## Q1.

### 1.

Select a.name from airports a
join routes r
On a.airport_id = r.airport_id
Limit 0;



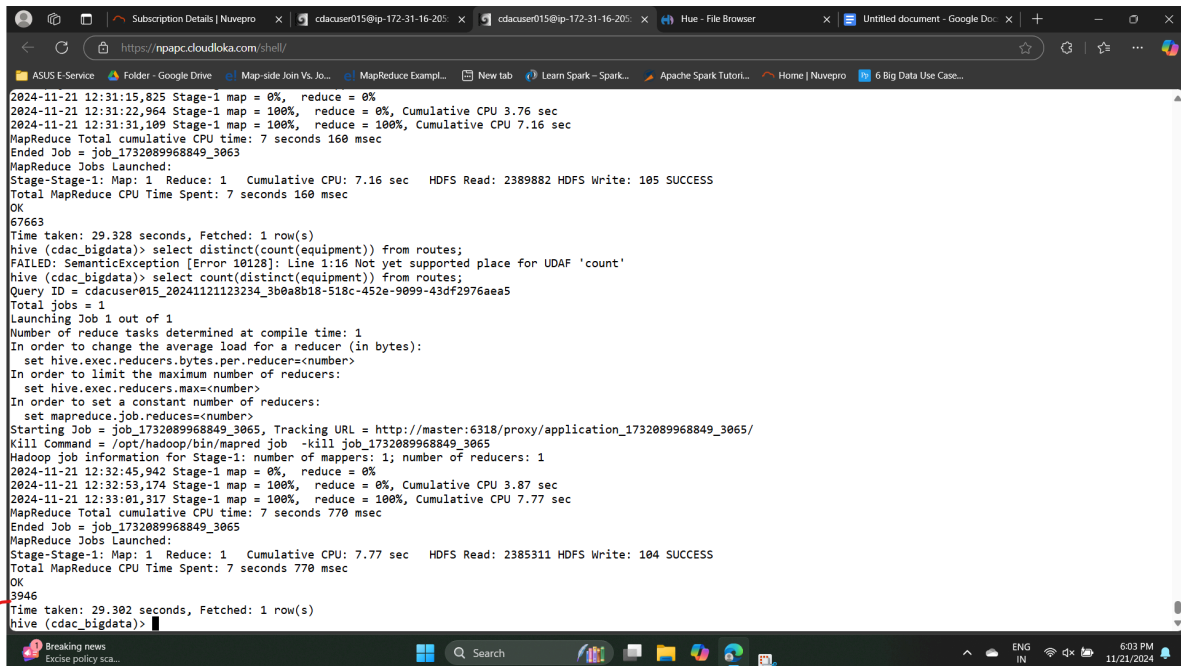### 2.

```
select a.name ,count(*)  from airlines a join routes r on
r.airline_id = a.airline_id   group by name limit 3;
```

Q3. select count(distinct(equipment)) from routes;



Q2.
1.
create table routepart(airline_iata string,airline_id int,src_airport_iata string,src_airport_id int,dest_airport_id int,codeshare string,stops int,equipment string) partitioned by (dest_airport_iata string) row format delimited  fields terminated by ',' stored as textfile;

```
Time taken: 0.047 seconds, Fetched: 9 row(s)
hive (cdac_bigdata)> create table routepart(airline_iata string,airline_id int,src_airport_iata string,src_airport_id int,dest_airport_id int,codeshare string,stops in
t,equipment string) partitioned by (dest_airport_iata string) row format delimited by fields terminated by ',' stored as textfile;
NoViableAltException(53@[2032:103: ( tableRowFormatMapKeysIdentifier )?])
        at org.antlr.runtime.DFA.noViableAlt(DFA.java:158)
        at org.antlr.runtime.DFA.predict(DFA.java:116)
        at org.apache.hadoop.hive.ql.parse.HiveParser.rowFormatDelimited(HiveParser.java:27784)
        at org.apache.hadoop.hive.ql.parse.HiveParser.tableRowFormat(HiveParser.java:28001)
        at org.apache.hadoop.hive.ql.parse.HiveParser.createTableStatement(HiveParser.java:6765)
        at org.apache.hadoop.hive.ql.parse.HiveParser.ddlStatement(HiveParser.java:4295)
        at org.apache.hadoop.hive.ql.parse.HiveParser.execStatement(HiveParser.java:2494)
        at org.apache.hadoop.hive.ql.parse.HiveParser.statement(HiveParser.java:1420)
        at org.apache.hadoop.hive.ql.parse.ParseDriver.parse(ParseDriver.java:220)
        at org.apache.hadoop.hive.ql.parse.ParseUtils.parse(ParseUtils.java:74)
        at org.apache.hadoop.hive.ql.parse.ParseUtils.parse(ParseUtils.java:67)
        at org.apache.hadoop.hive.ql.Driver.compile(Driver.java:616)
        at org.apache.hadoop.hive.ql.Driver.compileInternal(Driver.java:1826)
        at org.apache.hadoop.hive.ql.Driver.compileAndRespond(Driver.java:1773)
        at org.apache.hadoop.hive.ql.Driver.compileAndRespond(Driver.java:1768)
        at org.apache.hadoop.hive.ql.reexec.ReExecDriver.compileAndRespond(ReExecDriver.java:126)
        at org.apache.hadoop.hive.ql.reexec.ReExecDriver.run(ReExecDriver.java:214)
        at org.apache.hadoop.hive.cli.CliDriver.processLocalCmd(CliDriver.java:239)
        at org.apache.hadoop.hive.cli.CliDriver.processCmd(CliDriver.java:188)
        at org.apache.hadoop.hive.cli.CliDriver.processLine(CliDriver.java:402)
        at org.apache.hadoop.hive.cli.CliDriver.executeDriver(CliDriver.java:821)
        at org.apache.hadoop.hive.cli.CliDriver.run(CliDriver.java:759)
        at org.apache.hadoop.hive.cli.CliDriver.main(CliDriver.java:683)
        at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
        at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
        at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
        at java.lang.reflect.Method.invoke(Method.java:498)
        at org.apache.hadoop.util.RunJar.run(RunJar.java:323)
        at org.apache.hadoop.util.RunJar.main(RunJar.java:236)
FAILED: ParseException line 1:229 cannot recognize input near 'by' 'fields' 'terminated' in serde properties specification
hive (cdac_bigdata)> create table routepart(airline_iata string,airline_id int,src_airport_iata string,src_airport_id int,dest_airport_id int,codeshare string,stops in
t,equipment string) partitioned by (dest_airport_iata string) row format delimited  fields terminated by ',' stored as textfile;
OK
Time taken: 0.075 seconds
hive (cdac_bigdata)>
```
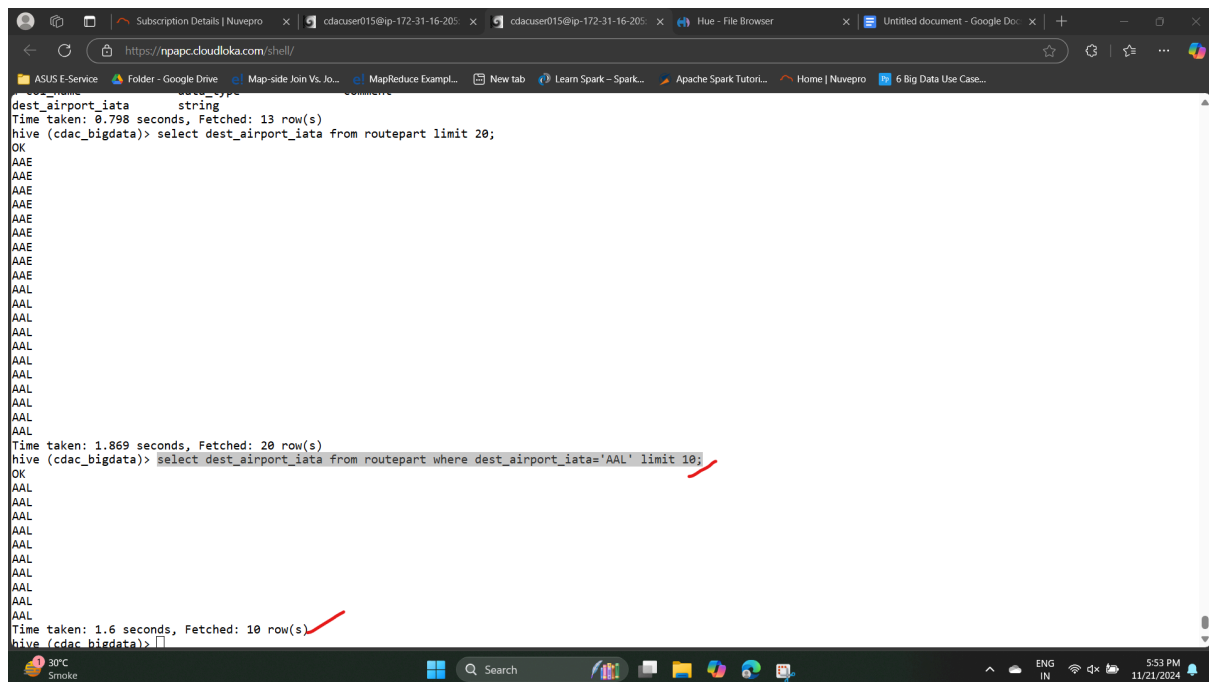
2.insert overwrite table routepart partition(dest_airport_iata) select
airline_iata,airline_id,src_airport_iata,src_airport_id,dest_airport_id,codes
hare,stops,equipment,dest_airport_iata from routes;



```
  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_2984, Tracking URL = http://master:6318/proxy/application_1732089968849_2984/
Kill Command = /opt/hadoop/bin/mapred job  -kill job_1732089968849_2984
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 4
2024-11-21 11:56:09,965 Stage-1 map = 0%,  reduce = 0%
2024-11-21 11:57:10,087 Stage-1 map = 0%,  reduce = 0%, Cumulative CPU 17.88 sec
2024-11-21 11:57:47,678 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 42.46 sec
2024-11-21 11:58:48,591 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 64.24 sec
2024-11-21 11:59:49,482 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 75.31 sec
2024-11-21 12:00:50,368 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 90.3 sec
2024-11-21 12:01:51,201 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 100.63 sec
2024-11-21 12:02:52,052 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 118.14 sec
2024-11-21 12:03:52,893 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 131.52 sec
2024-11-21 12:04:53,756 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 142.08 sec
2024-11-21 12:05:54,590 Stage-1 map = 67%,  reduce = 0%, Cumulative CPU 149.44 sec
2024-11-21 12:06:11,823 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 156.25 sec
2024-11-21 12:06:18,930 Stage-1 map = 100%,  reduce = 50%, Cumulative CPU 166.01 sec
2024-11-21 12:06:20,961 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU 176.34 sec
MapReduce Total cumulative CPU time: 2 minutes 56 seconds 340 msec
Ended Job = job_1732089968849_2984
Stage-4 is selected by condition resolver.
Stage-3 is filtered out by condition resolver.
Stage-5 is filtered out by condition resolver.
Moving data to directory hdfs://master:9000/user/hive/warehouse/cdac_bigdata.db/routepart/.hive-staging_hive_2024-11-21_11-55-57_553_5593767793457699168-1/-ext-10000
Loading data to table cdac_bigdata.routepart partition (dest_airport_iata=null)


        Time taken to load dynamic partitions: 144.604 seconds
        Time taken for adding to write entity : 0.044 seconds
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Reduce: 4   Cumulative CPU: 176.34 sec   HDFS Read: 2437026 HDFS Write: 11387971 SUCCESS
Total MapReduce CPU Time Spent: 2 minutes 56 seconds 340 msec
OK
Time taken: 954.127 seconds
hive (cdac_bigdata)>
```
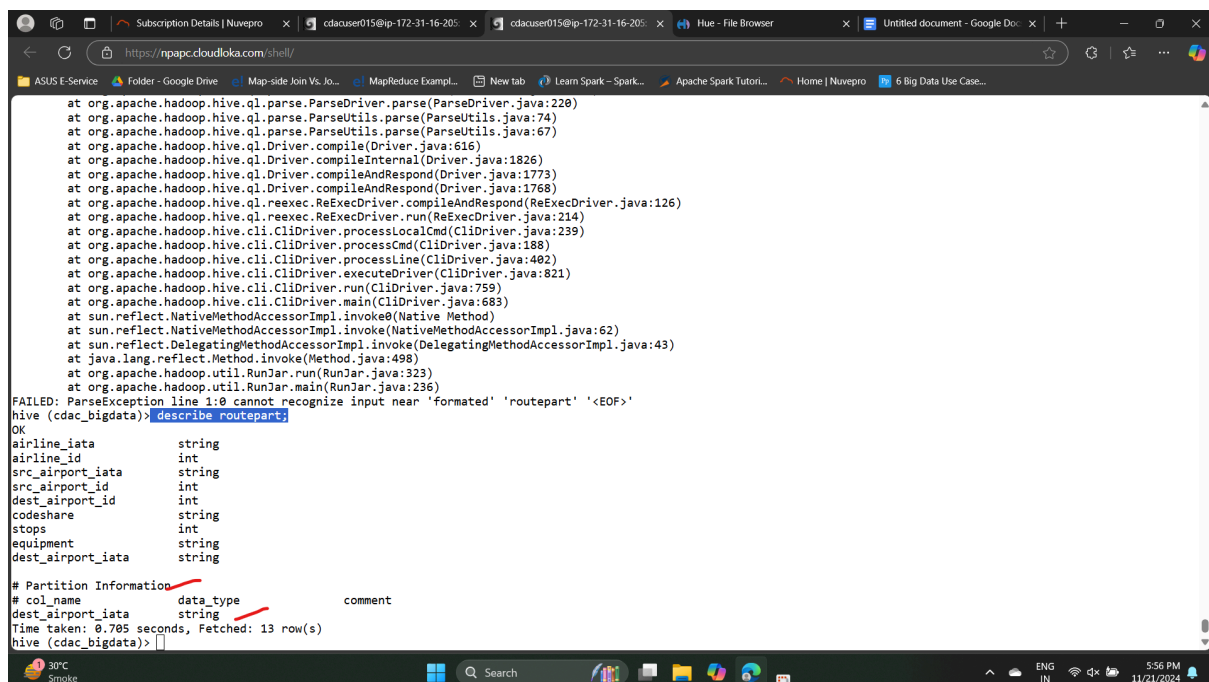
3.select dest_airport_iata from routepart where dest_airport_iata='AAL' limit 10;



4. describe routepart;



# Spark

Q2.

1.df.agg(max('booked_seats')).show()

```
NameError: name 'booked_seats' is not defined
>>> df.agg(max('booked_seats') & min('booked_seats') & avg('booked_seats')).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 '(max(`booked_seats`) AND min(`booked_seats`))' due to data type mismatch: '(max(`booked_seats`) AND min(`booked_se
ats`))' requires boolean type, not int;
'Aggregate [((max(booked_seats#19) AND min(booked_seats#19)) AND avg(cast(booked_seats#19 as bigint))) AS ((max(booked_seats) AND min(booked_seats)) AND avg(booked_sea
ts))#78]
+- Relation[Year#16,Quarter#17,Avg_rev_per_seat#18,booked_seats#19] csv

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

>>> df.agg(max('booked_seats')).show()
+----------------+
|max(booked_seats)|
+----------------+
|           49678|
+----------------+

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

>>>
```

**2.** `df.filter(col('avg_rev_per_seat') > 290).count()`



```
>>> df.filter('avg_rev_per_seat').count()
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
  File "/opt/spark-3.1.2/python/pyspark/sql/dataframe.py", line 1715, in filter
    jdf = self._jdf.filter(condition)
  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: filter expression '`avg_rev_per_seat`' of type double is not a boolean.;
Filter avg_rev_per_seat#18: double
+- Relation[Year#16,Quarter#17,Avg_rev_per_seat#18,booked_seats#19] csv

>>> df.filter('avg_rev_per_seat' > 290000).count()
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
TypeError: '>' not supported between instances of 'str' and 'int'
>>> df.select(distinct 'year').count('year')
  File "<stdin>", line 1
    df.select(distinct 'year').count('year')
                       ^
SyntaxError: invalid syntax
>>> df.printSchema()
root
 |-- Year: integer (nullable = true)
 |-- Quarter: integer (nullable = true)
 |-- Avg_rev_per_seat: double (nullable = true)
 |-- booked_seats: integer (nullable = true)

>>> df.filter('avg_rev_per_seat'> 290000).count()
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
TypeError: '>' not supported between instances of 'str' and 'int'
>>> df.filter(col('avg_rev_per_seat') > 290000).count()
0
>>> df.filter(col('avg_rev_per_seat') > 2900).count()
0
>>> df.filter(col('avg_rev_per_seat') > 290).count()
75
>>>
```

.

**3.**
`df.groupBy('quarter').agg(avg('booked_seats')).show()`

4. `df.select(distinct 'year').count().show()`

5.
`df.groupBy('year','quarter').agg(sum(col('avg_rev_per_seat')*col('booked_seats')).alias('Total_Rev')).show()`