

# CREDIT CARD FRAUD DETECTION – REAL TIME STREAMING

## PROBLEM APPROACH

1. Read the source card transactions data from Kafka topic (**transactions-topic-verified**) and define the proper data structure / schema on top of the streaming dataframe.
2. Derive the column “**status**” with two distinct values as “**GENUINE**” and “**FRAUD**” based on the following three rules:
  - a. **Credit Score:** - Get the credit score from card lookup HBase table based on card id fetched from the source transactions data and validate it. If the score is less than 200 then label it as “FRAUD” and move to next transaction. Otherwise proceed validating the second rule.
  - b. **Check Upper Control Limit (UCL):** - Get the amount from the transaction data and compare it with UCL value present in card lookup HBase table. If the amount is greater than UCL then label it as “FRAUD” and move to next transaction. Otherwise proceed validating the third and final rule.
  - c. **Zip code Distance:** - Here we calculate the distance travelled between two locations of the recent two transactions and its time difference and measure speed in terms of Kilometers/Second. We set the threshold value to **0.25 Kilometers/Second (900 Kilometers/Hour – average speed of commercial flights)**. If the speed measured is greater than threshold value, then label it as “FRAUD”. Otherwise, it can be treated as “GENUINE”.
3. Once the transaction status is derived, we can insert the card transactions source data along with its status in card transactions HBase table.
4. Update the postcode and transaction date in card lookup HBase table for the transactions happened only if the status is labelled as “GENUINE”.
5. After processing all the transactions within the same micro batch, display the card transactions information with its status in the console for live analysis.

## CODE IMPLEMENTATION

1. Import the necessary libraries related to PySpark, Spark Streaming, classes and functions definitions provided as boilerplate code (**dao.py, geo\_map.py and rules.py**).

```
# Import necessary libraries
from datetime import datetime
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, from_json, to_timestamp, udf
from pyspark.sql.types import StructType, StructField, StringType, IntegerType, LongType, TimestampType
from db.dao import HBaseDao
from db.geo_map import GEO_Map
from rules.rules import is_score_rejected, is_ucl_exceeded, is_zipcode_invalid
```

2. Define the required utility functions and UDF on top of it.
  - a. **get\_distance():** - This function is to calculate the distance between current and previous locations where transactions happened based on its latitude and longitude. It instantiates and calls GEO\_Map class (geo\_map.py) methods to get the latitude and longitude attributes value and calculate the distance in Kilometers.

### geo\_map.py

```

import math
import pandas as pd

class GEO_Map():
    """
    It hold the map for zip code and its latitude and longitude
    """
    __instance = None

    @staticmethod
    def get_instance():
        """ Static access method. """
        if GEO_Map.__instance == None:
            GEO_Map()
        return GEO_Map.__instance

    def __init__(self):
        """ Virtually private constructor. """
        if GEO_Map.__instance != None:
            raise Exception("This class is a singleton!")
        else:
            GEO_Map.__instance = self
            self.map = pd.read_csv("uszipsv.csv", header=None, names=['A','B','C','D','E'])
            self.map['A'] = self.map['A'].astype(str)

    def get_lat(self, pos_id):
        return self.map[self.map.A == pos_id ].B

    def get_long(self, pos_id):
        return self.map[self.map.A == pos_id ].C

    def distance(self, lat1, long1, lat2, long2):
        theta = long1 - long2
        dist = math.sin(self.deg2rad(lat1)) * math.sin(self.deg2rad(lat2)) + math.cos(self.deg2rad(lat1)) * math.cos(self.deg2rad(lat2)) * math.cos(self.deg2rad(theta))
        dist = math.acos(dist)
        dist = self.rad2deg(dist)
        dist = dist * 60 * 1.1515 * 1.609344
        return dist

    def rad2deg(self, rad):
        return rad * 180.0 / math.pi

```

### driver.py

```

# Utility function to calculate distance between last two locations where transactions happened
def get_distance(trans_postcode, lkp_postcode):
    """
    Calculate the distance between current and previous locations where transactions happened
    Input Arguments: trans_postcode (postcode from card transactions real time data - Kafka topic in JSON format)
                    lkp_postcode (postcode from card_lookup table (Hbase) of the same card_id from card_transactions real time data).
    """
    # Instantiate GEO_Map class object
    geo_obj = GEO_Map.get_instance()
    lat1 = float(geo_obj.get_lat(trans_postcode).values[0])
    long1 = float(geo_obj.get_long(trans_postcode).values[0])
    lat2 = float(geo_obj.get_lat(lkp_postcode).values[0])
    long2 = float(geo_obj.get_long(lkp_postcode).values[0])
    return geo_obj.distance(lat1, long1, lat2, long2)

```

- b. **get\_time\_diff\_secs():** - This function calculates the time difference between current and previous transactions timestamp.

### driver.py

```

# Utility function to calculate time difference between last two transactions happened
def get_time_diff_secs(trans_transactiondt, lkp_transactiondt):
    """
    Calculate the time difference between current and previous transactions timestamp.
    Input Arguments: trans_transactiondt (transaction_dt from card transactions real time data - Kafka topic in JSON format)
                    lkp_transactiondt (transaction_dt from card_lookup table (Hbase) of the same card_id from card_transactions real time data).
    """
    return (trans_transactiondt - lkp_transactiondt).total_seconds()

```

- c. **insert card transactions hbase():** - This function inserts card transactions real time data (Kafka topic in JSON format) into card\_transactions HBase table. It first instantiates HbaseDao class object to connect to HBase database and then call required methods to insert the data in card transactions table.

#### dao.py

```
import happybase

class HBaseDao:
    """
    Dao class for operation on HBase
    """
    __instance = None

    @staticmethod
    def get_instance():
        """ Static access method. """
        if HBaseDao.__instance == None:
            HBaseDao()
        return HBaseDao.__instance

    def __init__(self):
        if HBaseDao.__instance != None:
            raise Exception("This class is a singleton!")
        else:
            HBaseDao.__instance = self
            self.host = 'localhost'
            self.connect()

    def connect(self):
        for i in range(3):
            try:
                self.pool = happybase.ConnectionPool(size=3, host=self.host, port=9090)
                break
            except:
                print("Exception in connecting HBase")

    def get_data(self, key, table):
        for i in range(2):
            try:
                with self.pool.connection() as connection:
                    t = connection.table(table)
                    row = t.row(key)
                    return row
            except:
                self.reconnect()

"dao.py" 60L, 1168B
```

```
def write_data(self, key, row, table):
    error = None
    for i in range(2):
        try:
            with self.pool.connection() as connection:
                t = connection.table(table)
                return t.put(key, row)
        except Exception as e:
            error = e
            self.reconnect()
    raise Exception(key + str(e))

def reconnect(self):
    self.connect()
```

## driver.py

```
# Write card transactions final data to HBase table
def insert_card_transactions_hbase(transaction_data):
    """
    Insert card transactions real time data (Kafka topic in JSON format) into card transactions hbase table.
    Input Arguments: transaction_data (list object of card transactions real time data - Kafka topic in JSON format).
    """
    # Instantiate HBaseDao class object to connect to HBase
    hbase_obj = HBaseDao.get_instance()
    trans_key = '|'.join([str(v) for v in transaction_data])
    data = [b'transaction_data:card_id':bytes(str(transaction_data[0]), encoding='utf8'),
            b'transaction_data:member_id':bytes(str(transaction_data[1]), encoding='utf8'),
            b'transaction_data:amount':bytes(str(transaction_data[2]), encoding='utf8'),
            b'transaction_data:pos_id':bytes(str(transaction_data[3]), encoding='utf8'),
            b'transaction_data:postcode':bytes(str(transaction_data[4]), encoding='utf8'),
            b'transaction_data:transaction_dt':bytes(str(transaction_data[5]), encoding='utf8'),
            b'transaction_data:status':bytes(str(transaction_data[6]), encoding='utf8')]
    hbase_obj.write_data(bytes(trans_key, encoding='utf8'), data, 'card_transactions')
```

- d. **update card lookup hbase():** - This function updates postcode and transaction\_dt attributes in card lookup table HBase for which transactions are treated as GENUINE. It first instantiates HbaseDao class object to connect to HBase database and then call required methods to update the card lookup data.

## driver.py

```
# Update latest postcode and transaction date in card lookup HBase table
def update_card_lookup_hbase(trans_card_id, status_flag, lkp_data):
    """
    Update postcode and transaction_dt attributes in card lookup table HBase for which transactions are treated as GENUINE.
    Input Arguments: trans_card_id
                    status_flag
                    lkp_data (list object of postcode and transaction_dt card transactions real time data - Kafka topic in JSON format).
    """
    # Instantiate HBaseDao class object to connect to HBase
    hbase_obj = HBaseDao.get_instance()
    if status_flag == 'GENUINE':
        key = bytes(str(trans_card_id), encoding='utf8')
        data = [b'lkp_info:postcode':bytes(str(lkp_data[0]), encoding='utf8'),
                b'lkp_info:transaction_dt':bytes(str(lkp_data[1]), encoding='utf8')]
        hbase_obj.write_data(key, data, 'card_lookup')
```

- e. **check status():** - This function is to detect the fraudulent of the live transactions based on the rules defined (**mentioned in the problem approach**) and lookup data prepared as part of batch processing.

## rules.py

```
def is_score_rejected(score):
    """
    This function is to get the credit score of the particular card which transaction happened and see whether it has greater than or equal to 200. If not then it is labelled as Fraud.
    Input Arguments: score (credit_score from card_transactions real time data - Kafka topic in JSON format)
    """
    if score < 200:
        return True
    return False

def is_ucl_exceeded(amount, ucl):
    """
    This function is to check whether the amount exceeds the upper control limit of the transaction to be done. If it is exceeded then it is labelled as Fraud.
    Input Arguments: amount (card_transactions real time data - Kafka topic in JSON format)
                    ucl (UCL attribute value from card_lookup table - HBase)
    """
    if amount > ucl:
        return True
    return False

def is_zipcode_invalid(distance, time_diff_secs):
    """
    This function is to measure the speed in Kilometers/second travelled by the customer between two locations where customer made the recent two transactions (current and previous).
    Usually commercial airplanes travel with an average speed of around 800-950 Kilometers/hour. We considered the speed threshold value as 900 Kilometers/hour and converted it to seconds which is (900/3600) = 0.25 Kilometers/second. If the measure speed between two locations is greater than 0.25 Kilometers/second then it is labelled as Fraud.
    Input Arguments: distance (distance covered between two locations of recent two transactions)
                    time_diff_secs (Time difference between last two transactions happened)
    """
    if time_diff_secs == 0:
        if distance != 0:
            return True
        else:
            speed_kmps = distance / time_diff_secs
            if speed_kmps > 0.25:
                return True
            return False
    return False
```

## driver.py

```
def check_status(trans_card_id, trans_member_id, trans_amount, trans_pos_id, trans_postcode, trans_transactiondt):
    """
    This function is to detect the fraudulent of the live transactions based on the rules defined and lookup data prepared as part of batch processing.
    The following three rules to be validated:
    1. Check credit score
    2. Check UCL threshold with amount
    3. Measure the speed travelled in Kilometers/second between the locations of last two transactions.
    If any of the above rule fails then the transaction is labelled as FRAUD. Otherwise it is labelled as GENUINE.
    Input Arguments: trans_card_id (card id attribute from card_transactions real time data - Kafka topic in JSON format)
                    trans_member_id (member id attribute from Card_transactions real time data - Kafka topic in JSON format)
                    trans_amount (transaction amount from card_transactions real time data - Kafka topic in JSON format)
                    trans_pos_id (pos_id from card_transactions real time data - Kafka topic in JSON format)
                    trans_postcode (postcode from card_transactions real time data - Kafka topic in JSON format)
                    trans_transactiondt (transaction_dt from card_transactions real time data - Kafka topic in JSON format)
    """
    status_flag = 'GENUINE'
    # Instantiate HBaseDao class object to connect to HBase
    hbase_obj = HBaseDao.get_instance()
    lkp_row = hbase_obj.get_data(str(trans_card_id), 'card_lookup')
    if lkp_row == {}:
        pass
    else:
        lkp_score, lkp_postcode, lkp_transactiondt, lkp_ucl = [val.decode('utf8') for val in list(lkp_row.values())]
        lkp_score = int(lkp_score)
        lkp_postcode = int(lkp_postcode)
        lkp_transactiondt = datetime.strptime(lkp_transactiondt, '%Y-%m-%d %H:%M:%S')
        lkp_ucl = float(lkp_ucl)

        # Check transactions rules to determine credit score, UCL and zipcode distance.
        if is_score_rejected(lkp_score):
            status_flag = 'FRAUD'
        elif is_ucl_exceeded(trans_amount, lkp_ucl):
            status_flag = 'FRAUD'
        elif is_zipcode_invalid(get_distance(str(trans_postcode), str(lkp_postcode)), get_time_diff_secs(trans_transactiondt, lkp_transactiondt)):
            status_flag = 'FRAUD'
        else:
            pass

    # Call insert_card_transactions_hbase() function to insert the card transaction data with status derived
    insert_card_transactions_hbase(trans_card_id, trans_member_id, trans_amount, trans_pos_id, trans_postcode, trans_transactiondt, status_flag)
    # Call update_card_lookup_hbase() to Insert/update postcode and transaction_dt attributes in card_lookup HBase table
    update_card_lookup_hbase(trans_card_id, status_flag, (trans_postcode, trans_transactiondt))

    return status_flag
```

- f. Define UDF (fraud\_detection) on top of check\_status() utility function that to be used in spark streaming dataframe for deriving “status” attribute during real time. As part of this UDF, it insert the records into card transactions HBase table and update the postcode and transaction date in card lookup HBase table (as shown in check\_status() driver.py screenshot above).

## driver.py

```
# Create fraud detection UDF on top of status utility function
fraud_detection = udf(check_status, StringType())
```

3. Create a Spark Session which is an entry point to spark application and set the log level to 'ERROR'.

## driver.py

```
# Create a spark session
spark = SparkSession \
    .builder \
    .appName("Fraud Detection - Card Transactions") \
    .getOrCreate()

# Set log level to ERROR
spark.sparkContext.setLogLevel('ERROR')
```

- Read the card transactions real time data from Kafka source topic (transactions-topic-verified).

#### driver.py

```
# Read card transactions from Kafka topic : transactions-topic-verified
card_trans_source = spark \
    .readStream \
    .format("kafka") \
    .option("kafka.bootstrap.servers", "18.211.252.152:9092") \
    .option("subscribe", "transactions-topic-verified") \
    .option("startingOffsets", "earliest") \
    .load()
```

- Define the proper data structure or schema and assign it to the source data read from Kafka topic.

#### driver.py

```
# Defining the Data Structure for source card transactions (json format).
data_structure = StructType([
    StructField("card_id", LongType()),
    StructField("member_id", LongType()),
    StructField("amount", IntegerType()),
    StructField("pos_id", LongType()),
    StructField("postcode", IntegerType()),
    StructField("transaction_dt", StringType())
])

# Converting source data from json format to dataframe (table).
card_trans_formatted = card_trans_source \
    .select(from_json(col("value").cast("string"), data_structure).alias("card_transactions")) \
    .select("card_transactions.*") \
    .withColumn("transaction_dt", to_timestamp(col('transaction_dt'), 'dd-MM-yyyy HH:mm:ss')) \
    .select("card_id",
        "member_id",
        "amount",
        "pos_id",
        "postcode",
        "transaction_dt"
    )
```

- Add a new column “status” by calling the fraud\_detection UDF defined above. This UDF will also take of inserting and updating the records in card transactions and card lookup HBase tables.

#### driver.py

```
# Derive status attribute: GENUINE or FRAUD based on rules defined
card_trans_transformed = card_trans_formatted \
    .withColumn("status", fraud_detection(card_trans_formatted.card_id, card_trans_formatted.member_id, card_trans_formatted.amount, card_trans_
formatted.pos_id, card_trans_formatted.postcode, card_trans_formatted.transaction_dt)) \
    .select("card_id",
        "member_id",
        "amount",
        "pos_id",
        "postcode",
        "transaction_dt",
        "status"
    )
```

- Write the final transformed card transactions data with its status in the console for live analysis. It runs as a micro-batch with processing time of every “1 Minute” until it is interrupted/terminated.

driver.py

```
# Write the final transformed data to sink
card_trans_sink = card_trans_transformed \
    .writeStream \
    .outputMode("append") \
    .format("console") \
    .option("truncate", "false") \
    .trigger(processingTime="1 minute") \
    .start()

# Await termination
card_trans_sink.awaitTermination()
```

## DEPLOYMENT PLAN

### PRE-REQUISITES:

- Create an EMR cluster with **Hive, HBase, Sqoop, Spark, Zookeeper** applications installed in it as shown below:

Cluster: Credit\_Card\_Fraud\_Detection Waiting Cluster ready to run steps.

Summary	Configuration details
<p>ID: j-1QIUC79D8YE0A</p> <p>Creation date: 2023-01-08 06:16 (UTC+5:30)</p> <p>Elapsed time: 19 minutes</p> <p>After last step completes: Cluster waits</p> <p>Termination protection: Off <a href="#">Change</a></p> <p>Tags: -- <a href="#">View All / Edit</a></p> <p>Master public DNS: ec2-3-85-43-185.compute-1.amazonaws.com <a href="#">Connect to the Master Node Using SSH</a></p>	<p>Release label: emr-5.30.1</p> <p>Hadoop distribution: Amazon 2.8.5</p> <p><b>Applications: Hive 2.3.6, Hue 4.6.0, JupyterHub 1.1.0, Zeppelin 0.8.2, HBase 1.4.13, Sqoop 1.4.7, Spark 2.4.5, ZooKeeper 3.4.14, Livy 0.7.0</b></p> <p>Log URI: s3://aws-logs-145313842850-us-east-1/elasticmapreduce/</p> <p>EMRFS consistent view: Disabled</p> <p>Custom AMI ID: --</p>
<p>Application user interfaces</p> <p>Persistent user interfaces <a href="#">Spark history server, YARN timeline server, Tez UI</a></p> <p>On-cluster user interfaces <a href="#">Not Enabled</a> <a href="#">Enable an SSH Connection</a></p>	<p>Network and hardware</p> <p>Availability zone: us-east-1c</p> <p>Subnet ID: <a href="#">subnet-0f80dc2e</a></p> <p>Master: <span>Running</span> 1 m5.xlarge</p> <p>Core: --</p> <p>Task: --</p> <p>Cluster scaling: Not enabled</p> <p>Auto-termination: Not enabled</p>

Applications such as Hue, JupyterHub, Zeppelin, Livy are optional and depends based on your preference.

- Download and install MySQL connector in AWS EMR for performing Sqoop data ingestion from AWS RDS to HDFS.

**cd /home/hadoop**

**wget https://de-mysql-connector.s3.amazonaws.com/mysql-connector-java-8.0.25.tar.gz**

**tar -xvf mysql-connector-java-8.0.25.tar.gz**

```
cd mysql-connector-java-8.0.25/
sudo cp mysql-connector-java-8.0.25.jar /usr/lib/sqoop/lib/
```

- Download and install Hive-HBase handler jar file which will be used for card lookup data preparation for Hive and Hbase tables using PySpark. and place it in spark libraries folder. Also place some other required HBase jar files from HBase libraries to spark libraries folder.

```
cd /home/hadoop
wget https://repo1.maven.org/maven2/org/apache/hive/hive-hbase-
handler/3.1.3/hive-hbase-handler-3.1.3.jar
sudo cp hive-hbase-handler-3.1.3.jar /usr/lib/spark/jars/
sudo cp /usr/lib/hbase/*.jar /usr/lib/spark/jars/
sudo cp /usr/lib/hbase/lib/htrace-core-3.1.0-incubating.jar /usr/lib/spark/jars/
sudo cp /usr/lib/hbase/lib/guava-12.0.1.jar /usr/lib/spark/jars/
sudo cp /usr/lib/hbase/lib/metrics-core-2.2.0.jar /usr/lib/spark/jars/
```

- Create all the required tables in Hive and HBase such as **card\_member**, **member\_score**, **card\_transactions (historical data)** and **card\_lookup** tables and load the pre-processed data in it as part of batch processing (mid submission).

HBase:

```
hbase(main):001:0> list
TABLE
card_lookup
card_transactions
2 row(s) in 0.1870 seconds

=> ["card_lookup", "card_transactions"]
hbase(main):002:0> █
```

Hive:

```
[hadoop@ip-172-31-84-212 ~]$ hive

Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: false
hive> show databases;
OK
cred_financials_data
default
Time taken: 0.499 seconds, Fetched: 2 row(s)
hive> use cred_financials_data;
OK
Time taken: 0.019 seconds
hive> show tables;
OK
card_lookup
card_member
card_transactions
member_score
Time taken: 0.024 seconds, Fetched: 4 row(s)
hive> █
```



- Switch to root user in EMR cluster and install pandas and happybase python packages.

**sudo su root**

**yum install gcc (Press 'Y' and ENTER after it prompts to download)**

**yum install python3-devel (Press 'Y' and ENTER after it prompts to download)**

**pip install cython**

**pip install numpy**

**pip install pandas**

**pip install happybase**

**exit**

## DEPLOYMENT STEPS:

- Login to AWS EMR and create the python directory under home directory and place all the executable python files inside it (either using scp command or WinSCP tool) and change its permissions.

**cd /home/hadoop**

**mkdir -p /home/hadoop/python/src**

**Place the files using scp command from local to EMR**

```
C:\Users\Ganesh Jalakam>scp -r -P 22 -i "C:\Users\Ganesh Jalakam\Documents\AWS\Instance_Keys\RHEL.pem" "D:\IIITB_Data_Science\PG_DS_Course\Capstone_Project\Final_Submission\python\*" hadoop@ec2-3-85-43-185.compute-1.amazonaws.com:/home/hadoop/python/
dao.py 100% 1168 4.9KB/s 00:00
geo_map.py 100% 1205 4.3KB/s 00:00
__init__.py 100% 0 0.0KB/s 00:00
dao.cpython-37.pyc 100% 1732 6.7KB/s 00:00
__init__.cpython-37.pyc 100% 125 0.4KB/s 00:00
driver.py 100% 10KB 44.5KB/s 00:00
rules.py 100% 2760 11.9KB/s 00:00
__init__.py 100% 0 0.0KB/s 00:00
uszipsv.csv 100% 735KB 474.0KB/s 00:01
__init__.py 100% 0 0.0KB/s 00:00
C:\Users\Ganesh Jalakam>
```

**Change the executable files permissions:**

**cd /home/hadoop/python/src**

**chmod 755 \*.py**

**chmod 755 db/\*.py**

**chmod 755 rules/\*.py**

**driver.py - /home/hadoop/python/src**

```
[hadoop@ip-172-31-84-212 src]$ pwd
/home/hadoop/python/src
[hadoop@ip-172-31-84-212 src]$ ls -ltr
total 748
drwxrwxrwx 3 hadoop hadoop 76 Jan 8 01:35 db
-rwxr-xr-x 1 hadoop hadoop 10319 Jan 8 01:36 driver.py
drwxrwxrwx 2 hadoop hadoop 41 Jan 8 01:36 rules
-rw-rw-r-- 1 hadoop hadoop 752688 Jan 8 01:36 uszipsv.csv
-rwxr-xr-x 1 hadoop hadoop 0 Jan 8 01:36 __init__.py
[hadoop@ip-172-31-84-212 src]$
```

## dao.py and geo\_map.py – /home/hadoop/python/src/db

```
/home/hadoop/python/src/db
[hadoop@ip-172-31-84-212 db]$ ls -ltr
total 8
-rwxr-xr-x 1 hadoop hadoop 1168 Jan  8 01:35 dao.py
-rwxr-xr-x 1 hadoop hadoop 1205 Jan  8 01:35 geo_map.py
-rwxr-xr-x 1 hadoop hadoop    0 Jan  8 01:35 __init__.py
drwxrwxrwx 2 hadoop hadoop  63 Jan  8 01:35 __pycache__
[hadoop@ip-172-31-84-212 db]$
```

## rules.py - /home/hadoop/python/src/rules

```
/home/hadoop/python/src/rules
[hadoop@ip-172-31-84-212 rules]$ ls -ltr
total 4
-rwxr-xr-x 1 hadoop hadoop 2760 Jan  8 01:36 rules.py
-rwxr-xr-x 1 hadoop hadoop    0 Jan  8 01:36 __init__.py
[hadoop@ip-172-31-84-212 rules]$
```

2. Zip all the python dependencies into one file with name src.zip (containing dao.py, geo\_map.py and rules.py in it).

**cd /home/hadoop/python/src**

**zip src.zip \_\_init\_\_.py rules/\* db/\***

```
[hadoop@ip-172-31-84-212 src]$ zip src.zip __init__.py rules/* db/*
adding: __init__.py (stored 0%)
adding: rules/__init__.py (stored 0%)
adding: rules/rules.py (deflated 68%)
adding: db/dao.py (deflated 61%)
adding: db/geo_map.py (deflated 56%)
adding: db/__init__.py (stored 0%)
adding: db/__pycache__/ (stored 0%)
[hadoop@ip-172-31-84-212 src]$ ls -ltr
total 752
drwxrwxrwx 3 hadoop hadoop    76 Jan  8 01:35 db
-rwxr-xr-x 1 hadoop hadoop 10319 Jan  8 01:36 driver.py
drwxrwxrwx 2 hadoop hadoop   41 Jan  8 01:36 rules
-rw-rw-r-- 1 hadoop hadoop 752688 Jan  8 01:36 uszipsv.csv
-rwxr-xr-x 1 hadoop hadoop    0 Jan  8 01:36 __init__.py
-rw-rw-r-- 1 hadoop hadoop  2968 Jan  8 01:53 src.zip
```

3. Run the following command to enable Kafka Integration with Apache Spark.  
**export SPARK\_KAFKA\_VERSION=0.10**

- Run driver.py script to start spark streaming job which consumes real time transactions data and detects the status (FRAUD or GENUINE) and display the final output of each micro batches in the console.

**spark-submit --py-files src.zip --packages org.apache.spark:spark-sql-kafka-0-10\_2.11:2.4.5 --files uszipsv.csv driver.py**

### Console Output:

```
23/01/08 01:57:57 INFO JettyUtils: Adding filter org.apache.hadoop.yarn.server.webproxy.amfilter.AmipFilter to /static/sql/
23/01/08 01:57:58 INFO StateStoreCoordinatorRef: Registered StateStoreCoordinator endpoint

Batch: 0

+-----+-----+-----+-----+-----+-----+-----+
|card_id|member_id|amount|pos_id|postcode|transaction_dt|status|
+-----+-----+-----+-----+-----+-----+-----+
|348702330256514|37495066290|4380912|248063406800722|96774|2018-03-01 08:24:29|GENUINE|
|348702330256514|37495066290|6703385|78656277140812|84758|2018-06-02 04:15:03|GENUINE|
|348702330256514|37495066290|7454328|466952571393508|93645|2018-02-12 09:56:42|GENUINE|
|348702330256514|37495066290|4013428|45845320330319|15868|2018-06-13 05:38:54|GENUINE|
|348702330256514|37495066290|5495353|545499621965697|79033|2018-06-16 21:51:54|GENUINE|
|348702330256514|37495066290|3966214|369266342272501|22832|2018-10-21 03:52:51|GENUINE|
|348702330256514|37495066290|1753644|9475029292671|17923|2018-08-23 00:11:30|GENUINE|
|348702330256514|37495066290|1692115|27647525195860|55708|2018-11-23 17:02:39|GENUINE|
|5189563368503974|117826301530|9222134|525701337355194|64002|2018-03-01 20:22:10|GENUINE|
|5189563368503974|117826301530|4133848|182031383443115|26346|2018-09-09 01:52:32|GENUINE|
|5189563368503974|117826301530|8938921|799748246411019|76934|2018-12-09 05:20:53|GENUINE|
|5189563368503974|117826301530|1786366|131276818071265|63431|2018-08-12 14:29:38|GENUINE|
|5189563368503974|117826301530|9142237|564240259678903|50635|2018-06-16 19:37:19|GENUINE|
|540707334486464|1147922084344|6885448|887913906711117|59031|2018-05-05 07:53:53|GENUINE|
|540707334486464|1147922084344|4028209|11626605118182|80118|2018-08-11 01:06:50|GENUINE|
|540707334486464|1147922084344|3858369|896105817613325|53820|2018-07-12 17:37:26|GENUINE|
|540707334486464|1147922084344|9307733|729374116016479|14898|2018-07-13 04:50:16|GENUINE|
|540707334486464|1147922084344|4011296|543373367319647|44028|2018-10-17 13:09:34|GENUINE|
|540707334486464|1147922084344|9492531|211980095659371|49453|2018-04-21 14:12:26|GENUINE|
|540707334486464|1147922084344|7550074|345533088112099|15030|2018-09-29 02:34:52|GENUINE|
+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows
```

Job runs continuously as micro batches with 1 Minute time interval until it is terminated/interrupted.

## VALIDATION

- Check the card transactions count before and after streaming the data.

**HBase: - count 'card\_transactions'**

**BEFORE STREAMING – 53292 records**

```
Current count: 13000, row: 3781180267915971966020431295704|6203511|16127|922095772385354|2017-08-28 03:42:28|GENUINE
Current count: 14000, row: 4009218272111551|423517919211603|3454802|17968|949799330094128|2016-04-12 02:27:30|GENUINE
Current count: 15000, row: 4067184730430984|62171609446245|27995581|34142|989725801185837|2017-03-11 00:11:10|GENUINE
Current count: 16000, row: 4132381122041426|658327670442931|4640125|224301|222694174540297|2016-06-06 18:28:07|GENUINE
Current count: 17000, row: 4264553579186587|470021900536700|9146367|97492|119738629872810|2017-07-26 19:44:00|GENUINE
Current count: 18000, row: 4373464339970856|353188612655228|70707191|53820|896105817613325|2018-01-05 17:27:03|GENUINE
Current count: 19000, row: 4439501833420177|64000396834270|7241995|58212|430627506969079|2017-06-06 18:54:15|GENUINE
Current count: 20000, row: 4500499651579063|9821713214770331|5918030|49254|240254570257538|2017-11-19 18:21:24|GENUINE
Current count: 21000, row: 4566853351752365|208309336476264|5007068|42712|407204159700285|2016-03-31 06:03:05|GENUINE
Current count: 22000, row: 4600996294769125|641610955526739|2700042|62028|325976740802859|2017-10-07 23:19:53|GENUINE
Current count: 23000, row: 4689314809377828|260262056535812|9627742|24554|617477207911488|2017-08-21 03:10:13|GENUINE
Current count: 24000, row: 4766789897935106|689962742364435|2625284|48617|999365901769840|2017-08-20 06:06:12|GENUINE
Current count: 25000, row: 4863127030291206|788334823140096|4041377|23177|758727423991059|2016-11-24 20:48:46|GENUINE
Current count: 26000, row: 4907253800863053|723132659200637|8283600|49874|523030539477115|2017-03-15 13:01:59|GENUINE
Current count: 27000, row: 4995478705000641|733662898760334|353412|24557|696903855738500|2017-04-10 21:02:55|GENUINE
Current count: 28000, row: 5134479292018417|209432998940681|1693879|46713|994865202495162|2017-11-01 09:42:05|GENUINE
Current count: 29000, row: 516280888745682|990571384922260|9357712|82932|272510889999130|2016-08-16 13:53:15|GENUINE
Current count: 30000, row: 5186615811954262|85672266618984|1162144|24554|617477207911488|2018-01-19 22:02:13|GENUINE
Current count: 31000, row: 5221229593682054|762245749302598|2225338|46076|710963435453749|2017-12-19 09:22:27|GENUINE
Current count: 32000, row: 5258095619226135|479916651076477|4662643|78931|3188678414995|2017-12-20 06:26:28|GENUINE
Current count: 33000, row: 5294552751808411|809518123791925|6415146|22847|145809819892139|2017-04-20 09:16:35|GENUINE
Current count: 34000, row: 5342400871435050|61932267886721|641395125660|989260468134936|2016-03-30 16:00:23|GENUINE
Current count: 35000, row: 5380978184175608|291607862803202|6916930|78933|432050913975202|2017-02-17 19:03:46|GENUINE
Current count: 36000, row: 5414439899219272|948599790329037|2975624|17353|994050899536534|2017-09-28 12:51:34|GENUINE
Current count: 37000, row: 5481808794715436|689827807258904|5088559|79331|362455149134266|2016-12-20 09:31:00|GENUINE
Current count: 38000, row: 5534323829711423|527608877473344|733037|62313|808277142703095|2018-01-14 20:54:02|GENUINE
Current count: 39000, row: 5584977018792504|765899764008585|9068245|46510|189152062429368|2016-09-13 01:59:34|GENUINE
Current count: 40000, row: 6011139413319542|582288628480057|9799683|65604|470565180434292|2017-03-19 18:51:29|GENUINE
Current count: 41000, row: 6011525010455848|538555213501198|9400038|28719|146198474945328|2017-11-17 22:28:00|GENUINE
Current count: 42000, row: 6011782857327719|969966222627715|19828561|17061|841171694848680|2017-09-26 11:42:09|GENUINE
Current count: 43000, row: 6221796595498904|617313952879129|1551821|35674|834075578964661|2017-01-10 18:09:06|GENUINE
Current count: 44000, row: 6224271253849917|815187737253702|140177|83849|713116509563777|2018-01-25 18:29:05|GENUINE
Current count: 45000, row: 6225606551069826|284643419452603|7500507|22901|606202448444893|2017-12-05 03:01:57|GENUINE
Current count: 46000, row: 6228733641419063|610330639594612|3051758|39169|309527290827592|2016-07-14 03:34:45|GENUINE
Current count: 47000, row: 6447877814927926|907972949998745|923252|56685|982747668210744|2018-01-11 00:00:00|GENUINE
Current count: 48000, row: 6461366425954109|739325968607596|8515416|10800|153182020204393|2017-12-25 04:03:59|GENUINE
Current count: 49000, row: 6480152634975473|439083998526821|9304601|71047|181711798421306|2018-01-08 19:11:10|GENUINE
Current count: 50000, row: 6505080237250161|615754567307150|8801408|27341|486982167852820|2017-12-05 17:04:37|GENUINE
Current count: 51000, row: 6544876671165176|269098610760255|4604408|14510|536497882467098|2018-01-21 07:23:36|GENUINE
Current count: 52000, row: 6574255180086418|891702243060747|5991679|33097|891506971848956|2016-04-22 01:13:11|GENUINE
Current count: 53000, row: 6595814135933981|236664426408837|3243159|56162|834307885260185|2016-10-08 23:28:28|GENUINE
53292 row(s) in 2.8260 seconds
=> 53292
hbase(main):003:0>
```

## AFTER STREAMING – 59367 records

```

Current count: 19000, row: 4264553579186587|470021900536700|4272858|47578|552792990977827|2017-02-13 04:27:36|GENUINE
Current count: 20000, row: 4367125470626536|179997837474045|5811475|32535|945113329039084|2017-09-28 11:06:14|GENUINE
Current count: 21000, row: 4411296946890133|229073492285200|5784600|46986|931311742932204|2018-01-18 16:11:30|GENUINE
Current count: 22000, row: 4492848090158931|613838524128130|3001807|71043|366296280375812|2017-05-04 01:34:23|GENUINE
Current count: 23000, row: 4539744177029050|642560648641430|1999812|16921|198567387970489|2016-10-14 16:24:41|GENUINE
Current count: 24000, row: 4595507001684561|447917467655871|1078211|69136|68992652445515|2016-06-15 16:39:06|GENUINE
Current count: 25000, row: 462537189972166|306036466176496|4422540|61064|639432351795127|2017-10-07 12:47:53|GENUINE
Current count: 26000, row: 4719418574936017|500004043659347|1912413|54859|765832549139249|2016-10-31 19:13:16|GENUINE
Current count: 27000, row: 4782879464621468|257134899293254|3321180|30738|599036413189891|2017-06-22 22:53:28|GENUINE
Current count: 28000, row: 4865556006187980|454102358826314|3037500|88337|40012443036810|2017-10-04 00:31:46|GENUINE
Current count: 29000, row: 4907253800863053|723132659200637|3515852|83436|482389848742790|2017-04-11 00:33:35|GENUINE
Current count: 30000, row: 4981548424137027|996411635289270|5381248|98377|332218533399343|2017-01-12 08:46:29|GENUINE
Current count: 31000, row: 5127318999406559|391603008295007|1100805|27503|45264155049093|2017-04-12 22:12:19|GENUINE
Current count: 32000, row: 5155360614315751|403808871228232|5720247|38847|55900530339132|2016-04-21 00:38:29|GENUINE
Current count: 33000, row: 5175735819607028|227686059795799|7781043|52054|215945504043719|2018-02-11 00:00:00|GENUINE
Current count: 34000, row: 5204123273729700|133086195760104|9026991|13319|215972630915917|2017-11-01 07:56:45|GENUINE
Current count: 35000, row: 5221456036333304|77956963353858|3673522|73569|681145934981754|2016-11-04 23:17:24|GENUINE
Current count: 36000, row: 5263774491002052|839657476191015|7642957|985389|44524492|61050|2018-05-11 18:39:03|GENUINE
Current count: 37000, row: 5300868089771757|374629449245268|9885475|117986020341542|39421|2018-08-23 09:54:37|GENUINE
Current count: 38000, row: 5343422029247948|910294764593739|9686191|45850|753629127321217|2016-03-18 03:47:58|GENUINE
Current count: 39000, row: 5380978184175608|291607862803202|5072606|96744|996436980380873|2017-08-13 07:19:17|GENUINE
Current count: 40000, row: 5411856842735919|558289402758917|521850|65270|650284784307856|2016-02-16 17:28:13|GENUINE
Current count: 41000, row: 5462457589398710|198169089087246|8455206|23085|783487405458889|2017-12-21 14:02:30|GENUINE
Current count: 42000, row: 5517057937158439|825723134759441|1828956|50261|313053798939914|2017-06-30 08:12:24|GENUINE
Current count: 43000, row: 5572427538311236|976740397894598|2801376|62297|360207232746344|2017-04-15 04:37:44|GENUINE
Current count: 44000, row: 5595273277507573|262824903276197|9786992|16240|623747780727948|2017-12-11 00:00:00|GENUINE
Current count: 45000, row: 6011289432864968|726015889725417|6860143|46301|602179229683466|2016-09-18 04:53:05|GENUINE
Current count: 46000, row: 6011589374869018|430132783620523|8045248|27814|38466918460742|2017-07-08 00:32:41|GENUINE
Current count: 47000, row: 6011799606614876|557509021984037|4722088|24317|82504960837654|2018-01-30 21:22:36|GENUINE
Current count: 48000, row: 6221801089095872|457481752664480|9172465|512639120983433|42222|2018-04-08 01:44:58|GENUINE
Current count: 49000, row: 6223938739934916|821936471422885|6186915|42056|847187293322830|2016-05-18 02:12:46|GENUINE
Current count: 50000, row: 6225472681819758|718958735965781|2297129|48005|811847495441590|2017-02-22 14:55:55|GENUINE
Current count: 51000, row: 6227925299479230|565790839432279|9384231|513431311148007|22511|2018-06-18 14:14:06|GENUINE
Current count: 52000, row: 6443783023069597|523325937872991|2070595|13154|606469468197708|2017-07-16 17:27:38|GENUINE
Current count: 53000, row: 6457507104644777|683307404961691|434978|46392|234052426105135|2017-09-11 00:00:00|GENUINE
Current count: 54000, row: 6474282068752847|182802011203906|1511781|35646|191149080747748|2017-06-21 06:30:45|GENUINE
Current count: 55000, row: 6490380638769885|471193501554421|12117|39851|150112068926146|2017-10-28 03:52:37|GENUINE
Current count: 56000, row: 6515567258324915|203259382349255|3888023|74450|409035176270258|2017-02-12 02:07:33|GENUINE
Current count: 57000, row: 654595232279984|600768538287001|4283702|15956|325720495071260|2017-12-06 00:05:58|GENUINE
Current count: 58000, row: 6574456060675651|553127039477442|3607958|70368|94661058625308|2016-07-26 14:01:52|GENUINE
Current count: 59000, row: 6595638658736751|377918213614639|7015530|24464|45248890479808|2017-11-19 07:49:23|GENUINE
59367 row(s) in 1.6260 seconds
=> 59367
hbase(main):005:0>

```

**Hive: - select status, count(\*) from card\_transactions group by status;**

## BEFORE STREAMING

Status	_c1
FRAUD	82
GENUINE	53210

```

hive> select status, count(*) from card_transactions group by status;
Query ID = hadoop_20230108012501_db1b7cd1-2a03-468e-820a-1d2ba880cd36
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1673139411187_0014)

-----
VERTICES      MODE           STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED  1      1            0        0        0      0      0
Reducer 2 ..... container  SUCCEEDED  2      2            0        0        0      0      0
-----
VERTICES: 02/02 [=====] 100% ELAPSED TIME: 7.27 s
-----
OK
status _c1
FRAUD 82
GENUINE 53210
Time taken: 12.494 seconds, Fetched: 2 row(s)
hive>

```

## AFTER STREAMING

status	_c1
FRAUD	212
GENUINE	59155

```
hive> select status, count(*) from card_transactions group by status;
Query ID = hadoop_20230108020621_9e72e43d-2d80-49a8-ab62-0304aab4204e
Total jobs = 1
Launching Job 1 out of 1
Tez session was closed. Reopening...
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1673139411187_0016)
```

	VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	.....	container	SUCCEEDED	1	1	0	0	0	0
Reducer 2	.....	container	SUCCEEDED	2	2	0	0	0	0

```
VERTICES: 02/02 [=====>>>] 100% ELAPSED TIME: 7.30 s
OK
status _c1
FRAUD 212
GENUINE 59155
Time taken: 19.15 seconds, Fetched: 2 row(s)
hive>
```

2. Check the latest transaction date available in card\_transactions table.

**Hive: - select max(transaction\_dt) from card\_transactions;  
BEFORE STREAMING - 2018-02-11 00:00:09**

```
hive> select max(transaction_dt) from card_transactions;
Query ID = hadoop_20230108012636_d113f307-2408-4f8f-be2c-fb70bc3170a4
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1673139411187_0014)
```

	VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	.....	container	SUCCEEDED	1	1	0	0	0	0
Reducer 2	.....	container	SUCCEEDED	1	1	0	0	0	0

```
VERTICES: 02/02 [=====>>>] 100% ELAPSED TIME: 7.17 s
OK
_c0
2018-02-11 00:00:09
Time taken: 9.303 seconds, Fetched: 1 row(s)
hive>
```

**AFTER STREAMING - 2018-12-10 09:00:34**

```
hive> select max(transaction_dt) from card_transactions;
Query ID = hadoop_20230108020942_ed66881c-69e6-482c-9197-3eba6068c19b
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1673139411187_0016)
```

	VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	.....	container	SUCCEEDED	1	1	0	0	0	0
Reducer 2	.....	container	SUCCEEDED	1	1	0	0	0	0

```
VERTICES: 02/02 [=====>>>] 100% ELAPSED TIME: 6.36 s
OK
_c0
2018-12-10 09:00:34
Time taken: 8.877 seconds, Fetched: 1 row(s)
hive>
```