

Design of data pipeline for processing bookings and click stream data for YourOwnCabs Inc.

Atharv Darekar

10th September 2022

Contents

1	Read Data From Apache Kafka Topic	2
1.1	Analyze Incoming Data	2
1.2	Define Schema	2
1.3	Creating A Streaming DataFrame	3
1.4	Store The Streaming DataFrame	3
2	Import Table From Amazon Relational Database Service MySQL Server	3
2.1	Table Description	3
2.2	Import The Data Into HDFS	4
3	Ingest The Data Into Apache Hive	5
3.1	Ingest The Streaming Data	5
3.2	Ingest The Batch Data	5
4	Aggregation Over The Hive Tables Using Apache Spark	6
5	Run The Data Pipeline	7
5.1	Setup Directories To Stage Data	7
5.2	Execute Spark Structured Streaming Application	7
5.3	Setup Feeds Path To Ingest The Data Into Hive	7
5.4	Ingest The Data Into Hive	7
5.5	Aggregate The Bookings Data	8

1 Read Data From Apache Kafka Topic

1.1 Analyze Incoming Data

The data read from the Kafka topic is in JSON format. A sample data looks as given below:

```
{
  "customer_id": "62605529",
  "app_version": "2.1.36",
  "OS_version": "Android",
  "lat": "44.196239",
  "lon": "-15.354515",
  "page_id": "e7bc5fb2-1231-11eb-adc1-0242ac120002",
  "button_id": "a95dd57b-779f-49db-819d-b6960483e554",
  "is_button_click": "Yes",
  "is_page_view": "Yes",
  "is_scroll_up": "No",
  "is_scroll_down": "No",
  "timestamp\n": "2020-06-24 01:55:32\n"
}
```

1.2 Define Schema

From the sample data, it can be inferred that the data needs to be clean before imposing a proper schema. Hence the data is read as string type:

```
12 schema = StructType(
13     [
14         StructField(name="customer_id", dataType=StringType(),
15             ↪ nullable=True),
16         StructField(name="app_version", dataType=StringType(),
17             ↪ nullable=True),
18         StructField(name="OS_version", dataType=StringType(),
19             ↪ nullable=True),
20         StructField(name="lat", dataType=StringType(),
21             ↪ nullable=True),
22         StructField(name="lon", dataType=StringType(),
23             ↪ nullable=True),
24         StructField(name="page_id", dataType=StringType(),
25             ↪ nullable=True),
26         StructField(name="button_id", dataType=StringType(),
27             ↪ nullable=True),
28         StructField(name="is_button_click",
29             ↪ dataType=StringType(), nullable=True),
30         StructField(name="is_page_view", dataType=StringType(),
31             ↪ nullable=True),
32         StructField(name="is_scroll_up", dataType=StringType(),
33             ↪ nullable=True),
34         StructField(name="is_scroll_down", dataType=StringType(),
35             ↪ nullable=True),
36     ]
37 )
```

```

25         StructField(name="timestamp\n", dataType=StringType(),
26           ↪ nullable=True)
27     ]
28 )

```

1.3 Creating A Streaming DataFrame

The data is flatten, cleaned and stored in a dataframe.

```

29 click_stream_data = spark \
30     .readStream \
31     .format("kafka") \
32     .option("kafka.bootstrap.servers", "18.211.252.152:9092") \
33     .option("subscribe", "de-capstone3") \
34     .option("startingOffsets", "earliest") \
35     .load() \
36     .select(from_json(col("value").cast("string"),
37       ↪ schema).alias("click_stream_data")) \
38     .select("click_stream_data.*") \
39     .withColumnRenamed("OS_version", "os_version") \
40     .withColumnRenamed("timestamp\n", "timestamp")

```

1.4 Store The Streaming DataFrame

The streaming dataframe is written to the Apache Hadoop Distributed File System (HDFS) in CSV format.

```

41 write_click_stream_data = click_stream_data \
42     .writeStream \
43     .outputMode("append") \
44     .format("csv") \
45     .option("truncate", "false") \
46     .option("path", "/user/livy/click_stream_data") \
47     .option("checkpointLocation",
48       ↪ "hdfs:///user/livy/checkpoints/click_stream_data") \
49     .trigger(once=True) \
50     .start()

```

2 Import Table From Amazon Relational Database Service MySQL Server

2.1 Table Description

The Amazon RDS MySQL server stores the bookings data of the YourOwnCabs. The schema of the table is given below:

BOOKING_ID	STRING
CUSTOMER_ID	INT
DRIVER_ID	INT
CUSTOMER_APP_VERSION	STRING
CUSTOMER_PHONE_OS_VERSION	STRING
PICKUP_LAT	FLOAT
PICKUP_LON	FLOAT
DROP_LAT	FLOAT
DROP_LON	FLOAT
PICKUP_TIMESTAMP	TIMESTAMP
DROP_TIMESTAMP	TIMESTAMP
TRIP_FARE	FLOAT
TIP_AMOUNT	FLOAT
CURRENCY_CODE	STRING
CAB_COLOR	STRING
CAB_REGISTRATION_NO	STRING
CUSTOMER_RATING_BY_DRIVER	INT
RATING_BY_CUSTOMER	INT
PASSENGER_COUNT	INT

2.2 Import The Data Into HDFS

The bookings data hosted by the MySQL server is imported into the HDFS and stored in CSV format using the Apache Sqoop.

To import data from MySQL server using the Sqoop, MySQL connector needs to be setup on the Amazon Elastic Map Reduce (EMR) instance.

```

12 echo ">>> Download and install MySQL connector in Amazon Elastic
   ↪ Map Reduce (EMR) master node"
13
14 wget
   ↪ "https://de-mysql-connector.s3.amazonaws.com/mysql-connector-java-8.0.25.tar.gz"
15 sudo rm /usr/lib/sqoop/lib/mysql-connector-java-8.0.25.jar
16 sudo tar -xvf mysql-connector-java-8.0.25.tar.gz
   ↪ --strip-components=1 --directory /usr/lib/sqoop/lib/
   ↪ mysql-connector-java-8.0.25/mysql-connector-java-8.0.25.jar
17
18
19 echo ">>> Run the Apache Sqoop job to import bookings data from
   ↪ Amazon Relational Database Service (RDS) into HDFS"
20
21 sqoop import \
22 --connect
   ↪ jdbc:mysql://upgraddetest.cyaie1c9bmnf.us-east-1.rds.amazonaws.com/testdatabase
   ↪ \
23 --table bookings \
24 --username student --password STUDENT123 \
25 --target-dir /user/livy/bookings \
26 -m 1
27

```

3 Ingest The Data Into Apache Hive

3.1 Ingest The Streaming Data

The streaming data stored in the HDFS is ingested into Hive table using the serializing and de-serializing ensuring proper consumption of the data into the table.

```
1 CREATE TABLE IF NOT EXISTS TBL_CLICK_STREAMS
2 (
3     CUSTOMER_ID          INT,
4     APP_VERSION           STRING,
5     OS_VERSION            STRING,
6     LAT                   FLOAT,
7     LON                   FLOAT,
8     PAGE_ID               STRING,
9     BUTTON_ID             STRING,
10    IS_BUTTON_CLICK        STRING,
11    IS_PAGE_VIEW           STRING,
12    IS_SCROLL_UP           STRING,
13    IS_SCROLL_DOWN         STRING,
14    `TIMESTAMP`            TIMESTAMP
15 )
16 ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
17 STORED AS TEXTFILE;
18
19 LOAD DATA INPATH "hdfs:///user/livy/feeds_path/click_stream_data"
20 OVERWRITE INTO TABLE TBL_CLICK_STREAMS;
```

3.2 Ingest The Batch Data

The batch bookings data is ingested in a similar fashion as that of the streaming data.

```
1 CREATE TABLE IF NOT EXISTS TBL_BOOKINGS
2 (
3     BOOKING_ID            STRING,
4     CUSTOMER_ID           INT,
5     DRIVER_ID             INT,
6     CUSTOMER_APP_VERSION  STRING,
7     CUSTOMER_PHONE_OS_VERSION STRING,
8     PICKUP_LAT             FLOAT,
9     PICKUP_LON             FLOAT,
10    DROP_LAT               FLOAT,
11    DROP_LON               FLOAT,
12    PICKUP_TIMESTAMP        TIMESTAMP,
13    DROP_TIMESTAMP         TIMESTAMP,
14    TRIP_FARE               FLOAT,
15    TIP_AMOUNT              FLOAT,
16    CURRENCY_CODE           STRING,
17    CAB_COLOR               STRING,
```

```

18     CAB_REGISTRATION_NO          STRING,
19     CUSTOMER_RATING_BY_DRIVER    INT,
20     RATING_BY_CUSTOMER           INT,
21     PASSENGER_COUNT              INT
22 )
23 ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
24 STORED AS TEXTFILE;
25
26 LOAD DATA INPATH "hdfs:///user/livy/feeds_path/bookings"
27 OVERWRITE INTO TABLE TBL_BOOKINGS;

```

4 Aggregation Over The Hive Tables Using Apache Spark

The bookings data is aggregated date-wise using the PySpark, the aggregated data is stored as a Hive table as well as in CSV format in HDFS.

```

1  from pyspark.sql import SparkSession
2  from pyspark.sql.functions import col, to_date, count, sum as
   ↪ pyspark_sum
3
4  spark = SparkSession \
5      .builder \
6      .appName('aggregate_bookings_data') \
7      .enableHiveSupport() \
8      .getOrCreate()
9
10 bookings = spark.sql("SELECT * FROM TBL_BOOKINGS")
11
12 bookings_aggregated = bookings.withColumn("trip_date",
   ↪ to_date(col("pickup_timestamp"), "yyyy-MM-dd")) \
13     .groupBy(col("trip_date")) \
14     .agg(
15         count("booking_id").alias("number_of_trips"),
16         pyspark_sum("trip_fare").alias("total_trip_amount"),
17         pyspark_sum("tip_amount").alias("total_tip_amount"),
18
   ↪ pyspark_sum("passenger_count").alias("total_passenger_counts")
19     ) \
20     .orderBy("trip_date")
21
22 bookings_aggregated.createOrReplaceTempView("VW_AGGREGATED_BOOKINGS")
23 spark.sql("CREATE TABLE IF NOT EXISTS TBL_AGGREGATED_BOOKINGS \
24     AS SELECT * FROM VW_AGGREGATED_BOOKINGS")
25
26 bookings_aggregated \
27     .write \
28     .csv("/user/livy/aggregate_batch_data", header="false")

```

5 Run The Data Pipeline

The complete process from data ingestion, transformation and loading is automated using a simple bash script.

To execute the script, change the working directory to the root of the repository.¹

5.1 Setup Directories To Stage Data

```
1  #!/usr/bin/env bash
2
3  echo ">>> Setup Apache Hadoop Distributed File System (HDFS)
   ↳ directories to import click stream data"
4
5  hdfs dfs -rm -r /user/livy
6  hdfs dfs -mkdir -p /user/livy/click_stream_data
7
```

5.2 Execute Spark Structured Streaming Application

```
8  echo ">>> Run the Apache Spark Structured Streaming application
   ↳ to read data from Apache Kafka topic and store as CSV files
   ↳ in HDFS"
9
10 spark-submit --packages
   ↳ org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.5
   ↳ data_sourcing/import_streaming_data.py
11
```

5.3 Setup Feeds Path To Ingest The Data Into Hive

```
28 echo ">>> Setup feeds path and copy data into the feeds path for
   ↳ ingestion of the imported data into Apache Hive"
29
30 hdfs dfs -mkdir -p
   ↳ {/user/livy/feeds_path/click_stream_data,/user/livy/feeds_path/bookings}
31 hdfs dfs -cp /user/livy/click_stream_data/*.csv
   ↳ /user/livy/feeds_path/click_stream_data
32 hdfs dfs -cp /user/livy/bookings/part*
   ↳ /user/livy/feeds_path/bookings
33
```

5.4 Ingest The Data Into Hive

```
34 echo ">>> Ingest the click stream data and bookings data into
   ↳ Hive tables"
35
36 hive -f data_sourcing/ingest_streaming_data.sql
```

¹Refer to pipeline run video.

```
37  hive -f data_sourcing/ingest_batch_data.sql
38
```

5.5 Aggregate The Bookings Data

```
39  echo ">>> Run Apache Spark application to aggregate bookings
    ↪ data, load the aggregated data into a Hive table and store as
    ↪ CSV files in HDFS"
40
41  spark-submit data_transformation/aggregate_batch_data.py
```