

Questions

Q1. To emulate a live-stream of the traffic counter dataset, you are required to write a separate Python script that reads 10 records (of 2 each counter site - ignore the test site) every 5 seconds from traffic counter data file and stores them as separate files (countdata1, countdata2, countdata3, etc.) in the streaming directory on which your application is listening. (20 marks)

Ans: Python Code

```
import pandas as pd
import time

input_file = "per-vehicle-records-2021-01-15.csv"
output_file_prefix = "counterdata"

def read_vehicle_data():
    print("Reading data")
    reader = pd.read_csv(input_file, sep=',', chunksize=10, iterator=True)
    for i, j in enumerate(reader):
        df = next(reader)
        output_file = output_file_prefix + str(i) + ".csv"
        df.to_csv(output_file, index=False)
        time.sleep(5)
        print("Finished chunk")

if __name__ == "__main__":
    read_vehicle_data()
```

Explanation : This python script will read the file as pandas dataframe but in chunks size of 10 that will read 10 lines of file at time period of 5 seconds (used time.sleep(5)). Once the dataframe is generated it is converted to a CSV file. The script basically simulates the streaming data from a batch data.

Q2. Prepare the streaming application to read the data streams from the streaming directory using a batch length of 5 seconds. (10 marks)

Ans: Spark Code + Output

```
1  from pyspark.streaming import StreamingContext
2  ssc = StreamingContext(sc, 5)
```

```

1 # define the schema
2 TrafSchema = StructType([ StructField("cosit", StringType(), True), StructField("year", StringType(), True), StructField("month", StringType(), True), StructField("day", StringType(), True), StructField("hour",
StringType(), True), StructField("minute", StringType(), True), StructField("second", StringType(), True), StructField("millisecond", StringType(), True), StructField("minuteofday", StringType(), True),
StructField("lane", StringType(), True), StructField("lanename", StringType(), True), StructField("straddlelane", StringType(), True), StructField("straddlelanename", StringType(), True), StructField("class",
StringType(), True), StructField("classname", StringType(), True), StructField("length", StringType(), True), StructField("headway", StringType(), True), StructField("gap", StringType(), True),
StructField("speed", StringType(), True), StructField("weight", StringType(), True), StructField("temperature", StringType(), True), StructField("duration", StringType(), True), StructField("validitycode",
StringType(), True),StructField("numberofaxles", StringType(), True),StructField("axleweights", StringType(), True),StructField("axlespacings", StringType(), True)])
3

```

```

1 streamingInputDF = (
2     spark
3     .readStream
4     .schema(TrafSchema)
5     .option("maxFilesPerTrigger", 1) # Treat a sequence of files as a stream by picking one file at a time
6     .format("csv")
7     .load("dbfs:/FileStore/shared_uploads/email@gmail.com/streaming-files/")
8 )

```

Explanation : The code here is setting a streaming context in spark using a batch length of 5 seconds that means the streams will be read as a batch for 5 seconds for doing aggregations. It is also needed to set a schema for reading the input streams and setting the maxFilesPerTrigger as 1 to pick one file at a time as a stream.

Q3. Show total number of counts (on each site of M50) by vehicle class.

Ans: Spark code + Cassandra output

```

1 # Total number of counts by vehicle class
2 streamingVehicleClassCountsDF = (
3     streamingInputDF
4     .groupBy(
5         streamingInputDF.classname)
6     .count()
7 )

```

```

1 # Writing the Dataframe to Cassandra
2 %scala
3 //Writing the dataframe directly to cassandra
4 import org.apache.spark.sql.cassandra._
5
6
7 streamingVehicleClassCountsDF.write
8   .format("org.apache.spark.sql.cassandra")
9   .mode("overwrite")
10  .options(Map( "table" -> "streamingVehicleClassCountsDF", "keyspace" -> "test_keyspace"))
11  .save()

```

Explanation : Here we are grouping by vehicle class name to count the number of vehicles.

Stream 1 Output

Table ▾ +

	classname ▲	count ▲	
1	CAR	5	
2	HGV_RIG	4	
3	LGV	1	
4	classname	1	

Stream 2 Output

Table ▾ +

	classname ▲	count ▲	
1	CAR	8	
2	HGV_ART	2	
3	HGV_RIG	9	
4	LGV	1	

Stream 3 Output

Table ▾ +

	classname ▲	count ▲	
1	CAR	11	
2	HGV_ART	3	
3	HGV_RIG	15	
4	LGV	1	

Stream 4 Output

Table ▾ +

	classname ▲	count ▲	
1	CAR	15	
2	HGV_ART	4	
3	HGV_RIG	17	
4	LGV	4	

Q4. Compute the average speed (on each site on M50) by vehicle class.

Ans: Spark + Cassandra code + output

```
1 import pyspark.sql.functions as F
2 streamingVehicleClassAvgSpeedCountsDF = (
3     streamingInputDF
4     .groupBy(
5         streamingInputDF.classname)
6     .agg(F.round(F.avg(streamingInputDF.speed), 3))
7 )
```

```
1 # Writing the Dataframe to Cassandra
2 %scala
3 //Writing the dataframe directly to cassandra
4 import org.apache.spark.sql.cassandra._
5
6
7 streamingVehicleClassAvgSpeedCountsDF.write
8   .format("org.apache.spark.sql.cassandra")
9   .mode("overwrite")
10  .options(Map( "table" -> "streamingVehicleClassAvgSpeedCountsDF", "keyspace" -> "test_keyspace"))
11  .save()
```

Explanation : Here we are grouping by vehicle class name and calculating the average speed

Stream 1 Output

Table ▾ +

	classname ▲	round(avg(speed), 3) ▲
1	CAR	70.625
2	HGV_ART	70
3	HGV_RIG	70.333
4	LGV	76

Stream 2 Output



Q5: Find the top 3 busiest counter sites on M50.

Ans: Spark + Cassandra code + Output

```

1 streamingCountCositDF = (
2   streamingInputDF
3     .groupBy(
4       streamingInputDF.cosit)
5     .count()
6     .withColumnRenamed('count', 'total_count')
7     .orderBy()
8 )

```

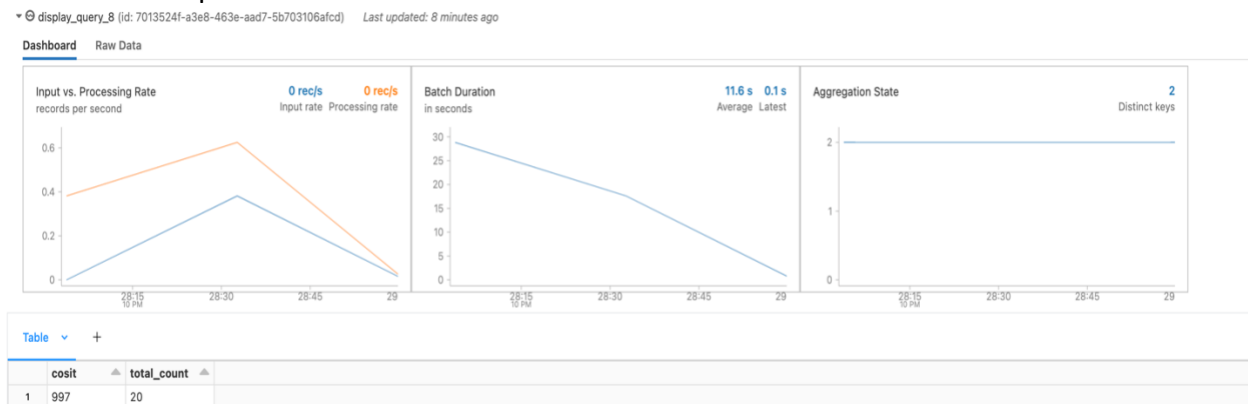
```
1 streamingCountCositDF.createTempView('count_cosit')
```

```
1 streamingTop3CositDF = spark.sql('select cosit, total_count from count_cosit order by total_count desc LIMIT 3')
2
```

```
1 # Writing the Dataframe to Cassandra
2 %scala
3 //Writing the dataframe directly to cassandra
4 import org.apache.spark.sql.cassandra._
5
6
7 streamingTop3CositDF.write
8   .format("org.apache.spark.sql.cassandra")
9   .mode("overwrite")
10  .options(Map( "table" -> "streamingTop3CositDF", "keyspace" -> "test_keyspace"))
11  .save()
```

Explanation : As streaming doesn't support ranking function yet, so using the Spark SQL to find the top 3 busiest counter sites, firstly the total count of cosit is calculated and then using that in descending order to calculate top 3 by using LIMIT 3 in the query.

Stream 1 Output



Stream 2 Output

	cosit	total_count	
1	997	10	

Stream 3 Output

	cosit ▲	total_count ▲	
1	997	30	

Q.6 Find total number of counts for HGVs on M50.

Ans:

```

1 streamingHGVCountsDF = (
2   streamingInputDF
3     .where((streamingInputDF.classname == 'HGV_RIG') | (streamingInputDF.classname == 'HGV_ART'))
4     .groupBy(
5       streamingInputDF.classname
6     )
7     .count()
8 )

```

```

1 |# Writing the Dataframe to Cassandra
2 %scala
3 //Writing the dataframe directly to cassandra
4 import org.apache.spark.sql.cassandra._
5
6
7 streamingHGVCountsDF.write
8   .format("org.apache.spark.sql.cassandra")
9   .mode("overwrite")
10  .options(Map( "table" -> "streamingHGVCountsDF", "keyspace" -> "test_keyspace"))
11  .save()

```

Explanation: Here we are first filtering the vehicles by HGV and then grouping by vehicle classname to find the counts of HGV

Stream 1 Output

	classname ▲	count ▲	
1	HGV_RIG	4	

Stream 2 Output

	classname ▲	count ▲	
1	HGV_ART	2	

Stream 3 Output

	classname ▲	count ▲	
1	HGV_ART	3	