# Assignment 20.1 Spark SQL 1

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

Solution:

> *val FileRDD = sc.textFile("/user/S20_Dataset_Holidays.txt")*
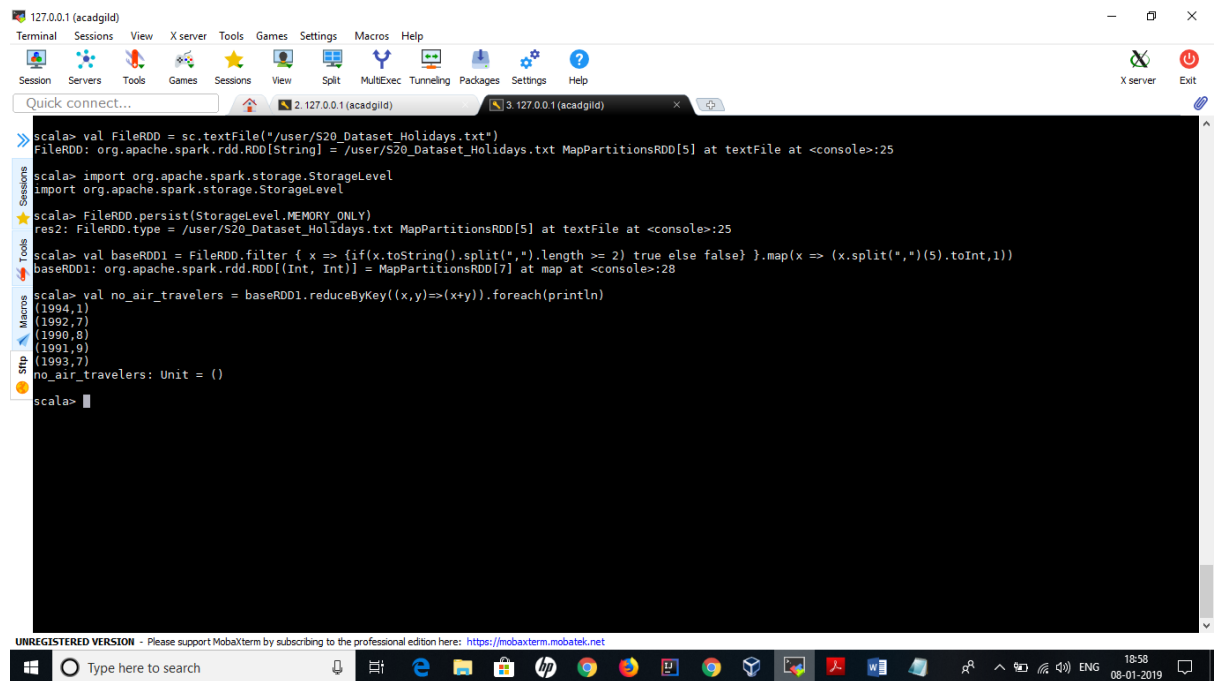>
> *import org.apache.spark.storage.StorageLevel*
>
> *FileRDD.persist(StorageLevel.MEMORY_ONLY)*
>
> *val baseRDD1 = FileRDD.filter { x => {if(x.toString().split(",").length >= 2) true else false}*
> *}.map(x => (x.split(",")(5).toInt,1))*
>
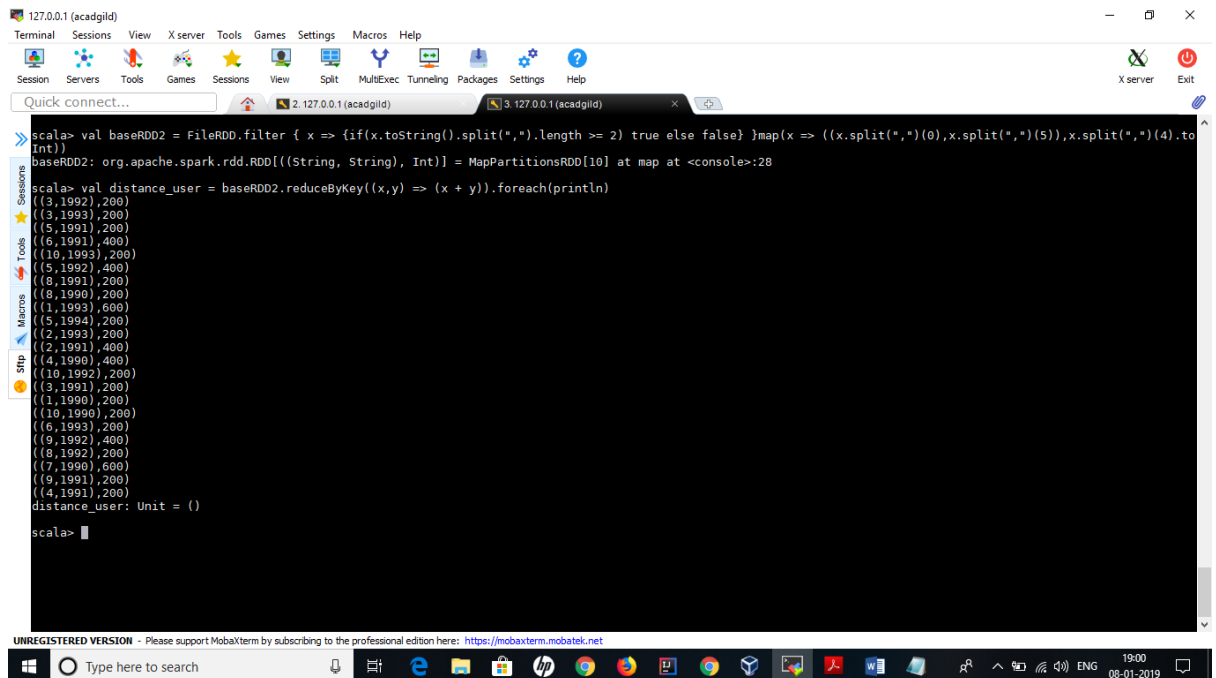> *val no_air_travelers = baseRDD1.reduceByKey((x,y)=>(x+y)).foreach(println)*

Output:



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

Solution:

> *val baseRDD2 = FileRDD.filter { x => {if(x.toString().split(",").length >= 2) true else false}*
> *}.map(x => ((x.split(",")(0),x.split(",")(5)),x.split(",")(4).toInt))*
>
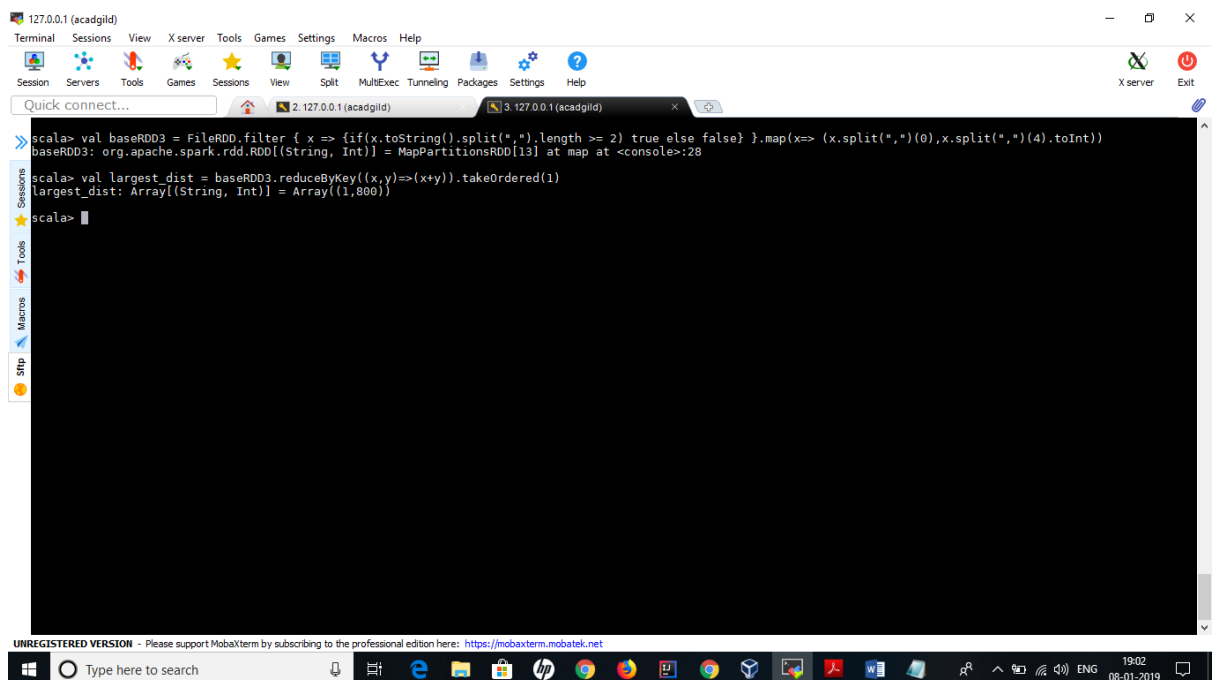> *val distance_user = baseRDD2.reduceByKey((x,y) => (x + y)).foreach(println)*

Output:



3. Which user has travelled the largest distance till date.

Solution:

*val baseRDD3 = FileRDD.filter { x => {if(x.toString().split(",").length >= 2) true else false}*
*}.map(x=> (x.split(",")(0),x.split(",")(4).toInt))*

*val largest_dist = baseRDD3.reduceByKey((x,y)=>(x+y)).takeOrdered(1)*

Output:

4. What is the most preferred destination for all users.

Solution:

*val baseRDD4 = FileRDD.filter { x => {if(x.toString().split(",").length >= 2) true else false}*
*}.map(x => (x.split(",")(2),1))*

*val dest = baseRDD4.reduceByKey((x,y)=>(x+y))*

*val dest =*
*baseRDD4.reduceByKey((x,y)=>(x+y)).takeOrdered(1)(Ordering[Int].reverse.on(_._2))*

Output:



5. Which route is generating the most revenue per year.

Solution:

*val FileRDD1 = sc.textFile("/user/S20_Dataset_Transport.txt")*

*val FileRDD2 = sc.textFile("/user/S20_Dataset_User_details.txt")*

*import org.apache.spark.storage.StorageLevel*

*FileRDD1.persist(StorageLevel.MEMORY_ONLY)*

*FileRDD2.persist(StorageLevel.MEMORY_ONLY)*


*val holidays =*
*FileRDD.map(x=>(x.split(",")(0).toInt,x.split(",")(1),x.split(",")(2),x.split(",")(3),x.split(",")(4)*
*.toInt,x.split(",")(5).toInt))*


*val transport = FileRDD1.map(x=> (x.split(",")(0),x.split(",")(1).toInt))*

*val user = FileRDD2.map(x=>(x.split(",")(0).toInt,x.split(",")(1),x.split(",")(2).toInt))*

*val holidaysmap = holidays.map(x=>x._4->(x._2,x._5,x._6))*
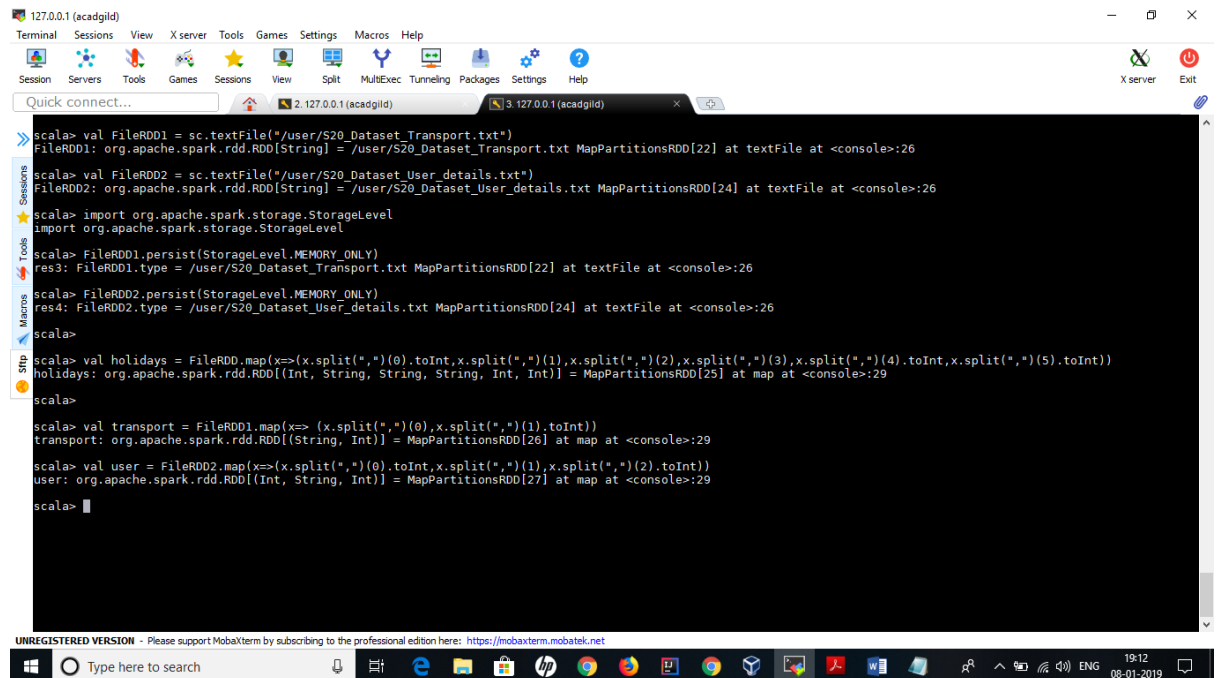
*val transportmap = transport.map(x=>x._1->x._2)*

*val join1 = holidaysmap.join(transportmap)*

*val route = join1.map(x=>(x._2._1._1->x._2._1._3)->(x._2._1._2*x._2._2))*

*val revenue = route.groupByKey().map(x=>x._2.sum->x._1)*

*val routemostrevenue = revenue.sortByKey(false).first()*

Output:

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

Solution:

*val userMap = holidays.map(x => x._4 -> (x._1,x._5,x._6))*

*val amount = userMap.join(transportmap)*

*val spend = amount.map(x => (x._2._1._1, x._2._1._3) -> (x._2._1._2 * x._2._2))*

*val total = spend.groupByKey().map(x => x._1 -> x._2.sum)*

*total.foreach(println)*

Output:

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

Solution:

*val AgeMap = user.map(x=>x._1->*

*{*

*if(x._3<20)*

*"20"*

*else if(x._3>35)*

*"35"*

*else "20-35"*

*})*

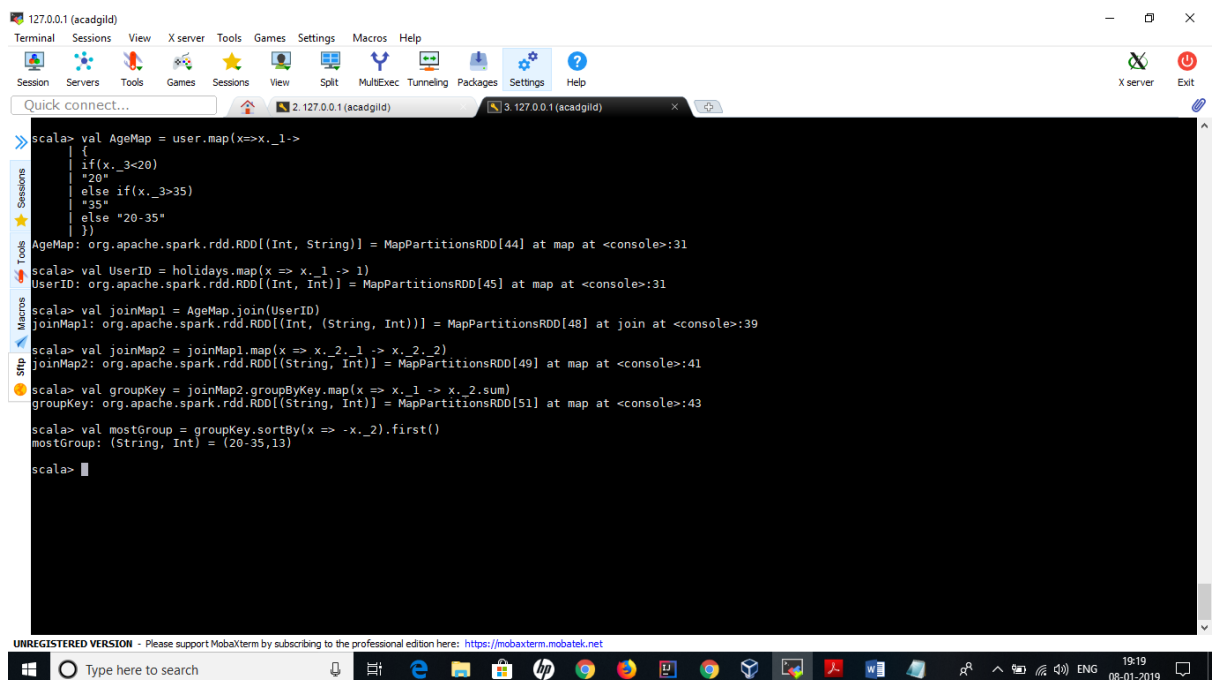*val UserID = holidays.map(x => x._1 -> 1)*

*val joinMap1 = AgeMap.join(UserID)*

*val joinMap2 = joinMap1.map(x => x._2._1 -> x._2._2)*

*val groupKey = joinMap2.groupByKey.map(x => x._1 -> x._2.sum)*

*val mostGroup = groupKey.sortBy(x => -x._2).first()*

Output: