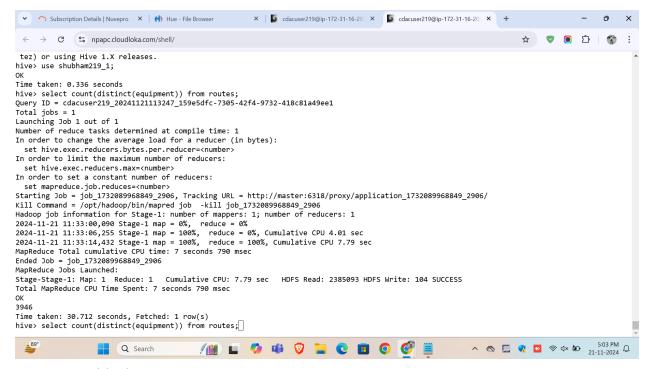
QUESTION 1;

1) select * from airport_1 a join routes r on a.iata = r.source_iata;

-					_										_		
· ^	Subscription Detail	× 5	Shell In A	Вох	×	19@ip-1	×	Hue - File Brow	ser X	Hue - File Browser ×	G cdacus	er219@ip-1	× +		-	0	×
\leftarrow \rightarrow	C º≅ npa	pc.cloudlo	oka.com/	shell/									☆	V	Ď	•	:
505	Phoenix-Mesa			esa	United States	AZA			307833	-111.655	1382	-7	N	Δ	merica	/Phoe	
ix G4		AZA	6505	PIA			0	M80									
505	Phoenix-Mesa			esa	United States	AZA			307833	-111.655	1382	-7	N	Δ	merica	/Phoe	
ix G4		AZA	6505	OGD			0	319		444 655	4200	_				/Dl	
505	Phoenix-Mesa			esa	United States	AZA			307833	-111.655	1382	-7	N	Д	merica	/Pnoe	
ix G4		AZA	6505	OAK			0	M80		444 655	4000	_				/Bl	
505 ix G4	Phoenix-Mesa 35	Gateway AZA	y M6 6505	esa MSO	United States 4216	AZA	0	KIWA 33.	307833	-111.655	1382	-7	N	Д	merica	/Pnoe	
1X 04 505	Phoenix-Mesa				United States				307833	111 (55	1382	-7	N		merica	/Dhaa	
ix G4		AZA	y ™ 6505	esa MOT		AZA	0	KIWA 33.: M80	30/833	-111.655	1382	-/	IV	μ.	merica	/ Pride	
505	Phoenix-Mesa			esa	United States	AZA	-		307833	-111.655	1382	-7	N		merica	/Dhoo	
ix G4		AZA	9 M	esa MLI		AZA	0	319	30/633	-111.055	1362	-/	IN	-	mer.rca	/ Piloe	
505	Phoenix-Mesa			esa	United States	AZA	-		307833	-111.655	1382	-7	N		merica	/Dhoo	
ix G4		AZA	6505	LAS		AZA	0	M80	30/633	-111.655	1302	-/	IV	-	uner ica	/ PIIOE	
505	Phoenix-Mesa			esa	United States	AZA			307833	-111.655	1382	-7	N	^	merica	/Dhoe	
ix G4		AZA	6505	IDA		, AZA	0	M80	,67655	-111.055	1362	-,	14	-	uner ica	i) Filoe	
505	Phoenix-Mesa			esa	United States	AZA	-		307833	-111.655	1382	-7	N	Δ	merica	/Phoe	
ix G4		AZA	6505	ICT		, A <u>-</u> A	0	M80	,0,055	111.055	1302	,			unci ico	,,,,,,,,	
505	Phoenix-Mesa			esa .	United States	AZA	-		307833	-111.655	1382	-7	N	Δ	merica	/Phoe	
ix G4		AZA	6505	GTF		, ,,_,,	0	M80		111.055	1302	,				,	
505	Phoenix-Mesa			esa	United States	AZA	-		307833	-111.655	1382	-7	N	Δ	merica	/Phoe	
ix G4		AZA	6505	RFD			0	319								,	
505	Phoenix-Mesa			esa	United States	AZA			307833	-111.655	1382	-7	N	Δ	merica	/Phoe	
ix G4		AZA	6505	GRR			0	319								,	
505	Phoenix-Mesa	Gateway	v Me	esa	United States	AZA		KIWA 33.	307833	-111.655	1382	-7	N	Δ	merica	/Phoe	
ix G4	35	AZA .	6505	RST	4048		0	319								-	
505	Phoenix-Mesa	Gateway	y Me	esa	United States	AZA		KIWA 33.	307833	-111.655	1382	-7	N	Δ	merica	/Phoe	
ix G4		AZA .	6505	SBN	4359		0	M80									
505	Phoenix-Mesa	Gateway	y Me	esa	United States	AZA		KIWA 33.	307833	-111.655	1382	-7	N	Δ	merica	/Phoe	
ix G4	35																
89°		2 - 1		-	- A	-21	60			6 🖺 🖫		~ =		0.4	4-	6:09 PM	0
-		Q Search		9		Щ	V			o 🚱 🗐 👶	^	Ø 🗀	AV N	d×	21	-11-2024	4

2) select count(airlne) from routes group by
(source_iata, destination_iata) order by
(source_iata, destination_iata) limit 3;

3) hive> select count(distinct(equipment)) from routes;



OUTPUT : 3946

Question 2 :

- 1) hive (shubham219_1) > Create table route_parittion(airline string , airline_id int , source_iata string , source_id string , destination_id string , codeshare string , stops int) partitioned by(destination_iata string) row format delimited fields terminated by "," stored as textfile;
- 2) insert overwrite table route_parittion select
 rp.airlne,rp.airline_id ,rp.source_iata, rp.source_id ,
 rp.destination_id,rp.codeshare , rp.stops , rp.destination_iata
 from routes rp;
- 3) select * from route_parittion where destination_iata =
 "ORD";

4)

SPARK

```
Question 1 :
>>> data = sc.textFile("/user/cdacuser219/airline 1")
>>> header = data.first()
>>> eliminate = data.filter(lambda x: x!=header)
>>> data map = eliminate.map(lambda x:
(int(x.split(",")[0]), int(x.split(",")[3])))
>>> for x in data map.collect():
          print(x)
. . .
>>> eliminate reduce = data map.reduceByKey(lambda x,y: x+y)
>>> for x in eliminate reduce.collect():
          print(x)
>>> filter row = eliminate reduce.filter(lambda x: x[1]>20000
and x[1] < 50000).count()
>>> print(filter row)
  uata = sc.textrile( /user/cuacuserzi9/airiine i )
IndentationError: unexpected indent
>>> data = sc.textFile("/user/cdacuser219/airline_1")
>>> header = data.first()
>>> eliminate = data.filter(lambda x: x!=header)^M
>>> eliminate = data.filter(lambda x: x!=header)
>>> data_map = eliminate.map(lambda x: (int(x.split(",")[0]),int(x.split(",")[3])))
>>> eliminate_reduce = data_map.reduceByKey(lambda x,y: x+y)
>>> filter_row = eliminate_reduce.filter(lambda x: x[1]>20000 and x[1]<500000).count()
>>> print(filter_row)
21
>>>
                         ^ 🖎 🛄 💎 🎅 d× 🕭 5:56 PM Q
         Q Search
>>> data = sc.textFile("/user/cdacuser219/airline 1")
>>> header = data.first()
>>> eliminate = data.filter(lambda x: x!=header)^M
>>> eliminate = data.filter(lambda x: x!=header)
>>> data map = eliminate.map(lambda x:
(int(x.split(",")[0]),int(x.split(",")[3])))
>>> eliminate reduce = data map.reduceByKey(lambda x,y: x+y)
```

```
(1, 1996)
Part 2
Question 1
>>> spark = SparkSession.builder.appName("exam").getOrCreate()
>>> data = spark.read.format("csv").option("header",
True).option("inferSchema", True).load("/user/cdacuser219/airline 1")
>>> data.show()
+---+
|Year|Quarter|Avg rev per seat|booked seats|
+---+
|1995|
         1 |
                   296.9|
                              46561|
|1995|
         2 |
                   296.8|
                              374431
        3 |
|1995|
                  287.51|
                              34128|
|1995|
        4 |
                  287.78|
                              303881
119961
         1 |
                   283.971
                             478081
|1996|
        2 |
                  275.78|
                             43020|
|1996|
        3 |
                  269.491
                              389521
|1996|
        4 |
                  278.33|
                             37443|
        1 |
|1997|
                   283.4|
                              350671
|1997|
        21
                  289.44|
                             465651
         3 |
|1997|
                  282.27|
                             388861
|1997|
         4 |
                  293.51
                              37454|
        1 |
                  304.74|
|1998|
                              31315|
|1998|
        2 |
                  300.97|
                             308521
|1998|
        3 |
                  315.25|
                              38118|
|1998|
        4 |
                   316.18|
                              35393|
                  331.74|
|1999|
        1 |
                             47453|
|1999|
        2 |
                   329.34|
                              38243|
|1999|
         3 |
                   317.22|
                              330481
|1999|
                   317.93|
        4 |
                              312561
+---+
only showing top 20 rows
>>> data.describe()
DataFrame[summary: string, Year: string, Quarter: string,
Avg rev per seat: string, booked seats: string]
>>> data.describe().show()
```

(1, 1995)(2, 1995)(3, 1995)(4, 1995)

summary	Year	Quarter	Avg_rev_per_seat								
booked_seats											
+											
+											
count	84	84	84								
84											
mean	2005.0										
2.5 329.7475000000006 39640.70238095238											
stddev 6.091669207609634 1.124748967976346											
32.64232664586615	5424.069182884	482									
min	1995	1	269.49								
30103											
max	2015	4	396.37								
49678											
+		+-	+								
+											

³⁾ data.groupBy('booked_seat').agg({'avg','booked_seats'}).show()