KafkaSparkStreaming.Java is he driver class for a streaming data processing framework which ingests real-time POS transaction data from Kafka.

NOTE:We already have the details on how data is getting populated in the tables which were created as part of project mid submission.

Below are the steps performed as part of final submission:

## Getting the stream of data(JSON format)  from Kafka stream:

This was achieved in driver class(KafkaSparkStreaming) using below code [Provided by Upgrad]

```
Logger.getLogger("org").setLevel(Level.OFF);
      Logger.getLogger("akka").setLevel(Level.OFF);
      SparkConf sparkConf = new
SparkConf().setAppName("KafkaSparkStreamingDemo").setMaster("lo
cal");
      JavaStreamingContext jssc = new
JavaStreamingContext(sparkConf, Durations.seconds(1));
      Map<String, Object> kafkaParams = new HashMap<>();
      kafkaParams.put("bootstrap.servers", "100.24.223.181:9092");
      kafkaParams.put("key.deserializer", StringDeserializer.class);
      kafkaParams.put("value.deserializer", StringDeserializer.class);
      kafkaParams.put("group.id", "upgraduserkafkasparknavjot77"); //
change here
      kafkaParams.put("auto.offset.reset", "earliest");
      kafkaParams.put("enable.auto.commit", true);
      Collection<String> topics = Arrays.asList("transactions-topic-
verified");

      JavaInputDStream<ConsumerRecord<String, String>> stream =
KafkaUtils.createDirectStream(jssc,
            LocationStrategies.PreferConsistent(),
            ConsumerStrategies.<String, String>Subscribe(topics,
kafkaParams));
      JavaDStream<String> jds = stream.map(x -> x.value());
```

**Cleaning the data:** Since we can get some corrupt data in stream there is an additional logic in code to determine if the JSON message is complete or not:

```
JavaDStream<String> jdsResult = jds.filter(y -> (y!=null &&
y.toString().startsWith("{") && y.toString().endsWith("}")));
```

***Connecting to HBase database to get data from Look up table:***

This was achieved in class HBaseConnection class where below code was written to connect to HBase DB:

```
org.apache.hadoop.conf.Configuration conf =
(org.apache.hadoop.conf.Configuration) HBaseConfiguration.create();

conf.setInt("timeout", 1200);

conf.set("hbase.master",
"ec2-54-86-225-139.compute-1.amazonaws.com:60000");

conf.set("hbase.zookeeper.quorum",
"ec2-54-86-225-139.compute-1.amazonaws.com");

conf.set("hbase.zookeeper.property.clientPort", "2181");

conf.set("zookeeper.znode.parent", "/hbase");

Connection con = ConnectionFactory.createConnection(conf);

//hbaseAdmin = new HBaseAdmin(conf);
hbaseAdmin = con.getAdmin();
```

Also the EC2 name was updated in hosts/etc file while running the code

## Getting the data from HBase DB

Getting the Now once we are able to connect to DB we need to fetch the data for a particular card ID to get the UCL score,Postal Code , Member ID and last transaction date.

NOTE: The table Users_lookup_hb was already created in NoSQL db Hbase(Explaination already provided in Mid Submission)

This was achieved using getStats function where all the fields in the particular record  are parsed ,converted to string and stored in variables in class names HBaseDAO

```
Admin hBaseAdmin1 = HbaseConnection.getHbaseAdmin();

        HTable table = null;
        try {
        Get g = new Get(Bytes.toBytes(card_id));
        table = new HTable(hBaseAdmin1.getConfiguration(),
"Users_lookup_hb");

        Result result = table.get(g);

        byte[] uclByte= result.getValue(Bytes.toBytes("cf"),
Bytes.toBytes("ucl"));
        byte[] scoreByte = result.getValue(Bytes.toBytes("cf"),
Bytes.toBytes("score"));
        byte[] postCodeByte= result.getValue(Bytes.toBytes("cf"),
Bytes.toBytes("postcode"));
        byte[] transactionDateByte= result.getValue(Bytes.toBytes("cf"),
Bytes.toBytes("transaction_dt"));

        String ucl = (uclByte != null) ? Bytes.toString(uclByte):"0";
        String score = (scoreByte != null) ?
Bytes.toString(scoreByte):"-1";
        String postCode = (postCodeByte != null) ?
Bytes.toString(postCodeByte):"0";
        String transactionDate = (transactionDateByte != null) ?
Bytes.toString(transactionDateByte):"0";
```

```
    // reading data from look up table
        return new
CardAnalysis(ucl+","+score+","+postCode+","+transactionDate);
```

**Updating the look up table:**This class also had logic to Update the card look up table based on the incoming stream.Function updatePostCodeTranDate was created where postCode and transactionDate were updated:

```
public static boolean updatePostCodeTranDate(String card_id,String postCode, String transactionDate) throws IOException {

        public static boolean updatePostCodeTranDate(String card_id,String postCode, String transactionDate) throws IOException {

        Admin hBaseAdmin1 = HbaseConnection.getHbaseAdmin();

        HTable table = null;
        try {
        Put p = new Put(Bytes.toBytes(card_id));
        table = new HTable(hBaseAdmin1.getConfiguration(), "Users_lookup_hb");
        p.add(Bytes.toBytes("cf"),
                Bytes.toBytes("postcode"),Bytes.toBytes(postCode));
        p.add(Bytes.toBytes("cf"),

Bytes.toBytes("transaction_dt"),Bytes.toBytes(transactionDate));
        table.put(p);
        return true;

        } catch (Exception e) {

        e.printStackTrace();
```

```
    } finally {

    try {

    if (table != null)

    table.close();

    } catch (Exception e) {

    e.printStackTrace();

    }

    }
    return false;

    }
```

**Updating the card transactions:** This I am assuming that we will be uploading card_transactions.csv to HDFS and load the data into Hive as per part of Mid submission.


## Validating if the transaction is fraud or not:


Data Validation: Logic is contained in CardAnalysis.java class. Once we are able to get all the three values(UCL,Score and postal code) , transaction data is then validated based on the three rules' parameters:

1) **Check UCL :** Here we are comparing the current amount to the UCL amount and determining if the current transaction is fraud or not:

```
   public boolean checkUCL(double newAmount) {
        if(this.upperControlLimit > 0)
        {
```

```
                    return this.upperControlLimit > newAmount ?
true:false;
            }
            else
                return true;
    }
```

2) **Check Member Score:** Below is the method used to determine the transaction based on member score:

```
public boolean checkScore() {
            if(this.score > -1)
                    return this.score > 200 ? true:false;
            else
                    return true;
    }
```

3) **Check Postal code:** This method is utilising DistanceUtility to calculate the distance between two postal codes -- This current postal code vs last transaction postal code fetched from the card look up table.Once the distance is calculated we calculate the speed based on the time difference and if the speed is more than 1000 KM/hr we treat that transaction as fraud:

```
public boolean checkPostCode(String newPostCode,String
newTransactionDate) throws ParseException,
NumberFormatException, IOException {
            SimpleDateFormat formatter = new SimpleDateFormat("dd-
MM-yyyy HH:mm:ss");
            DistanceUtility distUtility = new DistanceUtility();
            if(this.transactionDate != "0")
            {
                    double dist =
distUtility.getDistanceViaZipCode(newPostCode, this.postCode);
```

```
                double time =
(Math.abs(formatter.parse(newTransactionDate).getTime() -
formatter.parse(this.transactionDate).getTime())/ (1000*60*60));
                // if the speed < 1000 km/hr then GENUNIE else it is
FRAUD
                return dist/time < 1000 ? true:false;
        }
        else if(newTransactionDate != null)
                return true;
        else
                return false;

    }
```

**Displaying the output:** Once the transaction has been classified as fraud or genuine we display the output along with the driver class in driver class:

```
JavaDStream<String> jdsResult1 = jdsResult.map(new
Function<String,String>(){
        private static final long serialVersionUID = 1L;
        @Override
        public String call(String t) throws IOException,
NumberFormatException, ParseException {
                ObjectMapper maper = new ObjectMapper();
                TranData transactionData = maper.readValue(t,
TranData.class);
                PrintWriter writer = new PrintWriter(new File("data/
test.csv")) ;
            String status=null;
                CardAnalysis cardStats =
HbaseDAO.getStats(transactionData.getCard_id());
                if(cardStats != null &&
cardStats.checkUCL(transactionData.getAmount())
                        && cardStats.checkScore() &&
cardStats.checkPostCode(transactionData.getPostcode(),
transactionData.getTransaction_dt()))
                {
                    // Transaction validated successfully on all 3 rules
                      status="GENUINE";
```
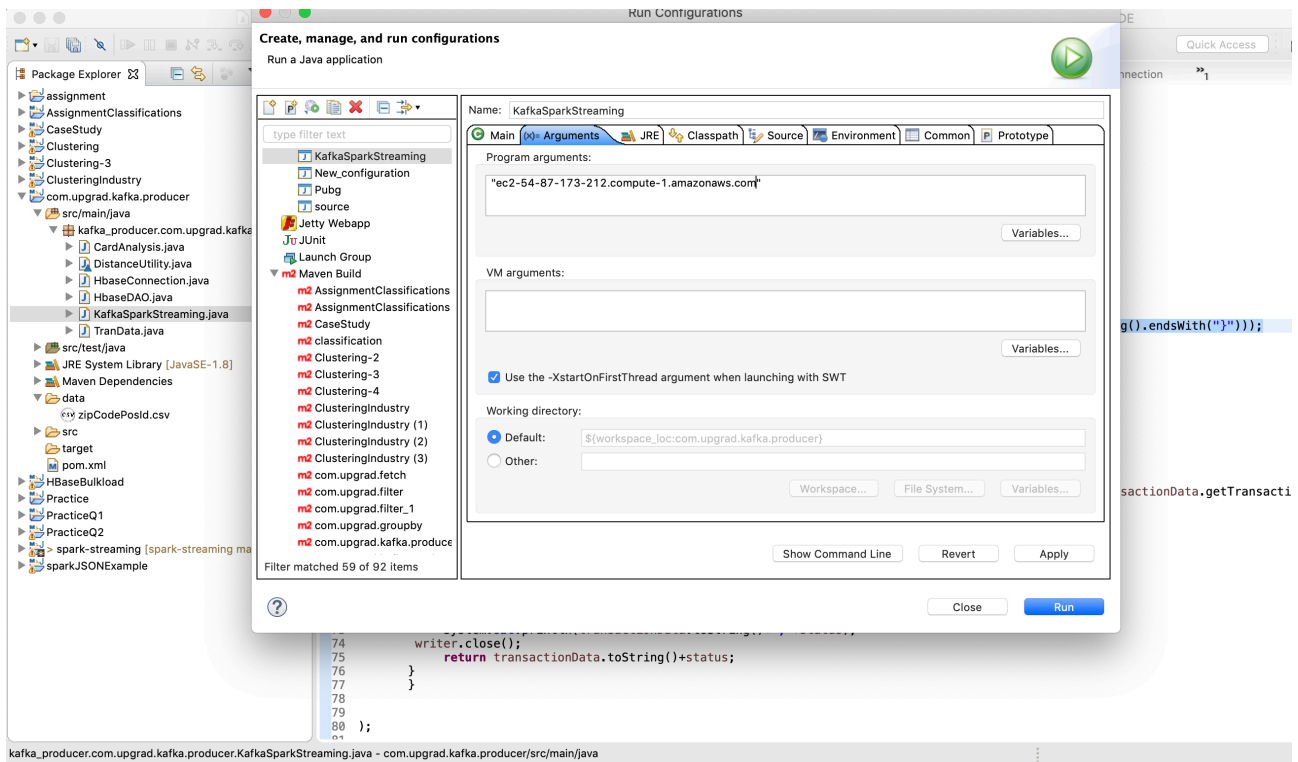
```
            }
            else
            {
                    // fraud transaction
                    status="FRAUD";
            }
        System.out.println(transactionData.toString()+","+status);

        writer.close();
            return transactionData.toString()+status;
        }
        }
```
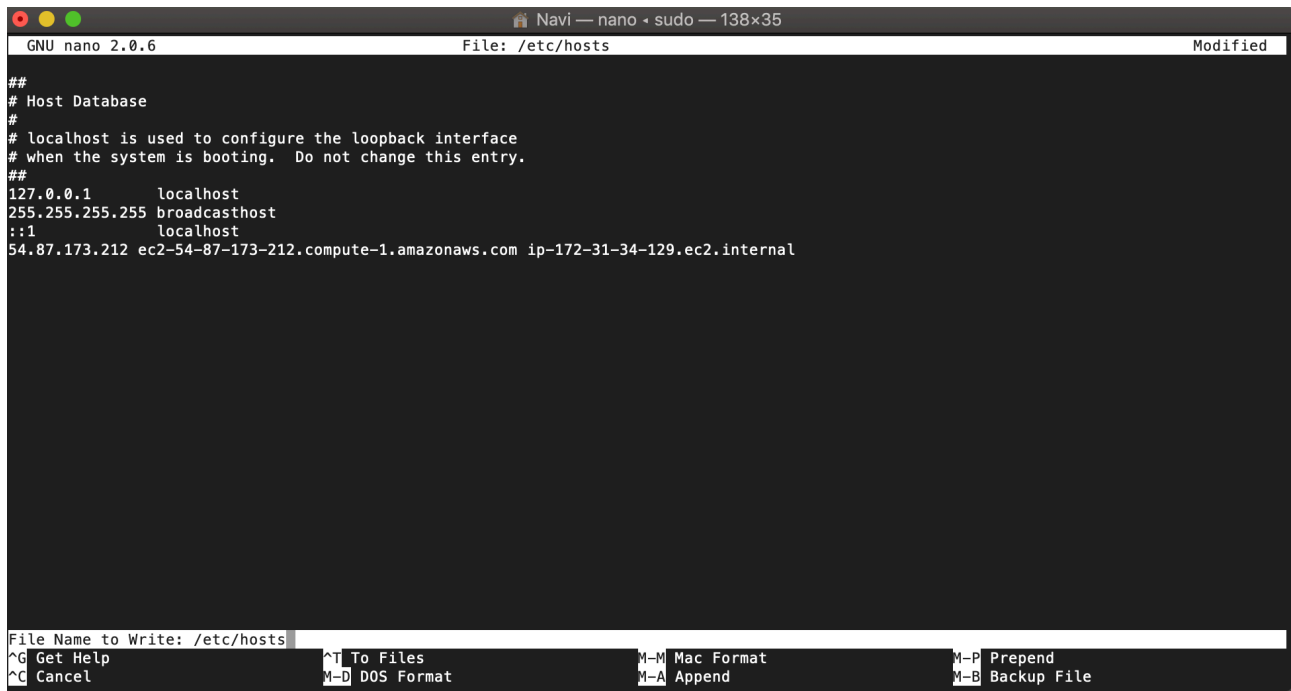
**Console output screenshot:** Please find below the initial screenshot where i have passed my EC2 IP as parameter(Full logs are attached separately)
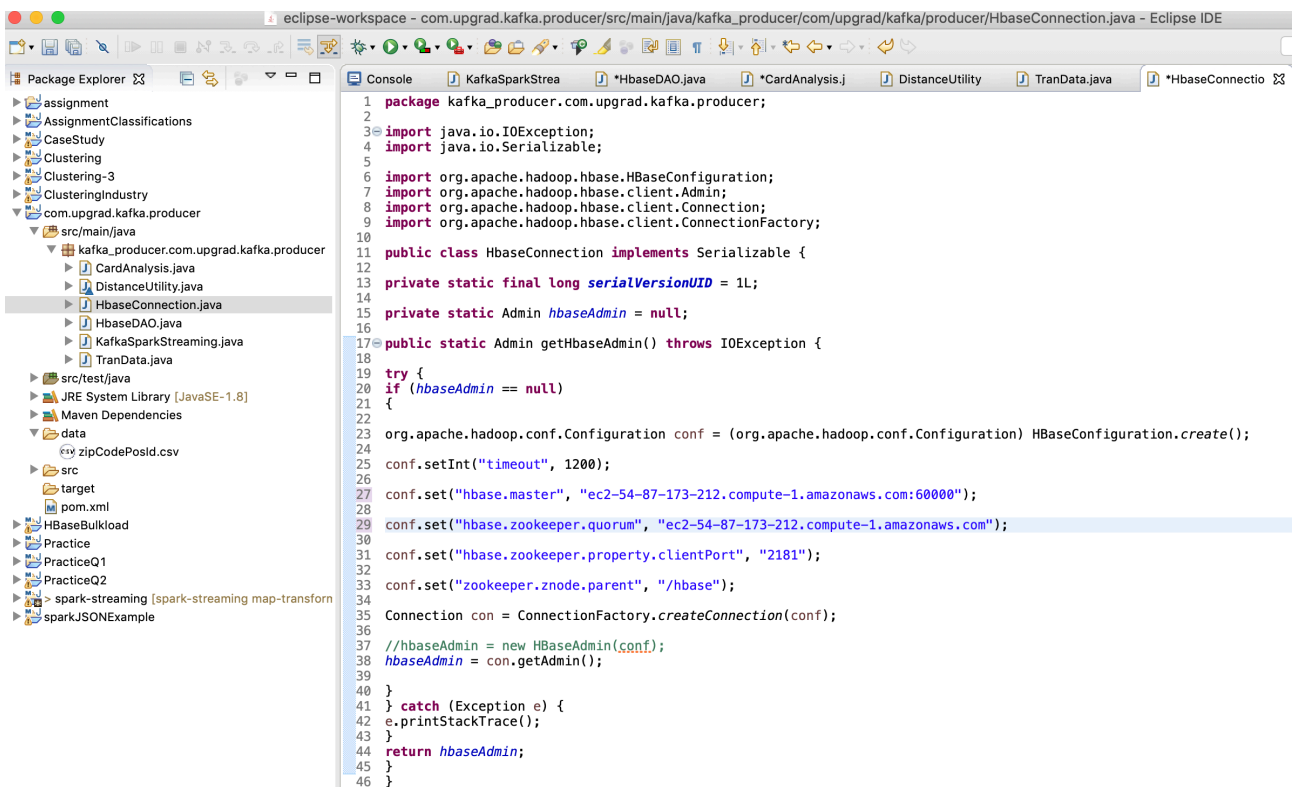
Argument settings:

## Etc/host settings:

```
GNU nano 2.0.6                    File: /etc/hosts                                    Modified

##
# Host Database
#
# localhost is used to configure the loopback interface
# when the system is booting.  Do not change this entry.
##
127.0.0.1       localhost
255.255.255.255 broadcasthost
::1             localhost
54.87.173.212 ec2-54-87-173-212.compute-1.amazonaws.com ip-172-31-34-129.ec2.internal




File Name to Write: /etc/hosts
^G Get Help        ^T To Files       M-M Mac Format    M-P Prepend
^C Cancel          M-D DOS Format    M-A Append        M-B Backup File
```

## HBaseConnection settings:

```java
package kafka_producer.com.upgrad.kafka.producer;

import java.io.IOException;
import java.io.Serializable;

import org.apache.hadoop.hbase.HBaseConfiguration;
import org.apache.hadoop.hbase.client.Admin;
import org.apache.hadoop.hbase.client.Connection;
import org.apache.hadoop.hbase.client.ConnectionFactory;

public class HbaseConnection implements Serializable {

    private static final long serialVersionUID = 1L;

    private static Admin hbaseAdmin = null;

    public static Admin getHbaseAdmin() throws IOException {

        try {
            if (hbaseAdmin == null)
            {

                org.apache.hadoop.conf.Configuration conf = (org.apache.hadoop.conf.Configuration) HBaseConfiguration.create();

                conf.setInt("timeout", 1200);

                conf.set("hbase.master", "ec2-54-87-173-212.compute-1.amazonaws.com:60000");

                conf.set("hbase.zookeeper.quorum", "ec2-54-87-173-212.compute-1.amazonaws.com");

                conf.set("hbase.zookeeper.property.clientPort", "2181");

                conf.set("zookeeper.znode.parent", "/hbase");

                Connection con = ConnectionFactory.createConnection(conf);

                //hbaseAdmin = new HBaseAdmin(conf);
                hbaseAdmin = con.getAdmin();

            }
        } catch (Exception e) {
            e.printStackTrace();
        }
        return hbaseAdmin;
    }
}
```

# Saving the console logs in an external file:



# Console Output: