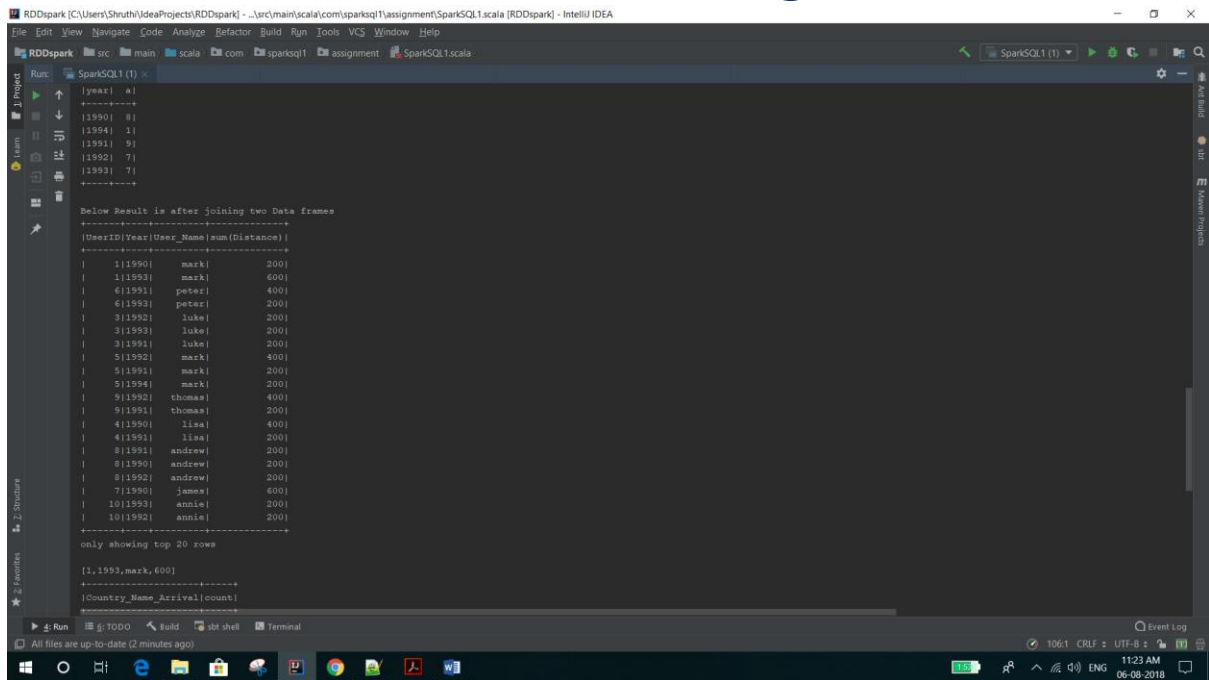


```
Run: SparkSQL1 (1) x
+-----+
|Year|count|
+-----+
|1990| 81|
|1994| 11|
|1991| 91|
|1992| 71|
|1993| 71|
+-----+

Using SQL & Temp View
+-----+
|year| a|
+-----+
|1990| 81|
|1994| 11|
|1991| 91|
|1992| 71|
|1993| 71|
+-----+

Below Result is after joining two Data frames
+-----+
|UserID|Year|User_Name|sum(Distance)|
+-----+
|1|1990| mark| 200|
|1|1993| mark| 600|
|4|1991| peter| 400|
|4|1993| peter| 200|
|3|1992| luke| 200|
|3|1993| luke| 200|
|3|1991| luke| 200|
|5|1992| mark| 400|
|3|1991| mark| 200|
|5|1994| mark| 200|
|9|1992| thomas| 400|
|5|1991| thomas| 200|
+-----+
```

SPARK SQL 1

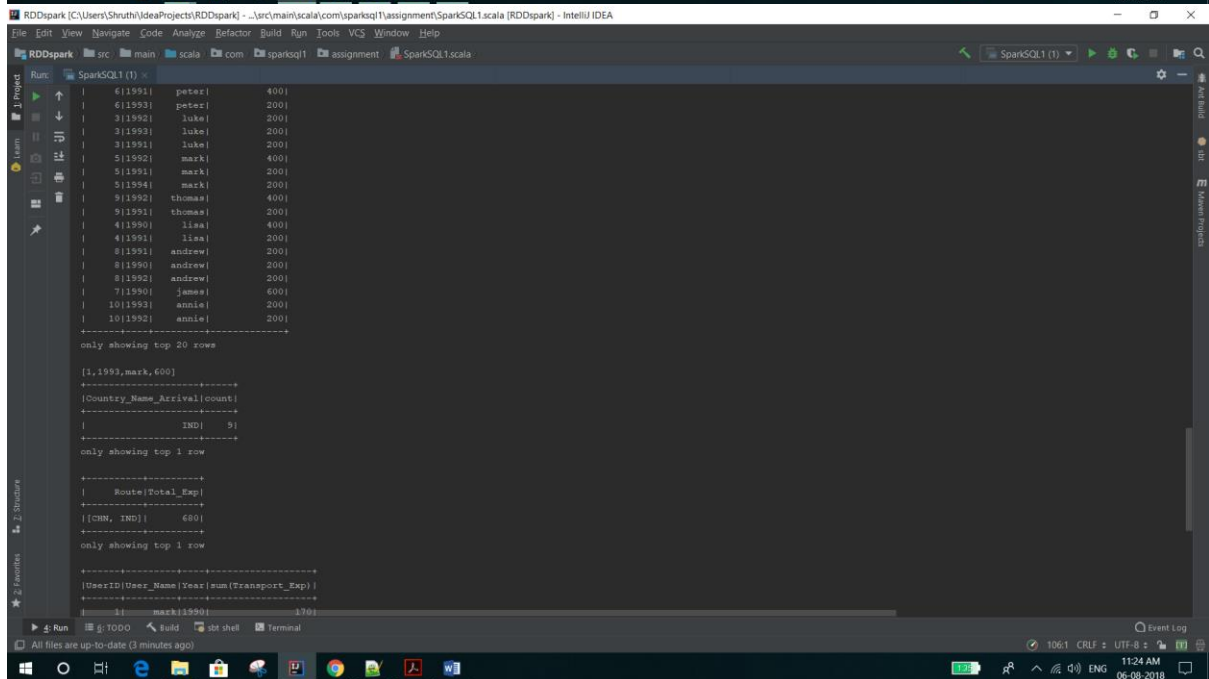
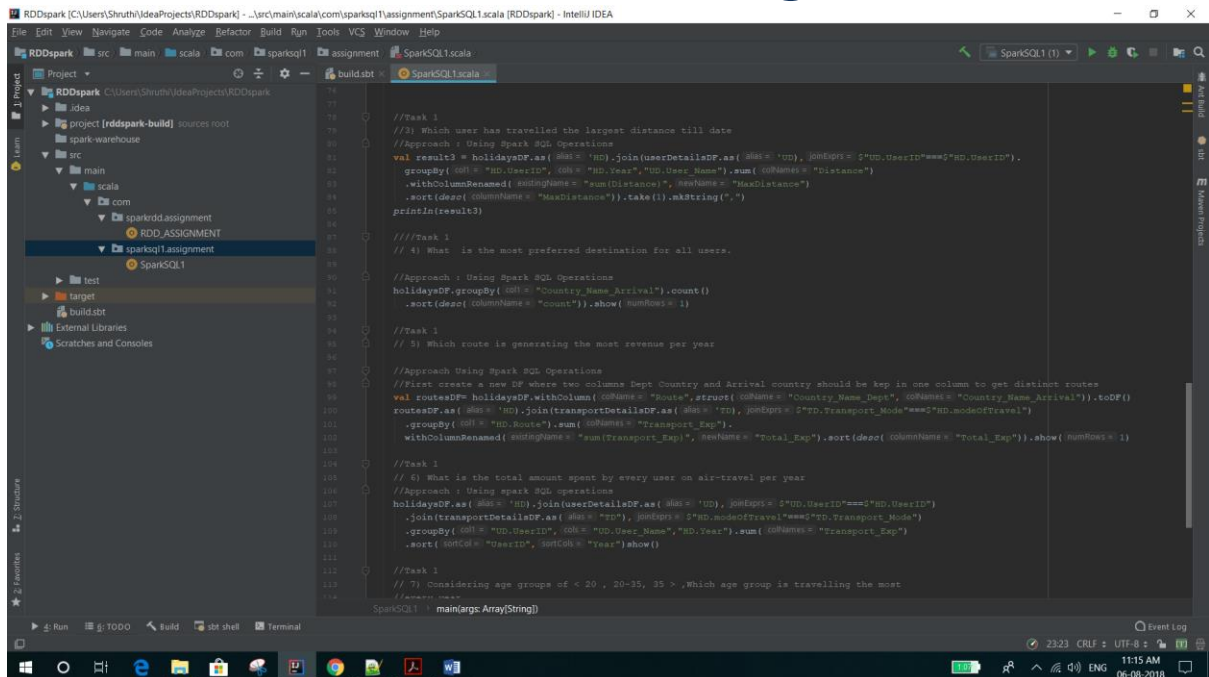


The screenshot shows the IntelliJ IDEA IDE with a SparkSQL query being executed. The query is a SELECT statement with a subquery in the FROM clause. The results are displayed in a table format, showing columns for User ID, Year, User Name, and the sum of Distance. The results are sorted by Year and then by User Name. The first row is [1, 1993, mark, 600]. The results are limited to the top 20 rows.

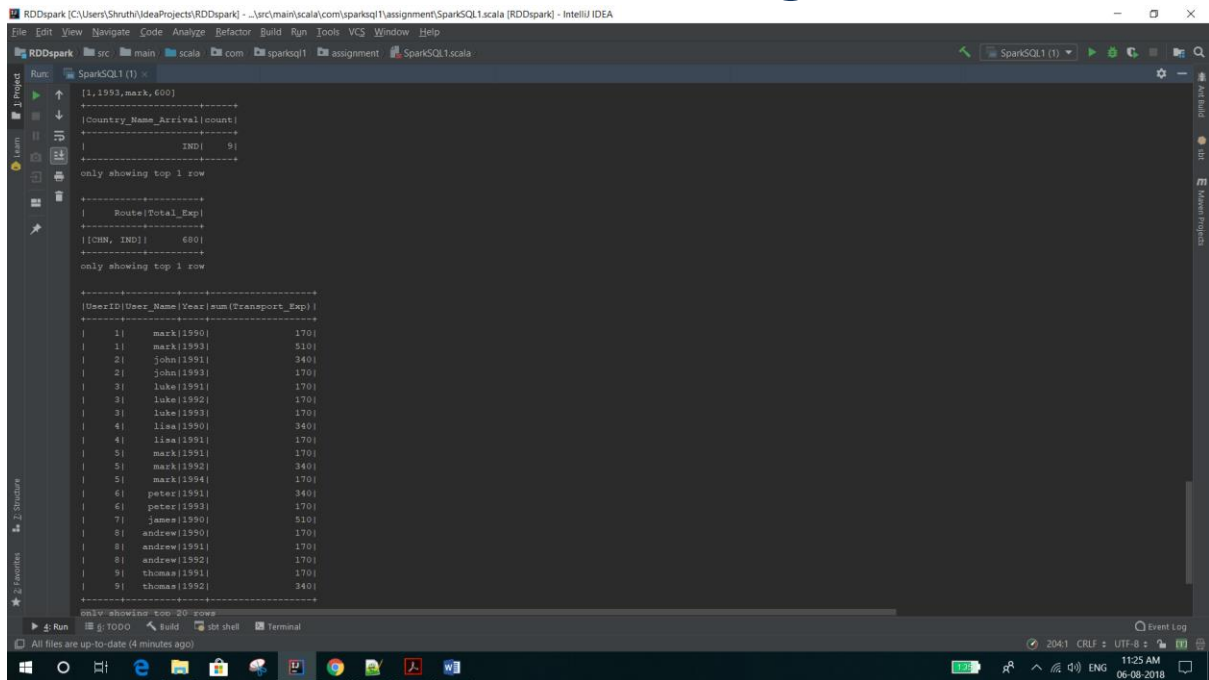
```
Run: SparkSQL (1) x
[year] a]
-----
[1990] 8]
[1994] 1]
[1991] 5]
[1992] 7]
[1993] 7]
-----

Below Result is after joining two Data frames
-----
[UserID|Year|User_Name|sum(Distance)]
-----
[ 1|1990| mark| 200]
[ 1|1993| mark| 600]
[ 4|1991| peter| 400]
[ 6|1993| peter| 200]
[ 3|1992| luke| 200]
[ 3|1993| luke| 200]
[ 3|1991| luke| 200]
[ 5|1992| mark| 400]
[ 5|1991| mark| 200]
[ 5|1994| mark| 200]
[ 9|1992| thomas| 400]
[ 9|1991| thomas| 200]
[ 4|1990| lisa| 400]
[ 4|1991| lisa| 200]
[ 8|1991| andrew| 200]
[ 8|1990| andrew| 200]
[ 8|1992| andrew| 200]
[ 7|1990| james| 600]
[10|1993| annie| 200]
[10|1992| annie| 200]
-----
only showing top 20 rows
[1,1993,mark,600]
-----
[Country_Name_Arrival|count]
-----
```

SPARK SQL 1



SPARK SQL 1



The screenshot shows the IntelliJ IDEA IDE with a SparkSQL query being executed. The query is as follows:

```
[1,1993,mark,600]
+-----+
|Country_Name_Arrival|count|
+-----+
|IND|          9|
+-----+
only showing top 1 row

+-----+
|Route|Total_Exp|
+-----+
|[OHM, IND]|      600|
+-----+
only showing top 1 row

+-----+
|UserID|User_Name|Year|sum(Transport_Exp)|
+-----+
| 1|    mark|1990|          170|
| 1|    mark|1993|          510|
| 2|   John|1991|          340|
| 2|   John|1993|          170|
| 3|   Luke|1991|          170|
| 3|   Luke|1992|          170|
| 3|   Luke|1993|          170|
| 4|   Lisa|1990|          340|
| 4|   Lisa|1991|          170|
| 5|   mark|1991|          170|
| 5|   mark|1992|          340|
| 5|   mark|1994|          170|
| 6|  Peter|1991|          340|
| 6|  Peter|1993|          170|
| 7|   James|1990|          510|
| 8| Andrew|1990|          170|
| 8| Andrew|1991|          170|
| 8| Andrew|1992|          170|
| 9| Thomas|1991|          170|
| 9| Thomas|1992|          340|
+-----+
```

The IDE interface includes a menu bar (File, Edit, View, Navigate, Code, Analyze, Refactor, Build, Run, Tools, VCS, Window, Help), a toolbar with icons for Run, Stop, and other actions, and a status bar at the bottom showing the current file is up-to-date (4 minutes ago) and the system clock (11:25 AM, 06-08-2018).

SPARK SQL 1

```
RDDEspark [C:\Users\Shruthi\IdeaProjects\RDDEspark] - ...src\main\scala\com\sparksql\assignment\SparkSQL1.scala [RDDEspark] - IntelliJ IDEA
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
RDDEspark | src | main | scala | com | sparksql | assignment | SparkSQL1.scala
Project | build.sbt | SparkSQL1.scala
//Task 1
// 5) Which route is generating the most revenue per year
//Approach Using Spark SQL Operations
//First create a new DF where two columns Dept Country and Arrival country should be kept in one column to get distinct routes
val routesDF= holidaysDF.withColumn(colName="Route",struct(colName="Country_Name_Dept", colName="Country_Name_Arrival")).coalesce(1)
routesDF.as( alias="RD").join(transportDetailsDF.as( alias="TD", joinHints="TD.Transport_Mode"==>"RD.ModeOfTravel"))
.groupBy(colName="Route").sum(colNames="Transport_Exp").
.withColumnRenamed(sum(Transport_Exp), newName="Total_Exp").sort(desc(colName="Total_Exp")).show(numRows=1)

//Task 1
// 6) What is the total amount spent by every user on air-travel per year
//Approach : Using spark SQL operations
holidaysDF.as( alias="RD").join(userDetailsDF.as( alias="UD", joinHints="UD.UserID"==>"RD.UserID"))
.join(transportDetailsDF.as( alias="TD", joinHints="TD.ModeOfTravel"==>"RD.Transport_Mode"))
.groupBy(colName="UD.UserID", colName="UD.User Name", "RD.Year").sum(colNames="Transport_Exp")
.sort(sortCol="UserID", sortCol="Year").show()

//Task 1
// 7) Grouping age groups of < 20 , 20-35 , 35 > , Which age group is travelling the most
//every year.
//another Approach of Case Statement
holidaysDF.as( alias="RD").join(userDetailsDF.as( alias="UD", joinHints="UD.UserID"==>"RD.UserID"))
.select(when(colName="UD.Age"<20, value="LessThan20"),
when(colName="UD.Age">20 && "UD.Age"<35, value="Between20And35"),
when(colName="UD.Age">35, value="Above35").alias( alias="AgeGroup"))
.groupBy(colName="AgeGroup").count().sort(desc(colName="count")).show(numRows=1)

SparkSQL1 | main(args: Array[String])
```

```
RDDEspark [C:\Users\Shruthi\IdeaProjects\RDDEspark] - ...src\main\scala\com\sparksql\assignment\SparkSQL1.scala [RDDEspark] - IntelliJ IDEA
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
RDDEspark | src | main | scala | com | sparksql | assignment | SparkSQL1.scala
Run | SparkSQL1 (1) |
only showing top 1 row
+-----+
|UserID|User Name|Year|sum(Transport_Exp)|
+-----+
| 11 | mark |1990 | 170 |
| 11 | mark |1993 | 510 |
| 21 | john |1991 | 340 |
| 21 | john |1993 | 170 |
| 31 | luke |1991 | 170 |
| 31 | luke |1992 | 170 |
| 31 | luke |1993 | 170 |
| 41 | lisa |1990 | 340 |
| 41 | lisa |1991 | 170 |
| 51 | mark |1991 | 170 |
| 51 | mark |1992 | 340 |
| 51 | mark |1994 | 170 |
| 61 | peter |1991 | 340 |
| 61 | peter |1993 | 170 |
| 71 | james |1990 | 510 |
| 81 | andrew |1990 | 170 |
| 81 | andrew |1991 | 170 |
| 81 | andrew |1992 | 170 |
| 91 | thomas |1991 | 170 |
| 91 | thomas |1992 | 340 |
+-----+
only showing top 20 rows
+-----+
| AgeGroup |count|
+-----+
|Between20And35| 13 |
+-----+
only showing top 1 row
Process finished with exit code 0
All files are up-to-date (4 minutes ago)
```