

## Final-Submission - Scripts Execution

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

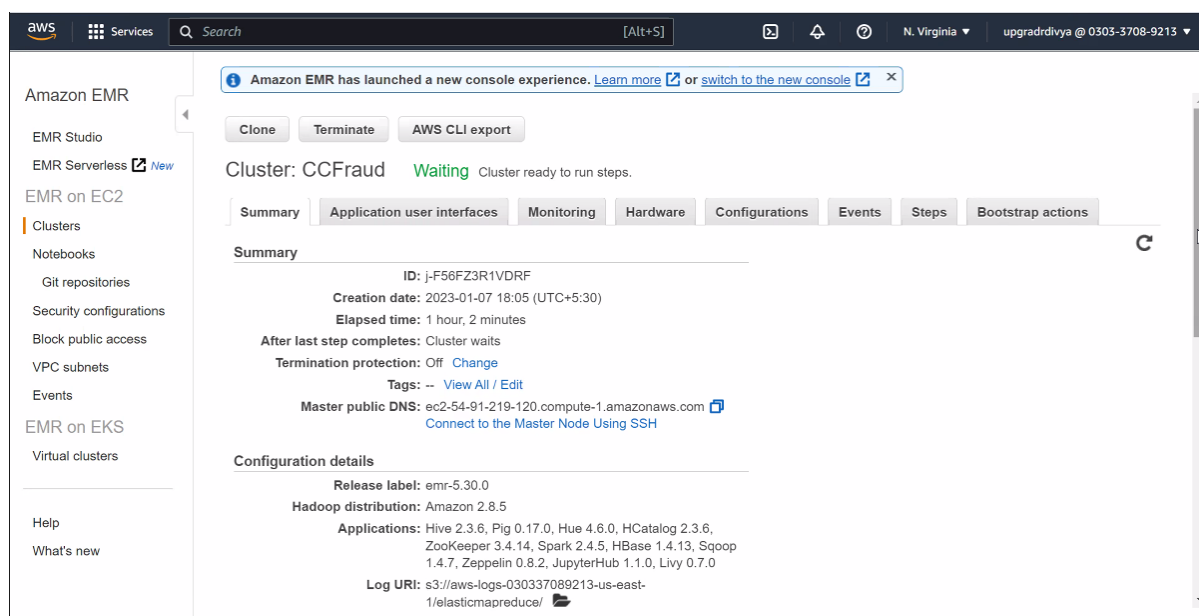
- For the below mentioned tasks, an EMR cluster with **Hadoop, Sqoop, Hive, HBase and Spark**, Root device EBS volume size **as 20 GB** was created.
- Tasks summary –

• **Task 5:** A streaming data processing framework that ingests real-time POS transaction data from Kafka was created. Further, the transaction data was validated based on the three rules' parameters (stored in the NoSQL database) discussed previously.

• **Task 6:** Updating the transactions data along with the status (fraud/genuine) in the card\_transactions table.

• **Task 7:** Storing the 'postcode' and 'transaction\_dt' of the current transaction in the look-up table of the NoSQL database if the transaction was classified as genuine.

### EMR Cluster Configuration:



The screenshot displays the AWS Management Console for an Amazon EMR cluster. The cluster is named 'CCFraud' and is currently in a 'Waiting' state, indicating it is ready to run steps. The console shows various tabs for cluster management, including Summary, Application user interfaces, Monitoring, Hardware, Configurations, Events, Steps, and Bootstrap actions. The Summary tab is selected, showing the following details:

- ID:** j-F56FZ3R1VDRF
- Creation date:** 2023-01-07 18:05 (UTC+5:30)
- Elapsed time:** 1 hour, 2 minutes
- After last step completes:** Cluster waits
- Termination protection:** Off (with a 'Change' link)
- Tags:** -- (with 'View All / Edit' link)
- Master public DNS:** ec2-54-91-219-120.compute-1.amazonaws.com (with a 'Connect to the Master Node Using SSH' link)

The Configuration details section shows:

- Release label:** emr-5.30.0
- Hadoop distribution:** Amazon 2.8.5
- Applications:** Hive 2.3.6, Pig 0.17.0, Hue 4.6.0, HCatalog 2.3.6, ZooKeeper 3.4.14, Spark 2.4.5, HBase 1.4.13, Sqoop 1.4.7, Zeppelin 0.8.2, JupyterHub 1.1.0, Livy 0.7.0
- Log URI:** s3://aws-logs-030337089213-us-east-1/elasticmapreduce/

- ### 1. Logging into EMR instance as “hadoop”:

```

hadoop@ip-172-31-28-134:~
 _ | _ | _ )
 _ | ( _ | /   Amazon Linux 2 AMI
 _ | \ _ | _ |

https://aws.amazon.com/amazon-linux-2/
104 package(s) needed for security, out of 170 available
Run "sudo yum update" to apply all updates.

EEEEEEEEEEEEEEEEEEEE MMMMMMM                      MMMMMMM RRRRRRRRRRRRRRRR
E::::::::::::::::::::E M::::::::M                      M::::::::M R::::::::::::R
EE::::::::EEEEEEEE::E M::::::::M                      M::::::::M R::::RRRRRR::::R
 E::::E          EEEE M::::::::M                      M::::::::M RR::::R          R::::R
 E::::E          M::::M M::::M M::::M M::::M          R:::R          R::::R
 E::::EEEEEEEEEE M::::M M::::M M::::M M::::M          R:::RRRRRR::::R
 E::::::::::::E M::::M M::::M M::::M M::::M          R::::::::RR
 E::::EEEEEEEEEE M::::M M::::M M::::M M::::M          R:::RRRRRR::::R
 E::::E          M::::M M::::M M::::M M::::M          R:::R          R::::R
 E::::E          EEEE M::::M          MMM M::::M          R:::R          R::::R
EE::::::::EEEEEEEE::E M::::M                      M::::M          R:::R          R::::R
E::::::::::::E M::::M                      M::::M RR::::R          R::::R
EEEEEEEEEEEEEEEEEEEE MMMMMMM                      MMMMMMM RRRRRRR          RRRRRR

[hadoop@ip-172-31-28-134 ~]$ █

```

2. Switching to root user and running pip install kafka-python and further running "sudo -i -u hadoop" to be a hadoop user

[illegible]

- Running the following commands in order to Install Happy base and start thrift server

- **sudo yum update**
- **sudo yum install python3-devel**
- **pip install happybase**
- **/usr/lib/hbase/bin/hbase-daemon.sh start thrift -p 9090**

```
hadoop@ip-172-31-28-134:~$
rubygem-psych.x86_64 0:2.0.0-36.amzn2.0.3
rubygem-rdoc.noarch 0:4.0.0-36.amzn2.0.3
rubygems.noarch 0:2.0.14.1-36.amzn2.0.3
selinux-policy.noarch 0:3.13.1-192.amzn2.6.8
selinux-policy-targeted.noarch 0:3.13.1-192.amzn2.6.8
strace.x86_64 0:4.26-1.amzn2.0.1
sysctl-defaults.noarch 0:1.0-3.amzn2
system-lsb.x86_64 0:4.1-27.amzn2.3.6
system-lsb-core.x86_64 0:4.1-27.amzn2.3.6
system-lsb-cxx.x86_64 0:4.1-27.amzn2.3.6
system-lsb-desktop.x86_64 0:4.1-27.amzn2.3.6
system-lsb-languages.x86_64 0:4.1-27.amzn2.3.6
system-lsb-submod-multimedia.x86_64 0:4.1-27.amzn2.3.6
system-lsb-submod-security.x86_64 0:4.1-27.amzn2.3.6
system-release.x86_64 1:2-14.amzn2
systemtap-runtime.x86_64 0:4.5-1.amzn2.0.1
systemtap-sdt-devel.x86_64 0:4.5-1.amzn2.0.1
texlive-base.noarch 2:2012-38.20130427_r30134.amzn2.0.5
texlive-dvipng.noarch 2:svn26689.1.14-38.amzn2.0.5
texlive-kpathsea.noarch 2:svn28792.0-38.amzn2.0.5
tkinter.x86_64 0:2.7.18-1.amzn2.0.5
tzdata.noarch 0:2022f-1.amzn2.0.1
update-motd.noarch 0:1.1.2-2.amzn2.0.2
xfsprogs.x86_64 0:5.0.0-10.amzn2.0.1
yum.noarch 0:3.4.3-158.amzn2.0.6

Replaced:
grub2.x86_64 1:2.02-35.amzn2.0.4
python-colorama.noarch 0:0.3.2-3.amzn2
system-lsb-printing.x86_64 0:4.1-27.amzn2.3.5
grub2-tools.x86_64 1:2.02-35.amzn2.0.4
python-six.noarch 0:1.9.0-2.amzn2

Complete!
[hadoop@ip-172-31-28-134 ~]$
```

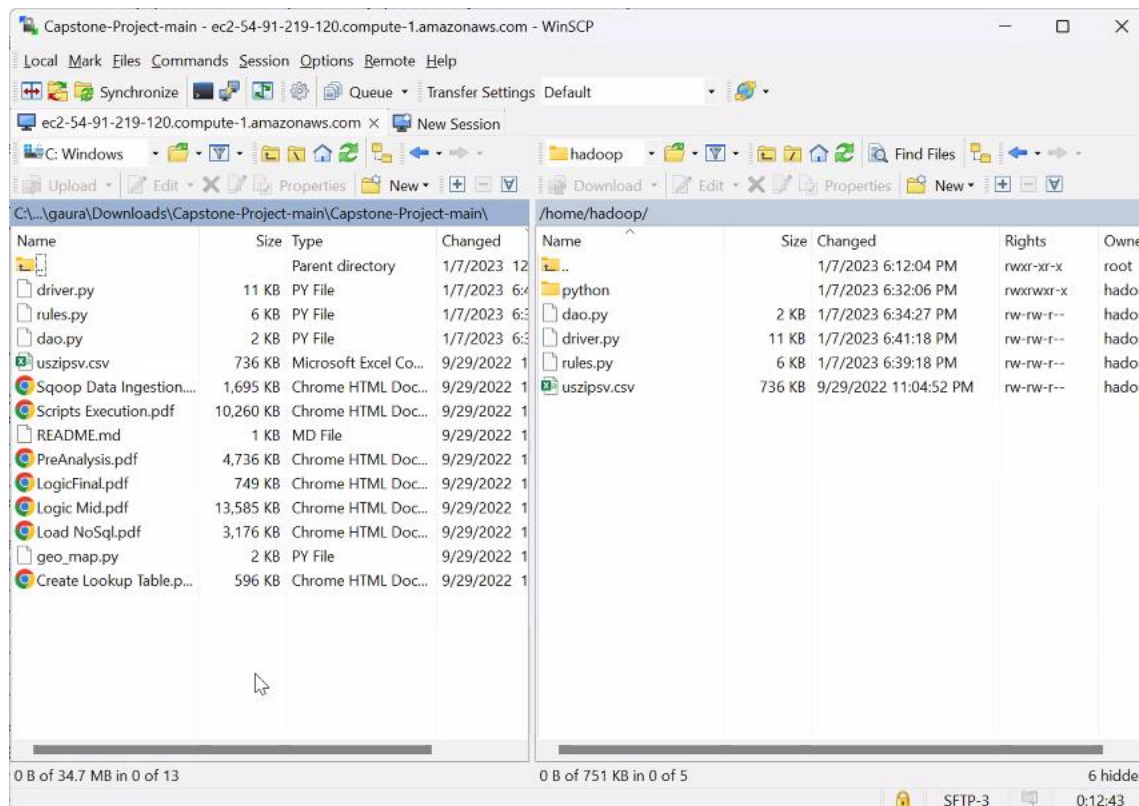
```
hadoop@ip-172-31-28-134:~$
Total download size: 288 k
Installed size: 753 k
Is this ok [y/d/N]: y
Downloading packages:
(1/4): python-srpm-macros-3-60.amzn2.0.1.noarch.rpm | 18 kB 00:00:00
(2/4): python-rpm-macros-3-60.amzn2.0.1.noarch.rpm | 14 kB 00:00:00
(3/4): python3-rpm-macros-3-60.amzn2.0.1.noarch.rpm | 12 kB 00:00:00
(4/4): python3-devel-3.7.15-1.amzn2.0.2.x86_64.rpm | 244 kB 00:00:00
-----
Total | 2.6 MB/s | 288 kB 00:00:00
Running transaction check
Running transaction test
Transaction test succeeded
Running transaction
Installing : python-srpm-macros-3-60.amzn2.0.1.noarch 1/4
Installing : python-rpm-macros-3-60.amzn2.0.1.noarch 2/4
Installing : python3-rpm-macros-3-60.amzn2.0.1.noarch 3/4
Installing : python3-devel-3.7.15-1.amzn2.0.2.x86_64 4/4
Verifying : python-rpm-macros-3-60.amzn2.0.1.noarch 1/4
Verifying : python3-rpm-macros-3-60.amzn2.0.1.noarch 2/4
Verifying : python3-devel-3.7.15-1.amzn2.0.2.x86_64 3/4
Verifying : python-srpm-macros-3-60.amzn2.0.1.noarch 4/4

Installed:
python3-devel.x86_64 0:3.7.15-1.amzn2.0.2

Dependency Installed:
python-rpm-macros.noarch 0:3-60.amzn2.0.1 python-srpm-macros.noarch 0:3-60.amzn2.0.1
python3-rpm-macros.noarch 0:3-60.amzn2.0.1

Complete!
[hadoop@ip-172-31-28-134 ~]$
```

- Downloaded files **db-> dao.py , geomap.py ,rules-> rules.py ,driver.py ,unzipsv.csv** from the resource section of the capstone project and transferred it to hadoop instance via WinSCP.



```
[hadoop@ip-172-31-28-134 ~]$ ls
dao.py driver.py python rules.py uszipsv.csv
[hadoop@ip-172-31-28-134 ~]$
```

- Updated the Public IP of the EC2 Instance "10.100.84.25" (self.host) in **dao.py** file.

```
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 = '10.100.84.25'
            for i in range(2):
                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:
```

- Updated **rules.py** with following parameters:

```
lookup_table = 'CC_LOOKUP_DATA_HBASE'
master_table = 'CC_transactions_hbase'
```

```
# List all the functions to check for the rules
from db.dao import HBaseDao
from db.geo_map import GEO_Map
from datetime import datetime
import uuid

# Create UDF functions
lookup_table = 'CC_LOOKUP_DATA_HBASE'
master_table = 'CC_transactions_hbase'
speed_threshold = 0.25 # km/sec - Average speed of flight 900 km/hr
```

- Created Python functions, containing the logic for the UDFs (**rules.py**)

**verify\_ucl\_data** : Function to verify the UCL rule that the Transaction amount should be less than Upper control limit (UCL)

```
def verify_ucl_data(card_id, amount):
    try:
        hbasedao = HBaseDao.get_instance()

        card_row = hbasedao.get_data(key=str(card_id), table=lookup_table)
        card_ucl = (card_row[b'card_data:ucl']).decode("utf-8")

        if amount < float(card_ucl):
            return True
        else:
            return False
    except Exception as e:
        raise Exception(e)
```

**verify\_credit\_score\_data**: Function to verify the credit score rule. Credit score for each member should be greater than 200

```
def verify_credit_score_data(card_id):
    try:
        hbasedao = HBaseDao.get_instance()

        card_row = hbasedao.get_data(key=str(card_id), table=lookup_table)
        card_score = (card_row[b'card_data:score']).decode("utf-8")

        if int(card_score) > 200:
            return True
        else:
            return False
    except Exception as e:
        raise Exception(e)
```



**verify\_postcode\_data:** Function to verify the following zipcode rules. ZIP code distance

```
def verify_postcode_data(card_id, postcode, transaction_dt):
    try:
        hbasedao = HBaseDao.get_instance()
        geo_map = GEO_Map.get_instance()

        card_row = hbasedao.get_data(key=str(card_id), table=lookup_table)
        last_postcode = (card_row[b'card_data:postcode']).decode("utf-8")
        last_transaction_dt = (card_row[b'card_data:transaction_dt']).decode("utf-8")

        current_lat = geo_map.get_lat(str(postcode))
        current_lon = geo_map.get_long(str(postcode))
        previous_lat = geo_map.get_lat(last_postcode)
        previous_lon = geo_map.get_long(last_postcode)

        dist = geo_map.distance(lat1=current_lat, long1=current_lon, lat2=previous_lat, long2=previous_lon)

        speed = calculate_speed(dist, transaction_dt, last_transaction_dt)

        if speed < speed_threshold:
            return True
        else:
            return False

    except Exception as e:
        raise Exception(e)
```

**calculate\_speed :** Function to calculate the speed from distance and transaction timestamp differentials.

```
except Exception as e:
    raise Exception(e)

"""
A function to calculate the speed from distance and transaction timestamp differentials
:param dist: (Float) Distance between postcodes
:param transaction_dt1: Transaction timestamp from the table
:param transaction_dt2: Transaction timestamp from the POS
:return: (Float) Speed
"""
def calculate_speed(dist, transaction_dt1, transaction_dt2):
    transaction_dt1 = datetime.strptime(transaction_dt1, '%d-%m-%Y %H:%M:%S')
    transaction_dt2 = datetime.strptime(transaction_dt2, '%d-%m-%Y %H:%M:%S')

    elapsed_time = transaction_dt1 - transaction_dt2
    elapsed_time = elapsed_time.total_seconds()

    try:
        return dist / elapsed_time
    except ZeroDivisionError:
        return 299792.458
# (Speed of light)
"""
```

**verify\_rules\_status:** A function to verify all the three rules - ucl, credit score and speed

```
def verify_rules_status(card_id, member_id, amount, pos_id, postcode, transaction_dt):
    hbasedao = HBaseDao.get_instance()

    # Check if the POS transaction passes all rules.
    # If yes, update the lookup table and insert data in master table as genuine.
    # Else insert the transaction in master table as Fraud.

    rule1 = verify_ucl_data(card_id, amount)
    rule2 = verify_credit_score_data(card_id)
    rule3 = verify_postcode_data(card_id, postcode, transaction_dt)

    if all([rule1, rule2, rule3]):
        status = 'GENUINE'
        hbasedao.write_data(key=str(card_id),
                           row=('card_data:postcode': str(postcode), 'card_data:transaction_dt': str(transaction_dt)),
                           table=lookup_table)
    else:
        status = 'FRAUD'

    new_id = str(uuid.uuid4()).replace('-', '')
    hbasedao.write_data(key=new_id,
                       row=('cardDetail:card_id': str(card_id), 'cardDetail:member_id': str(member_id),
                           'transactionDetail:amount': str(amount), 'transactionDetail:pos_id': str(pos_id),
                           'transactionDetail:postcode': str(postcode), 'transactionDetail:status': str(status),
                           'transactionDetail:transaction_dt': str(transaction_dt)),
                       table=master_table)
    return status[hadoop@ip-172-31-28-134 ~]$
```

8. Further, updating the 'driver.py' file with the following code  
Setting up the system dependencies and importing necessary libraries and modules

```
hadoop@ip-172-31-28-134:~$
#importing necessary libraries
import sys
from pyspark.sql import SparkSession
from pyspark.sql.functions import *
from pyspark.sql.types import *
import happybase
import math
import pandas as pd
from datetime import datetime
import uuid
```

9. Initializing the Spark session and reading input data from Kafka. Details of the Kafka broker, such as bootstrap server, port and topic name

- Connect to kafka topic using

**Bootstrap-server: 18.211.252.152**

**Port Number: 9092**

**Topic: transactions-topic-verified**

```
#initialising Spark session
spark = SparkSession \
    .builder \
    .appName("CreditCardFraud") \
    .getOrCreate()
spark.sparkContext.setLogLevel('ERROR')

# Reading input from Kafka
credit_data = spark.readStream \
    .format("kafka") \
    .option("kafka.bootstrap.servers", "18.211.252.152:9092") \
    .option("startingOffsets", "earliest") \
    .option("failOnDataLoss", "false") \
    .option("subscribe", "transactions-topic-verified") \
    .load()
```

10. Defining the JSON schema of each transaction

```
# Defining schema for transaction
dataSchema = StructType() \
    .add("card_id", LongType()) \
    .add("member_id", LongType()) \
    .add("amount", DoubleType()) \
    .add("pos_id", LongType()) \
    .add("postcode", IntegerType()) \
    .add("transaction_dt", StringType())
```

11. Reading the raw JSON data from Kafka as 'credit\_data\_stream' and defining UDF's to verify rules

```
# Casting raw data as string and aliasing
credit_data = credit_data.selectExpr("cast(value as string)")
credit_data_stream = credit_data.select(from_json(col="value", schema=dataSchema).alias("credit_data")).select(
    "credit_data.*")

# Define UDF which verifies all the rules for each transaction and updates the lookup and master tables
verify_all_rules = udf(verify_rules_status, StringType())

Final_data = credit_data_stream \
    .withColumn('status', verify_all_rules(credit_data_stream['card_id'],
                                         credit_data_stream['member_id'],
                                         credit_data_stream['amount'],
                                         credit_data_stream['pos_id'],
                                         credit_data_stream['postcode'],
                                         credit_data_stream['transaction_dt']))
```

## 12. Code to display output in console

```
# Write output to console as well
output_data = Final data \
    .select("card_id", "member_id", "amount", "pos_id", "postcode", "transaction_dt") \
    .writeStream \
    .trigger(processingTime = '1 seconds') \
    .outputMode("append") \
    .format("console") \
    .option("truncate", "false") \
    .start()
```

## 13. Define spark termination

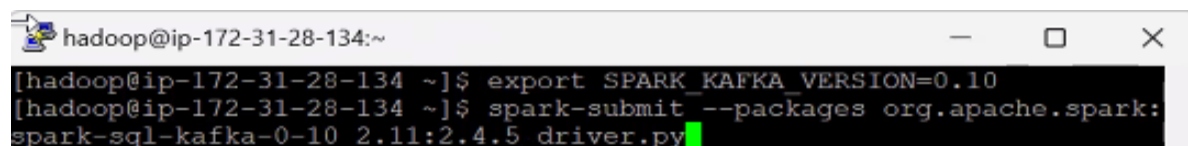
```
#indicating Spark to await termination
output_data.awaitTermination()
```

## 14. Setting the Kafka Version using the following command

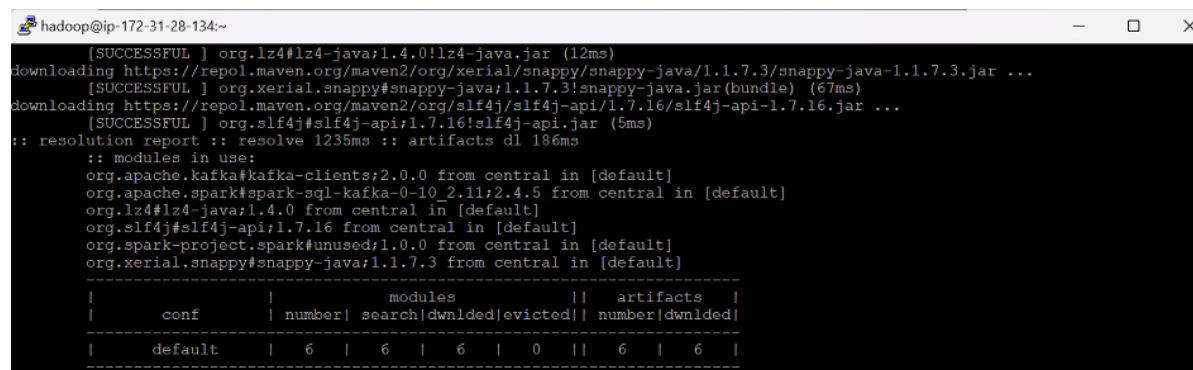
**export SPARK\_KAFKA\_VERSION=0.10**

## 15. Running the spark-submit command, specifying the Spark-SQL-Kafka package and python file

**spark-submit --packages org.apache.spark:spark-sql-kafka-0-10\_2.11:2.4.5 driver.py**



```
hadoop@ip-172-31-28-134:~
[hadoop@ip-172-31-28-134 ~]$ export SPARK_KAFKA_VERSION=0.10
[hadoop@ip-172-31-28-134 ~]$ spark-submit --packages org.apache.spark:
spark-sql-kafka-0-10_2.11:2.4.5 driver.py
```



```
hadoop@ip-172-31-28-134:~
[SUCCESSFUL ] org.lz4#lz4-java;1.4.0!lz4-java.jar (12ms)
downloading https://repo1.maven.org/maven2/org/xerial/snappy/snappy-java/1.1.7.3/snappy-java-1.1.7.3.jar ...
[SUCCESSFUL ] org.xerial.snappy#snappy-java;1.1.7.3!snappy-java.jar (bundle) (67ms)
downloading https://repo1.maven.org/maven2/org/slf4j/slf4j-api/1.7.16/slf4j-api-1.7.16.jar ...
[SUCCESSFUL ] org.slf4j#slf4j-api;1.7.16!slf4j-api.jar (5ms)
:: resolution report :: resolve 1235ms :: artifacts dl 186ms
:: modules in use:
org.apache.kafka#kafka-clients;2.0.0 from central in [default]
org.apache.spark#spark-sql-kafka-0-10_2.11:2.4.5 from central in [default]
org.lz4#lz4-java;1.4.0 from central in [default]
org.slf4j#slf4j-api;1.7.16 from central in [default]
org.spark-project.spark#unused;1.0.0 from central in [default]
org.xerial.snappy#snappy-java;1.1.7.3 from central in [default]
-----
|               |             modules              | artifacts |
|      conf     | number| search|dwnlded|evicted|| number|dwnlded|
-----
|      default  |    6  |    6  |    6  |    0  ||    6  |    6  |
-----
```



## 16. Checking the output in console:

```

-----
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|6885648|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  |
-----
nly showing top 20 rows

```

### 17. Count Data in Hbase: `count 'lookup_data_hive'`

[illegible]

Total number of records is **59367** which is matching with given requirement.