# Logic For First Submission

**<Properly explain the code, list the steps to run the code provided by you and attach screenshots of code execution>**


**->**
The main project is regarding capturing the cab rides data from different sources(Kafka and AWS RDS) and making them by consumable by creating aggregated data and storing the aggregated and non-aggregated data into hive tables.

Breaking the project down to multiple steps :


**Step 1** : Ingesting Click stream data from a Kafka topic into Hadoop
- ➔ We take the help of spark to read the data from the Kafka Topic and then store the data into the hadoop
- ➔ The data read from the Kafka topic is in the Json format so we will pull the Value of the Key Value in Json format and store them in hadoop
- ➔ For the above point we will use readstream and dataframes to store the data which is then written in the form of Json file again into Hadoop
- ➔ The Spark code is as follows:

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


# Initializing Spark session
spark = SparkSession \
    .builder \
    .appName("Kafka-to-local") \
    .getOrCreate()

#Reading Streaming data from de-capstone3 kafka topic[Shared by Upgrad]
# we are using starting offsets and setting it as earliest to pull all the data
df = spark.readStream \
    .format("kafka") \
    .option("kafka.bootstrap.servers", "18.211.252.152:9092") \
    .option("startingOffsets", "earliest") \
    .option("subscribe", "de-capstone3") \
    .load()
```

```
# We only copy the value part from the json format data
# we transform the same to string and then drop all the remaining value pulled
df= df \
    .withColumn('value_str',df['value'].cast('string').drop('value') \
    .drop('key','topic','partition','offset','timestamp','timestampType')

#Writing data from kakfa to Hadoop
df.writeStream \
  .format("json") \
  .outputMode("append") \
  .option("path", "/user/root/clickstream_data_dump") \
  .option("checkpointLocation", "/user/root/clickstream_data_dump_cp") \
  .start() \
  .awaitTermination()
```

➔ Copy this spark file to hadoop using winscp
➔ Then Copy this file to /user/livy for the code to be accesible to be run for spark submit job
➔ The File attached by the name "spark_kafka_to_local.py" contains the spark code which would need to be run as spark job using spark submit using the following code:
"

**export SPARK_KAFKA_VERSION=0.10**
**spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.1.2**
**spark_kafka_to_local.py**
"

The above code transfers the spark kafka dependecies for the spark job to interact with the kafka topic
➔ Screenshot for running the code

➔ Verifying if the json file was generated at the target location "/user/root/clickstream_data_dump"
  o Code: hadoop fs -ls /user/root/clickstream_data_dump

Screen shot for the same :



The screen shot shows that json file is generated with data copied from the kafka Topic

Screenshot of the data stored in the part file :



**Step 2:** The second step is to actually flatten the data
  ➔ The current data is string of data of the form of nested Json
  ➔ We write a Spark query to read this json file and then flatten the data -> make it in the form of a csv file [tabular format having all the data ]
  ➔ Spark query:

```
# importing the important classes/libraries
from pyspark.sql import SparkSession
from pyspark.sql.functions import *
from pyspark.sql.types import *
```

```
# Initializing Spark session
spark = SparkSession.builder \
    .appName("Flatten") \
    .master("local")\
    .getOrCreate()

# Reading the json file using spark.read.json function
df=spark.read.json("/user/root/clickstream_data_dump/part-00000-c09da594-6a92-4afe-
81ef-351b785d6e49-c000.json")

#Do not forget to change the name for the part file


#We take help of "get_json_object" to deal with nested json format to refer to the nested
#columns and providing alias at the same time
df_final=df.select(get_json_object(df['value_str'],"$.customer_id").alias("customer_id"),
        get_json_object(df['value_str'],"$.app_version").alias("app_version"),
        get_json_object(df['value_str'],"$.OS_version").alias("OS_version"),
        get_json_object(df['value_str'],"$.lat").alias("lat"),
        get_json_object(df['value_str'],"$.lon").alias("lon"),
        get_json_object(df['value_str'],"$.page_id").alias("page_id"),
        get_json_object(df['value_str'],"$.button_id").alias("button_id"),
        get_json_object(df['value_str'],"$.is_button_click").alias("is_button_click"),
        get_json_object(df['value_str'],"$.is_page_view").alias("is_page_view"),
        get_json_object(df['value_str'],"$.is_scroll_up").alias("is_scroll_up"),
        get_json_object(df['value_str'],"$.is_scroll_down").alias("is_scroll_down"),
        get_json_object(df['value_str'],"$.timestamp").alias("timestamp")
        )

# Finally writing the data into csv file and storing at the below given location
df_final.write.format('csv').mode('overwrite').save('/user/root/clickstream_flattened',heade
r='true')
```

➔ In order to run the spark submit program we need to run the following code( ensure that
the spark file name "spark_local_flatten.py" has been copied to /user/livy:
"

**export SPARK_KAFKA_VERSION=0.10**
**spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.1.2**
**spark_local_flatten.py**
"

➔ Screenshot of running the code :



Verifying that the csv file is generated at the target location
Code: hadoop fs -ls /user/root/clickstream_flattened

Screenshot :



Screenshot of the flattened data in the above csv file:
Code: hadoop fs -cat /user/root/clickstream_flattened/part-00000-01ef8efc-8ab9-4f9d-b037-0f398112fb84-c000.csv

##change the name of the csv file as per the new name generated when run

**Step 3:** The next steps involves ingesting data from AWS RDS table to the hadoop
- ➔ We take help of sqoop
- ➔ First we need to install Mysql connectors in our cluster so that we can use the jdbc connectors to connect to the database and pull the data for the same
- ➔ Code:

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/

Screenshot



- ➔ Now we can connect to the RDS using the jdbc connector and the user id and password shared
- ➔ Code:

```
sqoop import \
--connect jdbc:mysql://upgraddetest.cyaielc9bmnf.us-east-
1.rds.amazonaws.com/testdatabase  \
--table bookings \
--username student --password STUDENT123 \
--target-dir /user/root/bookings_1 \
--m 1
```

-> Screen shot

Screen shot showing 1000 records pulled:



Verifying the file generated at the target location:
Code: hadoop fs -ls /user/root/bookings_1



The image shows a part file was generated with the data pulled form the RDS

Screenshot of the data stored in the part file :

**Step 4:** The next step involves creating a aggregated file-> the number of bookings needs to be aggregated at daywise level.

The code is as follows:

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

Spark=
SparkSession.builder.appName("Aggregate_by_pick_up_date").master("local").g
etOrCreate()

df=spark.read.csv("/user/root/bookings_1/part-m-00000")
new_col =
["booking_id","customer_id","driver_id","customer_app_version","customer_phon
e_os_version","pickup_lat","pickup_lon","drop_lat",

"drop_lon","pickup_timestamp","drop_timestamp","trip_fare","tip_amount","curren
cy_code","cab_color","cab_registration_no","customer_rating_by_driver",
      "rating_by_customer","passenger_count"]

#Creating a new Dataframe form the old data frame and assigning a new
schema
new_df = df.toDF(*new_col)

new_df.show(5)

# Crearing a new column which has only the date part from the
pickup_timestamp column
df_new= new_df.withColumn("Pickupdate",col("pickup_timestamp").cast("date"))
df_new.show(5)
```
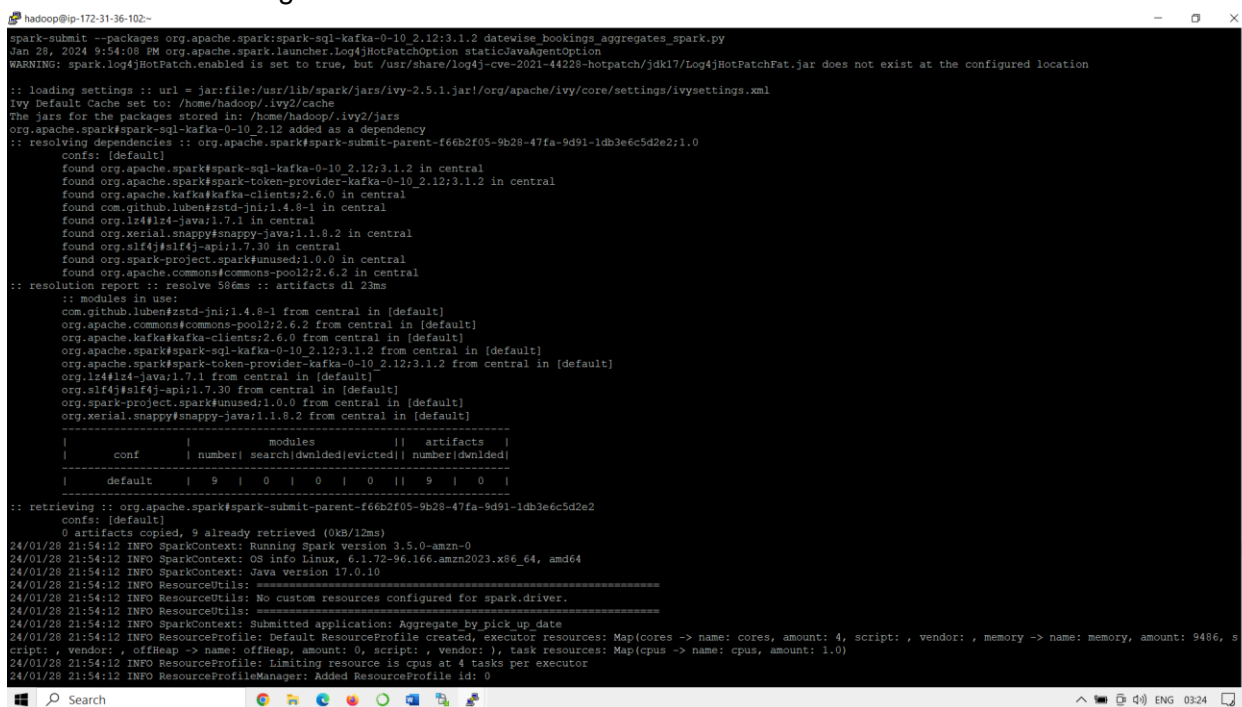
```
# grouping the data by the new column Pickupdate
agg_df =
df_new.groupBy("Pickupdate").agg(count("booking_id").alias("booking_count")).o
rderBy("Pickupdate")


# Copying the data to a csv file
agg_df.coalesce(1).write.format('csv').mode('overwrite').save('/user/root/datewise
_bookings_agg',header='true')
```

➔ Ensure the spark file is copied to /user/livy folder
➔ The attached file name is "datewise_bookings_aggregates_spark.py"
➔ Code to run the spark job -> "export SPARK_KAFKA_VERSION=0.10
　　　　　　　spark-submit --packages org.apache.spark:spark-sql-
　　　　　　　kafka-0-10_2.12:3.1.2
　　　　　　　datewise_bookings_aggregates_spark.py"
➔ Screenshot of running the file :



➔ Verifying the csv file got generated at the target location:
　　o Code: hadoop fs -ls /user/root/datewise_bookings_agg



➔ Screenshot of the data stored in the csv file

- o Code: hadoop fs -cat /user/root/datewise_bookings_agg/part-00000-c958e876-71d1-46ca-a84d-cd2273731394-c000.csv

## change the csv file name

➔ Screenshot of the above point :



**Step 5:** The next step involves generating the hive tables on the ingested data

Before creating the hive tables, Follow the following commands:

Type Hive to get into Hive CLI
Then Create a Database:
- ➔ Create Database CabRideProject;
- ➔ Use CabRideProject;

1. **Creating the Hive table to store the Booking data**
   ➔

   CREATE TABLE IF NOT EXISTS Bookings
   (
   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 DECIMAL(10, 2),
tip_amount DECIMAL(10, 2),
 currency_code STRING,
cab_color STRING,
cab_registration_no STRING,
customer_rating_by_driver INT,
rating_by_customer INT,
passenger_count INT
);
```

**Screen shot of the table creation**



2. **Creating table for storing then aggregated booking values [daywise]**
    ➔

   ```
   CREATE TABLE IF NOT EXISTS datewise_total_bookings
   ( pickup_date DATE,
   total_bookings INT );
   ```

   Screenshot for the same :



3. **Creating table for storing data for the click stream**
    ➔

   ```
   CREATE TABLE IF NOT EXISTS clickstream_data
   ( 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,
time_stamp TIMESTAMP
)  ;
```

Screenshot for the table creation command execution:



**<Command to load the data into Hive tables>**

Loading the data from hadoop to the hive tables;

1. Uploading the data from hadoop to Hive table for storing the bookings data ingested from AWS RDS to hadoop using sqoop

   ➔

   **LOAD DATA INPATH '/user/root/bookings_1/part-m-00000' OVERWRITE INTO TABLE Bookings**

   Screenshot for the running the same command and verifying the number of rows in the hive table

2. Uploading the data from hadoop to Hive table for the aggregated data created as a datewise aggregation of bookings

➔

**LOAD DATA INPATH '/user/root/datewise_bookings_agg' OVERWRITE INTO TABLE datewise_total_bookings;**

The image also contains verification of rows present in the hive table



3. Uploading the data form hadoop to hive -> This is for the data pulled form kafka and stored in the hadoop

➔

**LOAD DATA INPATH '/user/root/clickstream_flattened' into table clickstream_data;**



Below image also holds the verification query to run the numbers of rows