

Final-Submission – Logic Explanation

- Explanation of the solution to the streaming layer problem

***Please zoom to 180% or 200% to see screenshots with better clarity

- In order to complete below tasks, I have created EMR cluster with **Hadoop, Sqoop, Hive, HBase and Spark**, Root device EBS volume size as 20 GB
 - Task 5:** Create a streaming data processing framework that ingests real-time POS transaction data from Kafka. The transaction data is then validated based on the three rules' parameters (stored in the NoSQL database) discussed previously.
 - Task 6:** Update the transactions data along with the status (fraud/genuine) in the card_transactions table.
 - Task 7:** Store the 'postcode' and 'transaction_dt' of the current transaction in the look-up table in the NoSQL database if the transaction was classified as genuine.

EMR Cluster Configuration:

Step 1: Software and Steps
Step 2: Hardware
Step 3: General Cluster Settings
Step 4: Security

Software Configuration

Release
emr-5.30.1

☒ Hadoop 2.8.5
☐ JupyterHub 1.1.0
☐ Ganglia 3.7.2
☒ Hive 2.3.6
☐ MXNet 1.5.1
☒ Hue 4.6.0
☒ Spark 2.4.5

☐ Zeppelin 0.8.2
☐ Tez 0.9.2
☒ HBase 1.4.13
☐ Presto 0.232
☒ Sqoop 1.4.7
☐ Phoenix 4.14.3
☐ HCatalog 2.3.6

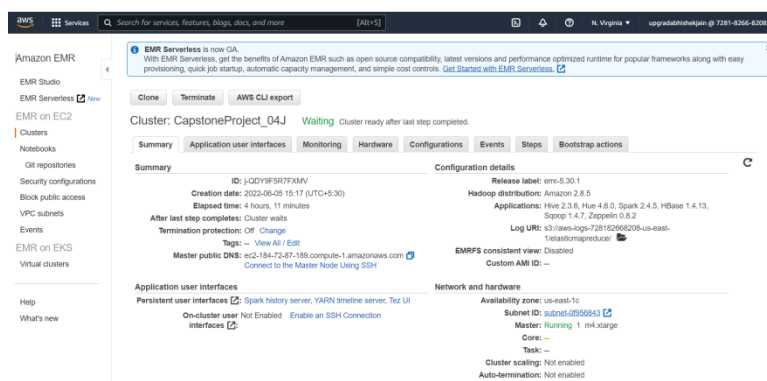
☐ Livy 0.7.0
☐ Flink 1.10.0
☐ Pig 0.17.0
☐ ZooKeeper 3.4.14
☐ Mahout 0.13.0
☐ Oozie 5.2.0
☐ TensorFlow 1.14.0

EBS Root Volume

Specify the root device volume size up to 100 GiB. This sizing applies to all instances in the cluster. [Learn more](#)

Root device EBS volume size GiB

Cancel
Previous
Next



The screenshot shows the AWS EMR console interface. The left sidebar contains navigation links for Amazon EMR, EMR Studio, EMR Serverless, Clusters, Notebooks, Git repositories, Security configurations, Block public access, VPC subnets, Events, EMR on EKS, Virtual clusters, Help, and What's new. The main content area displays the configuration details for a cluster named 'CapstoneProject_04J' in the 'Waiting' state. The configuration details include:

- Summary:** ID: j-QDY9F8R7FMV, Creation date: 2022-06-05 15:17 (UTC+5:30), Elapsed time: 4 hours, 11 minutes, After last step completes: Cluster waits, Termination protection: Off, Tags: View All / Edit, Master public DNS: ec2-184-72-87-188.compute-1.amazonaws.com, Connect to the Master Node Using SSH.
- Configuration details:** Release label: emr-5.30.1, Hadoop distribution: Amazon 2.8.5, Applications: Hive 2.3.6, Hue 4.6.0, Spark 2.4.5, HBase 1.4.13, Sqoop 1.4.7, Zeppelin 0.8.2, Log URI: s3://aws-logs-728162068208-us-east-1/elasticmapreduce/, EMRFS consistent view: Disabled, Custom AMI ID: --.
- Application user interfaces:** Persistent user interfaces: Spark History server, YARN timeline server, Tez UI, On-cluster user: Not Enabled, Enable an SSH Connection, Interfaces: --.
- Network and hardware:** Availability zone: us-east-1c, Subnet ID: subnet-d8568a53, Master: Running 1 m4.xlarge, Core: --, Task: --, Cluster scaling: Not enabled, Auto-termination: Not enabled.

2. Logged into EMR instance as “hadoop”:

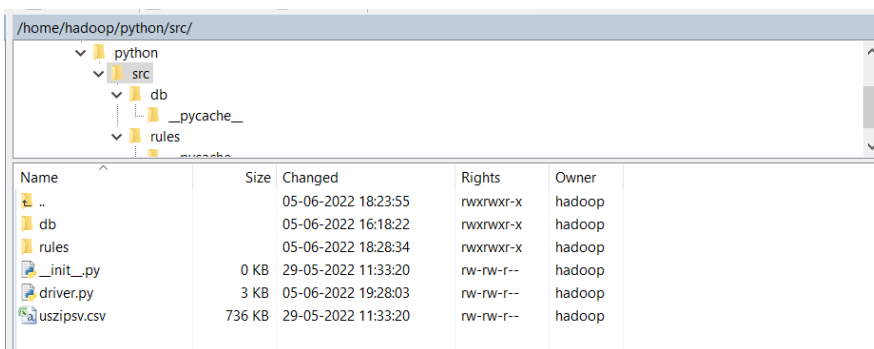
```
hadoop@ip-172-31-23-33:~  
Authenticating with public key "imported-openssh-key"  
Last login: Sun Jun  5 12:39:23 2022  
  
      _|_  _|_  )  
      _|_  ( _|_ /  Amazon Linux 2 AMI  
      __|_\__|__|_  
  
https://aws.amazon.com/amazon-linux-2/  
  
EEEEEEEEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRRRRRRRRRRRR  
E::::::::::::::::::::E M::::::::M          M::::::::M R::::::::::::R  
EE::::EEEEEEEEEE::E M::::::::M          M::::::::M R::::RRRRRR::::R  
 E::::E          EEEEE M::::::::M          M::::::::M RR::::R          R::::R  
 E::::E          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          R::RRRRRR::::R  
 E::::E          M::::M M::M M::::M          R:::R          R::::R  
 E::::E          EEEEE M::::M          MMM M::::M          R:::R          R::::R  
EE::::EEEEEEEEEE::E M::::M          M::::M          R:::R          R::::R  
E::::::::::::::::::::E M::::M          M::::M RR::::R          R::::R  
EEEEEEEEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRRR          RRRRRR  
  
[hadoop@ip-172-31-23-33 ~]$
```

3. Switch to root user and run `pip install kafka-python` and then again use `sudo -i -u hadoop` to be a hadoop user

```
[root@ip-172-31-23-33 ~]# pip install kafka-python
WARNING: Running pip install with root privileges is generally not a good idea.
Try `pip3 install --user` instead.
Requirement already satisfied: kafka-python in /usr/local/lib/python3.7/site-packages (2.0.2)
```

- Run 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`

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



```
hadoop@ip-172-31-23-33:~/python/src
[hadoop@ip-172-31-23-33 src]$ ls
db driver.py _init_.py rules unzipsv.csv
[hadoop@ip-172-31-23-33 src]$
```

- Updated the Public IP of your EC2 Instance **"184.72.87.189"**(self.host) in **dao.py** file

```
hadoop@ip-172-31-23-33:~/python/src/db
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 = '184.72.87.189'
            #self.host = 'localhost'
            for i in range(2):
```

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

lookup_table = 'lookup_data_hbase'

master_table = 'card_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 = 'lookup_data_hbase'
master_table = 'card_transactions_hbase'
```

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

verify_ucl_data : Function to verify the UCL rule 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 of 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 : A function to calculate the speed from distance and transaction timestamp differentials

```
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
```

9. Next, I updated the 'driver.py' file with the following code
Setting up the system dependencies and importing necessary libraries and modules

```
#importing necessary libraries
import os
import sys
from pyspark.sql import SparkSession
from pyspark.sql.functions import *
from pyspark.sql.types import *
from rules.rules import *
```

10. Initializing the Spark session and reading input data from Kafka mentioning the details of the Kafka broker, such as bootstrap server, port and topic name
1. 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()
```

11. Define JSON schema of each transactions

```
# 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())
```

12. Read the raw JSON data from Kafka as 'credit_data_stream' and Define 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']))
```

13. 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 \
    .outputMode("append") \
    .format("console") \
    .option("truncate", False) \
    .start()
```

14. Define spark termination

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

15. Set the Kafka Version using the following command

```
export SPARK_KAFKA_VERSION=0.10
```

16. Run 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-172-31-23-33:~/python/src
```

```
[hadoop@ip-172-31-23-33 src]$ export SPARK_KAFKA_VERSION=0.10
[hadoop@ip-172-31-23-33 src]$ spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.5 driver.py
```



```

hadoop@ip-172-31-23-33:~/python/src
[hadoop@ip-172-31-23-33 src]$ export SPARK_KAFKA_VERSION=0.10
[hadoop@ip-172-31-23-33 src]$ spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.5 driver.py
ivy Default Cache set to: /home/hadoop/.ivy2/cache
The jars for the packages stored in: /home/hadoop/.ivy2/jars
:: loading settings :: url = jar:file:/usr/lib/spark/jars/ivy-2.4.0.jar!/org/apache/ivy/core/settings/ivysettings.xml
org.apache.spark#spark-sql-kafka-0-10_2.11 added as a dependency
:: resolving dependencies :: org.apache.spark#spark-submit-parent-0cc3580a-528f-4d96-aled-f668602ac5ce;1.0
  confs: [default]
  found org.apache.spark#spark-sql-kafka-0-10_2.11:2.4.5 in central
  found org.apache.kafka#kafka-clients;2.0.0 in central
  found org.lz4#lz4-java;1.4.0 in central
  found org.xerial.snappy#snappy-java;1.1.7.3 in central
  found org.slf4j#slf4j-api;1.7.16 in central
  found org.spark-project.spark#unused;1.0.0 in central
:: resolution report :: resolve 528ms :: artifacts dl 21ms
  :: 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]
-----
  | conf | number | search | dwnlded | evicted | number | dwnlded |
-----+-----+-----+-----+-----+-----+-----+
  | default | 6 | 0 | 0 | 0 | 6 | 0 |
-----
:: retrieving :: org.apache.spark#spark-submit-parent-0cc3580a-528f-4d96-aled-f668602ac5ce
  confs: [default]
  0 artifacts copied, 6 already retrieved (0KB/12ms)
22/06/05 12:59:23 INFO SparkContext: Running Spark version 2.4.5-amzn-0
22/06/05 12:59:23 INFO SparkContext: Submitted application: CreditFraud
22/06/05 12:59:23 INFO SecurityManager: Changing view acls to: hadoop
22/06/05 12:59:23 INFO SecurityManager: Changing modify acls to: hadoop
22/06/05 12:59:23 INFO SecurityManager: Changing view acls groups to:
22/06/05 12:59:23 INFO SecurityManager: Changing modify acls groups to:
22/06/05 12:59:23 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions:
22/06/05 12:59:24 INFO Utils: Successfully started service 'sparkDriver' on port 33835.
22/06/05 12:59:24 INFO SparkEnv: Registering MapOutputTracker
22/06/05 12:59:24 INFO SparkEnv: Registering BlockManagerMaster
22/06/05 12:59:24 INFO BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology info
22/06/05 12:59:24 INFO BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
22/06/05 12:59:24 INFO DiskBlockManager: Created local directory at /mnt/tmp/blockmgr-bd558144-04a9-4075-a87f-96fa8c568453
22/06/05 12:59:24 INFO MemoryStore: MemoryStore started with capacity 1028.8 MB
22/06/05 12:59:24 INFO SparkEnv: Registering OutputCommitCoordinator
22/06/05 12:59:25 INFO Utils: Successfully started service 'SparkUI' on port 4040.
22/06/05 12:59:25 INFO SparkUI: Bound SparkUI to 0.0.0.0, and started at http://ip-172-31-23-33.ec2.internal:4040
22/06/05 12:59:25 INFO Utils: Using initial executors = 50, max of spark.dynamicAllocation.initialExecutors, spark.dynamicAllocation

```

17. Check Output in console :

```

atch: 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 |
-----+-----+-----+-----+-----+-----+
nly showing top 20 rows

```


18. Count Data in Hbase: `count 'lookup_data_hive'`

```

current count: 20000, row: 27899
current count: 21000, row: 28599
current count: 22000, row: 29799
current count: 23000, row: 30699
current count: 24000, row: 31599
current count: 25000, row: 32499
current count: 26000, row: 33399
current count: 27000, row: 341724964658347.21078177559189.12-04-2018152438.2021-01-04171328.398477
current count: 28000, row: 346618652451437.540752175694215.29-04-2018001259.2021-01-04171400.227023
current count: 29000, row: 35244
current count: 30000, row: 36164
current count: 31000, row: 370582035946789.433446648625434.08-07-2018034337.2021-01-04171349.489439
current count: 32000, row: 375804375521605.880937166605449.26-05-2018130045.2021-01-04171430.733012
current count: 33000, row: 38174
current count: 34000, row: 39074
current count: 35000, row: 39977
current count: 36000, row: 40768
current count: 37000, row: 41560
current count: 38000, row: 42357
current count: 39000, row: 4318541450494035.496412742732167.12-02-2018145807.2021-01-04171356.009418
current count: 40000, row: 43999
current count: 41000, row: 44794
current count: 42000, row: 45544
current count: 43000, row: 46304
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: 5257502990314019.205172644964018.14-07-2018070014.2021-01-04171327.667742
current count: 52000, row: 53290
current count: 53000, row: 54020
current count: 54000, row: 4211
current count: 55000, row: 6478888641720566.273286841077378.06-10-2018212851.2021-01-04171333.585477
current count: 56000, row: 6965
current count: 57000, row: 7865
current count: 58000, row: 8768
current count: 59000, row: 9668
5367 row(s) in 3.8145 seconds
> 55367

```

Total number for record is **59367** which is matching with given requirement