# Assignment Day 20

## Task 1:
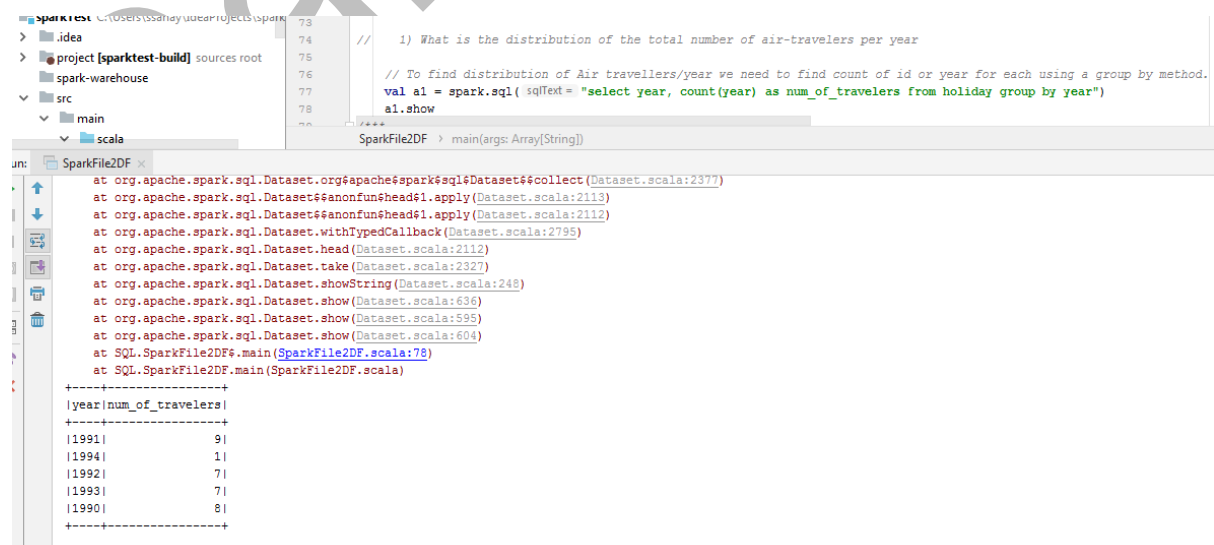
1) What is the distribution of the total number of air-travellers per year ?

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

3) Which user has travelled the largest distance till date?

4) What is the most preferred destination for all users?

5)Which route is generating the most revenue per year?

6) What is the 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?

## Ans:

**Note: Program files are properly documented for a detailed description of each instruction used within the program.**

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

**ScreenShot**:

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

**ScreenShot:**



**3)** Which user has travelled the largest distance till date?

**ScreenShot:**

## 4) What is the most preferred destination for all users?

## ScreenShot:



```scala
//Selecting destination & count of destination & providing an alias Dest_Visited from holiday view & creating an Alias
val a4_top = spark.sql( sqlText = "select destination, count(destination) as Dest_Visited from holiday group by destination")
a4_top.show

//Temp view created for saving the result of the query above
a4_top.createOrReplaceTempView( viewName = "a4_top")

//  4) What is the most preferred destination for all users.
//Selecting the maximum of all the destination grouped above in the view.
val a4 = spark.sql( sqlText = "select destination as `Most Visited Destination` from a4_top where Dest_Visited IN(select max(Dest_Visited
a4.show
```

```
+-----------+------------+
|destination|Dest_Visited|
+-----------+------------+
|        AUS|           5|
|        PAK|           5|
|        RUS|           6|
|        IND|           9|
|        CHN|           7|
+-----------+------------+

+------------------------+
|Most Visited Destination|
+------------------------+
|                     IND|
+------------------------+
```

## 5) Which route is generating the most revenue per year?

## ScreenShot:



```scala
holi_trans.createOrReplaceTempView( viewName = "holi_trans")

//From the holi_trans view above selecting the route & mutiplying the count of trnsport_mode with it's respective cost/unit
//to fetch the revenue generated by a particular route each year.
val max_revenue = spark.sql( sqlText = "select source,destination,year, (count(transport_mode) * cost_per_unit) as Revenue from holi_tran

//Temp view created for saving the result of the query above
max_revenue.createOrReplaceTempView( viewName = "max_revenue")

//  5)Which route is generating the most revenue per year
//selecting the route that genarates maximum revenue each year from max_revenue view & storing the result in a5.
val a5 = spark.sql( sqlText = "select source,destination,year,Revenue from max_revenue group by source,destination,year,Revenue HAVING R
a5.show
```

```
        at java.lang.Thread.run(Thread.java:748)
+------+-----------+----+-------+
|source|destination|year|Revenue|
+------+-----------+----+-------+
|   CHN|        IND|1990|    340|
|   CHN|        RUS|1992|    340|
|   IND|        RUS|1991|    340|
|   AUS|        CHN|1993|    340|
|   CHN|        IND|1993|    340|
|   IND|        AUS|1991|    340|
|   RUS|        IND|1992|    340|
+------+-----------+----+-------+
```

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

**ScreenShot:**



```
6) What is the total amount spent by every user on air-travel per year
//To fetch the total amount spent by a particular user each year grouping by id,year & counting by each id,year
// & multiplying it by cost.
val a6 = spark.sql( sqlText = "select id,year,(count(id,year) * cost_per_unit) as `Total Amount Spent` from holi_trans1 group by id,yea
a6.show
```

```
    at SQL.SparkFile2DF.main(SparkFile2DF.scala)
+---+----+------------------+
| id|year|Total Amount Spent|
+---+----+------------------+
|  3|1993|               170|
|  6|1991|               340|
|  2|1991|               340|
|  3|1991|               170|
| 10|1990|               170|
|  9|1991|               170|
|  7|1990|               510|
|  8|1990|               170|
|  4|1991|               170|
|  8|1992|               170|
|  9|1992|               340|
|  5|1991|               170|
|  3|1992|               170|
|  5|1994|               170|
| 10|1992|               170|
|  4|1990|               340|
| 10|1993|               170|
|  2|1993|               170|
|  8|1991|               170|
|  1|1993|               510|
```

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

**ScreenShot:**



```
//Creating an Schema with column names as "Year","ageGroup","Travelled"
val colName = Seq("Year","ageGroup","Travelled")

//Schema of toDF_data is (._1,._2,._3),convert it into (year,ageGroup,Distance) in yearGroupSortNew by specying
//above three columns as column names.
val tab_data = toDF_data.toDF(colName:_*)

// Register the DataFrame as a temporary view tab_data
tab_data.createOrReplaceTempView( viewName = "tab_data")

//The final SQL statements to get the desired result from view tab_data. Selecting all from tab_data & inner join to self
//on year & where travelled distance is max by giving a self alias to tab_data view as a & b respectively.
val a7 = spark.sql( sqlText = "select a.* from tab_data a inner join(select Year,max(Travelled) as Max from tab_data group by Year) b on
//Displaying final output as the age group that is travelling the most every year.
//    7) Considering age groups of < 20 , 20-35, 35 > ,Which age group is travelling the most every year.
a7.show
```

```
    at SQL.SparkFile2DF.main(SparkFile2DF.scala)
+----+--------+---------+
|Year|ageGroup|Travelled|
+----+--------+---------+
|1993|      20|     1000|
|1992|      35|      800|
|1991|   20-35|      800|
|1994|   20-35|      200|
|1990|   20-35|     1000|
+----+--------+---------+

Process finished with exit code 0
```

****************************** **End** ******************************