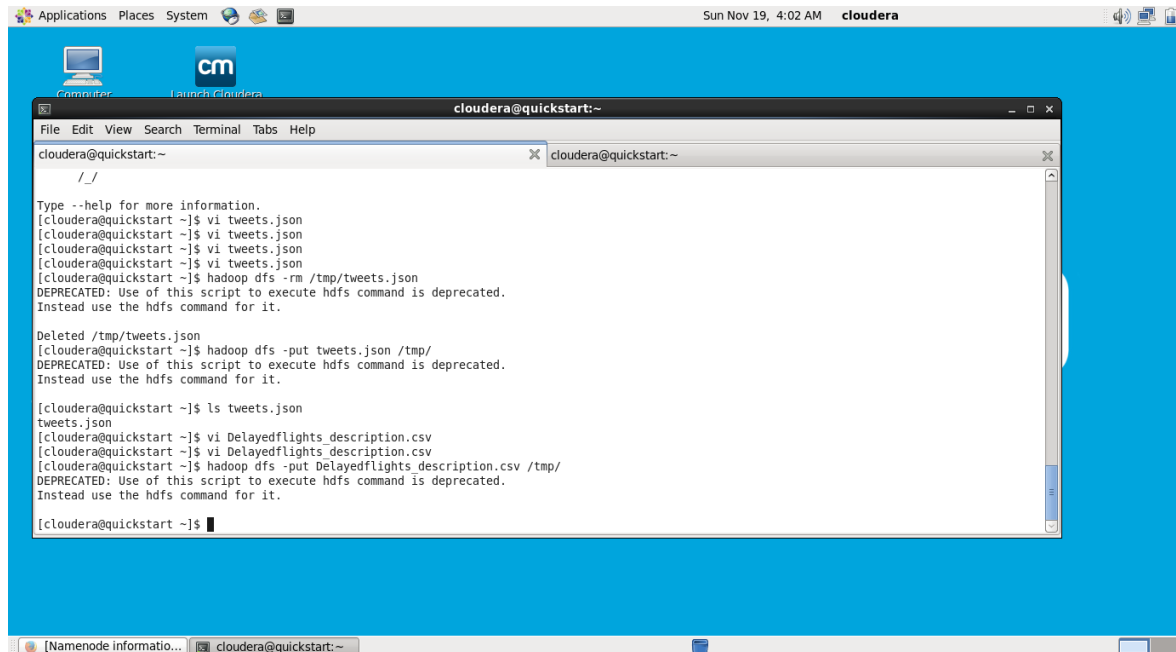1. Copy delayed_flights.csv to HDFS using
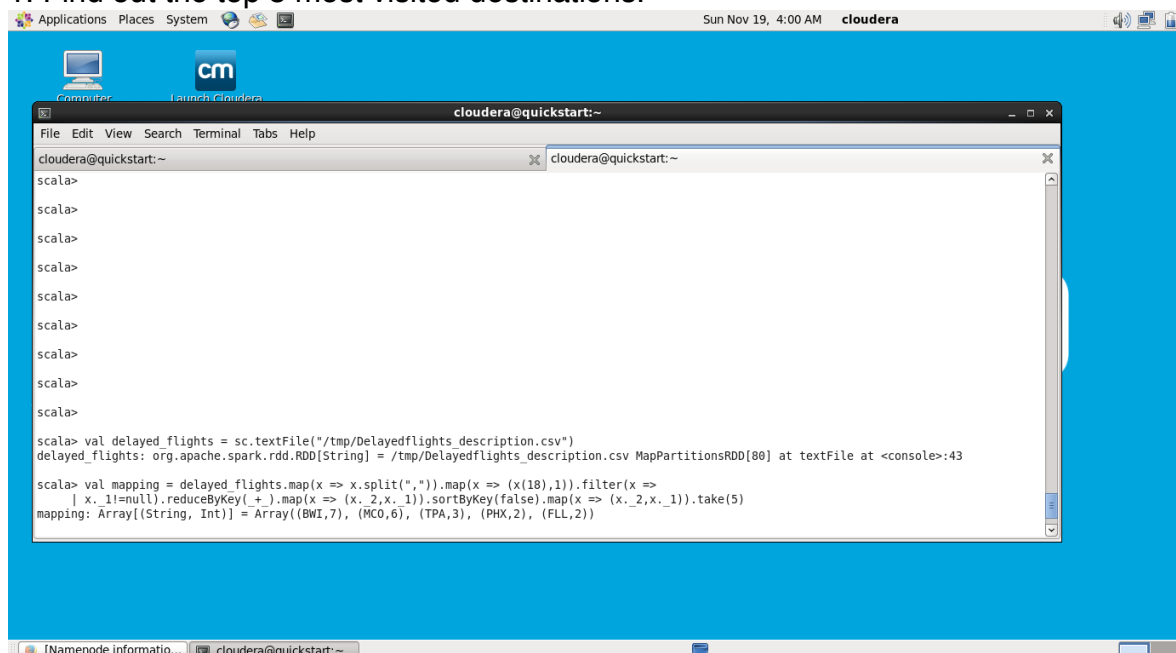
# NOTE: Copied only half the file as not able to copy full file from local system to cloudera, so some results are empty as there is no data.



Problem statement 1:
1. Find out the top 5 most visited destinations.



Problem statement 2

Which month has seen the most number of cancellations due to bad weather?



Problem statement 3:
Top ten origins with the highest AVG departure delay



Problem statement 4
Which route (origin & destination) has seen the maximum diversion?

cloudera@quickstart:~

File  Edit  View  Search  Terminal  Tabs  Help

cloudera@quickstart:~                                                    cloudera@quickstart:~

```
scala>

scala>

scala>

scala>

scala>

scala>

scala>

scala>

scala> val delayed_flights = sc.textFile("/tmp/Delayedflights_description.csv")
delayed_flights: org.apache.spark.rdd.RDD[String] = /tmp/Delayedflights_description.csv MapPartitionsRDD[108] at textFile at <console>:43

scala> val diversion = delayed_flights.map(x => x.split(",")).filter(x => ((x(24).equals("1")))).map(x =>
     | ((x(17)+","+x(18)),1)).reduceByKey(_+_).map(x => (x._2,x._1)).sortByKey(false).map(x =>
     | (x._2,x._1)).take(10).foreach(println)
diversion: Unit = ()

scala>
```

[Namenode informatio...]   cloudera@quickstart:~