Pyspark

```
>>> airline = sc.textFile("/user/cdacuser87123/airlines.csv")
>>> airline.count()
85
>>> header = airline.first()
>>> airline = airline.filter(lambda line: line!=header)
>>> airline.count()
84
print(airline.take(5))
['1995,1,296.9,46561', '1995,2,296.8,37443', '1995,3,287.51,34128', '1995,4,287.78,30388', '1996,1,283.97,47808']
>>> split = airline.map(lambda a : (a.split(",")[0],a.split(",")[1],float(a.split(",")[2]),int(a.split(",")[3])))
>>> print(split.take(5))
[('1995', '1', 296.9, 46561), ('1995', '2', 296.8, 37443), ('1995', '3', 287.51, 34128), ('1995', '4', 287.78, 30388), ('1996', '1', 283.97, 47808)]
```
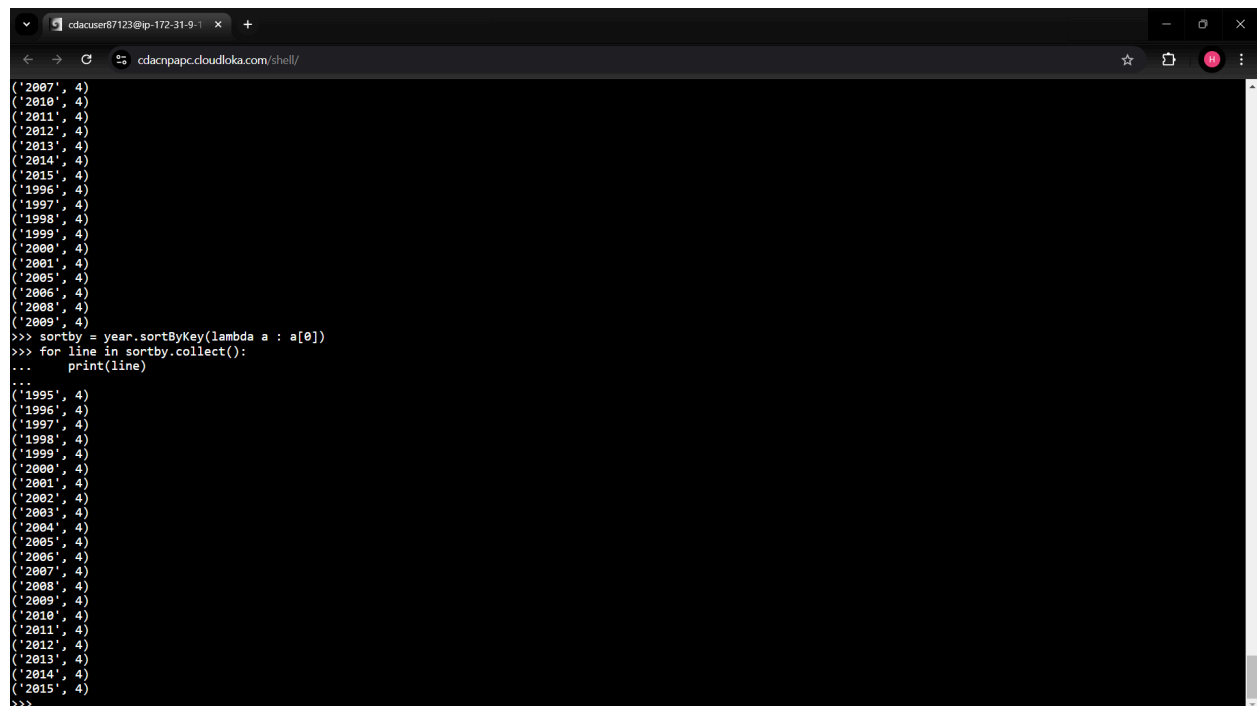
## Q1

**a**

```
max_seats = split.filter(lambda a : a[3]>40000)
>>> max_seats.count()
```

## Output



```
        merger.mergeValues(iterator)
  File "/opt/spark-3.1.2/python/pyspark/shuffle.py", line 240, in mergeValues
    for k, v in iterator:
ValueError: too many values to unpack (expected 2)

        at org.apache.spark.api.python.BasePythonRunner$ReaderIterator.handlePythonException(PythonRunner.scala:517)
        at org.apache.spark.api.python.PythonRunner$$anon$3.read(PythonRunner.scala:652)
        at org.apache.spark.api.python.PythonRunner$$anon$3.read(PythonRunner.scala:635)
        at org.apache.spark.api.python.BasePythonRunner$ReaderIterator.hasNext(PythonRunner.scala:470)
        at org.apache.spark.InterruptibleIterator.hasNext(InterruptibleIterator.scala:37)
        at scala.collection.Iterator$GroupedIterator.fill(Iterator.scala:1209)
        at scala.collection.Iterator$GroupedIterator.hasNext(Iterator.scala:1215)
        at scala.collection.Iterator$$anon$10.hasNext(Iterator.scala:458)
        at org.apache.spark.shuffle.sort.BypassMergeSortShuffleWriter.write(BypassMergeSortShuffleWriter.java:132)
        at org.apache.spark.shuffle.ShuffleWriteProcessor.write(ShuffleWriteProcessor.scala:59)
        at org.apache.spark.scheduler.ShuffleMapTask.runTask(ShuffleMapTask.scala:99)
        at org.apache.spark.scheduler.ShuffleMapTask.runTask(ShuffleMapTask.scala:52)
        at org.apache.spark.scheduler.Task.run(Task.scala:131)
        at org.apache.spark.executor.Executor$TaskRunner.$anonfun$run$3(Executor.scala:497)
        at org.apache.spark.util.Utils$.tryWithSafeFinally(Utils.scala:1439)
        at org.apache.spark.executor.Executor$TaskRunner.run(Executor.scala:500)
        at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1149)
        at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:624)
        ... 1 more

>>> print(split.take(5))
[('1995', '1', 296.9, 46561), ('1995', '2', 296.8, 37443), ('1995', '3', 287.51, 34128), ('1995', '4', 287.78, 30388), ('1996', '1', 283.97, 47808)]
>>> split = airline.map(lambda a : (a.split(",")[0],a.split(",")[1],float(a.split(",")[2]),int(a.split(",")[3])))
>>> print(split.take(5))
[('1995', '1', 296.9, 46561), ('1995', '2', 296.8, 37443), ('1995', '3', 287.51, 34128), ('1995', '4', 287.78, 30388), ('1996', '1', 283.97, 47808)]
>>> combine = split.map(lambda a : (a[0],1),a[1])
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'a' is not defined
>>> combine = split.map(lambda a : (a[0],1),a[1]))
  File "<stdin>", line 1
    combine = split.map(lambda a : (a[0],1),a[1]))
                                                 ^
SyntaxError: unmatched ')'
>>> max_seats = split.filter(lambda a : a[3]>40000)
>>> max_seats.count()
38
>>>
```

**b**

```
>>> combine = split.map(lambda a : (a[0],1))
>>> year = combine.reduceByKey(lambda a,b: a+b)
```

```
sortby = year.sortByKey(lambda a : a[0])
>>> for line in sortby.collect():
...     print(line)
```

Output



Q2

a)
```
combine = split.map(lambda a : (a[0],a[2]))
```

```
>>> max_per_seat = combine.max(key = lambda a :
a[1])
max_per_seat = combine.min(key = lambda a : a[1])
>>> print(max_per_seat)
avg_per_seat = combine.map(lambda a :
a[1]).mean()
>>> print(avg_per_seat)
```
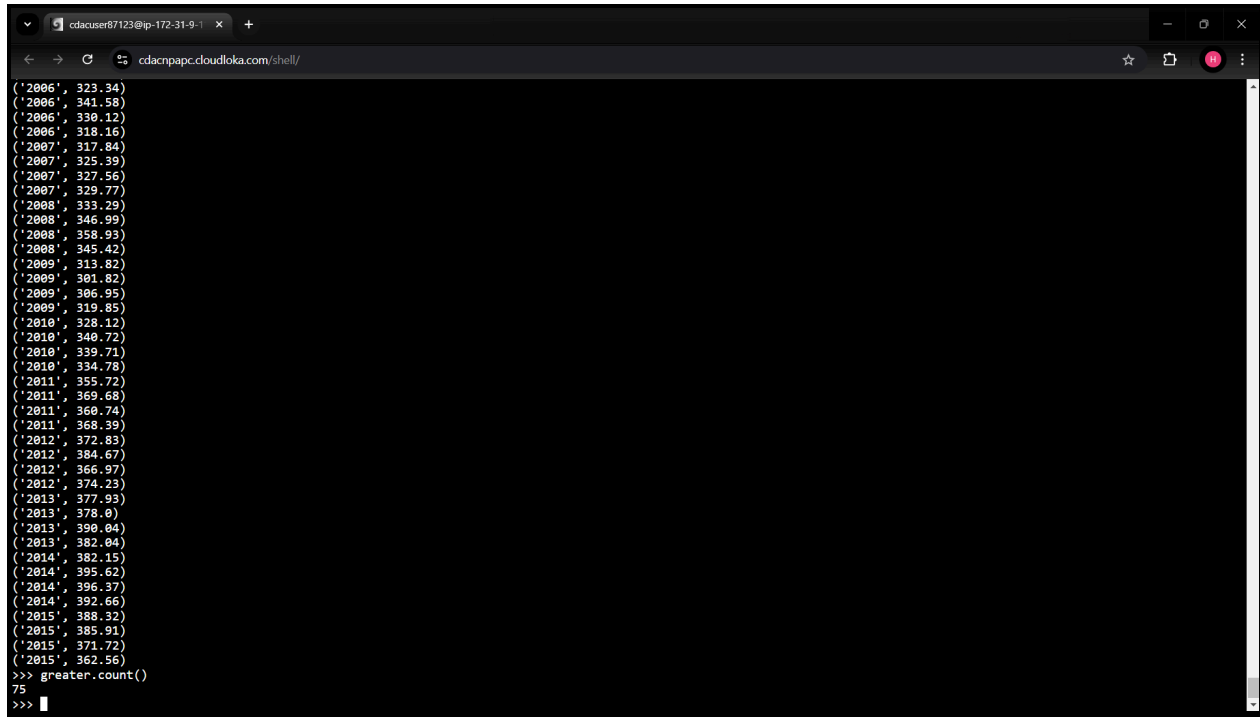
Output:



b)
```
combine = split.map(lambda a : (a[0],a[2]))
```

```
greater = combine.filter(lambda a : (a[1]>290))
greater.count()
```



```
('2006', 323.34)
('2006', 341.58)
('2006', 330.12)
('2006', 318.16)
('2007', 317.84)
('2007', 325.39)
('2007', 327.56)
('2007', 329.77)
('2008', 333.29)
('2008', 346.99)
('2008', 358.93)
('2008', 345.42)
('2009', 313.82)
('2009', 301.82)
('2009', 306.95)
('2009', 319.85)
('2010', 328.12)
('2010', 340.72)
('2010', 339.71)
('2010', 334.78)
('2011', 355.72)
('2011', 369.68)
('2011', 360.74)
('2011', 368.39)
('2012', 372.83)
('2012', 384.67)
('2012', 366.97)
('2012', 374.23)
('2013', 377.93)
('2013', 378.0)
('2013', 390.04)
('2013', 382.04)
('2014', 382.15)
('2014', 395.62)
('2014', 396.37)
('2014', 392.66)
('2015', 388.32)
('2015', 385.91)
('2015', 371.72)
('2015', 362.56)
>>> greater.count()
75
>>>
```

c)

```
combine = split.map(lambda a : (a[1],a[3]))
>>> print(combine)
PythonRDD[57] at RDD at PythonRDD.scala:53
>>> print(combine.take())
print(combine.take(4))
```

```
[('1', 46561), ('2', 37443), ('3', 34128), ('4',
30388)]
>>> co = combine.reduceByKey(lambda a,b : a+b)
>>> print(co.collect())
[('1', 873761), ('4', 821351), ('2', 807596),
('3', 827111)]
>>> sortby = co.sortBy(lambda a : a[0])
>>> for line in sortby.collect():
...      print(line)
```



d)

```
>>> combine = split.map(lambda a : (a[0],a[1]))
>>> co  = combine.reduceByKey(lambda a,b: a+b)
>>> for l in co.collect():
...         print(l)
```



e)

```
>>> combine = split(lambda a :
(a[0],(a[2]*a[3])))
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
TypeError: 'PipelinedRDD' object is not callable
>>> combine = split.map(lambda a :
(a[0],(a[2]*a[3])))
>>> print(combine.take(5))
[('1995', 13823960.899999999), ('1995',
11113082.4), ('1995', 9812141.28), ('1995',
8745058.639999999), ('1996', 13576037.760000002)]
>>> revenue = combine,reduceByKey(lambda a,b :
a+b)
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'reduceByKey' is not defined
>>> revenue = combine.reduceByKey(lambda a,b :
a+b)
>>> print(revenue.collect())
```

```
>>> print(revenue.collect())
[('1995', 43494243.22), ('2002', 47499146.5), ('2003', 49273210.83), ('2004', 50631364.949999996), ('2007', 57309216.07), ('2010', 54861521.29), ('2011', 51888286.22),
('2012', 62199127.28), ('2013', 66363208.71), ('2014', 62624175.85000001), ('2015', 62378990.57), ('1996', 46358778.03), ('1997', 45385236.16), ('1998', 42035717.78), (
'1999', 48757714.48), ('2000', 52342926.550000004), ('2001', 55533779.9999999), ('2005', 46376786.24), ('2006', 50437898.419999994), ('2008', 57653170.760000005), ('20
09', 46746446.59)]
```
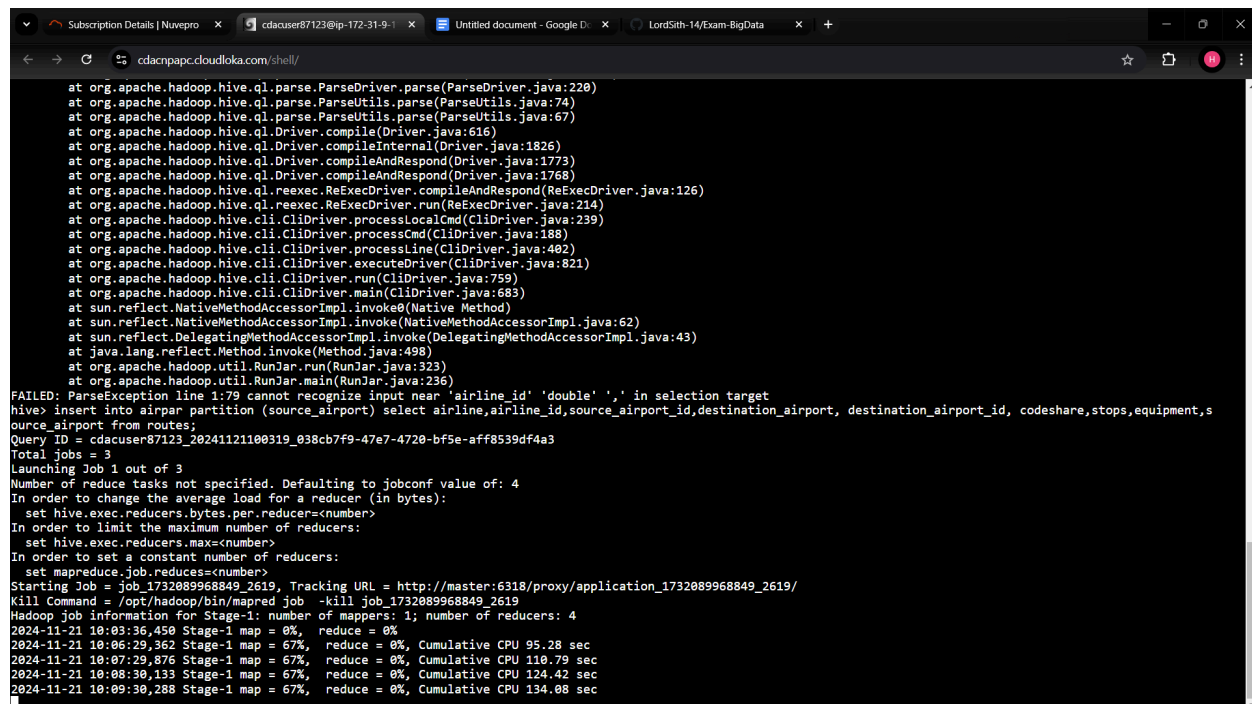
# Hive

## Q2)

### A
create table airpar(airline string,airline_id double,source_airport_id double,destination_airport string, destination_airport_id double, codeshare string,stops int,equipment) partitioned by (source_airport string) row format delimited fields terminated by ',' stored as textFile;

insert into airpar partition(source_airport) select airline ,airline_id ,source_airport_id ,destination_airport , destination_airport_id , codeshare,stops,equipment from routes



b)

```
insert overwrite airpar partition(source_airport) select airline
,airline_id ,source_airport_id ,destination_airport ,
destination_airport_id , codeshare,stops,equipment from routes
where source_airport = "JFK"
```

c)
```
select source_airport ,destination_airport from airpar  where
source_airport = "LAX"
```