# HIVE
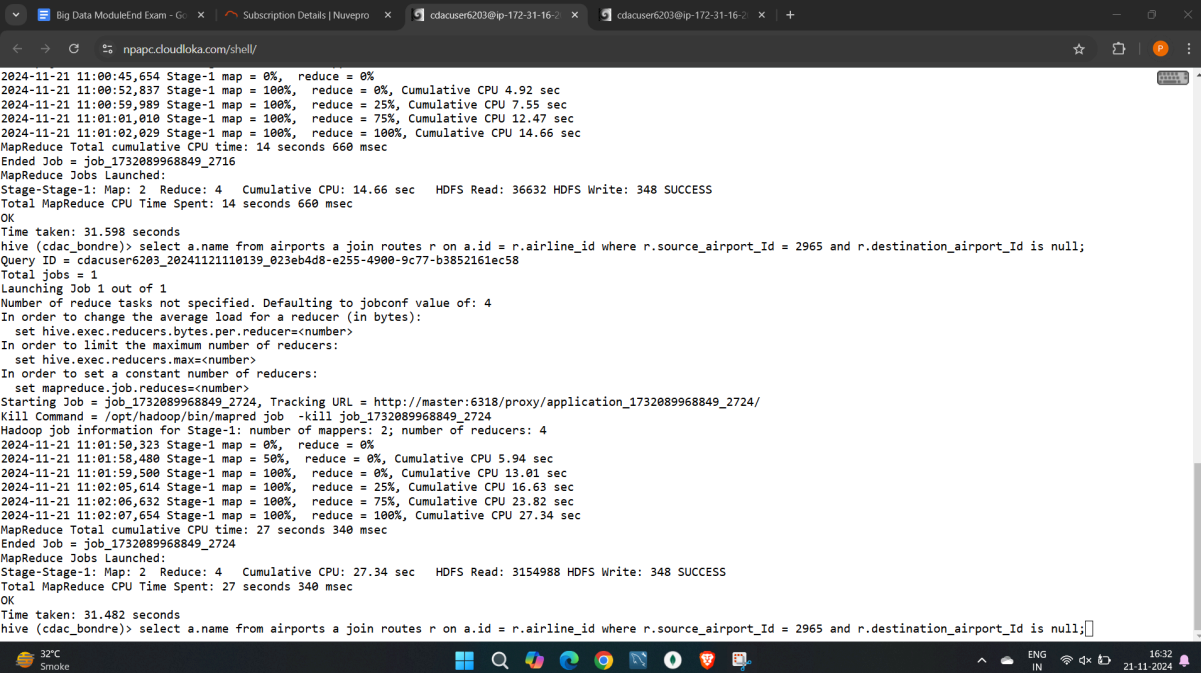
## Q1.

1.
→ select a.name from airports a join routes r on a.id = r.airline_id where r.source_airport_Id = 2965 and r.destination_airport_Id is null;



2.
→ select distinct(a.name) from airlines a join routes r on a.id = r.airline_id order by max(source_airport_id) desc limit 3 ;

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

Q2
1.
→ create table routes_partioned (airline string,airline_id int,source_airport string,source_airport_id int,destination_airport_id int, codeshare string,stops int, equipment string)
                    > partitioned by (destination_airport string)
                    > row format delimited
                    > fields terminated by ','
                    > stored as textfile;



2.
→  insert overwrite table routes_partioned select airline,airline_id,source_airport,source_airport_id,destination_airport_id ,codeshare,stops,equipments from routes;

**The query got stucked in between the i have shared the query with ScreenShot**

```
Time taken: 0.062 seconds
hive (cdac_bondre)> insert overwrite table routes_partioned select airline,airline_id,source_airport,source_airport_id,destination_airport_id ,codeshare,stops,equipment
s from routes;
FAILED: SemanticException [Error 10004]: Line 1:139 Invalid table alias or column reference 'equipments': (possible column names are: airline, airline_id, source_airpor
t, source_airport_id, destination_airport, destination_airport_id, codeshare, stops, equipment)
hive (cdac_bondre)> insert overwrite table routes_partioned select airline,airline_id,source_airport,source_airport_id,destination_airport_id ,codeshare,stops,equipment
 from routes;
FAILED: SemanticException [Error 10044]: Line 1:23 Cannot insert into target table because column number/types are different 'routes_partioned': Table insclause-0 has 9
 columns, but query has 8 columns.
hive (cdac_bondre)> insert overwrite table routes_partioned select airline,airline_id,source_airport,source_airport_id,destination_airport_id ,codeshare,stops,equipment
 destination_airport from routes;
FAILED: SemanticException [Error 10044]: Line 1:23 Cannot insert into target table because column number/types are different 'routes_partioned': Table insclause-0 has 9
 columns, but query has 8 columns.
hive (cdac_bondre)> insert overwrite table routes_partioned select airline,airline_id,source_airport,source_airport_id,destination_airport_id ,codeshare,stops,equipment
, destination_airport from routes;
Query ID = cdacuser6203_20241121113409_4191a38e-0c0d-4543-be0d-a28d7c0b0215
Total jobs = 3
Launching Job 1 out of 3
Number of reduce tasks not specified. Defaulting to jobconf value of: 4
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1732089968849_2912, Tracking URL = http://master:6318/proxy/application_1732089968849_2912/
Kill Command = /opt/hadoop/bin/mapred job  -kill job_1732089968849_2912
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 4
2024-11-21 11:34:19,780 Stage-1 map = 0%,  reduce = 0%
2024-11-21 11:35:19,790 Stage-1 map = 0%,  reduce = 0%, Cumulative CPU 19.61 sec
2024-11-21 11:35:49,231 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 55.15 sec
2024-11-21 11:36:50,104 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 80.58 sec
2024-11-21 11:37:50,986 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 92.34 sec
2024-11-21 11:38:51,856 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 107.29 sec
2024-11-21 11:39:52,712 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 117.07 sec
2024-11-21 11:40:53,566 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 150.93 sec
2024-11-21 11:41:54,416 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 159.96 sec
2024-11-21 11:42:55,257 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 168.72 sec
2024-11-21 11:43:56,087 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 176.88 sec
```

Q3
→

**SPARK**

**Q2**
1.
→

```
Max- df.agg(max('booked_seats')).show()
Min- df.agg(min('booked_seats')).show()
Avg- df.agg(avg('booked_seats')).show()
```

```
|1998|      2|         300.97|       30852|
|1998|      3|         315.25|       38118|
|1998|      4|         316.18|       35393|
|1999|      1|         331.74|       47453|
|1999|      2|         329.34|       38243|
|1999|      3|         317.22|       33048|
|1999|      4|         317.93|       31256|
+----+-------+---------------+------------+
only showing top 20 rows

>>> df.printSchema()
root
 |-- year: integer (nullable = true)
 |-- quarter: integer (nullable = true)
 |-- avg_rev_per_seat: float (nullable = true)
 |-- booked_seats: integer (nullable = true)

>>> from pyspark.sql.functions import max,min,avg,sum
>>> df.agg(max('booked_seats')).show()
+-----------------+
|max(booked_seats)|
+-----------------+
|            49678|
+-----------------+

>>> df.agg(min('booked_seats')).show()
+-----------------+
|min(booked_seats)|
+-----------------+
|            30103|
+-----------------+

>>> df.agg(avg('booked_seats')).show()
+-----------------+
|avg(booked_seats)|
+-----------------+
|39640.70238095238|
+-----------------+

>>>
```

2.
→
df.agg(count(col('avg_rev_per_seat'))).agg(col('avg_rev_per_seat')
<290).show()

3.
→ df.groupBy('quarter').agg(avg(col('booked_seats'))).show()

```
NameError: name 'row' is not defined
>>> from pyspark.sql.functions import count,col,row
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
ImportError: cannot import name 'row' from 'pyspark.sql.functions' (/opt/spark-3.1.2/python/pyspark/sql/functions.py)
>>> df.agg(count(col('avg_rev_per_seat'))).agg(col('avg_rev_per_seat')<290).show()
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
  File "/opt/spark-3.1.2/python/pyspark/sql/dataframe.py", line 1816, in agg
    return self.groupBy().agg(*exprs)
  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_seat`' given input columns: [count(avg_rev_per_seat)];
'Aggregate [('avg_rev_per_seat < 290) AS (avg_rev_per_seat < 290)#70]
+- Aggregate [count(avg_rev_per_seat#2) AS count(avg_rev_per_seat)#67L]
   +- Relation[year#0,quarter#1,avg_rev_per_seat#2,booked_seats#3] csv

>>> df.groupBy('quarter').agg(avg(col('booked_seats'))).show()
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'avg' is not defined
>>> from pyspark.sql.functions import count,col,row,avg
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
ImportError: cannot import name 'row' from 'pyspark.sql.functions' (/opt/spark-3.1.2/python/pyspark/sql/functions.py)
>>> from pyspark.sql.functions import count,col,avg
>>> df.groupBy('quarter').agg(avg(col('booked_seats'))).show()
+-------+------------------+
|quarter|  avg(booked_seats)|
+-------+------------------+
|      1|41607.666666666664|
|      3| 39386.23809523809|
|      4| 39111.95238095238|
|      2| 38456.95238095238|
+-------+------------------+

>>>
```

4.
→ df.groupBy("year").agg(count('quarter')).show()

```
ImportError: cannot import name 'row' from 'pyspark.sql.functions' (/opt/spark-3.1.2/python/pyspark/sql/functions.py)
>>> from pyspark.sql.functions import count,col,avg
>>> df.groupBy('quarter').agg(avg(col('booked_seats'))).show()
+-------+------------------+
|quarter|  avg(booked_seats)|
+-------+------------------+
|      1|41607.666666666664|
|      3| 39386.23809523809|
|      4| 39111.95238095238|
|      2| 38456.95238095238|
+-------+------------------+

>>> df.groupBy("year").agg(count('quarter')).show()
+----+--------------+
|year|count(quarter)|
+----+--------------+
|2003|             4|
|2007|             4|
|2015|             4|
|2006|             4|
|2013|             4|
|1997|             4|
|2014|             4|
|2004|             4|
|1996|             4|
|1998|             4|
|2012|             4|
|2009|             4|
|1995|             4|
|2001|             4|
|2005|             4|
|2000|             4|
|2010|             4|
|2011|             4|
|2008|             4|
|1999|             4|
+----+--------------+
only showing top 20 rows

>>>
```

5.

```
→ df.agg('quarter').agg(sum(col('avg_rev_per_seat') *
col('booked_seats')).alias("revenue")).agg(max(col("revenue"))).sh
ow()
```