

Code Logic - Retail Data Analysis

Below is the step by step explanation of the code logic and flow:

1. Imported all the required libraries for this project.

```
1 # Importing required libraries
2
3 from pyspark.sql import SparkSession
4 from pyspark.sql.functions import *
5 from pyspark.sql.types import *
```

2. Defined separate utility functions to determine total items in an order, total cost of an order and whether the order is a new buy order or a return order.
 - a. Utility function to determine total items in an order i.e. sum of the quantities of each item in the order.

```
8
9 def get_total_item_count(items):
10     """
11     This method is used to calculate the total number of items in an order.
12     :param items: The list of items in an order.
13     :return total_count: The sum of quantities of each item in the order.
14     """
15     total_count = 0
16     for item in items:
17         total_count = total_count + int(item['quantity'])
18     return total_count
```

- b. Utility function to determine total cost of an order i.e. sum of quantity multiplied by unit price for each item in the order. If the order is a new buy order, cost is returned else if the order is a return order, then cost is multiplied by -1 and then returned.

```
20 def get_total_order_cost(items, order_type):
21     """
22     This method is used to calculate the total cost of an order.
23     :param items: The list of items in an order.
24     :param order_type: This determines whether the order is a fresh one or a return.
25     :return cost: This return sum of quantity * unit_price for each item in the order. If the order is a return order, the total cost is multiplied by -1.
26     """
27     cost = 0
28     for item in items:
29         cost = cost + (int(item['quantity'])*float(item['unit_price']))
30     if order_type == 'ORDER':
31         return cost
32     else:
33         return -1*cost
```

- c. Utility function to determine the value if 'is_order' column. If the order is a new buy order, the value is set to 1 else 0.

```

35 def get_is_order(order):
36     """
37     This method is to determine whether a the order is a new buy order or not.
38     :param order: Detrmines the type of order on which the value of 'is_order' column is set.
39     :return 0/1: Return 1 if order is a new buy order else return 0.
40     """
41     if order == 'ORDER':
42         return 1
43     else:
44         return 0
45

```

- d. Utility function to determine the value if 'is_return' column. If the order is a return order, the value is set to 1 else 0.

```

46 def get_is_return(order):
47     """
48     This method is to determine whether a the order is a return order or not.
49     :param order: Detrmines the type of order on which the value of 'is_return' column is set.
50     :return 0/1: Return 1 if order is a returnn order else return 0.
51     """
52     if order == 'RETURN':
53         return 1
54     else:
55         return 0
56

```

3. Establishing the spark session and setting the spark log level to 'ERROR'.

```

57 # Setting up the spark session
58
59 spark = SparkSession \
60     .builder \
61     .appName('RetailDataAnalysis') \
62     .getOrCreate()
63
64 # Setting the spark log level to ERROR.
65
66 spark.sparkContext.setLogLevel('ERROR')
67

```

4. Connect to the given kafka server and topic using the spark session created above from where the data is to be read.

```

68 # Reading data from the given kafka server and topic
69
70 orderRaw = spark.readStream \
71     .format('kafka') \
72     .option("kafka.bootstrap.servers", "18.211.252.152:9092") \
73     .option('subscribe', 'real-time-project') \
74     .option("startingOffsets", "latest") \
75     .load()

```

5. Defining the JSON schema for the reading the data from kafka topic

```

77 # Defining the json Schema of the data to read it properly.
78
79 jsonSchema = StructType() \
80     .add("invoice_no", StringType()) \
81     .add("country", StringType()) \
82     .add("timestamp", TimestampType()) \
83     .add("type", StringType()) \
84     .add("items", ArrayType(StructType([
85         StructField("SKU", StringType()),
86         StructField("title", StringType()),
87         StructField("unit_price", StringType()),
88         StructField("quantity", StringType())
89     ])))

```

6. Casting the raw data as a string and aliasing it to data. In addition to this, defining the UDF's (user defined functions) to calculate additional columns 'total_items', 'total_cost', 'is_order' and 'is_return' using the utility functions defined in the starting.

```

90 # Converting raw data to a string and aliasing it as data.
91
92 orderStream = orderRaw.select(from_json(col('value').cast('string'), jsonSchema).alias('data')).select('data.*')
93
94 # Defining a UDF's with utility functions for calculating different column values.
95
96 # UDF to get total items in an order.
97 add_total_item_count = udf(get_total_item_count, IntegerType())
98
99 # UDF to get total cost of an order.
100 add_total_order_cost = udf(get_total_order_cost, FloatType())
101
102 # UDF to determine if an order is a new buy order.
103 add_is_order = udf(get_is_order, IntegerType())
104
105 # UDF to determine if an order is a return order.
106 add_is_return = udf(get_is_return, IntegerType())

```

7. Calculating the additional columns using UDF's defined above.

```

108 # Calculating additional columns using the UDF's defined above.
109
110 expandedOrderStream = orderStream \
111     .withColumn("total_items", add_total_item_count(orderStream.items)) \
112     .withColumn('total_cost', add_total_order_cost(orderStream.items, orderStream.type)) \
113     .withColumn('is_order', add_is_order(orderStream.type)) \
114     .withColumn('is_return', add_is_return(orderStream.type))
115

```

8. Writing the information collected for every 1 minutes time interval to the console.

```

116 # Writing the summarized information to console at a processing interval of 1 minute.
117
118 extendedOrderQuery = expandedOrderStream \
119     .select("invoice_no", "country", "timestamp", "total_cost", "total_items", "is_order", "is_return") \
120     .writeStream \
121     .outputMode("append") \
122     .format("console") \
123     .option("truncate", "false") \
124     .trigger(processingTime="1 minute") \
125     .start()
126

```

9. The next step is to calculate the time based KPI's (Key Performance Indicators) in a window of 1 minute. To calculate the below KPI's I'm grouping the data based on 1 minute time interval on timestamp column. Also, watermark is set to 1minute so that the data coming late up to 1 minute is also processed.

- a. Orders Per Minute (OPM): It is calculated as the count of total invoice numbers in that batch.
- b. Total Sale Volume: It is defined as the sum of 'total_cost' column for that batch.
- c. Rate of Return: It is defined as the average of 'is_return' column for that batch.
- d. Average Transaction Size: it is calculated as the average 'total_cost' column for that batch.

```

127 # Calculating Time Based Key Performance Indicators
128
129 aggStreamByTime = expandedOrderStream \
130     .withWatermark("timestamp", "1 minute") \
131     .groupBy(window("timestamp", "1 minute", "1 minute")) \
132     .agg(sum("total_cost").alias("total_sale_volume"),
133         count("invoice_no").alias("OPM"),
134         avg("is_return").alias('rate_of_return'),
135         avg("total_cost").alias('average_transaction_size'))
136

```

10. Following this, the next step is to write the time based KPI's output to a JSON file on HDFS.

```
137 # Writing Time Based KPI values to a json file.
138
139 queryByTime = aggStreamByTime \
140     .select("window", "OPM", 'total_sale_volume', 'rate_of_return', 'average_transaction_size') \
141     .writeStream \
142     .outputMode('append') \
143     .format('json') \
144     .option("truncate", "false") \
145     .option('path', '/output/timeBasedKPI') \
146     .option("checkpointLocation", "/output/timeBasedKPI") \
147     .trigger(processingTime="1 minute") \
148     .start()
```

11. The next step is to calculate KPI's based on time and country. To calculate these, we're grouping the data based on 1 minute interval on timestamp column and then the country column. Again, watermark is set to 1minute so that the data coming late up to 1 minute is also processed.

```
152 aggStreamByTimeCountry = expandedOrderStream \
153     .withWatermark('timestamp', '1 minute') \
154     .groupBy(window('timestamp', '1 minute', '1 minute'), 'country') \
155     .agg(sum("total_cost").alias("total_sale_volume"),
156         count("invoice_no").alias("OPM"),
157         avg("is_return").alias('rate_of_return'))
```

12. Finally, we need to write the time and country based API's to JSON files as well and then await for the termination of above three queries.

```
161 queryByTimeCountry = aggStreamByTimeCountry \
162     .select("window", 'country', "OPM", 'total_sale_volume', 'rate_of_return') \
163     .writeStream \
164     .outputMode('append') \
165     .format('json') \
166     .option('truncate', 'false') \
167     .option('path', '/output/timeCountryBasedKPI') \
168     .option("checkpointLocation", "/output/timeCountryBasedKPI") \
169     .trigger(processingTime = '1 minute') \
170     .start()
171
172 # Awaiting termination of the above written queries.
173
174 extendedOrderQuery.awaitTermination()
175 queryByTime.awaitTermination()
176 queryByTimeCountry.awaitTermination()
177
```

13. The next step is to run the code. The below command is to be executed on spark EMR cluster.

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

