## Question 1:

```
1 select ap.name ,r.source_airport , r.destination_airport from
route1 r join air
port1 ap on r.source_airport = ap.iata join airport1 src on
src.iata=r.source_airport join airport1
desc on desc.iata = r.destination_airport where r.source_airport =
r.destination_airport limit 10 ;
```

**2.**`select equipment , count(*) as total from routes group by equipment order by total desc limit 10;`



**3:**
`select airline ,count(*)as total from route1 group by airline order by total desc limit 10;`



Question 2:

**1:** 
```
create table route_partition (airline string,airline_id
int,source_airport_id int,destination_airport string,
destination_airport_id int,codeshare
string,  stops int   ,equipment string ) partitioned by
(source_airport string) row format delimited fields terminated by
"," stored as textfile;

insert overwrite table route_partition partition(source_airport)
select airline ,airline_id ,source_airport, source_airport_id ,
destination_airpor
t,destination_airport_id ,codeshare , stops , equipment from
route1 distribute by source_airport;
```

3.
Select * from routes_partition where source_airport = 'LAX'



**Spark:**

Question 1:

2

```
df=spark.read.format('csv').option("header",True).option("inferSch
ema",True).load("/user/cdacuser82313/airlinedata.csv")
>>> df.count()
```

84

```
df.select("year").distinct().orderBy("year",ascending=True).show()
```



```
df.select("year").distinct().count()
```

```
|2000|              154376|
|2010|              163741|
|2011|              142647|
|2008|              166897|
|1999|              150000|
+----+--------------------+
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

>>> df.select("year").distinct().count()
21
>>>
```

Question 2:

1
```
df.groupBy("year").agg((F.min("avg_rev_per_seat").alias("minimun")
),(F.max("avg_rev_per_seat").alias("maximun")),(F.avg("avg_rev_per
_seat").alias("average"))).order
By("year",ascending=True).show()
```



```
File "<stdin>", line 1, in <module>
  File "/opt/spark-3.1.2/python/pyspark/sql/group.py", line 118, in agg
    jdf = self._jgd.agg(exprs[0]._jc,
  File "/opt/spark-3.1.2/python/lib/py4j-0.10.9-src.zip/py4j/java_gateway.py", line 1304, in __call__
  File "/opt/spark-3.1.2/python/pyspark/sql/utils.py", line 117, in deco
    raise converted from None
pyspark.sql.utils.AnalysisException: cannot resolve '`avg_rev_per_year`' given input columns: [Avg_rev_per_seat, Quarter, Year, booked_seats];
'Aggregate [year#16], [year#16, min('avg_rev_per_year) AS minimun#60, max('avg_rev_per_year) AS maximun#62, avg('avg_rev_per_year) AS average#64]
+- Relation[Year#16,Quarter#17,Avg_rev_per_seat#18,booked_seats#19] csv

>>> df.groupBy("year").agg((F.min("avg_rev_per_seat").alias("minimun")),(F.max("avg_rev_per_seat").alias("maximun")),(F.avg("avg_rev_per_seat").alias("average"))).order
By("year",ascending=True).show()
+----+-------+-------+------------------+
|year|minimun|maximun|           average|
+----+-------+-------+------------------+
|1995| 287.51|  296.9|          292.2475|
|1996| 269.49| 283.97|          276.8925|
|1997| 282.27| 293.51|           287.155|
|1998| 300.97| 316.18|           309.285|
|1999| 317.22| 331.74|          324.0575|
|2000| 336.66| 340.23|          339.0325|
|2001| 299.81| 347.69|          319.7975|
|2002|  303.3| 320.02|           312.525|
|2003| 312.39| 319.19|          315.4675|
|2004| 296.54| 320.23|           305.875|
|2005| 301.39| 314.76|           307.185|
|2006| 318.16| 341.58|             328.3|
|2007| 317.84| 329.77|            325.14|
|2008| 333.29| 358.93|          346.1575|
|2009| 301.82| 319.85|            310.61|
|2010| 328.12| 340.72|          335.8325|
|2011| 355.72| 369.68|363.63250000000005|
|2012| 366.97| 384.67|           374.675|
|2013| 377.93| 390.04|          382.0025|
|2014| 382.15| 396.37|             391.7|
+----+-------+-------+------------------+
only showing top 20 rows

>>>
```
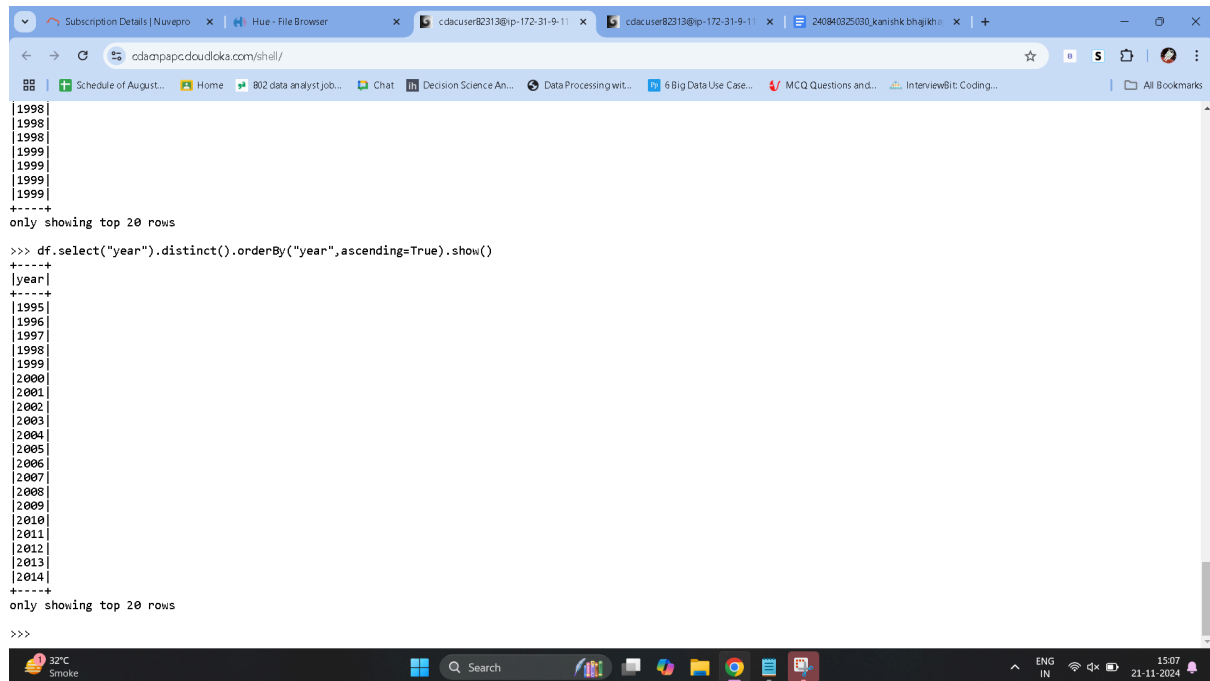
**3**

```
df.groupBy("year").agg(F.sum("booked_seats").alias("total_for_all_
quarter")).show()
```



**4**

```
>>>df.select("year").distinct().show()
```



**5**

```
>>> df.groupBy("year").agg(F.sum(F.col("avg_rev_per_seat")*
F.col("booked_seats")).alias("total_revenue")).orderBy("year",asce
```

```
nding=False).show()
```

```
>>> df.groupBy("year").agg(F.sum(F.col("avg_rev_per_seat")* F.col("booked_seats")).alias("total_revenue")).orderBy("year",ascending=False).show()
+----+--------------------+
|year|       total_revenue|
+----+--------------------+
|2015|       6.237899057E7|
|2014| 6.262417585000001E7|
|2013|       6.636320871E7|
|2012|       6.219912728E7|
|2011|       5.188828622E7|
|2010|       5.486152129E7|
|2009|       4.674644659E7|
|2008|5.765317076000005E7|
|2007|       5.730921607E7|
|2006|5.0437898419999994E7|
|2005|       4.637678624E7|
|2004|5.0631364949999996E7|
|2003|       4.927321083E7|
|2002|        4.74991465E7|
|2001| 5.553377999999999E7|
|2000|5.234292655000004E7|
|1999|       4.875771448E7|
|1998|       4.203571778E7|
|1997|       4.538523616E7|
|1996|       4.635877803E7|
+----+--------------------+
only showing top 20 rows

>>>
```