

# Retail Data Analysis Project

Code Logic Explanation

**Submitted by:**

Prerna Gupta (prerna.gupta1792@gmail.com)

# Code Logic- Retail Data Analysis

## 1. Importing necessary spark libraries and functions

```
4 from pyspark.sql import SparkSession
5 from pyspark.sql.functions import *
6 from pyspark.sql.types import *
7 from pyspark.sql.window import Window
8
```

The code starts by importing the necessary libraries.

- SparkSession is the entry point to programming Spark with the DataFrame API. It allows you to create a DataFrame and execute SQL queries on the data.
- Second line imports all functions from the pyspark.sql.functions module. These functions are used to manipulate data within DataFrames. Examples include functions for data aggregation, filtering, and transformation such as sum, avg, etc.
- Third line imports all types from the pyspark.sql.types module. These types are used to define the schema for DataFrames. Examples include StringType, IntegerType, StructType, StructField, etc.
- Fourth line imports the Window class from the pyspark.sql.window module. Window class is used to specify windowing functions (such as ranking, aggregations etc.) that operate over a specified range of rows relative to the current row.

## 2. Creating Spark Session

```
9 #creating SparkSession
10 spark = SparkSession \
11     .builder \
12     .appName("RetailDataKPI") \
13     .getOrCreate()
14 spark.sparkContext.setLogLevel('ERROR')
```

Starting the spark session with RetailDataKPI name to establish connection to the spark cluster. getOrCreate() method is used to either create a new SparkSession if one does not already exist or return the existing SparkSession. spark.sparkContext.setLogLevel('ERROR') sets the log level for the Spark context to ERROR, so that only error messages will be logged, which can help reduce the amount of log output and make it easier to spot issues.

### 3. Creating UDFs (User Defined Functions)

- **is\_Order:** Defining a function using if-else statement to check if the type of order is a new order, returning 1 if it's a new order.

```
16 # Creating UDFs
17 #Checking for a new order
18 def is_order(type):
19     if type=="ORDER":
20         return 1
21     else:
22         return 0
23
```

- **is\_return:** Defining a function using if-else statement to check if it's a return order, returning 1 if it's a return.

```
24 #Checking for return of an order
25 def is_return(type):
26     if type=="RETURN":
27         return 1
28     else:
29         return 0
```

- **total\_items:** Creating a function to calculate the total number of items in an order. Using for loop summed up the quantities of each product in a particular invoice to get the number of total items associated with a particular transaction.

```
31 # Calculating the total number of items
32 def total_items(items):
33     if items is not None:
34         item_count =0
35         for item in items:
36             item_count = item_count + item['quantity']
37         return item_count
```

- **total\_cost :** Defining a function to calculate the total cost of a particular transaction. Calculated the cost by multiplying the unit price with quantity and adding it for all products in the transaction. If it's a return transaction then multiplying the cost by -1 as it will not contribute to the revenue.

```
39 # Calculating the total cost of the order
40 def total_cost(items,type):
41     if items is not None:
42         total_cost =0
43         item_price =0
44         for item in items:
45             item_price = (item['quantity']*item['unit_price'])
46             total_cost = total_cost+ item_price
47             item_price=0
48         if type == 'RETURN':
49             return total_cost *(-1)
50         else:
51             return total_cost
```

#### 4. Registering UDFs

```
53 # Register UDFs
54 is_order = udf(is_order, IntegerType())
55 is_return = udf(is_return, IntegerType())
56 add_total_items = udf(total_items, IntegerType())
57 add_total_cost = udf(total_cost, FloatType())
```

Each line registers a Python function (is\_order, is\_return, total\_items, total\_cost) as a UDF in PySpark, specifying the expected return type for each function. These UDFs will be used to transform streaming data and adding columns to the batch output in console.

#### 5. Reading Data from Kafka topic

```
61 # Reading input data from Kafka topic
62 raw_data_stream = spark \
63     .readStream \
64     .format("kafka") \
65     .option("kafka.bootstrap.servers", "18.211.252.152:9092") \
66     .option("subscribe", "real-time-project") \
67     .option("startingOffsets", "latest") \
68     .load()
```

Above code sets up a PySpark streaming DataFrame (raw\_data\_stream) to read real-time data from a Kafka topic named **'real-time-project'**. It connects to a Kafka broker at **18.211.252.152:9092** and starts reading from the latest offset.

#### 6. Defining the Schema

```
70 # Defininig json Schema
71 json_schema = StructType() \
72     .add("invoice_no", LongType()) \
73     .add("country", StringType()) \
74     .add("timestamp", TimestampType()) \
75     .add("type", StringType()) \
76     .add("items", ArrayType(StructType([
77         StructField("SKU", StringType()),
78         StructField("title", StringType()),
79         StructField("unit_price", FloatType()),
80         StructField("quantity", IntegerType())
81     ])))
```

Defining the schema for the incoming json data from the kafka topic.

#### 7. Deserialize Kafka Messages

```
82
83 orders_stream_data = raw_data_stream.select(from_json(col("value").cast("string"), json_schema).alias("data")).select("data.*")
84
```

Above code Reads the value column from Kafka messages, converts it to a string, parses the JSON data using the provided schema, and aliases it as data after that it selects all fields from the parsed JSON data, making them individual columns in the resulting DataFrame.

## 8. Adding new columns

```
86 # Calculating additional columns from the stream
87 orders_stream_output = orders_stream_data \
88     .withColumn("total_cost", add_total_cost(orders_stream_data.items,orders_stream_data.type)) \
89     .withColumn("total_items", add_total_items(orders_stream_data.items)) \
90     .withColumn("is_order", is_order(orders_stream_data.type)) \
91     .withColumn("is_return", is_return(orders_stream_data.type))
```

Above code adds the 4 more columns to the dataframe using the UDFs created in the beginning of the code by taking the streaming data as input and transforming it.

## 9. Writing to Console

```
94 # Writing the summarised input table to the console
95 orders_batch = orders_stream_output \
96     .select("invoice_no", "country", "timestamp", "total_cost", "total_items", "is_order", "is_return") \
97     .writeStream \
98     .outputMode("append") \
99     .format("console") \
100     .option("truncate", "false") \
101     .option("path", "/Console_output") \
102     .option("checkpointLocation", "/Console_output_checkpoints") \
103     .trigger(processingTime="1 minute") \
104     .start()
```

Above code lines selects the specific columns from the orders\_stream\_output DataFrame and writes it to console in append mode every 1 minute and also maintains checkpoints to keep track of the streaming state and progress.

## 10. Calculating Time Based KPIs

```
106 # Calculating Time based KPIs
107 agg_on_time = orders_stream_output \
108     .withWatermark("timestamp", "1 minutes") \
109     .groupby(window("timestamp", "1 minute")) \
110     .agg(sum("total_cost").alias("total_volume_of_sales"),
111         avg("total_cost").alias("average_transaction_size"),
112         count("invoice_no").alias("OPM"),
113         avg("is_return").alias("rate_of_return")) \
114     .select("window.start", "window.end", "OPM", "total_volume_of_sales", "average_transaction_size", "rate_of_return")
115
```

Above processes streaming order data, aggregates it in 1-minute intervals, and calculates several metrics like total sales volume, average transaction size, number of transactions, and return rate for each 1-minute window, then selects the required columns and stores in agg\_on\_time DF.

## 11. Calculating Time and Country based KPIs

```
116 # Calculating Time and country based KPIs
117 agg_on_time_country = orders_stream_output \
118     .withWatermark("timestamp", "1 minutes") \
119     .groupBy(window("timestamp", "1 minutes"), "country") \
120     .agg(sum("total_cost").alias("total_volume_of_sales"),
121         count("invoice_no").alias("OPM"),
122         avg("is_return").alias("rate_of_return")) \
123     .select("window.start", "window.end", "country", "OPM", "total_volume_of_sales", "rate_of_return")
```

Above code lines processes streaming data and calculate the KPIs like Order per minute (OPM), total volume of sales and rate of return in a 1-minute interval and then aggregating the results on country in the 1-minute window.

## 12. Writing to HDFS

```
125 # Writing to the json files in hdfs : Time based KPI values
126 ByTime = agg_on_time.writeStream \
127     .format("json") \
128     .outputMode("append") \
129     .option("truncate", "false") \
130     .option("path", "timeKPIvalue") \
131     .option("checkpointLocation", "timeKPIvalue_cp") \
132     .trigger(processingTime="1 minutes") \
133     .start()
134
135 # Writing to the hdfs : Time and country based KPI values
136 ByTime_country = agg_on_time_country.writeStream \
137     .format("json") \
138     .outputMode("append") \
139     .option("truncate", "false") \
140     .option("path", "time_countryKPIvalue") \
141     .option("checkpointLocation", "time_countryKPIvalue_cp") \
142     .trigger(processingTime="1 minutes") \
143     .start()
```

Above code lines write the Kpis calculated based on time and based time-country into json files in append mode in HDFS and also maintains the checkpoint information for a processing time of 1 minute.

## 13. Await Termination

```
145 orders_batch.awaitTermination()
146 ByTime.awaitTermination()
147 ByTime_country.awaitTermination()
```

Above 3 lines are used to keep the streaming jobs running and continuously processing incoming data. They ensure that the Spark application does not exit immediately after starting but remains active to handle streaming data as it arrives.

#### 14. Running the pyspark script

```
[hadoop@ip-172-31-7-149 ~]$ export SPARK_KAFKA_VERSION=0.10
```

```
[hadoop@ip-172-31-7-149 ~]$ spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.5.1 spark-streaming.py
```

Above commands are used to submit the spark job and read data from kafka. In essence, this command runs a Spark Streaming application defined in **spark-streaming.py**, with the necessary Kafka connector (spark-sql-kafka-0-10\_2.12:3.5.1) to enable integration between Spark and Kafka. The export SPARK\_KAFKA\_VERSION=0.10 line ensures compatibility with **Kafka version 0.10**.