Question 1

1.

```
select name
from airport ap join routes r
on ap.iata=r.src_airport_iata
join airport ap2
on ap2.iata=r.dest_airport_iata limit 6;
```

2.

Question 2

**1.partition table**

```
hive
set hive.cli.print.current.db=true;
use cdac_gauri;

create table q2_routes_table
(airline_iata string,
airline_id int,
src_airport_id int,
dest_airport_iata string,
dest_airport_id int,
codeshare string,
stops int,
equipment string)
partitioned by(src_airport_iata string)
row format delimited
fields terminated by ','
stored as textfile;
```
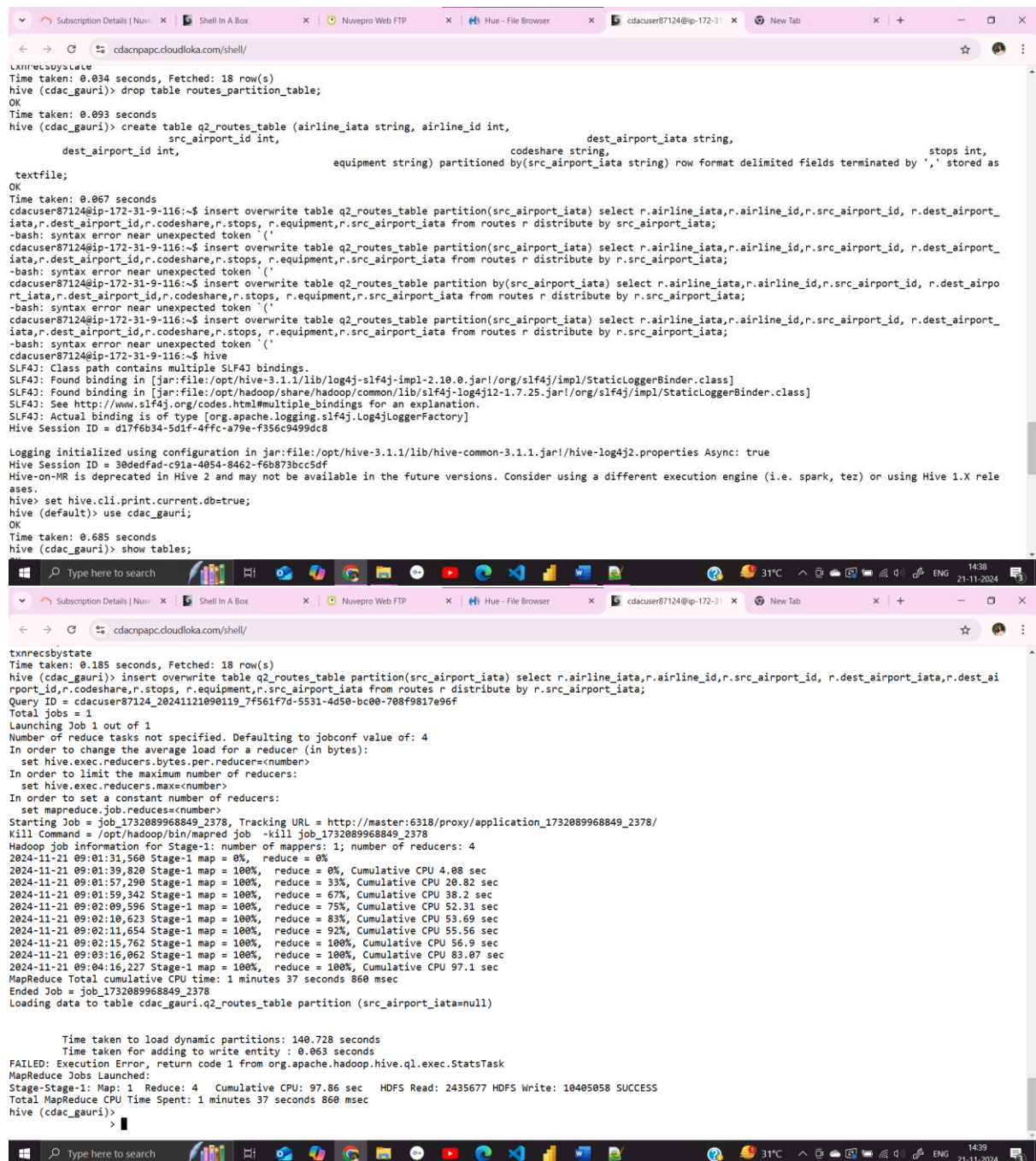
insert overwrite table q2_routes_table partition(src_airport_iata)

select r.airline_iata,r.airline_id,r.src_airport_id,

r.dest_airport_iata,r.dest_airport_id,r.codeshare,r.stops,

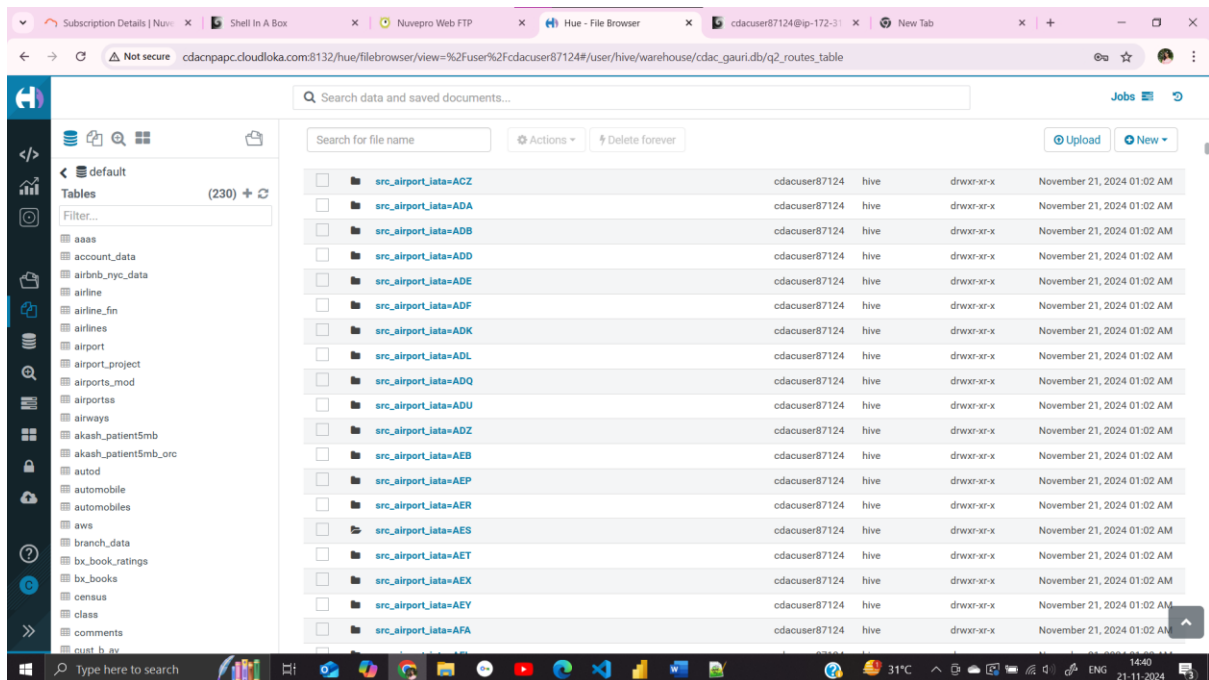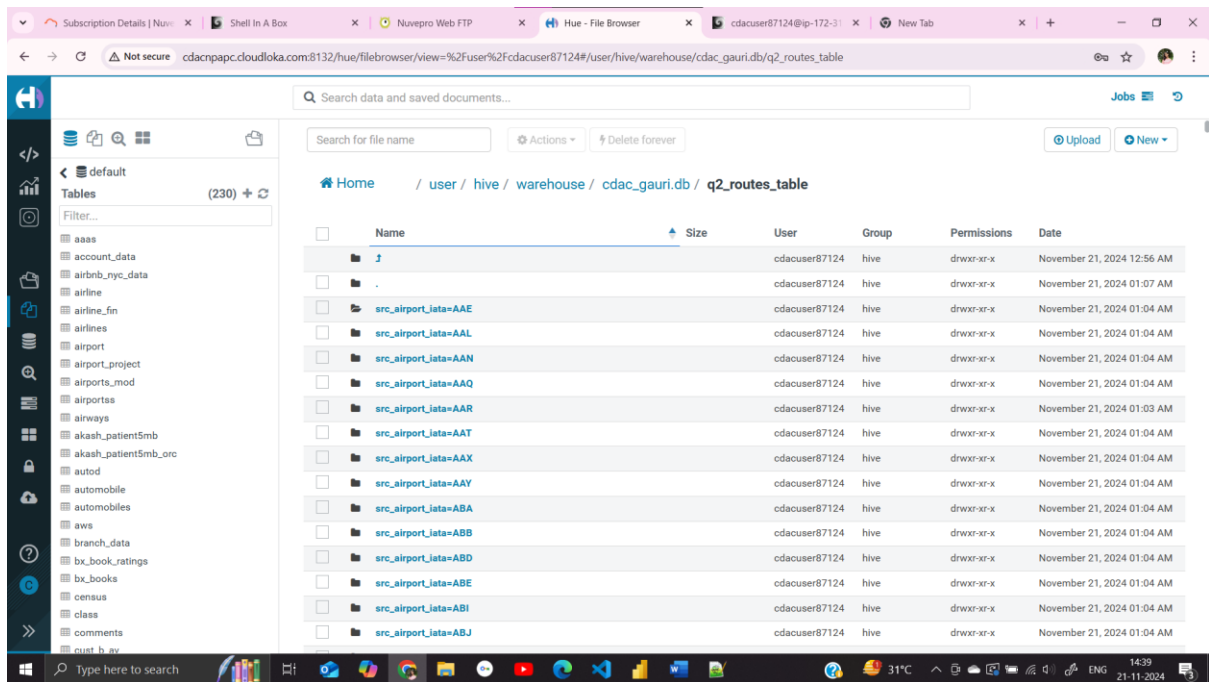r.equipment,r.src_airport_iata from routes r distribute by r.src_airport_iata;

**2.Insert data into partition table for source airport = "JFK"**

create table q2_routes_table2

(airline_iata string,

airline_id int,

src_airport_id int,

dest_airport_iata string,

dest_airport_id int,

codeshare string,

stops int,

equipment string)

partitioned by(src_airport_iata string)

row format delimited

fields terminated by ','

stored as textfile;


insert overwrite table q2_routes_table partition(src_airport_iata)

select r.airline_iata,r.airline_id,r.src_airport_id,

r.dest_airport_iata,r.dest_airport_id,r.codeshare,r.stops,

r.equipment,r.src_airport_iata from routes r distribute by r.src_airport_iata="JFK";



Shell is not working

3.

**select * from q2_routes_table where src_airport_iata="LAX";**

```
txnrecsbystate
Time taken: 0.165 seconds, Fetched: 19 row(s)
hive (cdac_gauri)> select * from q2_routes_table where src_airport_iata="LAX";
OK
HA      2688    3484    HNL     3728            0       332     LAX
HA      2688    3484    OGG     3456            0       763     LAX
HU      2660    3484    PVG     3406    Y       0       777     LAX
AZ      596     3484    FCO     1555            0       772     LAX
AZ      596     3484    CDG     1382    Y       0       388 772 LAX
Y4      5325    3484    AGU     1785            0       320     LAX
Y4      5325    3484    GDL     1804            0       320 319 LAX
Y4      5325    3484    MEX     1824            0       320     LAX
Y4      5325    3484    MLM     1821            0       320     LAX
Y4      5325    3484    UPN     1835            0       320 319 LAX
IB      2822    3484    LHR     507     Y       0       744 77W 388         LAX
IB      2822    3484    MAD     1229            0       346 342 LAX
Y4      5325    3484    ZCL     1855            0       320     LAX
AY      2350    3484    MAD     1229            0       346 343 LAX
AY      2350    3484    LHR     507             0       744 388 LAX
US      5265    3484    ABQ     4019    Y       0       CRJ CR7 LAX
JJ      4867    3484    LIM     2789    Y       0       763     LAX
US      5265    3484    AKL     2006    Y       0       773 772 LAX
US      5265    3484    AUS     3673            0       M83 738 LAX
US      5265    3484    BDL     3825            0       738     LAX
US      5265    3484    BNA     3690            0       738     LAX
US      5265    3484    BOS     3448            0       757 738 LAX
US      5265    3484    CLT     3876            0       321     LAX
US      5265    3484    CMH     3759            0       738     LAX
US      5265    3484    DCA     3520            0       738     LAX
US      5265    3484    DEN     3751            0       CR7     LAX
US      5265    3484    DFW     3670    Y       0       738 763 757 M83 LAX
US      5265    3484    ELP     3559    Y       0       CRJ     LAX
US      5265    3484    EUG     4099    Y       0       CRJ     LAX
US      5265    3484    FAT     3687    Y       0       CRJ     LAX
US      5265    3484    GRU     2564            0       777     LAX
US      5265    3484    GUA     1767            0       320     LAX
US      5265    3484    HNL     3728            0       757     LAX
JL      2987    3484    LIM     2789    Y       0       763     LAX
JL      2987    3484    NRT     2279            0       773     LAX
US      5265    3484    IAD     3714            0       738     LAX
US      5265    3484    TAH     3550            0       CR7     LAX
```

**SPARK**

**Q1.**

**q1**

data=spark.textFile("/user/cdacuser87124/spark/airlines_data.csv")

header=data.first()

eliminate=data.map(lambda a: a!= header )

map_data=eliminate.map(lambda a :a.split(","))

**Q2 dataframe**

spark=SparkSession.builder.appName("airline").getOrCreate()

df=spark.read.csv("/user/cdacuser87124/spark/airlines_data.csv",header=True,inferSchema=True)

from pyspark.sql.functions import *;

**1.**

df.agg(max("Avg_rev_per_seat").alias("max_avg_rev")).show()

df.agg(min("Avg_rev_per_seat").alias("min_avg_rev")).show()

df.agg(avg("Avg_rev_per_seat").alias("avg_rev")).show()x



**2.**

new_df=df.filter(df.Avg_rev_per_seat > "290")

new_df.count()  ----75

**3.**

```
df.groupBy("Quarter").agg(sum("booked_seats").alias("total_booked_seats")).show()
```



**4.**

**5.**

```
df.select("Year").distinct().show()
```

```
|2011|
|2008|
|1999|
+----+
only showing top 20 rows

>>>
>>>
>>>
>>>
>>>
>>>
>>>
>>> df.select("Year").distinct().show()
+----+
|Year|
+----+
|2003|
|2007|
|2015|
|2006|
|2013|
|1997|
|2014|
|2004|
|1996|
|1998|
|2012|
|2009|
|1995|
|2001|
|2005|
|2000|
|2010|
|2011|
|2008|
|1999|
+----+
only showing top 20 rows

>>>
```

5.

df.groupBy("Year").agg(sum(col("Avg_rev_per_seat")*col("booked_seats")).alias("cumulative_avg_rev")).show()

```
|2012|    6.2199127728E7|
|2009|    4.674644659E7|
|1995|    4.349424322E7|
|2001| 5.553377999999999E7|
|2005|    4.637678624E7|
|2000|5.234292655000004E7|
|2010|    5.486152129E7|
|2011|    5.188828622E7|
|2008|5.7653170760000005E7|
|1999|    4.875771448E7|
+----+--------------------+
only showing top 20 rows

>>> df.groupBy("Year").agg(sum(col("Avg_rev_per_seat")* col("booked_seats")).alias("cumulative_avg_rev")).show()
+----+--------------------+
|Year|  cumulative_avg_rev|
+----+--------------------+
|2003|    4.927321083E7|
|2007|    5.730921607E7|
|2015|    6.237899057E7|
|2006|5.0437898419999994E7|
|2013|    6.636320871E7|
|1997|    4.538523616E7|
|2014| 6.262417585000001E7|
|2004|5.0631364949999996E7|
|1996|    4.635877803E7|
|1998|    4.203571778E7|
|2012|    6.219912728E7|
|2009|    4.674644659E7|
|1995|    4.349424322E7|
|2001| 5.553377999999999E7|
|2005|    4.637678624E7|
|2000|5.234292655000004E7|
|2010|    5.486152129E7|
|2011|    5.188828622E7|
|2008|5.7653170760000005E7|
|1999|    4.875771448E7|
+----+--------------------+
only showing top 20 rows

>>>
```