# BIG DATA MODULE END PAPER

## ROLL NO.240840325072
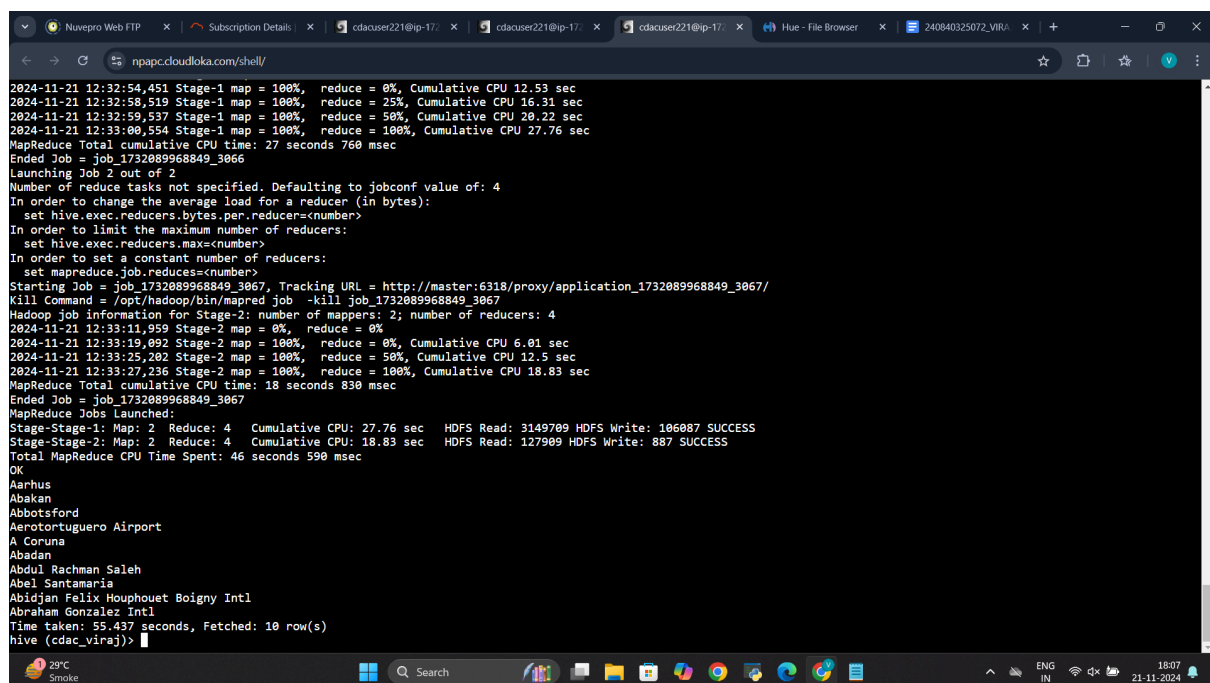## NAME:VIRAJ SHIRKE

HIVE:
QUESTION 1:

1.
select distinct a.name from airport a join routes r on a.iata = r.src_airport_iata limit 10;



2.select count(*) from airlines a join routes r on r.airline_id=a.airline_id group by a.airline_id order by count(*) desc limit 3;

3.select count(distinct(equipment)) from routes;



QUESTION 2:

1.create table routes_partitioned( airline_iata string ,airline_id int ,src_airport_iata string,src_airport_id int , dest_airport_id int , codeshare st
ring , stops int , equipment string)
partitioned by (dest_airport_iata)
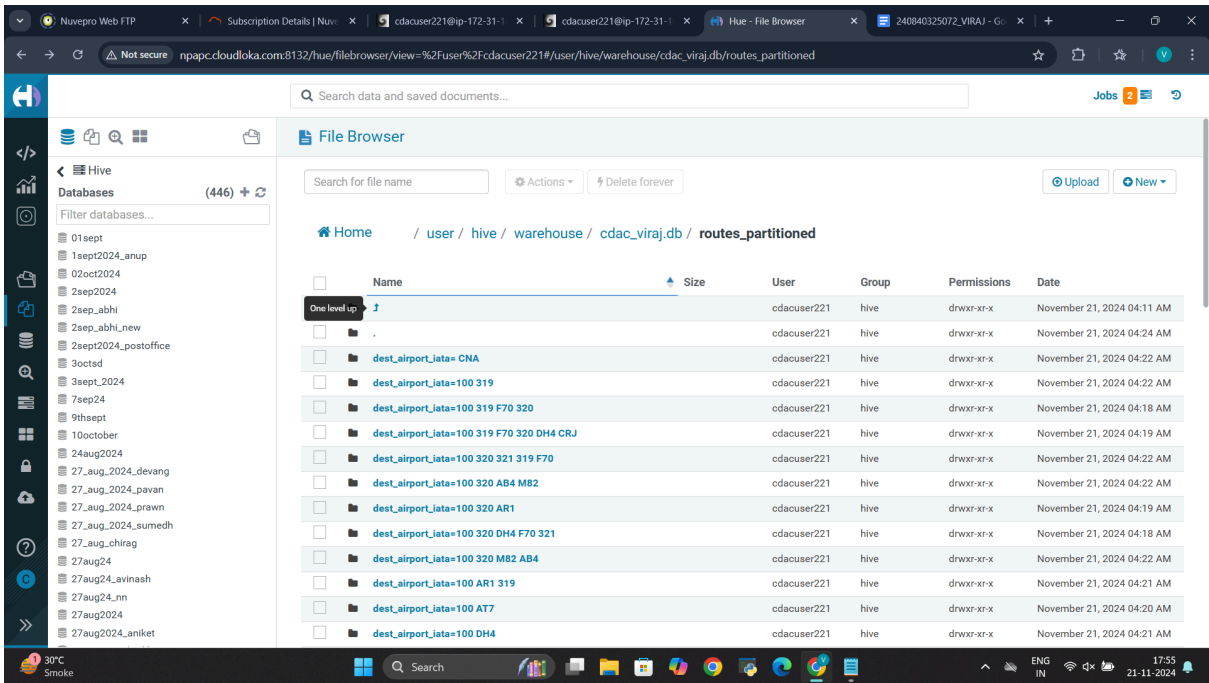row format delimited
fields terminated by ','

stored as textfile;



2.
Insert overwrite table routes_partitioned PARTITION(dest_airport_iata) select
r.airline_iata ,r.airline_id ,r.src_airport_iata ,r.src_airport_id ,r.dest_airport_iata,
r.dest_airport_id , r.codeshare ,r.stops,r.equipment from routes r DISTRIBUTE By
dest_airport_iata;

```
hive (cdac_viraj)> desc routes_partitioned;
OK
airline_iata          string
airline_id            int
src_airport_iata      string
src_airport_id        int
dest_airport_id       int
codeshare             string
stopd                 int
equipment             string
dest_airport_iata     string

# Partition Information
# col_name             data_type              comment
dest_airport_iata      string
Time taken: 0.323 seconds, Fetched: 13 row(s)
```

3.


PYSPARK :
QUESTION 1:
1.
filterRDD = splitRDD.map(lambda a:(int(a[3]) > 20000 and int(a[3]) < 50000)



2.
quarterRDD = splitRDD.map(lambda a: (str[0],str[1]))
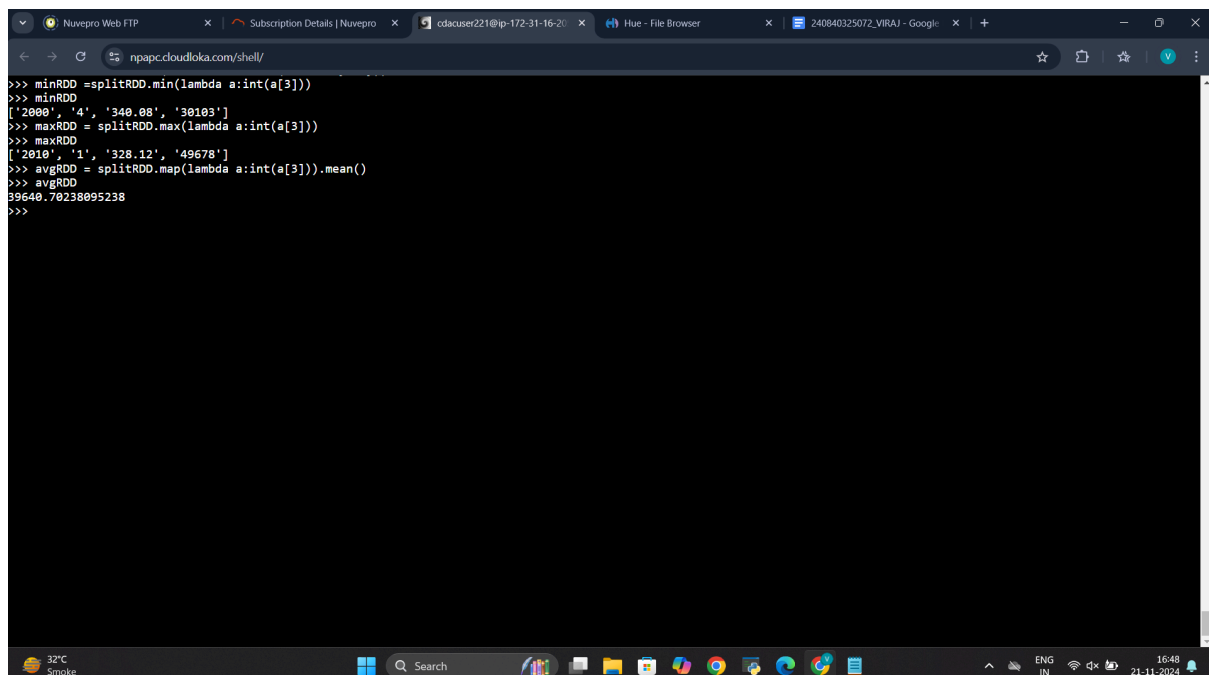for line in  quarterRDD .collect():

print(line)



QUESTION 2:

1.
Min,Max,Avg:
minRDD =splitRDD.min(lambda a:int(a[3]))
maxRDD = splitRDD.max(lambda a:int(a[3]))
avgRDD = splitRDD.max(lambda a:int(a[3])).mean()
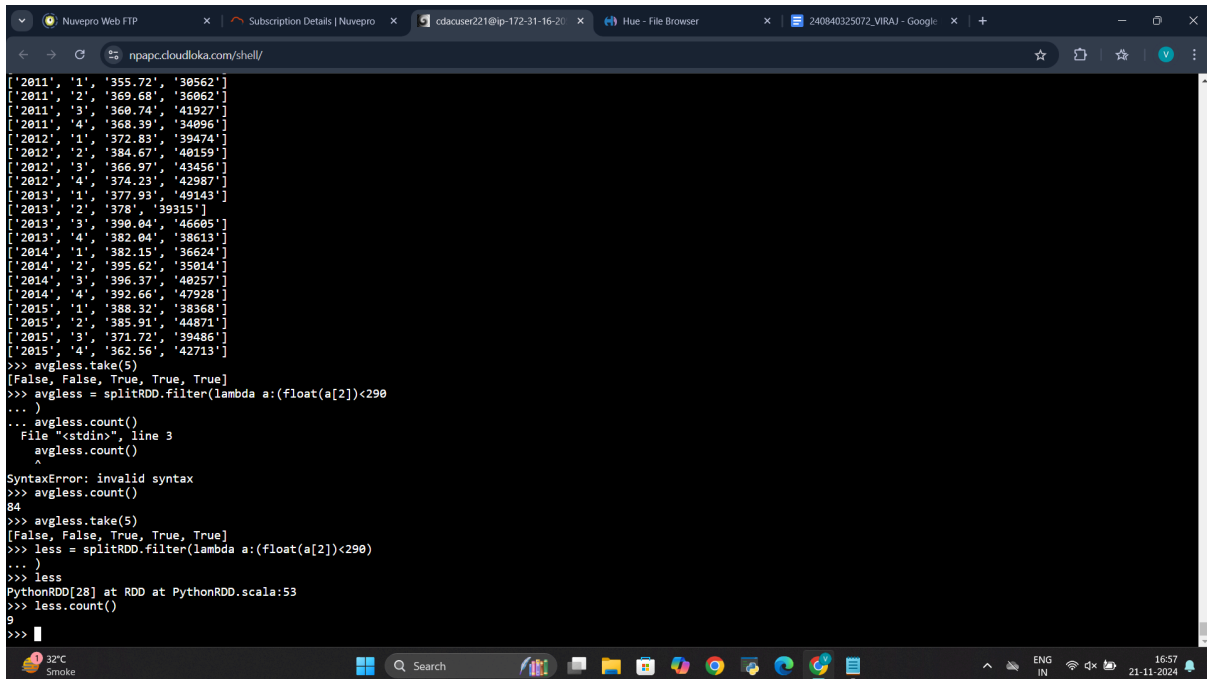
## 2.

```
less = splitRDD.filter(lambda a:(float(a[2])<290)
```



## 3.



## 4.

```
year = splitRDD.map(lambda a:(a[0],1))
```

countyear = year.reduceByKey(lambda a,b:a+b)



5.
totalrev = splitRDD.map(lambda a: ((int(a[1])), float(a[2])*int(a[3])))
high = totalrev.reduceByKey(lambda a,b:a+b)
max = high.sortBy(lambda a:-a[1])