

# Code Logic - Retail Data Analysis

```

1  #importing necessary libraries and python modules
2
3  import sys
4  import os
5  from pyspark.sql import SparkSession
6  from pyspark.sql.functions import *
7  from pyspark.sql.types import *
8  from ast import literal_eval

```

Required python modules and libraries are imported here. Literal\_eval will be used to convert string from items column into a proper python list of dictionary.

```

9
10 # get total cost. It will be arrived at by summing up the unit price and quantity of products.
11 def get_total_cost(items):
12     items = literal_eval(items)
13     total_cost = 0
14     for item in items:
15         total_cost += item["unit_price"] * item["quantity"]
16     return total_cost
17
18 # get total items. It will be arrived by summing up the total quantity of the products.
19 def get_total_items(items):
20     items = literal_eval(items)
21     total_items = 0
22     for item in items:
23         total_items += item["quantity"]
24     return total_items
25
26 # if that order is ORDER or RETURN. In case the category is ORDER return 1 else i.e., category is RETURN, return 0 for order type.
27 def type_order(category):
28     if category == "ORDER":
29         return 1
30     return 0
31
32 # if that order is ORDER or RETURN. In case the category is ORDER return 0 else i.e., category is RETURN, return 1 for return type.
33 def type_return(category):
34     if category == "RETURN":
35         return 1
36     return 0
37
38

```

Below are the details of the custom functions:

1. `get_total_cost(items)`: This function takes item as an argument and calculate the total cost by summing up the unit price and quantity of products. The formulae to calculate get total cost is :

$$\sum(\text{quantity} * \text{unitprice})$$

2. `get_total_items(items)`: This function takes item as an argument and used to retrieve the total items by summing up total quantity of the products. The formulae to calculate the total item is:

$$\Sigma(\text{quantity})$$

3. `type_order(category)`: This function takes category as an argument, and it is used to map type of order. If that order is ORDER or RETURN. In case the category is ORDER return 1 else i.e., category is RETURN, return 0 for order type.
4. `type_return(category)`: This function takes category as an argument, and it is used to map type of order. If that order is ORDER or RETURN. In case the category is ORDER return 0 else i.e., category is RETURN, return 1 for return type.

```

39
40     if len(sys.argv) != 4:
41         print("Usage: spark-submit spark-streaming.py <hostname> <port> <topic>")
42         exit(-1)
43
44     host = sys.argv[1]
45     port = sys.argv[2]
46     topic = sys.argv[3]
47
48     spark = SparkSession \
49 ———» .builder \
50 ———» .appName("RetailDataAnalysis") \
51 ———» .getOrCreate()
52     spark.sparkContext.setLogLevel('ERROR')
53
54     bootstrap_server = host + ":" + port
55
56     lines = spark \
57 ———» .readStream \
58 ———» .format("kafka") \
59 ———» .option("kafka.bootstrap.servers", bootstrap_server) \
60 ———» .option("subscribe", topic) \
61 ———» .load()
62

```

In the beginning, the host, port number and topic is received from the command line argument.  
 host: 18.211.252.152  
 port: 9092  
 topic: real-time-project

This is used to read the spark stream from kafka bootstrap server.

```

62 schema = StructType() \
63     .add("invoice_no", StringType()) \
64     .add("country", StringType()) \
65     .add("timestamp", TimestampType()) \
66     .add("type", StringType()) \
67     .add("items", StringType())
68
69 raw_data = lines.selectExpr("cast(value as string)").select(from_json("value", schema).alias("temp")).select("temp.*")
70
71 # create user-defined functions for each
72 total_cost = udf(lambda items: get_total_cost(items))
73 total_quantity = udf(lambda items: get_total_items(items))
74 is_order = udf(lambda types: type_order(types))
75 is_return = udf(lambda types: type_return(types))
76
77 new_df = raw_data
78 new_df = new_df.withColumn("total_cost", total_cost("items")) \
79     .withColumn("total_items", total_quantity("items")) \
80     .withColumn("is_order", is_order("type")) \
81     .withColumn("is_return", is_return("type"))
82
83 # create kafka dataframe with invoice number, country, timestamp, total cost, total items, is order and is return
84 kafkaDF = new_df.select(["invoice_no", "country", "timestamp", "total_cost", "total_items", "is_order", "is_return"])
85 kafkaDF = kafkaDF.withColumn("total_cost", when(kafkaDF.is_order == 1, kafkaDF.total_cost).otherwise(-kafkaDF.total_cost))
86
87
88

```

At first the schema is created. Schema consist of the following:

- i. invoice\_no – String type
- ii. country – String type
- iii. timestamp – Timestamp type
- iv. type – String type
- v. items – String type

```
raw_data = lines.selectExpr("cast(value as string)").select(from_json("value",
schema).alias("temp")).select("temp.*")
```

The above line reads the data in SQL data frame format.

```
# create user-defined functions for each
total_cost = udf(lambda items: get_total_cost(items))
total_quantity = udf(lambda items: get_total_items(items))
is_order = udf(lambda types: type_order(types))
is_return = udf(lambda types: type_return(types))
```

user defined functions are created and it will return custom function outputs.

The next lines will create a new data frame with columns total\_cost, total\_quantity, is\_order, is\_return. Total cost will become negative if the order is 0 else the total cost remains as is.

```

88
89 # streaming raw data
90 query0 = kafkaDF.select(["invoice_no", "country", "timestamp", "total_cost", "total_items", "is_order", "is_return"])
91
92
93 # create time-based KPI with tumbling window of one minute
94 query1 = kafkaDF.select(["timestamp", "invoice_no", "total_cost", "is_order", "is_return"])
95 query1 = query1.withWatermark("timestamp", "1 minute").groupBy(window("timestamp", "1 minute")) \
96     .agg(round(sum("total_cost"), 2).alias("total_sales_volume"), count("invoice_no").alias("OPM"), \
97         round(sum("is_return") / (sum("is_order") + sum("is_return")), 2).alias("rate_of_return"), \
98         round(sum("total_cost") / count("invoice_no"), 2).alias("average_transaction_size"))
99
100
101 # create time-and-country based KPI with tumbling window of one minute
102 query2 = kafkaDF.select(["timestamp", "invoice_no", "country", "total_cost", "is_order", "is_return"])
103 query2 = query2.withWatermark("timestamp", "1 minute").groupBy(window("timestamp", "1 minute"), "country") \
104     .agg(round(sum("total_cost"), 2).alias("total_sales_volume"), count("invoice_no").alias("OPM"), \
105         round(sum("is_return") / (sum("is_order") + sum("is_return")), 2).alias("rate_of_return"))
106

```

The next steps explain the batch SQL data frames with proper schema are created.

Query 0 gives us the output in console.

Query 1 to create time-based KPI with tumbling window of 1 minute.

Query 2 to create time-and-country based KPI with tumbling window of 1 minute.

To calculate the KPIs, it can be done by summing total costs because the values of total\_cost are already in decent positive and negative values that were transformed earlier.

```

107
108 # write stream data to write the time-based KPIs into one minute window each
109 query0 = query0.writeStream \
110     .format("console") \
111     .outputMode("append") \
112     .option("truncate", "false") \
113     .trigger(processingTime="1 minute") \
114     .start()
115
116 #write stream data into json format
117 query1 = query1.writeStream \
118     .format("json") \
119     .outputMode("append") \
120     .option("truncate", "false") \
121     .option("path", "/user/ec2-user/real-time-project/warehouse/op1") \
122     .option("checkpointLocation", "hdfs:///user/ec2-user/real-time-project/warehouse/checkpoints1") \
123     .trigger(processingTime="1 minute") \
124     .start()
125
126 query2 = query2.writeStream \
127     .format("json") \
128     .outputMode("append") \
129     .option("truncate", "false") \
130     .option("path", "/user/ec2-user/real-time-project/warehouse/op2") \
131     .option("checkpointLocation", "hdfs:///user/ec2-user/real-time-project/warehouse/checkpoints2") \
132     .trigger(processingTime="1 minute") \
133     .start()
134
135 query0.awaitTermination()
136 query1.awaitTermination()
137 query2.awaitTermination()

```

Here the dataframes will be written. Query 0 will write in console and query 1 and query 2 will print the JSON format.

Await termination wait for the stream to finish.