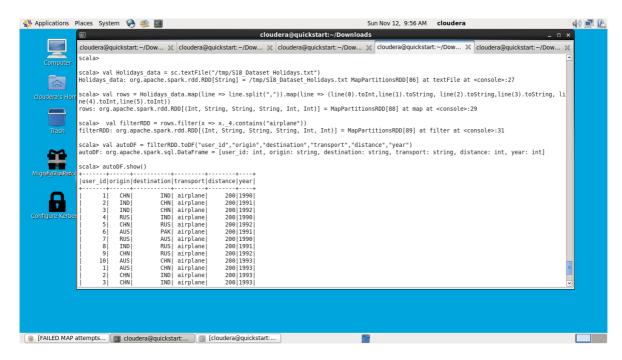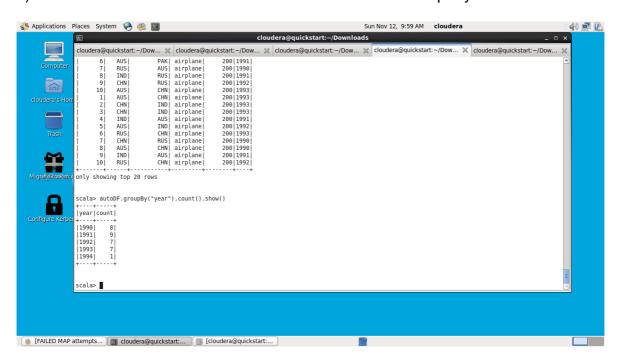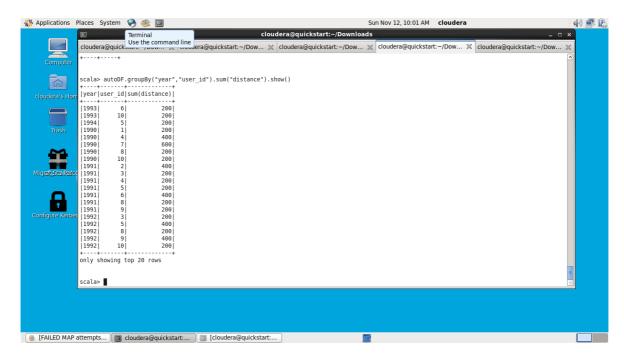Convert text file to data frame to run queries:
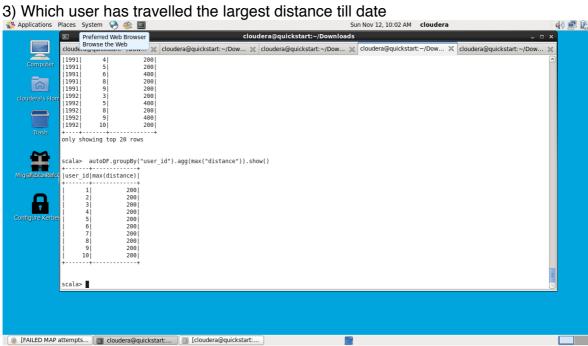


# 1) What is the distribution of the total number of air-travelers per year



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

```
scala> autoDF.groupBy("year","user_id").sum("distance").show()
+----+-------+-------------+
|year|user_id|sum(distance)|
+----+-------+-------------+
|1993|      6|          200|
|1993|     10|          200|
|1994|      5|          200|
|1990|      1|          200|
|1990|      4|          400|
|1990|      7|          600|
|1990|      8|          200|
|1990|     10|          200|
|1991|      2|          400|
|1991|      3|          200|
|1991|      4|          200|
|1991|      5|          200|
|1991|      6|          400|
|1991|      8|          200|
|1991|      9|          200|
|1992|      3|          200|
|1992|      5|          400|
|1992|      8|          200|
|1992|      9|          400|
|1992|     10|          200|
+----+-------+-------------+
only showing top 20 rows

scala>
```

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

```
|1991|      4|          200|
|1991|      5|          200|
|1991|      6|          400|
|1991|      8|          200|
|1991|      9|          200|
|1992|      3|          200|
|1992|      5|          400|
|1992|      8|          200|
|1992|      9|          400|
|1992|     10|          200|
+----+-------+-------------+
only showing top 20 rows

scala>  autoDF.groupBy("user_id").agg(max("distance")).show()
+-------+-------------+
|user_id|max(distance)|
+-------+-------------+
|      1|          200|
|      2|          200|
|      3|          200|
|      4|          200|
|      5|          200|
|      6|          200|
|      7|          200|
|      8|          200|
|      9|          200|
|     10|          200|
+-------+-------------+

scala>
```

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

```
                       |
                       |  ;
<console>:3: error: ')' expected but ';' found.
                       ;
                       ^

scala> autoDF.groupby("distance").agg(max("distance")).show
<console>:36: error: value groupby is not a member of org.apache.spark.sql.DataFrame
              autoDF.groupby("distance").agg(max("distance")).show
                     ^

scala> autoDF.groupBy("distance").agg(max("distance")).show
+--------+-------------+
|distance|max(distance)|
+--------+-------------+
|     200|          200|
+--------+-------------+


scala> autoDF.groupBy("destination").count().show()
+-----------+-----+
|destination|count|
+-----------+-----+
|        CHN|    7|
|        RUS|    6|
|        PAK|    5|
|        AUS|    5|
|        IND|    9|
+-----------+-----+


scala>
```