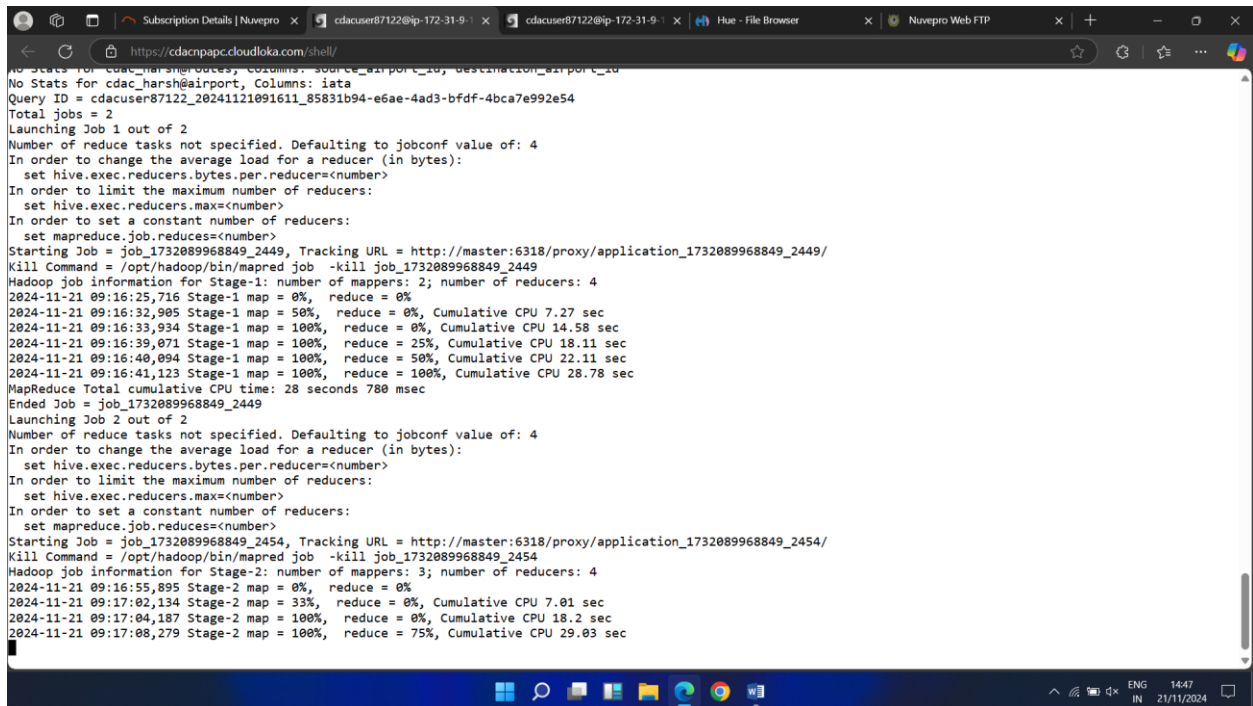Q1.

1.

```
select src.name from airport src  join routes r on src.iata =
r.source_airport_id join airport dest on
r.destination_airport_id=dest.iata wher
e src.iata=dest.iata limit 10;
```



Q1.

3.

```
select   count(*) as no_of_count , al.name   from routes r   join
airline al  on al.iata = r.source_airport_id    group by al.name
order by no_of_
count desc limit 1;
```

Q1.

2. select  count(*) as no_of_count  ,r.equipment   from routes r  join airline al  on al.iata = r.source_airport_id   group by al.name ,r

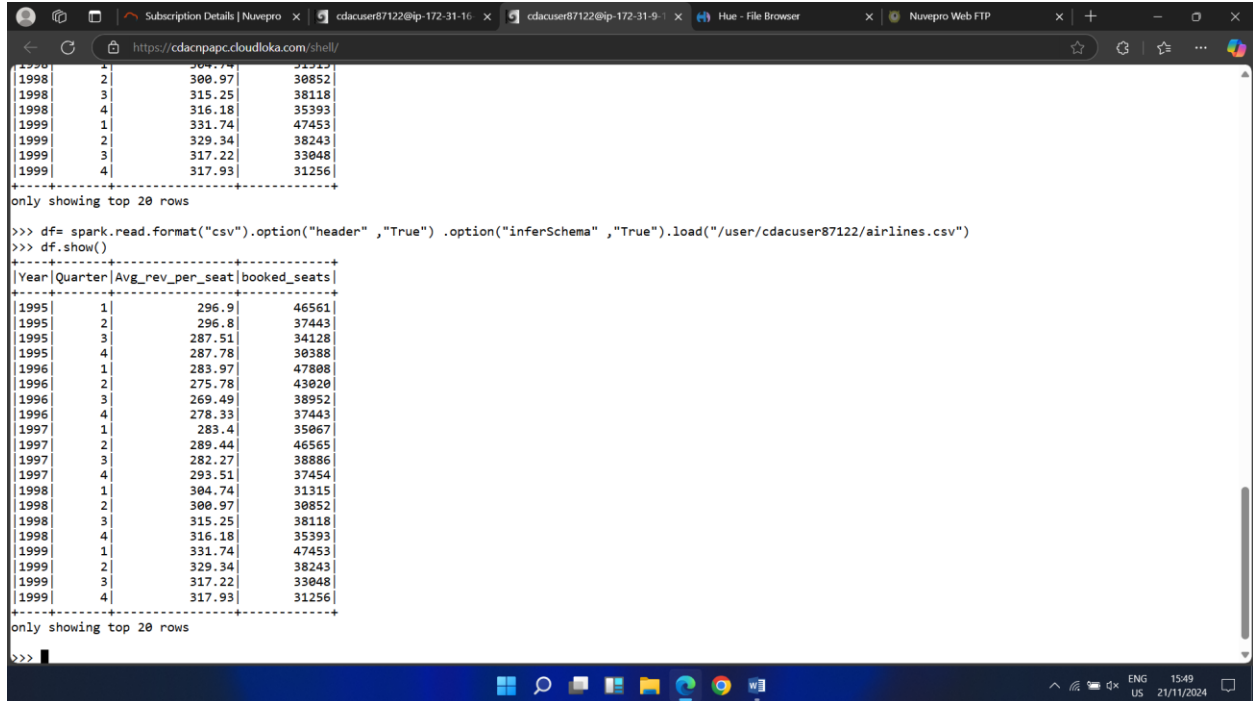.equipment order by no_of_count desc limit 1;

Data Frame

df= spark.read.format("csv").option("header" ,"True") .option("inferSchema" ,"True").load("/user/cdacuser87122/airlines.csv")

>>> df.show()



Q.2

1.

Find_insight =df.agg(min("avg_rev_per_seat") ,max("avg_rev_per_seat")avg("avg_rev_per_seat")")

Find_insight

2.

df.groupBy("avg_rev_per_seat">290).agg(count("avg_rev_per_seat")

3.

df.groupBy("Quarter")agg(sum("booked_seat").show()

4.

df.groupBy("Year").show()

5.

df.groupBy("Year").agg(sum("avg_rev_per_seat").limit(10)

RDD

1.