# Scripts Execution

**Explanation of the solution to the streaming layer problem**

- All the necessary libraries and functions have been imported:

```python
4  # Importing all the required function
5  from pyspark.sql import SparkSession
6  from pyspark.sql.functions import *
7  from pyspark.sql.types import *
8  from datetime import datetime
```

- Spark session is created and all the required python files are imported:

```python
9
10  # Creating Spark Session
11  spark = SparkSession  \
12          .builder  \
13          .appName("CapStone_Project")  \
14          .getOrCreate()
15  spark.sparkContext.setLogLevel('ERROR')
16  sc = spark.sparkContext
17
18  # Adding the required python files
19  sc.addPyFile('db/dao.py')
20  sc.addPyFile('db/geo_map.py')
21  sc.addFile('rules/rules.py')
22
23  # Importing all the required files
24  import dao
25  import geo_map
26  import rules
27
```

- The session is now connected to the server "18.211.252.152:9092" and to kafka topic "transactions-topic-verified"

```
27
28  # Reading from Kafka stream
29  lines = spark  \
30          .readStream  \
31          .format("kafka")  \
32          .option("kafka.bootstrap.servers","18.211.252.152:9092")  \
33          .option("subscribe","transactions-topic-verified")  \
34          .option("failOnDataLoss","false").option("startingOffsets","earliest")  \
35          .load()
36
```

- The schema is defined and parsed in the format:

```
36
37  # Defining Schema
38  schema =  StructType([
39                  StructField("card_id", StringType()),
40                  StructField("member_id", StringType()),
41                  StructField("amount", IntegerType()),
42                  StructField("pos_id", StringType()),
43                  StructField("postcode", StringType()),
44                  StructField("transaction_dt", StringType())
45              ])
46
47  # Parsing the data
48  parse = lines.select(from_json(col("value") \
49                                      .cast("string") \
50                                      ,schema).alias("parsed"))
51
52  # Parseing the dataframe
53  df_parsed = parse.select("parsed.*")
54
```

- The "look_up_table" and "card_transactions" table for card transaction details are used.
- A set of user defined functions are used in order to perform the required activates and to check whether the transactions are fraudulent or genuine.
- Function for fetching the Credit Scores from the look up tables using the "card_id".

```
54
55  # Adding Time stamp column
56  df_parsed = df_parsed.withColumn('transaction_dt_ts',unix_timestamp(df_parsed.transaction_dt, 'dd-MM-YYYY HH:mm:ss').cast(TimestampType()))
57
58  # Function for Credit Score
59  def score(a):
60      hdao = dao.HBaseDao.get_instance()
61      data_fetch = hdao.get_data(key=a, table='look_up_table')
62      return data_fetch['info:score']
63
64  # Defining UDF for Credit Score
65  score_udf = udf(score,StringType())
66
67  # Adding Score Column
68  df_parsed = df_parsed.withColumn("score",score_udf(df_parsed.card_id))
69
```

- A function for fetching postal code from the look up tables using the "card_id"

```python
69
70  # Function for Postal Code
71  def postcode(a):
72      hdao = dao.HBaseDao.get_instance()
73      data_fetch = hdao.get_data(key=a, table='look_up_table')
74      return data_fetch['info:postcode']
75
76  # Defining UDF for Postal Code
77  postcode_udf = udf(postcode,StringType())
78
79  # Adding Postal Code Column
80  df_parsed = df_parsed.withColumn("last_postcode", postcode_udf(df_parsed.card_id))
81
```

- Function to fetch Upper Control Limit from the look up tables using the "card_id"

```python
81  |
82  # Function for UCL
83  def ucl(a):
84      hdao = dao.HBaseDao.get_instance()
85      data_fetch = hdao.get_data(key=a, table='look_up_table')
86      return data_fetch['info:UCL']
87  `
88  # Defining UDF for UCL
89  ucl_udf = udf(ucl,StringType())
90
91  # Adding UCL Column
92  df_parsed = df_parsed.withColumn("UCL", ucl_udf(df_parsed.card_id))
93
```

- Function to calculate the distance between previous and current transaction postal codes from the look up tables and kafka stream:

```python
93
94  # Function for Calculating distance
95  def dist_cal(last_postcode,postcode):
96      gmap = geo_map.GEO_Map.get_instance()
97      last_lat = gmap.get_lat(last_postcode)
98      last_longg = gmap.get_long(last_postcode)
99      lat = gmap.get_lat(postcode)
100     longg = gmap.get_long(postcode)
101     d = gmap.distance(last_lat.values[0],last_longg.values[0],lat.values[0],longg.values[0])
102     return d
103
104 # Defining UDF for Distance
105 distance_udf = udf(dist_cal,DoubleType())
106
107 # Adding Distance Column
108 df_parsed = df_parsed.withColumn("distance",distance_udf(df_parsed.last_postcode,df_parsed.postcode))
109
```

- Function to fetch the last transactions date from the look up tables using the "card_id"

```python
109
110 # Function for Transaction Date
111 def Tdate(a):
112     hdao = dao.HBaseDao.get_instance()
113     data_fetch = hdao.get_data(key=a, table='look_up_table')
114     return data_fetch['info:transaction_date']
115
116 # Defining UDF for Transaction Date
117 Tdate_udf = udf(Tdate,StringType())
118
119 # Adding Transaction Date Column
120 df_parsed = df_parsed.withColumn("last_transaction_date",Tdate_udf(df_parsed.card_id))
121
```

- Function to calculate the time difference (in sec) between the previous and current transactions from the look up tables and kafka stream:

```python
121
122 # Adding Time stamp column
123 df_parsed = df_parsed.withColumn('last_transaction_date_ts',unix_timestamp(df_parsed.last_transaction_date, 'YYYY-MM-dd
    HH:mm:ss').cast(TimestampType()))
124
125 # Function to Calculate Time
126 def time_cal(last_date,curr_date):
127     d = curr_date - last_date
128     return d.total_seconds()
129
130 # Defining UDF for Calculating Time
131 time_udf = udf(time_cal,DoubleType())
132
133 # Adding Time diff column
134 df_parsed = df_parsed.withColumn('time_diff',time_udf(df_parsed.last_transaction_date_ts,df_parsed.transaction_dt_ts))
135
```

- The function to define the status of transaction is fraudulent or genuine:

```python
135
136 # Function to define the Status of the Transaction
137 def status_def(card_id,member_id,amount,pos_id,postcode,transaction_dt,transaction_dt_ts,last_transaction_date_ts,score,distance,time_diff):
138     hdao = dao.HBaseDao.get_instance()
139     geo = geo_map.GEO_Map.get_instance()
140     look_up = hdao.get_data(key=card_id, table='look_up_table')
141     status = 'FRAUD'
142     if rules.rules_check(data_fetch['info:UCL'],score,distance,time_diff,amount):
143         status = 'GENUINE'
144         data_fetch['info:transaction_date'] = str(transaction_dt_ts)
145         data_fetch['info:postcode'] = str(postcode)
146         hdao.write_data(card_id, data_fetch,'look_up_table')
147     row = {'info:postcode': bytes(postcode),'info:pos_id': bytes(pos_id),'info:card_id': bytes(card_id),'info:amount': bytes(amount),
148            'info:transaction_dt': bytes(transaction_dt),'info:member_id': bytes(member_id), 'info:status': bytes(status)}
149     key = '{0}.{1}.{2}.{3}'.format(card_id,member_id,str(transaction_dt),str(datetime.now())).replace(" ","").replace(":", "")
150     hdao.write_data(bytes(key),row,'card_transactions')
151     return status
152
153 # Defining UDF for Status
154 status_udf = udf(status_find,StringType())
155
156 # Adding Status Column
157 df_parsed = df_parsed.withColumn('status',status_udf(df_parsed.card_id,df_parsed.member_id,df_parsed.amount,df_parsed.pos_id,
158                                                       df_parsed.postcode,df_parsed.transaction_dt,df_parsed.transaction_dt_ts,
159
    df_parsed.last_transaction_date_ts,df_parsed.score,df_parsed.distance,df_parsed.time_diff))
160
```

- Function to define the rules to check if the transaction is genuine or fraudulent.
    1. Transaction amount < UCL
    2. Time difference in sec < 4 times the distance
    3. Credit Score < 200

```
1  # List all the functions to check for the rules
2
3  # Function to define rules
4  def rules_check(UCL,score,distance,time_diff,amount):
5      if amount < UCL:
6          if time_diff < (distance*4):
7              if score > 200:
8                  return True
9      return False
```

- Displaying the required columns in the console and the query has to terminated properly.

```
163
164  # Printing Output on Console
165  query1 = df_parsed \
166          .writeStream \
167          .outputMode("append") \
168          .format("console") \
169          .option("truncate", "False") \
170          .start()
171
172  # Terminating the Query
173  query1.awaitTermination()
```

- Using Putty connect to the EC2 initialized and logging in as root user.
- Connect to the Hbase Shell and look for the card transactions and look up tables which will be used in driver.py script and exit back to root user.
- Create a directory named CapStone and place the driver.py, geomap.py & dao.py (in db), rules.py (in rules) and uszipsv.csv
- In dao.py, the given self.host had to be replace to 'localhost' to run.
- Ensure thrift server is up and running.
- To make sure happybase is imported Give command - python -c 'import happybase'
- In root user the following commands are given:
  - wget https://ds-spark-sql-kafka-jar.s3.amazonaws.com/spark-sql-kafka-0-10_2.11-2.3.0.jar to create jar file required to run the spark code.
  - export SPARK_KAFKA_VERSION=0.10
  - spark2-submit --jars spark-sql-kafka-0-10_2.11-2.3.0.jar --files uszipsv.csv driver.py to execute the driver program.
- The program is executed and the desired output is displayed in the console with coulumns card_id, member_id, amount, pos_id, postcode, transaction_dt_ts and status.

```
-------------------------------------------
Batch: 0
-------------------------------------------
+----------------+------------+--------+----------------+--------+-------------------+-------+
|card_id         |member_id   |amount  |pos_id          |postcode|transaction_dt_ts  |status |
+----------------+------------+--------+----------------+--------+-------------------+-------+
|348702330256514 |37495066290 |4380912 |248063406800722 |96774   |2017-12-31 08:24:29|GENUINE|
|348702330256514 |37495066290 |6703385 |786562777140812 |84758   |2017-12-31 04:15:03|FRAUD  |
|348702330256514 |37495066290 |7454328 |466952571393508 |93645   |2017-12-31 09:56:42|GENUINE|
|348702330256514 |37495066290 |4013428 |45845320330319  |15868   |2017-12-31 05:38:54|GENUINE|
|348702330256514 |37495066290 |5495353 |545499621965697 |79033   |2017-12-31 21:51:54|GENUINE|
|348702330256514 |37495066290 |3966214 |369266342272501 |22832   |2017-12-31 03:52:51|GENUINE|
|348702330256514 |37495066290 |1753644 |9475029292671   |17923   |2017-12-31 00:11:30|FRAUD  |
|348702330256514 |37495066290 |1692115 |27647525195860  |55708   |2017-12-31 17:02:39|GENUINE|
|5189563368503974|117826301530|9222134 |525701337355194 |64002   |2017-12-31 20:22:10|GENUINE|
|5189563368503974|117826301530|4133848 |182031383443115 |26346   |2017-12-31 01:52:32|FRAUD  |
|5189563368503974|117826301530|8938921 |799748246411019 |76934   |2017-12-31 05:20:53|FRAUD  |
|5189563368503974|117826301530|1786366 |131276818071265 |63431   |2017-12-31 14:29:38|GENUINE|
|5189563368503974|117826301530|9142237 |564240259678903 |50635   |2017-12-31 19:37:19|GENUINE|
|5407073344486464|1147922084344|6885448 |887913906711117 |59031   |2017-12-31 07:53:53|FRAUD  |
|5407073344486464|1147922084344|4028209 |116266051118182 |80118   |2017-12-31 01:06:50|FRAUD  |
|5407073344486464|1147922084344|3858369 |896105817613325 |53820   |2017-12-31 17:37:26|GENUINE|
|5407073344486464|1147922084344|9307733 |729374116016479 |14898   |2017-12-31 04:50:16|FRAUD  |
|5407073344486464|1147922084344|4011296 |543373367319647 |44028   |2017-12-31 13:09:34|GENUINE|
|5407073344486464|1147922084344|9492531 |211980095659371 |49453   |2017-12-31 14:12:26|GENUINE|
|5407073344486464|1147922084344|7550074 |345533088112099 |15030   |2017-12-31 02:34:52|FRAUD  |
+----------------+------------+--------+----------------+--------+-------------------+-------+
only showing top 20 rows
```

- In Hbase Shell, scanning card_transactions table, the count of rows after execution is over 59000.

```
Current count: 30000, row: 36164
Current count: 31000, row: 370582035866789.433646648625434.08-07-2018034337.2021-01-04171349.489639
Current count: 32000, row: 375806375521605.880937166605469.26-05-2018130045.2021-01-04171430.733012
Current count: 33000, row: 38176
Current count: 34000, row: 39076
Current count: 35000, row: 39977
Current count: 36000, row: 40768
Current count: 37000, row: 41560
Current count: 38000, row: 42387
Current count: 39000, row: 4318541450654035.496612742732167.12-02-2018145807.2021-01-04171356.009418
Current count: 40000, row: 43999
Current count: 41000, row: 44784
Current count: 42000, row: 45546
Current count: 43000, row: 46306
Current count: 44000, row: 47134
Current count: 45000, row: 47925
Current count: 46000, row: 48730
Current count: 47000, row: 49500
Current count: 48000, row: 50351
Current count: 49000, row: 5120
Current count: 50000, row: 51888
Current count: 51000, row: 5257502990314019.205172644364018.14-07-2018070014.2021-01-04171327.867742
Current count: 52000, row: 53290
Current count: 53000, row: 5620
Current count: 54000, row: 6211
Current count: 55000, row: 6478888441720966.273246841077378.06-10-2018212851.2021-01-04171333.585477
Current count: 56000, row: 6968
Current count: 57000, row: 7868
Current count: 58000, row: 8768
Current count: 59000, row: 9668
59367 row(s) in 3.8140 seconds

=> 59367
```