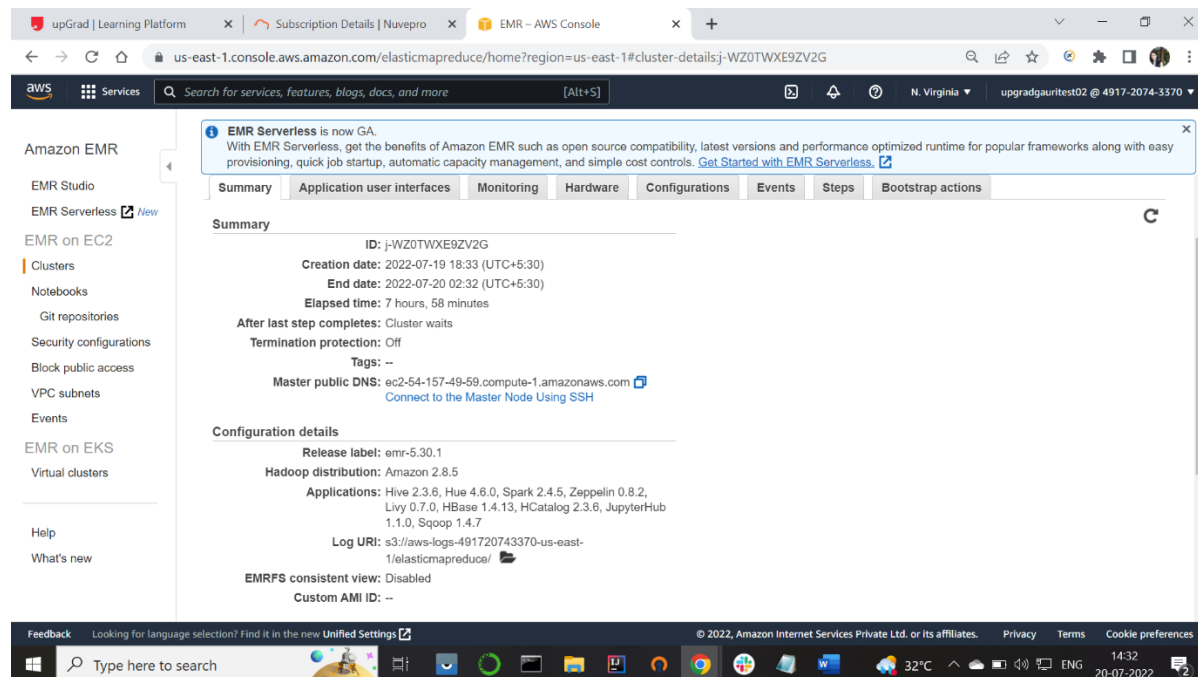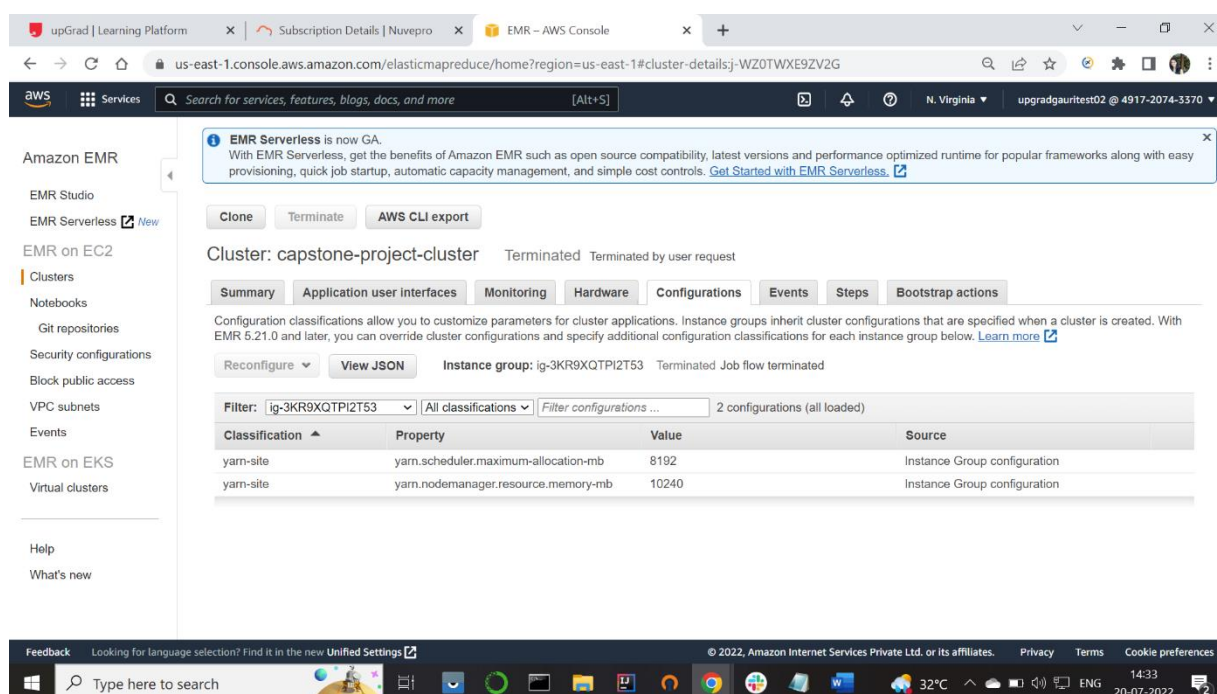# Logic For First Submission

For the capstone project, I have created EMR Cluster with applications as shown in the screenshot. Also configured YARN parameters for this cluster and login with hadoop user in PUTTY terminal

- **Task 1**: To write a job to consume clickstream data from Kafka and ingest to Hadoop.
1. Created a **spark_kafka_to_local.py** file and imported necessary libraries

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import *
from pyspark.sql.types import *
```

2. Established spark connection

```
spark =
SparkSession.builder.appName("KafkaRead").getOrCreate()
spark.sparkContext.setLogLevel('ERROR')
```

3. Read data from kafka server and topic given

```
lines = spark.readStream.format("kafka") \
        .option("kafka.bootstrap.servers","18.211.252.152:9092") \
        .option("subscribe","de-capstone3") \
        .option("failOnDataLoss","false") \
        .option("startingOffsets", "earliest") \
        .load()
```

4. Casted raw data as string

```
kafkaDF = lines.selectExpr("cast(key as string)","cast(value as string)")
```

5. Wrote Kafka data into json file

```
output = kafkaDF \
        .writeStream \
        .outputMode("append") \
        .format("json") \
        .option("truncate", "false") \
        .option("path","/user/hadoop/clickStreamData/") \
        .option("checkpointLocation", "/user/hadoop/clickstream_checkpoint/") \
        .start()

output.awaitTermination()
```

6. Logged in to the EMR instance and below command is executed to download Spark-SQL-Kafka jar file. This jar is used to run the Spark Streaming-Kafka codes
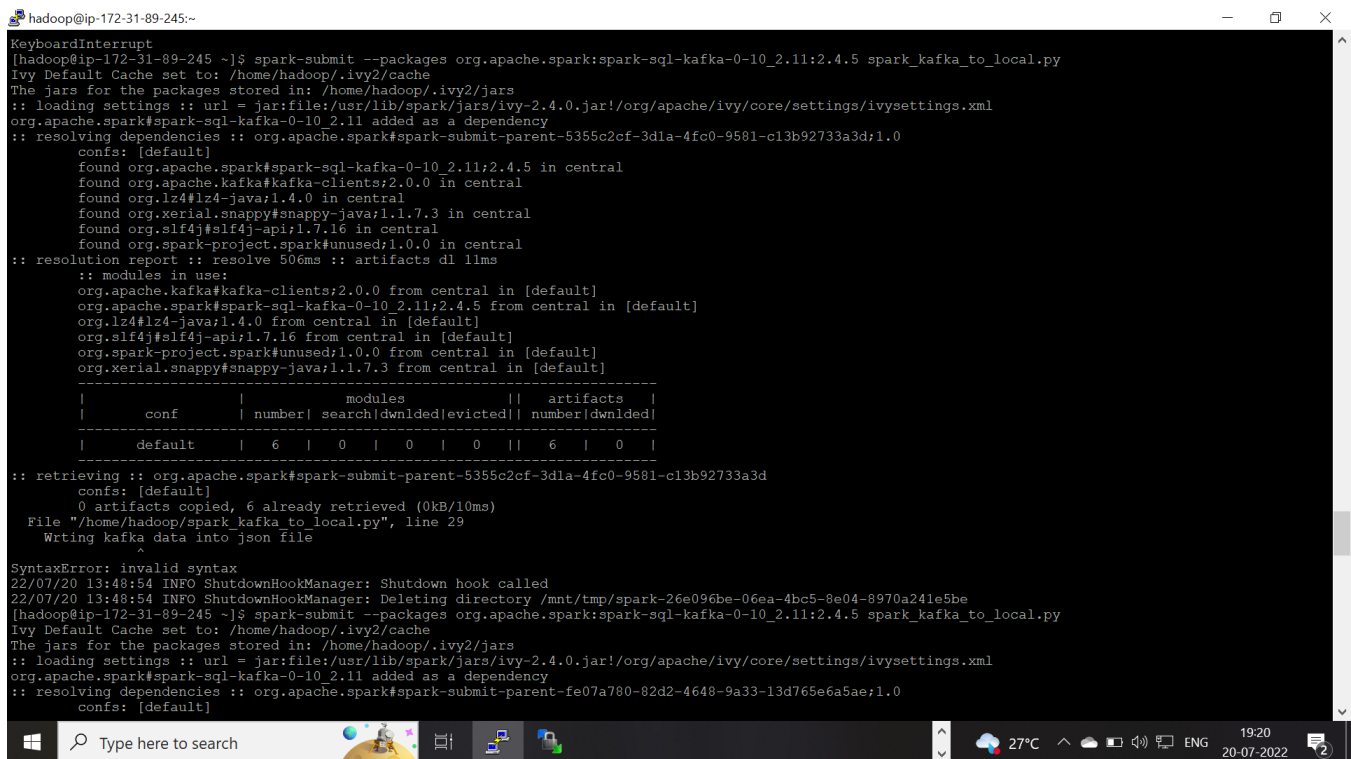
```
wget https://ds-spark-sql-kafka-jar.s3.amazonaws.com/spark-sql-
kafka-0- 10_2.11-2.3.0.jar
```

7. Kafka version is set using the following command:

```
export SPARK_KAFKA_VERSION=0.10
```

8. Submitted the spark job using command below:

```
spark-submit --packages org.apache.spark:spark-sql-kafka-0-
10_2.11:2.4.5 spark_kafka_to_local.py
```

9.  The data extracted from Kafka was in nested json format. Hence wrote a pyspark job in **spark_local_flatten.py** file to flatten the data and load into Hadoop

10. Imported necessary libraries

```
from pyspark.sql import SparkSession
from pyspark.sql import functions as F
from pyspark.sql.functions import col
from pyspark.sql.types import *
```

11. Established a spark connection

```
spark=SparkSession \
        .builder \
        .appName('transformKafkaData') \
        .master('yarn') \
        .getOrCreate()
```

12. Read extracted data stored in json format

```
df=spark.read.json('/user/hadoop/clickStreamData/')
```

13. Flattened raw data using regexp_replace function and stored the raw data into respective columns in a dataframe

```
flatten_df=df.withColumn("value",
F.split(F.regexp_replace(F.regexp_replace((F.regexp_replace("value",'\{|}',"")
),'\:',','),'\"|"',"").cast("string"),','))\
.withColumn("customer_id", F.element_at("value",2))\
withColumn("app_version", F.element_at("value",4))\
.withColumn("OS_version",F.element_at("value",6))\
.withColumn("lat",F.element_at("value",8))\
.withColumn("lon", F.element_at("value",10))\
.withColumn("page_id", F.element_at("value",12))\
.withColumn("button_id",F.element_at("value",14))\
.withColumn("is_button_click",F.element_at("value",16))\
.withColumn("is_page_view",F.element_at("value",18))\
.withColumn("is_scroll_up",F.element_at("value",20))\
.withColumn("is_scroll_down",F.element_at("value",22))\
.withColumn("date_hour",F.element_at("value",24))\
.withColumn("minutes",F.element_at("value",25))\
.withColumn("seconds",F.element_at("value",26))\
.drop("value")
```

14. Concatenated **date_hour, minutes and seconds** column to make it into timestamp format:

```
flatten_df=flatten_df.select("*",F.concat(col("date_hour"),F.lit(":")
,col("minutes"),F.lit(":"),col("seconds")).alias("timestamp"))
```

15. Removed extra character \n from timestamp column to make the data more structured

```
flatten_df =
flatten_df.select("*").withColumn("timestamp",F.expr("substring(timestamp, 1,
length(timestamp)-2)")).drop("date_hour").drop("minutes").drop("seconds")
```

16. Wrote the flattened dataframe in csv file

```
flatten_df.write.option("header","true").csv('/user/hadoop/
clickStream_flatten_data/')
```

## 17. Executed spark_local_flatten.py file using command:

```
spark-submit spark_local_flatten.py
```



```
         at java.lang.Thread.run(Thread.java:750)
[hadoop@ip-172-31-89-245 ~]$ spark-submit spark_local_flatten.py
22/07/20 13:52:56 INFO SparkContext: Running Spark version 2.4.5-amzn-0
22/07/20 13:52:56 INFO SparkContext: Submitted application: transformKafkaData
22/07/20 13:52:56 INFO SecurityManager: Changing view acls to: hadoop
22/07/20 13:52:56 INFO SecurityManager: Changing modify acls to: hadoop
22/07/20 13:52:56 INFO SecurityManager: Changing view acls groups to:
22/07/20 13:52:56 INFO SecurityManager: Changing modify acls groups to:
22/07/20 13:52:56 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users  with view permissions: Set(hadoop); groups with vi
ew permissions: Set(); users  with modify permissions: Set(hadoop); groups with modify permissions: Set()
22/07/20 13:52:56 INFO Utils: Successfully started service 'sparkDriver' on port 42075.
22/07/20 13:52:56 INFO SparkEnv: Registering MapOutputTracker
22/07/20 13:52:56 INFO SparkEnv: Registering BlockManagerMaster
22/07/20 13:52:56 INFO BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
22/07/20 13:52:56 INFO BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
22/07/20 13:52:56 INFO DiskBlockManager: Created local directory at /mnt/tmp/blockmgr-a2022508-4b7f-4efb-a144-925da52f7027
22/07/20 13:52:56 INFO MemoryStore: MemoryStore started with capacity 1028.8 MB
22/07/20 13:52:57 INFO SparkEnv: Registering OutputCommitCoordinator
22/07/20 13:52:57 INFO Utils: Successfully started service 'SparkUI' on port 4040.
22/07/20 13:52:57 INFO SparkUI: Bound SparkUI to 0.0.0.0, and started at http://ip-172-31-89-245.ec2.internal:4040
22/07/20 13:52:57 INFO Utils: Using initial executors = 50, max of spark.dynamicAllocation.initialExecutors, spark.dynamicAllocation.minExecutors and spark.e
xecutor.instances
22/07/20 13:52:58 INFO RMProxy: Connecting to ResourceManager at ip-172-31-89-245.ec2.internal/172.31.89.245:8032
22/07/20 13:52:58 INFO Client: Requesting a new application from cluster with 1 NodeManagers
22/07/20 13:52:58 INFO Client: Verifying our application has not requested more than the maximum memory capability of the cluster (8192 MB per container)
22/07/20 13:52:58 INFO Client: Will allocate AM container, with 896 MB memory including 384 MB overhead
22/07/20 13:52:58 INFO Client: Setting up container launch context for our AM
22/07/20 13:52:58 INFO Client: Setting up the launch environment for our AM container
22/07/20 13:52:58 INFO Client: Preparing resources for our AM container
22/07/20 13:52:58 WARN Client: Neither spark.yarn.jars nor spark.yarn.archive is set, falling back to uploading libraries under SPARK_HOME.
22/07/20 13:53:01 INFO Client: Uploading resource file:/mnt/tmp/spark-5ffcce5f-7f0e-4d19-b5f3-f3c4cdd7c691/__spark_libs__698612333779150397.zip -> hdfs://ip-
172-31-89-245.ec2.internal:8020/user/hadoop/.sparkStaging/application_1658324257852_0003/__spark_libs__698612333779150397.zip
22/07/20 13:53:02 INFO Client: Uploading resource file:/etc/spark/conf/hive-site.xml -> hdfs://ip-172-31-89-245.ec2.internal:8020/user/hadoop/.sparkStaging/a
pplication_1658324257852_0003/hive-site.xml
22/07/20 13:53:02 INFO Client: Uploading resource file:/usr/lib/spark/python/lib/pyspark.zip -> hdfs://ip-172-31-89-245.ec2.internal:8020/user/hadoop/.sparkS
taging/application_1658324257852_0003/pyspark.zip
22/07/20 13:53:02 INFO Client: Uploading resource file:/usr/lib/spark/python/lib/py4j-0.10.7-src.zip -> hdfs://ip-172-31-89-245.ec2.internal:8020/user/hadoop
/.sparkStaging/application_1658324257852_0003/py4j-0.10.7-src.zip
22/07/20 13:53:02 INFO Client: Uploading resource file:/mnt/tmp/spark-5ffcce5f-7f0e-4d19-b5f3-f3c4cdd7c691/__spark_conf__1486980818810018452.zip -> hdfs://ip
-172-31-89-245.ec2.internal:8020/user/hadoop/.sparkStaging/application_1658324257852_0003/__spark_conf__.zip
22/07/20 13:53:02 INFO SecurityManager: Changing view acls to: hadoop
22/07/20 13:53:02 INFO SecurityManager: Changing modify acls to: hadoop
22/07/20 13:53:02 INFO SecurityManager: Changing view acls groups to:
22/07/20 13:53:02 INFO SecurityManager: Changing modify acls groups to:
```



```
22/07/20 13:53:28 INFO ContextCleaner: Cleaned accumulator 43
22/07/20 13:53:28 INFO ContextCleaner: Cleaned accumulator 59
22/07/20 13:53:28 INFO ContextCleaner: Cleaned accumulator 58
22/07/20 13:53:28 INFO ContextCleaner: Cleaned accumulator 49
22/07/20 13:53:28 INFO ContextCleaner: Cleaned accumulator 57
22/07/20 13:53:28 INFO ContextCleaner: Cleaned accumulator 42
22/07/20 13:53:28 INFO ContextCleaner: Cleaned accumulator 41
22/07/20 13:53:28 INFO ContextCleaner: Cleaned accumulator 63
22/07/20 13:53:28 INFO ContextCleaner: Cleaned accumulator 46
22/07/20 13:53:28 INFO ContextCleaner: Cleaned accumulator 62
22/07/20 13:53:28 INFO ContextCleaner: Cleaned accumulator 54
22/07/20 13:53:28 INFO ContextCleaner: Cleaned accumulator 50
22/07/20 13:53:28 INFO SparkContext: Invoking stop() from shutdown hook
22/07/20 13:53:28 INFO SparkUI: Stopped Spark web UI at http://ip-172-31-89-245.ec2.internal:4040
22/07/20 13:53:28 INFO YarnClientSchedulerBackend: Interrupting monitor thread
22/07/20 13:53:28 INFO YarnClientSchedulerBackend: Shutting down all executors
22/07/20 13:53:28 INFO YarnSchedulerBackend$YarnDriverEndpoint: Asking each executor to shut down
22/07/20 13:53:28 INFO SchedulerExtensionServices: Stopping SchedulerExtensionServices
(serviceOption=None,
 services=List(),
 started=false)
22/07/20 13:53:28 INFO YarnClientSchedulerBackend: Stopped
22/07/20 13:53:28 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
22/07/20 13:53:28 INFO MemoryStore: MemoryStore cleared
22/07/20 13:53:28 INFO BlockManager: BlockManager stopped
22/07/20 13:53:28 INFO BlockManagerMaster: BlockManagerMaster stopped
22/07/20 13:53:28 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
22/07/20 13:53:28 INFO SparkContext: Successfully stopped SparkContext
22/07/20 13:53:28 INFO ShutdownHookManager: Shutdown hook called
22/07/20 13:53:28 INFO ShutdownHookManager: Deleting directory /mnt/tmp/spark-cf4ebecc-8f3c-4327-923a-2cdec1e6fcfa
22/07/20 13:53:28 INFO ShutdownHookManager: Deleting directory /mnt/tmp/spark-5ffcce5f-7f0e-4d19-b5f3-f3c4cdd7c691/pyspark-21801177-29af-48c4-9085-3408702681
35
22/07/20 13:53:28 INFO ShutdownHookManager: Deleting directory /mnt/tmp/spark-5ffcce5f-7f0e-4d19-b5f3-f3c4cdd7c691
[hadoop@ip-172-31-89-245 ~]$ hadoop fs -ls
Found 4 items
drwxr-xr-x   - hadoop hadoop          0 2022-07-20 13:53 .sparkStaging
drwxr-xr-x   - hadoop hadoop          0 2022-07-20 13:50 clickStreamData
drwxr-xr-x   - hadoop hadoop          0 2022-07-20 13:53 clickStream_flatten_data
drwxr-xr-x   - hadoop hadoop          0 2022-07-20 13:50 clickstream_checkpoint
[hadoop@ip-172-31-89-245 ~]$ hadoop fs -ls clickStream_flatten_data
Found 2 items
-rw-r--r--   1 hadoop hadoop          0 2022-07-20 13:53 clickStream_flatten_data/_SUCCESS
-rw-r--r--   1 hadoop hadoop     454733 2022-07-20 13:53 clickStream_flatten_data/part-00000-bd4158e2-11f9-4f5c-a72b-f21c1ced6e27-c000.csv
[hadoop@ip-172-31-89-245 ~]$
```
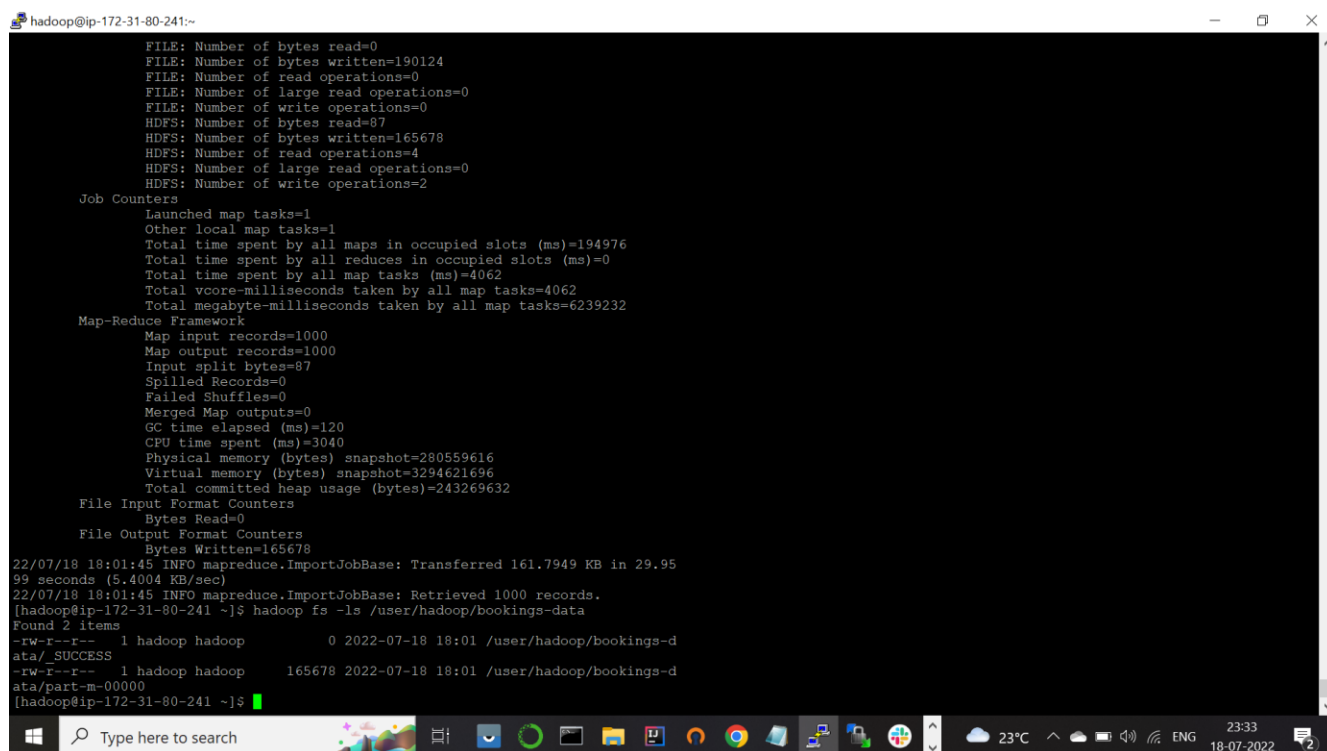
## 18. Screenshot of the flattened data

- **Task 2**: To write a script to ingest the relevant bookings data from AWS RDS to Hadoop.

1. For this task, first I installed mysql connector and MySQL on EMR cluster
2. Checked if the directory to load the data is already present in Hadoop

```
hadoop fs -rm -r /user/hadoop/bookings-data
```

3. Imported data from AWS RDS to Hadoop using command:

```
sqoop import \
--connect
jdbc:mysql://upgraddetest.cyaielc9bmnf.us-
east-1.rds.amazonaws.com/testdatabase \
--table bookings \
--username student --password STUDENT123 \
--target-dir /user/hadoop/bookings-data \
--m 1
```

4. Viewed imported data in hdfs using command

```
hadoop fs -ls /user/hadoop/bookings-data
```

```
hadoop fs -cat /user/hadoop/bookings-data/part-m-00000
```

5-4562,2,1,4
BK7038242108,16698708,27674043,4.4.4,Android,-11.9722805,-167.370596,70.401513,173.928091,2020-05-06 18:31:33.0,2020-07-14 02:15:47.0,467,94,INR,green,236-56
-9503,2,1,3
BK1805019792,72602005,80532337,3.2.26,iOS,-80.120948,151.351093,37.0961345,-143.290695,2020-08-01 21:33:09.0,2020-04-08 18:07:14.0,693,62,INR,olive,709-92-61
70,1,3,1
BK1718641612,67559017,61308810,4.4.30,iOS,35.814832,118.787533,36.3208175,-69.32082,2020-06-21 01:20:34.0,2020-03-25 12:11:42.0,932,63,INR,fuchsia,723-15-949
3,2,5,4
BK7168466183,15278180,42292257,3.1.40,iOS,50.214794,142.629899,19.394059,145.434588,2020-01-28 16:32:38.0,2020-05-28 10:25:29.0,263,21,INR,fuchsia,568-30-416
9,2,4,1
BK4797431513,28119681,42859819,1.1.8,Android,36.3909495,-69.675475,45.985664,-141.017147,2020-04-17 10:46:25.0,2020-10-17 00:56:14.0,481,1,INR,yellow,410-04-
8620,2,5,2
BK817041033,43759137,43883071,3.1.5,Android,82.2011535,53.591668,-43.3610485,-115.986783,2020-06-02 22:13:15.0,2020-10-18 23:20:14.0,260,81,INR,maroon,151-30
-6951,2,3,1
BK3091366811,78172665,13232332,1.1.19,iOS,-54.240295,-77.782015,-22.0201085,71.882567,2020-09-17 21:17:22.0,2020-07-25 03:14:58.0,751,89,INR,blue,062-69-8682
,5,3,1
BK2211742187,98294088,38321934,3.4.35,Android,-60.7373515,-21.294616,-25.7328075,-12.309869,2020-02-08 11:34:53.0,2020-04-17 04:14:10.0,295,22,INR,maroon,379
-54-1940,4,3,1
BK5718782371,91413624,50796177,1.2.17,iOS,12.212756,-49.495039,35.83743,21.181406,2020-08-02 15:24:44.0,2020-06-02 16:24:03.0,649,1,INR,olive,687-83-0160,4,1
,3
BK6925041861,16754182,65101252,2.3.4,iOS,-50.394933,-156.101704,81.0656725,-135.869812,2020-09-11 10:49:04.0,2020-03-05 22:44:59.0,168,69,INR,purple,711-68-5
952,4,2,3
BK1273485977,77394438,78242097,3.4.20,iOS,-81.8097455,-2.09345,-46.8590855,-69.444371,2020-01-24 03:41:51.0,2020-02-11 06:38:16.0,442,99,INR,green,210-68-582
5,3,4,4
BK1840869868,70959923,34539083,1.1.28,Android,-62.9599155,6.000395,-10.9973665,81.94683,2020-08-23 05:50:43.0,2020-07-31 14:24:55.0,875,40,INR,olive,401-09-0
475,5,3,2
BK997882598,58504148,46158041,4.3.35,iOS,29.9369725,-124.215858,-58.4375475,-89.123929,2020-01-12 08:48:05.0,2020-04-20 20:47:16.0,230,3,INR,green,266-60-742
4,3,1,1
BK5662087214,77440134,39470327,3.2.17,Android,13.965136,16.864746,-42.6375745,172.758412,2020-02-02 12:03:37.0,2020-09-26 06:29:56.0,941,88,INR,gray,279-59-5
349,5,4,1
BK8865887274,71215906,71671654,2.3.21,iOS,-85.653217,-6.88182,9.9079,-11.557289,2020-08-12 05:42:41.0,2020-06-15 18:24:22.0,894,72,INR,green,228-40-4173,2,2,
4
BK9708496297,29394176,20731780,2.4.7,Android,-45.18792,-33.406776,-67.895486,78.633127,2020-05-31 12:33:42.0,2020-04-12 21:21:13.0,938,15,INR,aqua,616-72-497
1,2,5,3
BK605227090,94946719,98071218,3.4.8,iOS,-4.8054875,-76.313234,76.0732515,-136.36527,2020-07-25 14:13:40.0,2020-07-04 14:08:20.0,716,19,INR,teal,112-71-5489,2
,5,3
BK8602928713,33696426,68228240,4.2.33,iOS,29.8689545,83.997957,-43.8161155,88.585796,2020-07-10 22:27:16.0,2020-06-07 20:32:40.0,763,65,INR,fuchsia,069-53-21
45,5,1,1
BK6243816121,74681879,25278350,2.3.24,iOS,-53.204624,-50.218077,85.7211825,-25.819898,2020-05-11 20:53:51.0,2020-08-26 11:52:47.0,196,78,INR,purple,874-93-24
59,4,1,4
BK9843664360,67564464,71913052,2.1.30,iOS,-10.031329,-49.100434,73.858805,157.577843,2020-09-20 22:51:48.0,2020-01-27 22:11:07.0,289,49,INR,maroon,876-76-759
8,3,3,3
BK6282269780,86988153,87739332,4.4.35,iOS,64.2850085,95.841271,-84.535181,-96.117739,2020-06-23 05:38:14.0,2020-07-11 16:40:57.0,196,89,INR,purple,740-96-793
1,4,5,4
[hadoop@ip-172-31-80-241 ~]$

- **Task 3**: To create aggregates for finding date-wise total bookings using the Spark script.

1. For this task, created a **datewise_bookings_aggregates_spark.py** file.
2. Imported necessary libraries

```
from pyspark.sql import SparkSession
from pyspark.sql import functions as F
from pyspark.sql.functions import col
from pyspark.sql.types import *
```

3. Established spark connection

```
spark = SparkSession \
        .builder \
        .appName('aggregateBatchData') \
        .master('yarn') \
        .getOrCreate()
```

4. Read data from csv filed extracted from AWS RDS and stored in HDFS

```
df=spark.read.csv('/user/hadoop/bookings-
data/',header=False,inferSchema = True)
```

5. Added column data according to given data

```
new_columns=["booking_id","customer_id","driver_id","customer_app_version","
customer_phone_os_version","pickup_lat","pickup_lon","drop_lat","drop_lon","
pickup_timestamp","drop_timestamp","trip_fare","tip_amount","currency_code",
"cab_color","cab_registration_no","customer_rating_by_driver","rating_by_cus
tomer","passenger_count"]

new_df = df.toDF(*new_columns)
```

6. Created a new column with date extracted from **pickup_timestamp** column

```
new_df = new_df.withColumn("date", F.to_date(F.col("pickup_timestamp")))
```

7. Datewise bookings aggregate using groupBy function

```
aggregate_df = new_df.groupby('date').count()
```

8. Wrote **aggregate_df** dataframe in csv files in HDFS

```
aggregate_df.write.csv('/user/hadoop/bookings_aggregate_data/')
```

9. Executed **datewise_bookings_aggregates_spark.py** using command:

```
spark-submit datewise_bookings_aggregates_spark.py
```

10. Screenshot of the csv files and data

- **Task 4**: To create a Hive-managed table for clickstream data, bookings data and aggregated data:

1. Below command is used to launch hive CLI

```
hive
```

2. Database **cab_rides_data** is created using command:

```
create database if not exists cab_rides_data;
```

3. Command to create clickStreamData table:
   As the clickStreamData has column header, I have used command to skip the first
row

```
create table if not exists clickStreamData(
customer_id int,
app_version string,
os_version string,
lat double,
lon double,
page_id string,
button_id string,
is_button_click string,
is_page_view string,
is_scroll_up string,
is_scroll_down string,
`timestamp` timestamp)
row format delimited fields terminated by ',' lines
terminated by '\n' stored as textfile
tblproperties("skip.header.line.count"="1");
```

4. Command to load data from HDFS to **clickStreamData** table:

```
load data inpath '/user/hadoop/clickStream_flatten_data/' into
table clickStreamData;
```

5. Screenshot of the loaded data in clickStreamData table:



6. To check the number of rows in **clickStreamData** table:

7. Command to create **bookingsData** table:

```
create table if not exists bookingsData(
booking_id string,
customer_id int,
driver_id int,
customer_app_version string,
customer_phone_os_version string,
pickup_lat double,
pickup_lon double,
drop_lat double,
drop_lon double,
pickup_timestamp timestamp,
drop_timestamp timestamp,
trip_fare double,
tip_amount double,
currency_code string,
cab_color string,
cab_registration_no string,
customer_rating_by_driver int,
rating_by_customer int,
passenger_count int)
row format delimited fields terminated by ',' lines terminated
by '\n' stored as textfile;
```

8. Command to load the data into **bookingsData** table:

```
load data inpath '/user/hadoop/bookings-
data/' into table bookingsData;
```

9. Screenshot of the loaded data and count of rows in **bookingsData** table:



10. Command to create **bookingsAggregateData** table and load the data:
In bookings_aggregate_data, we have **date** column with **date** as **datatype**. Hence to cast the column in date type, I have created a temporary table named as **testAggregateData** and loaded with data from HDFS:

Command to create testAggregateData table:

```
create table if not exists testAggregateData(
`date` string,
no_of_bookings int)
row format delimited fields terminated by ','
lines terminated by '\n' stored as textfile;
```

Command to load data from HDFS:

```
load data inpath '/user/hadoop/bookings_aggregate_data/'
into table testAggregateData;
```

Screenshot of the loaded data



In the next step, I have created a table named as **bookingsAggregateData** to cast the column date into date datatype.

Command to create **bookingsAggregateData** table**:**

```
create table bookingsAggregateData as select
cast(`date` as date),no_of_bookings from
testAggregateData;
```

Screenshot of the data loaded in bookingsAggregateData table and total count of rows: