# BIG DATA HADOOP & SPARK TRAINING
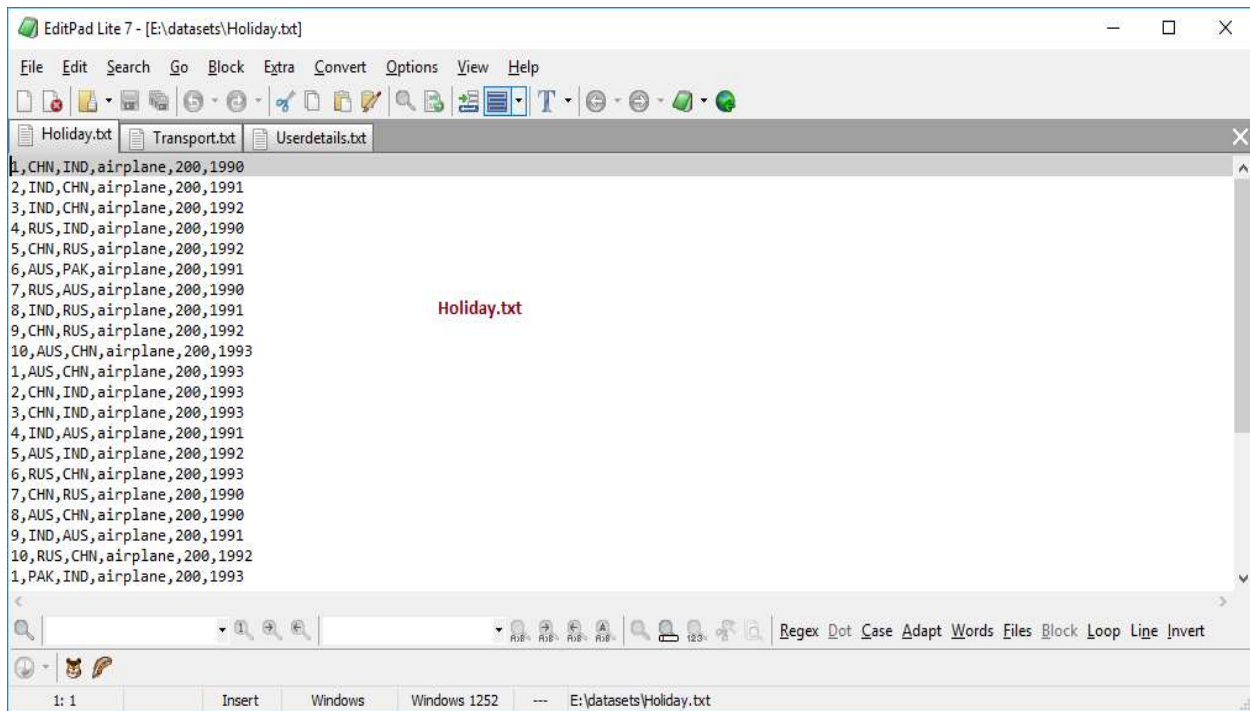
ASSIGNMENT ON SPARK SQL I

Rashmi Krishna

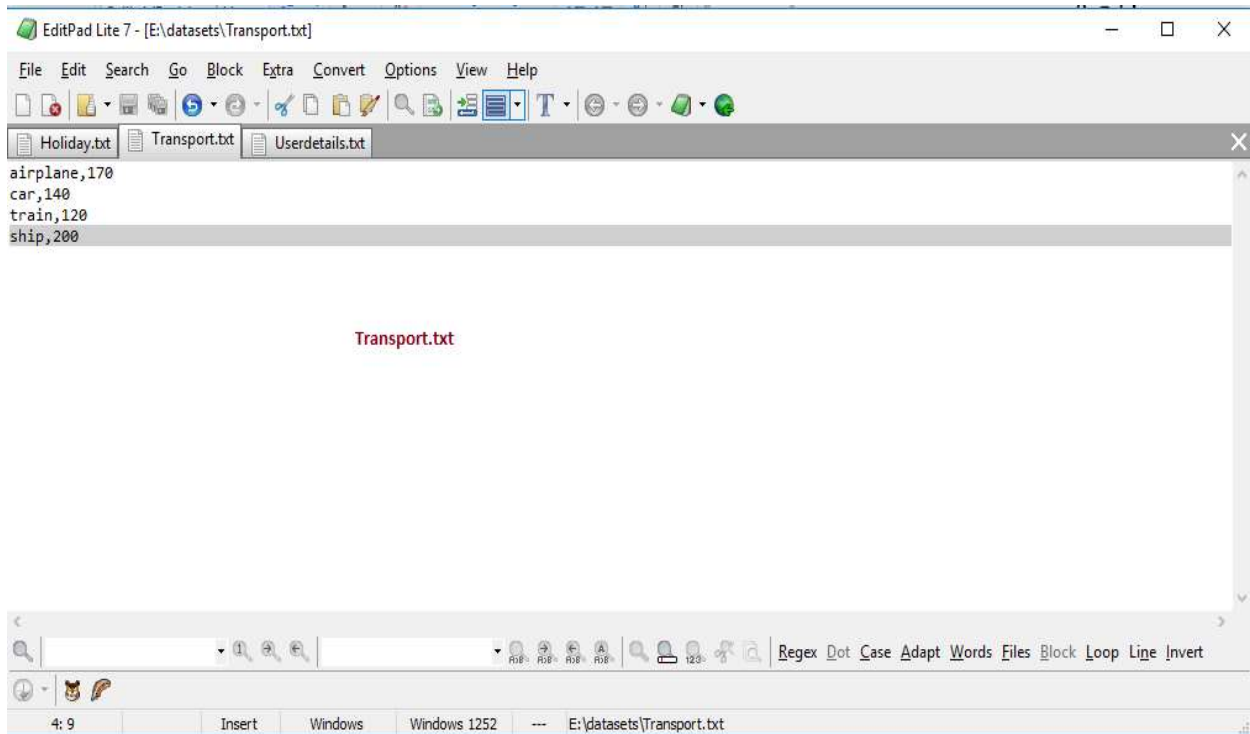# Task 1

Input files required to perform the below operations:



Holiday.txt

```
EditPad Lite 7 - [E:\datasets\Holiday.txt]
File  Edit  Search  Go  Block  Extra  Convert  Options  View  Help

Holiday.txt   Transport.txt   Userdetails.txt
1,CHN,IND,airplane,200,1990
2,IND,CHN,airplane,200,1991
3,IND,CHN,airplane,200,1992
4,RUS,IND,airplane,200,1990
5,CHN,RUS,airplane,200,1992
6,AUS,PAK,airplane,200,1991
7,RUS,AUS,airplane,200,1990
8,IND,RUS,airplane,200,1991
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
1,PAK,IND,airplane,200,1993
```
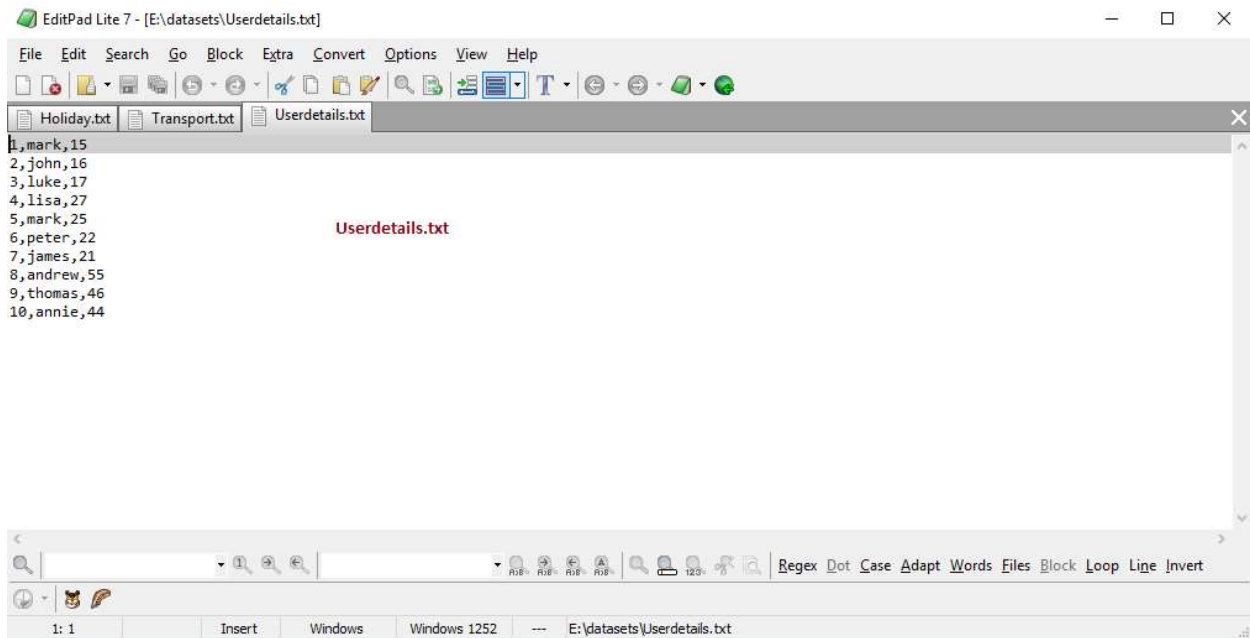
Holiday.txt

1: 1      Insert    Windows    Windows 1252   ---   E:\datasets\Holiday.txt

Transport.txt

```
EditPad Lite 7 - [E:\datasets\Transport.txt]
File  Edit  Search  Go  Block  Extra  Convert  Options  View  Help

Holiday.txt   Transport.txt   Userdetails.txt
airplane,170
car,140
train,120
ship,200
```

Transport.txt

4: 9      Insert    Windows    Windows 1252   ---   E:\datasets\Transport.txt

Below tasks are performed in **IntelliJ IDEA.**

Created a Spark-Scala Application and performed following steps to read the above mentioned text data to spark.

- ➢ Created case class for each text files which represents the schema for the respective text files.
- o **case class** Holidays_details(User_ID: Int, Country_name_Arr: String, Country_name_dest: String, Transport_mode: String, Distance: Int, year: Long)
- o **case class** Transport_mode(Trans_mode: String, Transportation_Exp: Int)
- o **case class** User_ID(Users_ID: Int, User_name: String, User_age: Int)
- ➢ created a spark session object for spark application
  - o val spark = SparkSession
  - o .builder()
  - o .master("local")
  - o .appName("Spark SQL basic example")
  - o .config("spark.some.config.option", "some-value")
  - o .getOrCreate()
  - o println("Spark Session Object created")

- ➢ Imported spark implicits to convert RDD's to DataFrames implicitly
  - o import spark.implicits._

- ➢ created a spark context object to read contents of text file to Spark
  - o val HolidaysFromFile = spark.sparkContext // create spark context object

- o    .textFile("E:\\datasets\\Holiday.txt")// path of the text file in the local file system
- o    .map(_.split(","))//splitting the input file base on ',' seperator
- o    .map(attributes => Holidays_details(attributes(0).toInt, attributes(1), attributes(2), attributes(3), attributes(4).toInt, attributes(5).toInt)) //assisgning the inputs to attributes of case class "Holidays_details
- o    .toDF()// converting RDD to Data Frame
- o    HolidaysFromFile.show() // displaying contents of the variable in which contents of text file is stored

➢ Similarly different DataFrames are created for other two text files as well, as shown in the below screen shot.

```scala
//loading Transport text file
val TransportFromFile = spark.sparkContext
  .textFile( path = "E:\\datasets\\Transport.txt")
  .map(_.split( regex = ","))
  .map(attributes => Transport_mode(attributes(0), attributes(1).toInt))
TransportFromFile.toDF().show()

//loading Userdetails text file
val UserIDFromFile = spark.sparkContext
  .textFile( path = "E:\\datasets\\Userdetails.txt")
  .map(_.split( regex = ","))
  .map(attributes => User_ID(attributes(0).toInt, attributes(1), attributes(2).toInt))
UserIDFromFile.toDF().show()
```

Loading/reading text file called "Transport.txt" from local file system to Spark using the case class "Transport_mode"

Loading/reading text file called "Userdetails.txt" from local file system to Spark using the case class "User_ID"



```
18/05/12 10:11:10 INFO TaskSetManager: Finished task 0.0 in stage 0.0 (TID 0) in 413 ms on localhost (executor driver) (1/1)
18/05/12 10:11:10 INFO TaskSchedulerImpl: Removed TaskSet 0.0, whose tasks have all completed, from pool
18/05/12 10:11:10 INFO DAGScheduler: ResultStage 0 (show at asnt20.scala:27) finished in 0.454 s
18/05/12 10:11:10 INFO DAGScheduler: Job 0 finished: show at asnt20.scala:27, took 0.818321 s
18/05/12 10:11:11 INFO CodeGenerator: Code generated in 152.101746 ms
+-------+---------------+----------------+--------------+--------+----+
|User_ID|Country_name_Arr|Country_name_dest|Transport_mode|Distance|year|
+-------+---------------+----------------+--------------+--------+----+
|      1|            CHN|             IND|      airplane|     200|1990|
|      2|            IND|             CHN|      airplane|     200|1991|
|      3|            IND|             CHN|      airplane|     200|1992|
|      4|            RUS|             IND|      airplane|     200|1990|
|      5|            CHN|             RUS|      airplane|     200|1992|
|      6|            AUS|             PAK|      airplane|     200|1991|
|      7|            RUS|             AUS|      airplane|     200|1990|
|      8|            IND|             RUS|      airplane|     200|1991|
|      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
```

Contents of Holiday.txt

**Screenshot 1: NewProj1 - IntelliJ IDEA**

```
NewProj1 [C:\Users\Laptop\IdeaProjects\NewProj1] - ...\src\main\scala\asnt20.scala [newproj1] - IntelliJ IDEA

File  Edit  View  Navigate  Code  Analyze  Refactor  Build  Run  Tools  VCS  Window  Help

NewProj1   src   main   scala   asnt20.scala                                                                    asnt20

Run:   asnt20

18/05/12 10:11:11 INFO DAGScheduler: Final stage: ResultStage 1 (show at asnt20.scala:39)
18/05/12 10:11:11 INFO DAGScheduler: Parents of final stage: List()
18/05/12 10:11:11 INFO DAGScheduler: Missing parents: List()
18/05/12 10:11:11 INFO DAGScheduler: Submitting ResultStage 1 (MapPartitionsRDD[13] at show at asnt20.scala:39), which has no missing parents
18/05/12 10:11:11 INFO MemoryStore: Block broadcast_3 stored as values in memory (estimated size 9.1 KB, free 902.4 MB)
18/05/12 10:11:11 INFO MemoryStore: Block broadcast_3_piece0 stored as bytes in memory (estimated size 4.4 KB, free 902.4 MB)
18/05/12 10:11:11 INFO BlockManagerInfo: Added broadcast_3_piece0 in memory on 192.168.56.1:54232 (size: 4.4 KB, free: 902.7 MB)
18/05/12 10:11:11 INFO SparkContext: Created broadcast 3 from broadcast at DAGScheduler.scala:996
18/05/12 10:11:11 INFO DAGScheduler: Submitting 1 missing tasks from ResultStage 1 (MapPartitionsRDD[13] at show at asnt20.scala:39)
18/05/12 10:11:11 INFO TaskSchedulerImpl: Adding task set 1.0 with 1 tasks
18/05/12 10:11:11 INFO TaskSetManager: Starting task 0.0 in stage 1.0 (TID 1, localhost, executor driver, partition 0, PROCESS_LOCAL, 5924 bytes)
18/05/12 10:11:11 INFO Executor: Running task 0.0 in stage 1.0 (TID 1)
18/05/12 10:11:12 INFO HadoopRDD: Input split: file:/E:/datasets/Transport.txt:0+42
18/05/12 10:11:12 INFO ContextCleaner: Cleaned accumulator 0
18/05/12 10:11:12 INFO ContextCleaner: Cleaned accumulator 1
18/05/12 10:11:12 INFO Executor: Finished task 0.0 in stage 1.0 (TID 1). 1499 bytes result sent to driver
18/05/12 10:11:12 INFO TaskSetManager: Finished task 0.0 in stage 1.0 (TID 1) in 317 ms on localhost (executor driver) (1/1)
18/05/12 10:11:12 INFO TaskSchedulerImpl: Removed TaskSet 1.0, whose tasks have all completed, from pool
18/05/12 10:11:12 INFO DAGScheduler: ResultStage 1 (show at asnt20.scala:39) finished in 0.319 s
18/05/12 10:11:12 INFO DAGScheduler: Job 1 finished: show at asnt20.scala:39, took 0.346637 s
18/05/12 10:11:12 INFO CodeGenerator: Code generated in 23.219237 ms
+----------+------------------+
|Trans_mode|Transportation_Exp|
+----------+------------------+
|  airplane|               170|
|       car|               140|
|     train|               120|
|      ship|               200|
+----------+------------------+

18/05/12 10:11:12 INFO BlockManagerInfo: Removed broadcast_1_piece0 on 192.168.56.1:54232 in memory (size: 5.2 KB, free: 902.7 MB)
18/05/12 10:11:12 INFO MemoryStore: Block broadcast_4 stored as values in memory (estimated size 127.2 KB, free 902.3 MB)
```

Contents of transport.txt

**Screenshot 2: NewProj1 - IntelliJ IDEA**

```
NewProj1 [C:\Users\Laptop\IdeaProjects\NewProj1] - ...\src\main\scala\asnt20.scala [newproj1] - IntelliJ IDEA

File  Edit  View  Navigate  Code  Analyze  Refactor  Build  Run  Tools  VCS  Window  Help

NewProj1   src   main   scala   asnt20.scala                                                                    asnt20

Run:   asnt20

18/05/12 10:11:12 INFO TaskSchedulerImpl: Adding task set 2.0 with 1 tasks
18/05/12 10:11:12 INFO TaskSetManager: Starting task 0.0 in stage 2.0 (TID 2, localhost, executor driver, partition 0, PROCESS_LOCAL, 5926 bytes)
18/05/12 10:11:12 INFO Executor: Running task 0.0 in stage 2.0 (TID 2)
18/05/12 10:11:12 INFO HadoopRDD: Input split: file:/E:/datasets/Userdetails.txt:0+116
18/05/12 10:11:12 INFO Executor: Finished task 0.0 in stage 2.0 (TID 2). 1589 bytes result sent to driver
18/05/12 10:11:12 INFO TaskSetManager: Finished task 0.0 in stage 2.0 (TID 2) in 27 ms on localhost (executor driver) (1/1)
18/05/12 10:11:12 INFO TaskSchedulerImpl: Removed TaskSet 2.0, whose tasks have all completed, from pool
18/05/12 10:11:12 INFO DAGScheduler: ResultStage 2 (show at asnt20.scala:46) finished in 0.028 s
18/05/12 10:11:12 INFO DAGScheduler: Job 2 finished: show at asnt20.scala:46, took 0.058412 s
18/05/12 10:11:12 INFO CodeGenerator: Code generated in 17.49142 ms
+--------+---------+--------+
|Users_ID|User_name|User_age|
+--------+---------+--------+
|       1|     mark|      15|
|       2|     john|      16|
|       3|     luke|      17|
|       4|     lisa|      27|
|       5|     mark|      25|
|       6|    peter|      22|
|       7|    james|      21|
|       8|   andrew|      55|
|       9|   thomas|      46|
|      10|    annie|      44|
+--------+---------+--------+

18/05/12 10:11:12 INFO SparkSqlParser: Parsing command: Holiday_data
18/05/12 10:11:13 INFO SparkSqlParser: Parsing command:  select year,count("year") from Holiday_data Group by year
18/05/12 10:11:13 INFO CodeGenerator: Code generated in 54.24509 ms
18/05/12 10:11:13 INFO CodeGenerator: Code generated in 99.082548 ms
18/05/12 10:11:13 INFO SparkContext: Starting job: show at asnt20.scala:50
18/05/12 10:11:13 INFO DAGScheduler: Registering RDD 24 (show at asnt20.scala:50)
18/05/12 10:11:13 INFO DAGScheduler: Got job 3 (show at asnt20.scala:50) with 1 output partitions
18/05/12 10:11:13 INFO DAGScheduler: Final stage: ResultStage 4 (show at asnt20.scala:50)
18/05/12 10:11:13 INFO DAGScheduler: Parents of final stage: List(ShuffleMapStage 3)
```
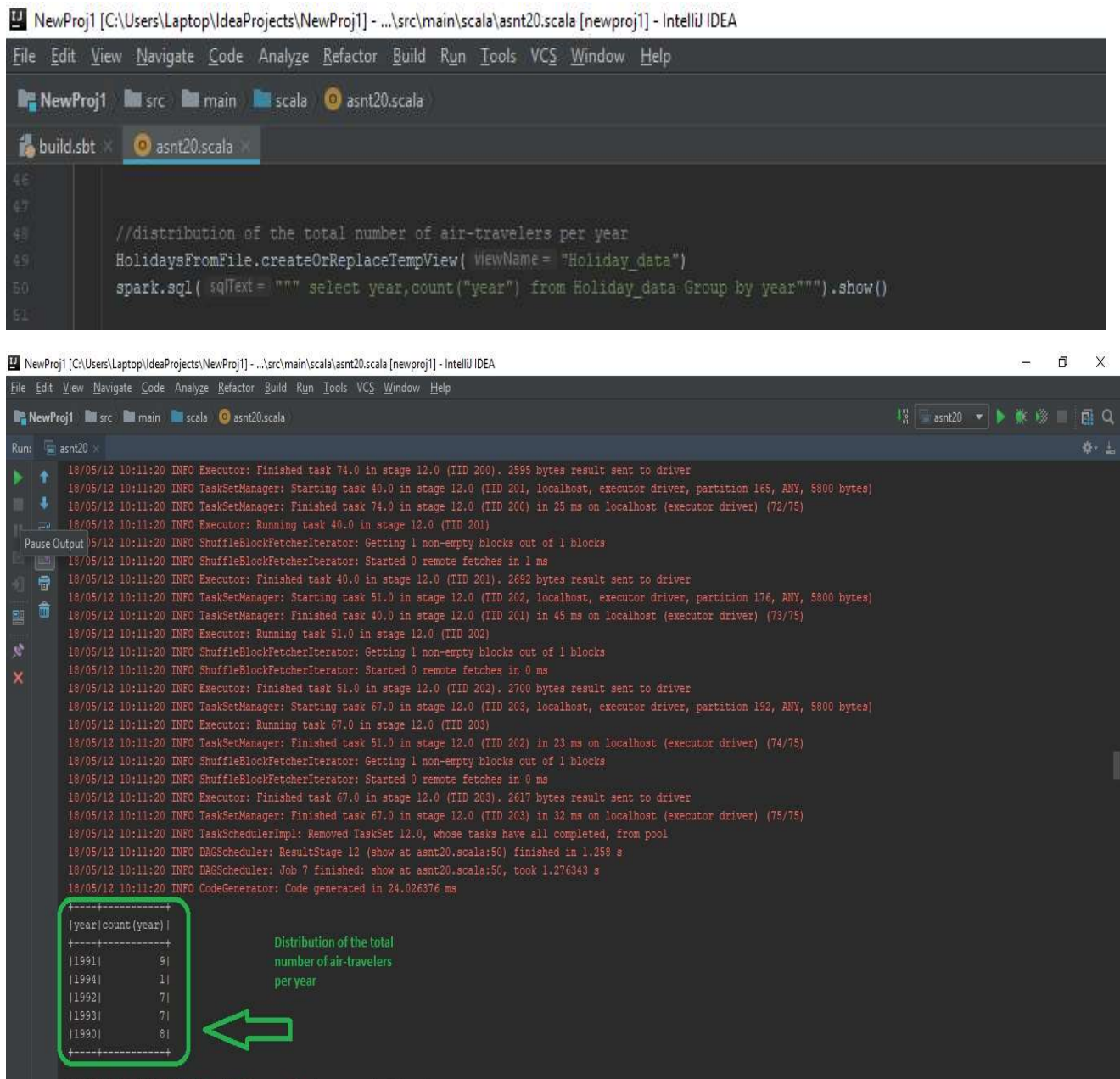
Contents of User Details.txt

1) What is the distribution of the total number of air-travelers per year?

Below code will generate the solution for this task:

➢ HolidaysFromFile.createOrReplaceTempView("Holiday_data")
   //create a view on the Data Frames called "Holiday_data"
➢ spark.sql(""" select year,count("year") from Holiday_data Group by year""").show()

   Write a query to select year and calculate count of year from the view created

2) What is the total air distance covered by each user per year

Below code will generate the solution for this task:

- val joindf = HolidaysFromFile.as ('a).join(UserIDFromFile.toDF().as('b), $"a.User_ID" === $"b.Users_ID")//join two data frames based on User_ID
- joindf.createOrReplaceTempView("join_view")//create a view on the joindf variable
- spark.sql(""" select User_ID,User_name,year,sum(Distance)as Total_distance from join_view Group by year,User_ID,User_name order by Total_distance desc""").show()                //Write a SQL query to find the total air distance covered by each user per year

3) Which user has travelled the largest distance till date

Below code will generate the solution for this task:

- val task1 = spark.sql(""" select User_ID,User_name,year,sum(Distance)as Total_distance from join_view Group by year,User_ID,User_name""")//saving the SQL query in a variable "task1"
- task1.toDF().createOrReplaceTempView("Distance_view")// creating view on task1 variable
- spark.sql("""select User_ID,User_name,year,max(Total_distance) as Maximum_distance from Distance_view Group by year, User_ID,User_name order by Maximum_distance desc""")
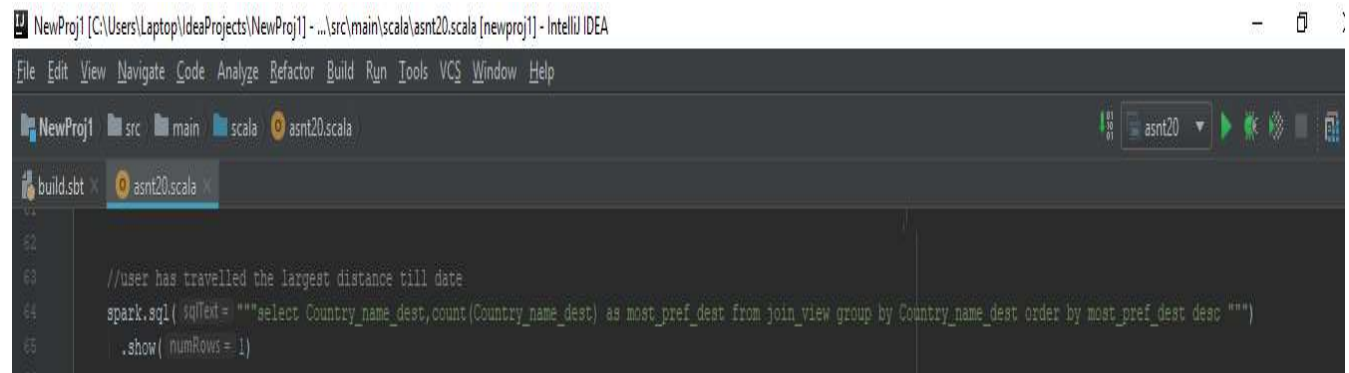- .show(1)// creating a SQL query to find which user has travelled the largest distance till date
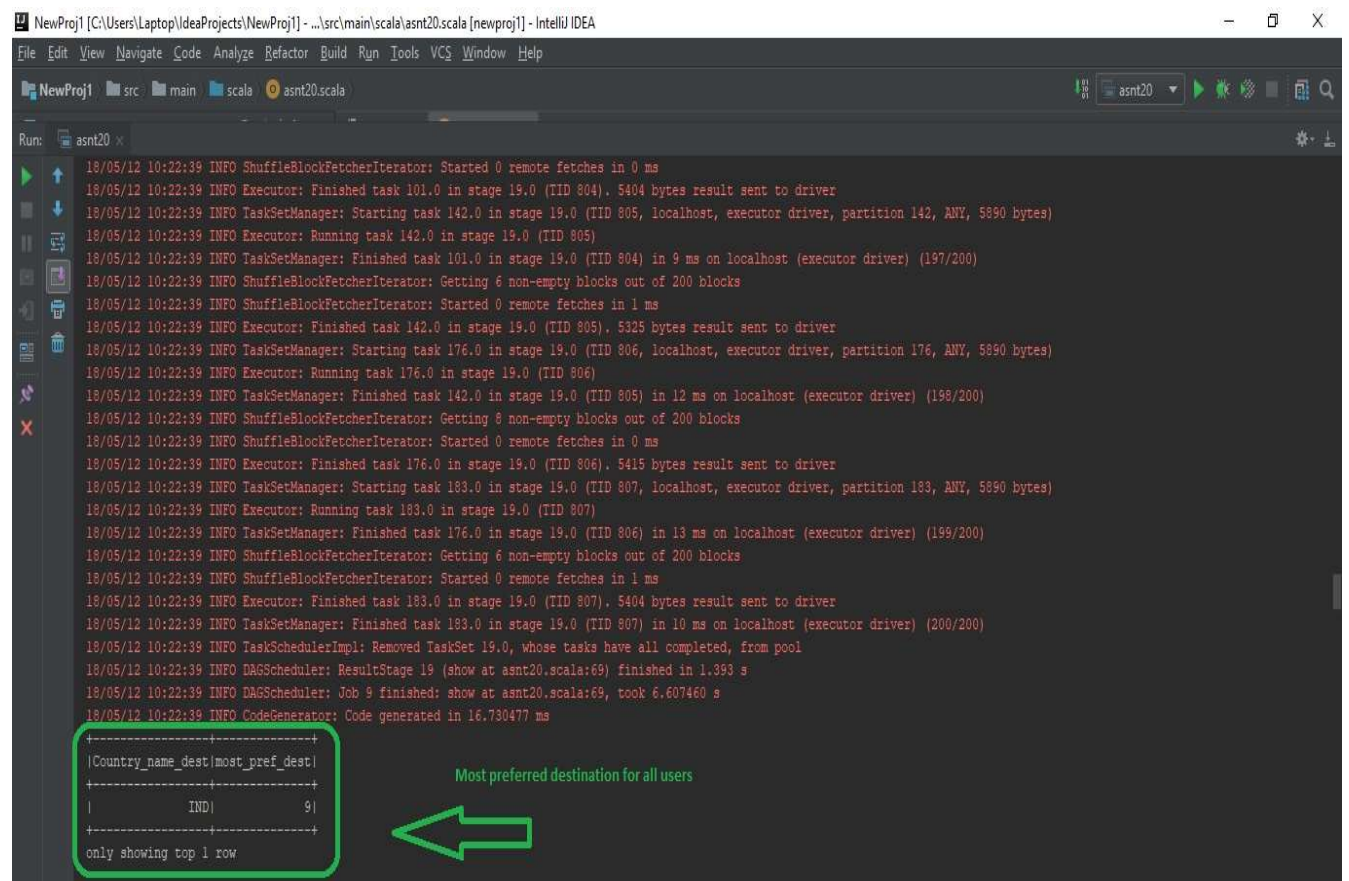
4) What is the most preferred destination for all users.

Below code will generate the solution for this task:

spark.sql("""select Country_name_dest,count(Country_name_dest) as most_pref_dest from join_view group by Country_name_dest order by most_pref_dest desc """) .show(1)

// write an sql query to find most preferred destination for all users from the view previously created //called "join_view"





Most preferred destination for all users

5)Which route is generating the most revenue per year

Below code will generate the solution for this task:

- val joineddf = HolidaysFromFile.as('c').join(TransportFromFile.toDF().as('d), $"c.Transport_mode" === $"d.Trans_mode", joinType = "left_outer")
  // to join two Dataframes based on Transportation mode
- joineddf.createOrReplaceTempView("revenue_view")
  //creating a view on the join previously created
- val revenue = spark.sql(   // creating an sql query to calculate the revenue
- """select User_ID,Country_name_dest, Transport_mode,year, count(Transport_mode) * sum(Transportation_Exp) as revenue_exp from revenue_view
- | group by User_ID,year,Country_name_dest,Transport_mode""".stripMargin)
- revenue.toDF().createOrReplaceTempView("max_revenue")// converting the previously evaluated variable to DataFrame and create a view on the same called"max_revenue"
- spark.sql(
- """select Country_name_dest, Transport_mode,year,max(revenue_exp) as maximum_revenue from max_revenue
- | group by Country_name_dest, Transport_mode, year
- | order by maximum_revenue desc""".stripMargin).show()

// write an sql query to find which route is generating the most revenue per year

```
18/05/12 10:22:44 INFO Executor: Finished task 87.0 in stage 22.0 (TID 1008). 5629 bytes result sent to driver
18/05/12 10:22:44 INFO TaskSetManager: Starting task 172.0 in stage 22.0 (TID 1009, localhost, executor driver, partition 172, ANY, 6218 bytes)
18/05/12 10:22:44 INFO Executor: Running task 172.0 in stage 22.0 (TID 1009)
18/05/12 10:22:44 INFO TaskSetManager: Finished task 87.0 in stage 22.0 (TID 1008) in 18 ms on localhost (executor driver) (199/200)
18/05/12 10:22:44 INFO ShuffleBlockFetcherIterator: Getting 0 non-empty blocks out of 1 blocks
18/05/12 10:22:44 INFO ShuffleBlockFetcherIterator: Started 0 remote fetches in 0 ms
18/05/12 10:22:44 INFO ShuffleBlockFetcherIterator: Getting 1 non-empty blocks out of 1 blocks
18/05/12 10:22:44 INFO ShuffleBlockFetcherIterator: Started 0 remote fetches in 0 ms
18/05/12 10:22:44 INFO Executor: Finished task 172.0 in stage 22.0 (TID 1009). 5618 bytes result sent to driver
18/05/12 10:22:44 INFO TaskSetManager: Finished task 172.0 in stage 22.0 (TID 1009) in 14 ms on localhost (executor driver) (200/200)
18/05/12 10:22:44 INFO TaskSchedulerImpl: Removed TaskSet 22.0, whose tasks have all completed, from pool
18/05/12 10:22:44 INFO DAGScheduler: ResultStage 22 (show at asnt20.scala:77) finished in 3.486 s
18/05/12 10:22:44 INFO DAGScheduler: Job 10 finished: show at asnt20.scala:77, took 3.863660 s
+----------------+--------------+----+---------------+
|Country_name_dest|Transport_mode|year|maximum_revenue|
+----------------+--------------+----+---------------+
|             IND|      airplane|1990|            170|
|             CHN|      airplane|1991|            170|
|             CHN|      airplane|1992|            170|
|             RUS|      airplane|1992|            170|
|             PAK|      airplane|1991|            170|
|             AUS|      airplane|1990|            170|
|             RUS|      airplane|1991|            170|
|             CHN|      airplane|1993|            170|
|             IND|      airplane|1993|            170|
|             AUS|      airplane|1991|            170|
|             IND|      airplane|1992|            170|
|             RUS|      airplane|1990|            170|
|             CHN|      airplane|1990|            170|
|             PAK|      airplane|1990|            170|
|             AUS|      airplane|1993|            170|
|             PAK|      airplane|1994|            170|
+----------------+--------------+----+---------------+
```

Route generating the most revenue per year

5) What is the total amount spent by every user on air-travel per year

Below code will generate the solution for this task:

- val expense = spark.sql(// create a sql query to calculate the transportation expenses from //previously created revenue_view
- """select User_ID,Country_name_dest, Transport_mode,year,sum(Transportation_Exp) as Total_exp from revenue_view
- | group by User_ID,year,Country_name_dest, Transport_mode """.stripMargin)
- .filter(col("Transport_mode") === "airplane")// filter transport mode as airplane
- val new_joindf = UserIDFromFile.toDF().as('e).join(expense.as('f), $"e.Users_ID" === $"f.User_ID")// join two data frames based on UserID
- new_joindf.toDF().createOrReplaceTempView("expense_view") //create a view on variable created
- spark.sql("""select User_ID, Transport_mode, year, Total_exp from expense_view group by User_ID, year ,Total_exp,Transport_mode """).show() //an sql query to find total expenses spent by user for air-travel

```scala
object asnt20 {

    val expense = spark.sql(
        """select User_ID,Country_name_dest, Transport_mode,year,sum(Transportation_Exp) as Total_exp from revenue_view
          | group by User_ID,year,Country_name_dest, Transport_mode """.stripMargin)
        .filter(col( colName = "Transport_mode") === "airplane")

    // val airplane_exp = expense.filter(col("Transport_mode") === "airplane").show()
    val new_joindf = UserIDFromFile.toDF().as( alias = 'e).join(expense.as( alias = 'f), joinExprs = $"e.Users_ID" === $"f.User_ID")
    new_joindf.toDF().createOrReplaceTempView( viewName = "expense_view")
    spark.sql( sqlText = """select User_ID, Transport_mode, year, Total_exp from expense_view group by User_ID, year ,Total_exp,Transport_mode """).show()
```

Run:  asnt20

```
18/05/12 10:22:51 INFO ShuffleBlockFetcherIterator: Getting 1 non-empty blocks out of 200 blocks
18/05/12 10:22:51 INFO ShuffleBlockFetcherIterator: Started 0 remote fetches in 0 ms
18/05/12 10:22:51 INFO Executor: Finished task 97.0 in stage 42.0 (TID 1337). 5714 bytes result sent to driver
18/05/12 10:22:51 INFO TaskSetManager: Finished task 97.0 in stage 42.0 (TID 1337) in 7 ms on localhost (executor driver) (100/100)
18/05/12 10:22:51 INFO TaskSchedulerImpl: Removed TaskSet 42.0, whose tasks have all completed, from pool
18/05/12 10:22:51 INFO DAGScheduler: ResultStage 42 (show at asnt20.scala:83) finished in 0.764 s
18/05/12 10:22:51 INFO DAGScheduler: Job 14 finished: show at asnt20.scala:83, took 0.777474 s
+-------+--------------+----+---------+
|User_ID|Transport_mode|year|Total_exp|
+-------+--------------+----+---------+
|      1|      airplane|1990|      170|
|      1|      airplane|1993|      170|
|      6|      airplane|1991|      170|
|      6|      airplane|1993|      170|
|      3|      airplane|1992|      170|
|      3|      airplane|1993|      170|
|      3|      airplane|1991|      170|
|      5|      airplane|1992|      170|
|      5|      airplane|1991|      170|
|      5|      airplane|1994|      170|
|      9|      airplane|1992|      170|
|      9|      airplane|1991|      170|
|      4|      airplane|1990|      170|
|      4|      airplane|1991|      170|
|      8|      airplane|1991|      170|
|      8|      airplane|1990|      170|
|      8|      airplane|1992|      170|
|      7|      airplane|1990|      170|
|     10|      airplane|1993|      170|
|     10|      airplane|1992|      170|
+-------+--------------+----+---------+
only showing top 20 rows

18/05/12 10:22:51 INFO SparkContext: Invoking stop() from shutdown hook
```

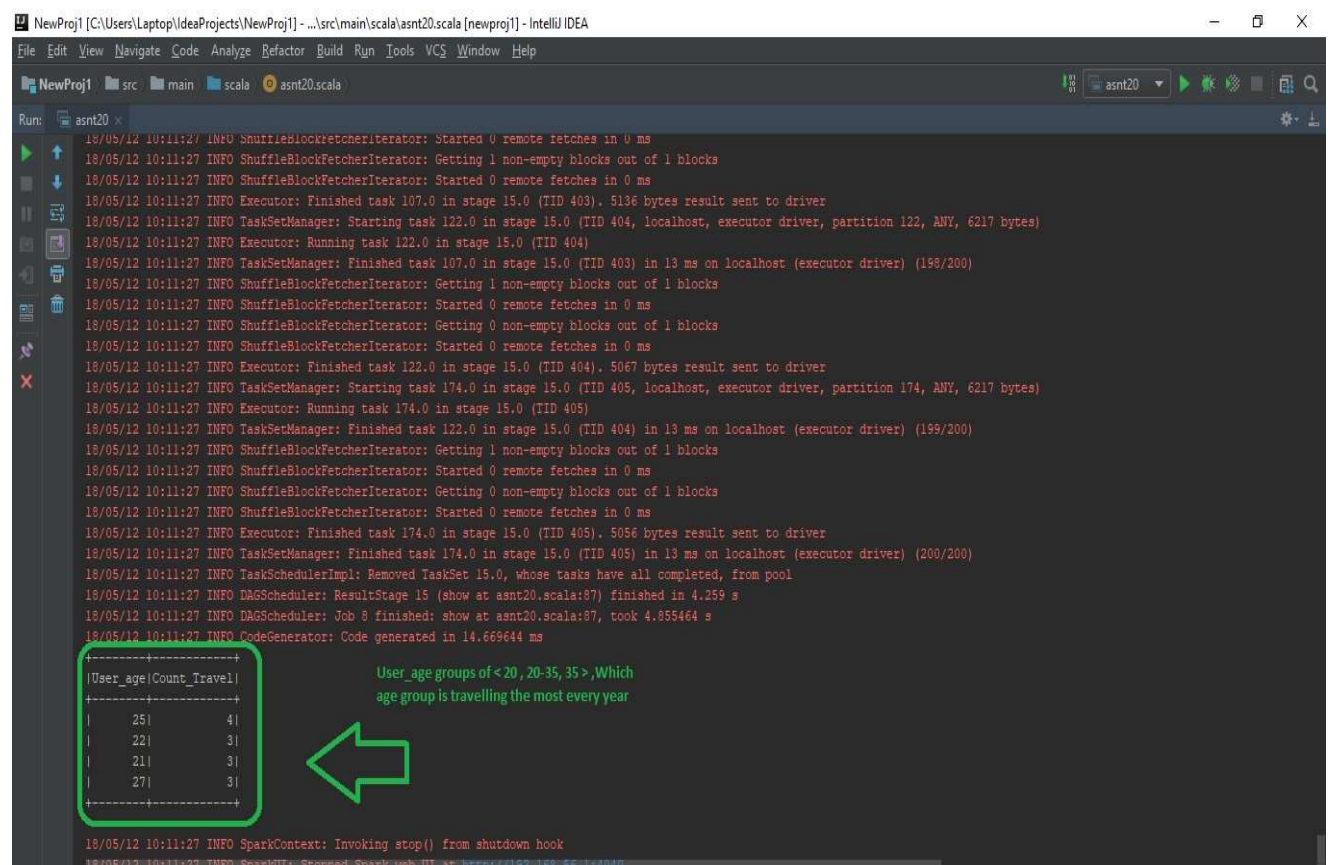Total amount spent by every user on air-travel per year

7) Considering age groups of < 20 , 20-35, 35 > ,Which age group is travelling the most every year.

spark.sql( """"select  User_age,count(Users_ID) as Count_Travel from join_view WHERE User_age >=20 AND User_age<=35 group by User_age,User_ID order by Count_Travel desc """").show()

// sql query to take count of number of users from previously created join_view and filter the records //based on Users_age between 20 to 35, to find which age group is travelling the most.