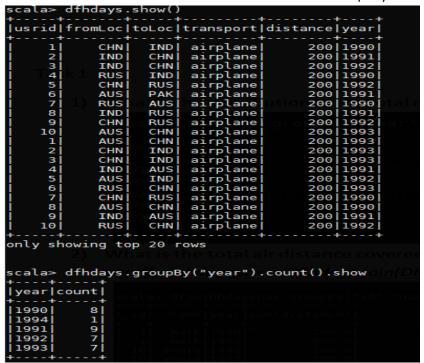
```
import org.apache.spark.sql.SparkSession
import org.apache.spark.sql.types.
val schmhdays = StructType(Seq(
StructField("usrid", IntegerType, true),
StructField("fromLoc", StringType, true),
StructField("toLoc", StringType, true),
StructField("transport", StringType, true),
StructField("distance", IntegerType, true),
StructField("year", IntegerType, true)))
val dfhdays = spark.read.format("csv").option("header",
"false").schema(schmhdays).load("resources/Holidays.csv")
val schmusers = StructType(Seq(
StructField("usrid", IntegerType, true),
StructField("username", StringType, true),
StructField("age", IntegerType, true)))
val dfusers = spark.read.format("csv").option("header",
"false").schema(schmusers).load("resources/Users.csv")
val schmtransp = StructType(Seq(
StructField("transport", StringType, true),
StructField("fare", IntegerType, true)))
val dftransp = spark.read.format("csv").option("header",
"false").schema(schmtransp).load("resources/Transport.csv")
dfusers.createOrReplaceTempView("Vdfusers")
val dfVusers = sql("select * from Vdfusers ")
dfVusers.show
dftransp.createOrReplaceTempView("Vdftransp")
val dfVtransp = sql("select * from Vdftransp")
dfVtransp.show
dfhdays.createOrReplaceTempView("Vdfhdays")
val dfVhdays = sql("select * from Vdfhdays")
dfVhdays.show
```

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



2) What is the total air distance covered by each user per year var dfjoinhdayspax = dfhdays.join(dfusers, "usrid")

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

```
scala> var dfjointransp = dfhdays.join(dftransp."transport")
dfjointransp: org.apache.spark.sql.DataFrame = [transport: string, usrid: int ... 5 more fields]

scala> var dfmostrevenue = dfjointransp.groupBy("transport", "year").sum("fare")
dfmostrevenue: org.apache.spark.sql.DataFrame = [transport: string, year: int ... 1 more field]

scala> dfmost
dfmostrevenue.show
transport|year|sum(fare)|
transport|year|sum(fare)|
iairplane|1993| 1190|
iairplane|1990| 1360|
iairplane|1991| 170|
iairplane|1992| 1190|

scala> var dfmostrevenue = dfjointransp.groupBy("transport", "year").sum("fare") as "mostrevenue_year"
dfmostrevenue: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [transport: string, year: int ... 1 more field]

scala> dfmostrevenue.show
transport|year|sum(fare)|
transport|year|sum(fare)|
iairplane|1993| 1190|
iairplane|1991| 1530|
iairplane|1991| 1530|
iairplane|1994| 170|
airplane|1994| 170|
airplane|1992| 1190|
```

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

var dfjointranspuser = dfhdays.join(dftransp,"transport").join(dfusers,"usrid")

```
scala> var dfmospentuser = dfjointranspuser.groupBy("transport", "year","usrid","username").sum("fare") as "SpentByUser_year"
dfmospentuser: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [transport: string, year: int ... 3 more fields]
 scala> dfmospentuser.show
|transport|year|usrid|username|sum(fare)|
     airplane | 1993|
airplane | 1992|
airplane | 1992|
airplane | 1992|
airplane | 1991|
airplane | 1991|
airplane | 1991|
airplane | 1993|
airplane | 1993|
airplane | 1993|
airplane | 1991|
airplane | 1993|
airplane | 1994|
airplane | 1994|
airplane | 1994|
airplane | 1996|
airplane | 1990|
airplane | 1990|
airplane | 1990|
airplane | 1990|
                                                                                                                            510
                                                               10
53
8
5
9
4
2
10
8
9
1
6
6
10
7
5
4
                                                                               john
annie
mark
luke
andrew
mark
                                                                                                                            340
                                                                                                                            170
340
170
170
170
340
170
170
170
170
                                                                               thomas
lisa
john
annie
andrew
thomas
                                                                                  mark
peter
                                                                                                                            170
170
340
                                                                                   peter
annie
james
                                                                                                                            170
510
       airplane|1994
airplane|1990
                                                                                       mark
lisa
                                                                                                                            170
340
    nly showing top 20 rows
```

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

```
scala> dfagegroup.show
|usrid|transport|fromLoc|toLoc|distance|year|fare|username|age|
      4 airplane
                           RUSI
                                              200 | 1990 | 170 |
                                                                     lisa| 27|
                                  IND
                                              200 | 1992 | 170 |
200 | 1991 | 170 |
      5 airplane
6 airplane
                                  RUS
                                                                     mark | 25
                           CHN
                                                                    peter 22
james 21
                           AUS
                                  PAK
                                              200 | 1990 | 170 |
200 | 1991 | 170 |
          airplane
                                  AUS
                           RUS
                                                                     lisa| 27|
mark| 25|
      4
          airplane
                                  AUS
                           IND
                                              200 1992 170
      5|
          airplane
                           AUS
                                  IND
                                             200|1992| 170|

200|1993| 170|

200|1990| 170|

200|1990| 170|

200|1991| 170|

200|1991| 170|

200|1990| 170|

200|1994| 170|
                                                                    peter| 22|
james| 21|
      6
          airplane
                           RUS
                                  CHN
          airplane
                                  RUS
                           CHN
                                                                     lisa| 27|
mark| 25|
      4
          airplane
                           CHN
                                  PAK
      5
                                  PAK
          airplane
                           IND
          airplane
                                                                    peter 22
      6
                           PAK
                                  RUS
                                                                    james| 21|
mark| 25|
          airplane
                           CHN
                                  IND
      5 airplane
                           CHN
                                  PAK
scala> dfagegroup.count()
res38: Long = 13
scala> dfjointransp
dfjointransp
                 dfjointranspuser
scala> dfjointranspuser.filter($"age" === 20).count()
res39: Long = 0
scala> dfjointranspuser.filter($"age" === 35).count()
res40: Long = 0
scala> dfjointranspuser.filter($"age" >= 35).count()
res41: Long = 9
scala> dfjointranspuser.filter($"age" <= 20).count()
res42: Long = 10
```