

Code Logic - Retail Data Analysis

The below Python Script is used to process the input data streams into the resultant JSON files

spark-streaming.py

Setting up the libraries and modules

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

The following functions are used to calculate the various attributes

1. Total Cost UDF - Total cost of an order arrived at by summing up the cost of all products in that invoice. The return cost is treated as a loss. Thus, for return orders, this value will be negative. To calculate the income from sale of each product, the unit price of the product is multiplied with the quantity of the product purchased.

```
def calculate_total_cost(items, trn_type):
    if items is not None:
        total_cost = 0
        item_price = 0
        for item in items:
            item_price = (item['quantity'] * item['unit_price'])
            total_cost = total_cost + item_price
            item_price = 0

        if trn_type == "RETURN":
            return total_cost * -1
        else:
            return total_cost
```

2. Total Items UDF - Total number of items present in an order. To calculate the number of products in every invoice the quantity ordered of each product in that invoice is summed up

```
def calculate_total_item_count(items):
    if items is not None:
        total_count = 0
        for item in items:
            total_count = total_count + item['quantity']
        return total_count
```

3. Is Order UDF - This flag denotes whether an order is a new order or not. If this invoice is for a return order, the value should be 0. To determine if invoice is for an order or not an if-else statement is used

```
def flag_isOrder(trn_type):
    if trn_type == "ORDER":
        return(1)
    else:
        return(0)
```

4. Is Return UDF - This flag denotes whether an order is a return order or not. If this invoice is for a new sales order, the value should be 0. To determine if invoice is for a return or not an if-else statement is used

```
def flag_isReturn(trn_type):
    if trn_type == "RETURN":
        return(1)
    else:
        return(0)
```

Initialising the Spark session and setting the log level to error

```
spark = SparkSession \
    .builder \
    .appName("spark-streaming") \
    .getOrCreate()
spark.sparkContext.setLogLevel('ERROR')
```

Reading input data from Kafka using sparkreadStream in 'retail corp orderData' mentioning the details of the Kafka broker, such as bootstrap server, port and topic name

```
retail_corp_orderData = spark.readStream \
    .format("kafka") \
    .option("kafka.bootstrap.servers", "18.211.252.152:9092") \
    .option("startingOffsets", "earliest") \
    .option("failOnDataLoss", "false") \
    .option("subscribe", "real-time-project") \
    .load()
```

Defining JSON schema of each order, using appropriate datatypes and StructField in the case of item attribute

```
jsonSchema = StructType() \
    .add("invoice_no", LongType()) \
    .add("country", StringType()) \
    .add("timestamp", TimestampType()) \
    .add("type", StringType()) \
    .add("items", ArrayType(StructType([
        StructField("SKU", StringType()),
        StructField("title", StringType()),
        StructField("unit_price", FloatType()),
        StructField("quantity", IntegerType()),
    ])))
```

Reading the raw JSON data from Kafka as 'order stream' by casting it to string and storing it with the alias 'data'

```
orderStream = retail_corp_orderData.select(from_json(col("value").cast("string"),
jsonSchema).alias("data")).select("data.*")
```

Defining the UDFs by Converting the Python functions defined earlier, and assigning the appropriate return datatype

```
sum_total_order_cost = udf(calculate_total_cost, FloatType())
sum_total_item_count = udf(calculate_total_item_count, IntegerType())
sum_isOrder = udf(flag_isOrder, IntegerType())
sum_isReturn = udf(flag_isReturn, IntegerType())
```

Calculating the additional columns according to the required input values

```
updatedOrderStream = orderStream \
    .withColumn("total_cost", sum_total_order_cost(orderStream.items, orderStream.type)) \
    .withColumn("total_items", sum_total_item_count(orderStream.items)) \
    .withColumn("is_order", sum_isOrder(orderStream.type)) \
    .withColumn("is_return", sum_isReturn(orderStream.type))
```

Writing the summarized input values to console, using 'append' output method and applying truncate as false and setting the processing time to 1 minute

```
extendedOrderQuery = updatedOrderStream \
    .select("invoice_no", "country", "timestamp", "total_cost", "total_items", "is_order", "is_return") \
    .writeStream \
    .outputMode("append") \
```

```
.format("console") \
.option("truncate", "false") \
.trigger(processingTime = "1 minute") \
.start()
```

Calculating time-based KPIs (Total sale volume, OPM, Rate of return, Average transaction size) having tumbling window of one minute and watermark of one minute.

```
timebasedKPI = updatedOrderStream \
.withWatermark("timestamp", "1 minute") \
.groupBy(window("timestamp", "1 minute", "1 minute")) \
.agg(sum("total_cost").alias("total_sale_volume"),
     count("invoice_no").alias("OPM"),
     avg("is_return").alias("rate_of_return"),
     avg("total_cost").alias("average_transaction_size")
    ) \
.select("window", "OPM", "total_sale_volume", "average_transaction_size", "rate_of_return" )
```

Writing the time-based KPIs data to HDFS - HDFS into JSON files for each one-minute window, using 'append' output mode, setting truncate as false, and specifying the HDFS output path for both the KPI files and for their checkpoints. Ten 1-minute window batches were taken.

```
queryByTime = timebasedKPI.writeStream \
.format("json") \
.outputMode("append") \
.option("truncate", "false") \
.option("path", "/user/ec2-user/time_kpi") \
.option("checkpointLocation", "/user/ec2-user/time_kpi_checkpoints") \
.trigger(processingTime="1 minute") \
.start()
```

Calculating time-and-country-based KPIs (Total sale volume, OPM, Rate of return) having tumbling window of one minute and watermark of one minute. Here the data is grouped by window and country both

```
timecountrybasedKPI = updatedOrderStream \
.withWatermark("timestamp", "1 minute") \
.groupBy(window("timestamp", "1 minute", "1 minute"), "country") \
.agg(sum("total_cost").alias("total_sale_volume"),
     count("invoice_no").alias("OPM"),
     avg("is_return").alias("rate_of_return")) \
.select("window", "country", "OPM", "total_sale_volume", "rate_of_return" )
```

Writing the time-and-country-based KPIs data to HDFS into JSON files for each one-minute window, using 'append' output mode, setting truncate as false, and specifying the HDFS output path for both the KPI files and for their checkpoints. Ten 1-minute window batches were taken.

```
queryByCountry = timecountrybasedKPI.writeStream \  
  .format("json") \  
  .outputMode("append") \  
  .option("truncate","false") \  
  .option("path","/user/ec2-user/country_kpi") \  
  .option("checkpointLocation","/user/ec2-user/country_kpi_checkpoints") \  
  .trigger(processingTime="1 minute") \  
  .start()
```

Indicating Spark to await termination

```
extendedOrderQuery.awaitTermination()  
queryByCountry.awaitTermination()  
queryByTime.awaitTermination()
```

Console Command

To execute this project, I have used an EMR cluster along with an EC2 Kafka Instance.

Step I – Creation and initiation of the EC2 Kafka Instance

The steps to set up the EC2 instance are as follows

- i. From the AWS Management Console, selecting the EC2 service and under Images selecting AMIs, from there selecting 'Public Images'. Select the AMI ID 'ami-0b47786b148f4dd16'
- ii. Selecting the configure security group details and adding ports for Kafka module
- iii. Reviewing and launching the instance
- iv. Selecting the Network and Security and creating Elastic IP and associate it with the EC2 Kafka Instance
- v. Logging in the EC2 instance as 'ec2-user' using Putty terminal and changing the directory to /home/ec2-user/downloads using 'cd' command
- vi. Directing to the kafka directory using `cd /home/ec2-user/downloads/kafka_2.12-2.3.0` and to go inside the config directory using `cd config/`
- vii. To make changes to the server.properties file, edit this file enter the command `vi server.properties`. In place of your.host.name enter the IPv4 Public IP of your EC2 instance.

- viii. To the Kafka directory using the `cd kafka_2.12-2.3.0/` command and then start the Zookeeper using `bin/zookeeper-server-start.sh config/zookeeper.properties`
- ix. In the new terminal go to the Kafka directory using the `cd kafka_2.12-2.3.0/` command and Start the Kafka server using `bin/kafka-server-start.sh config/server.properties`
- x. The producer is started using `bin/kafka-console-producer.sh --broker-list 18.211.252.152:9092 --topic real-time-project`

Step II – Creation and initiation of the EMR Cluster

- i. Creation of python script file using vi `spark-streaming.py` and adding put code in that
- ii. Exporting Spark Kafka version using `export SPARK_KAFKA_VERSION=0.10`
- iii. Hitting the below spark submit command to store the batches in a file

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

```
hadoop@ip-172-31-40-239:~$
22/08/13 18:58:12 INFO SharedState: Setting hive.metastore.warehouse.dir ('null') to the value of spark.sql.warehouse.dir ('hdfs:///user/spark/warehouse').
22/08/13 18:58:12 INFO SharedState: Warehouse path is 'hdfs:///user/spark/warehouse'.
22/08/13 18:58:12 INFO JettyUtils: Adding filter org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter to /SQL.
22/08/13 18:58:12 INFO JettyUtils: Adding filter org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter to /SQL/json.
22/08/13 18:58:12 INFO JettyUtils: Adding filter org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter to /SQL/execution.
22/08/13 18:58:12 INFO JettyUtils: Adding filter org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter to /SQL/execution/json.
22/08/13 18:58:12 INFO JettyUtils: Adding filter org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter to /static/sql.
22/08/13 18:58:13 INFO StateStoreCoordinatorRef: Registered StateStoreCoordinator endpoint
-----
Batch: 0
-----
+-----+-----+-----+-----+-----+-----+-----+
|invoice_no|country|timestamp|total_cost|total_items|is_order|is_return|
+-----+-----+-----+-----+-----+-----+-----+
|154132550699322|United Kingdom|2022-08-11 09:42:15|4.95|136|1|0|
|154132550699323|United Kingdom|2022-08-11 09:42:15|9.9|111|1|0|
|154132550699324|United Kingdom|2022-08-11 09:42:50|4.56|124|1|0|
|154132550699325|United Kingdom|2022-08-11 09:42:52|19.92|124|1|0|
|154132550699326|United Kingdom|2022-08-11 09:42:56|34.800003|138|1|0|
|154132550699327|United Kingdom|2022-08-11 09:43:13|17.400002|112|1|0|
|154132550699328|United Kingdom|2022-08-11 09:43:16|3.3|114|1|0|
|154132550699329|United Kingdom|2022-08-11 09:43:17|7.38|118|1|0|
|154132550699330|EIRE|2022-08-11 09:43:21|16.6|14|1|0|
|154132550699331|United Kingdom|2022-08-11 09:43:23|22.5|16|1|0|
|154132550699332|United Kingdom|2022-08-11 09:43:27|16.5|126|1|0|
|154132550699333|United Kingdom|2022-08-11 09:43:38|3.48|138|1|0|
|154132550699334|United Kingdom|2022-08-11 09:43:43|19.8|114|1|0|
|154132550699335|United Kingdom|2022-08-11 09:43:57|45.0|151|1|0|
|154132550699336|United Kingdom|2022-08-11 09:44:03|0.55|14|1|0|
|154132550699337|United Kingdom|2022-08-11 09:44:11|16.5|113|1|0|
|154132550699338|United Kingdom|2022-08-11 09:44:13|52.64|116|1|0|
|154132550699339|United Kingdom|2022-08-11 09:44:14|15.0|116|1|0|
|154132550699340|United Kingdom|2022-08-11 09:44:19|0.84|110|1|0|
|154132550699341|United Kingdom|2022-08-11 09:44:19|11.9|169|1|0|
+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows
```

iv. Checking the hdfs locations for both time and time-country based KPIs

For time based KPI

`hadoop fs -ls /user/ec2-user/time_kpi`

Using the cat command opening one of the json file for time KPI

`hadoop fs -cat /user/ec2-user/time_kpi/part-00004-dcac6af3-becf-4430-9f46-558c2d5dd9eb-c000.json`

```

-rw-r--r-- 1 hadoop hadoop 199 2022-08-13 19:14 /user/ec2-user/time_kpi/part-00004-dcac6af3-becf-4430-9f46-558c2d5dd9eb-c000.
json
[hadoop@ip-172-31-40-239 ~]$ hadoop fs -cat /user/ec2-user/time_kpi/part-00004-dcac6af3-becf-4430-9f46-558c2d5dd9eb-c000.json
{"window":{"start":"2022-08-13T19:06:00.000Z","end":"2022-08-13T19:07:00.000Z"},"OPM":16,"total_sale_volume":169.73000061511993,"avera
ge_transaction_size":10.608125038444996,"rate_of_return":0.0625}
[hadoop@ip-172-31-40-239 ~]$ hadoop fs -cat /user/ec2-user/time_kpi/part-00197-9883bf9c-f368-4519-ba78-11712b1b0320-c000.json
{"window":{"start":"2022-08-11T19:52:00.000Z","end":"2022-08-11T19:53:00.000Z"},"OPM":12,"total_sale_volume":145.2400015592575,"avera
ge_transaction_size":12.10333346327146,"rate_of_return":0.0}
{"window":{"start":"2022-08-12T20:29:00.000Z","end":"2022-08-12T20:30:00.000Z"},"OPM":13,"total_sale_volume":161.77999639511108,"avera
ge_transaction_size":12.444615107316237,"rate_of_return":0.07692307692307693}
{"window":{"start":"2022-08-12T08:03:00.000Z","end":"2022-08-12T08:04:00.000Z"},"OPM":12,"total_sale_volume":111.8899998664856,"avera
ge_transaction_size":9.324166655540466,"rate_of_return":0.0}
{"window":{"start":"2022-08-13T10:27:00.000Z","end":"2022-08-13T10:28:00.000Z"},"OPM":8,"total_sale_volume":124.45999801158905,"avera
ge_transaction_size":15.557499751448631,"rate_of_return":0.0}
{"window":{"start":"2022-08-13T16:00:00.000Z","end":"2022-08-13T16:01:00.000Z"},"OPM":13,"total_sale_volume":345.30000257492065,"avera
ge_transaction_size":26.56153865960928,"rate_of_return":0.07692307692307693}
{"window":{"start":"2022-08-13T04:45:00.000Z","end":"2022-08-13T04:46:00.000Z"},"OPM":11,"total_sale_volume":123.20999896526337,"avera
ge_transaction_size":11.200908996842124,"rate_of_return":0.0}
{"window":{"start":"2022-08-13T14:20:00.000Z","end":"2022-08-13T14:21:00.000Z"},"OPM":5,"total_sale_volume":86.9200029373169,"average_
transaction_size":17.38400058746338,"rate_of_return":0.2}
{"window":{"start":"2022-08-11T12:32:00.000Z","end":"2022-08-11T12:33:00.000Z"},"OPM":4,"total_sale_volume":53.29999923706055,"average
transaction_size":13.324999809265137,"rate_of_return":0.0}
{"window":{"start":"2022-08-13T07:48:00.000Z","end":"2022-08-13T07:49:00.000Z"},"OPM":11,"total_sale_volume":102.34000039100647,"avera
ge_transaction_size":9.303636399182407,"rate_of_return":0.0}
{"window":{"start":"2022-08-13T09:51:00.000Z","end":"2022-08-13T09:52:00.000Z"},"OPM":11,"total_sale_volume":168.9899983406067,"avera
ge_transaction_size":15.362727121873336,"rate_of_return":0.09090909090909091}
{"window":{"start":"2022-08-11T12:40:00.000Z","end":"2022-08-11T12:41:00.000Z"},"OPM":6,"total_sale_volume":66.56000077724457,"average
transaction_size":11.093333462874094,"rate_of_return":0.0}
{"window":{"start":"2022-08-13T10:12:00.000Z","end":"2022-08-13T10:13:00.000Z"},"OPM":16,"total_sale_volume":324.5100042819977,"avera
ge_transaction_size":20.281875267624855,"rate_of_return":0.0}
{"window":{"start":"2022-08-12T13:24:00.000Z","end":"2022-08-12T13:25:00.000Z"},"OPM":8,"total_sale_volume":25.180001497268677,"avera
ge_transaction_size":3.1475001871585846,"rate_of_return":0.25}
{"window":{"start":"2022-08-11T21:48:00.000Z","end":"2022-08-11T21:49:00.000Z"},"OPM":9,"total_sale_volume":112.89000499248505,"avera
ge_transaction_size":12.543333888053894,"rate_of_return":0.11111111111111111}
{"window":{"start":"2022-08-11T12:42:00.000Z","end":"2022-08-11T12:43:00.000Z"},"OPM":6,"total_sale_volume":244.30000352859497,"avera
ge_transaction_size":40.71666725476583,"rate_of_return":0.16666666666666666}
{"window":{"start":"2022-08-11T11:51:00.000Z","end":"2022-08-11T11:52:00.000Z"},"OPM":9,"total_sale_volume":209.2899982906342,"avera
ge_transaction_size":23.254443314340378,"rate_of_return":0.0}
[hadoop@ip-172-31-40-239 ~]$

```

For time-country based KPI

`hadoop fs -ls /user/ec2-user/country_kpi`

Using the cat command opening for one of the json file for a country KPI

`hadoop fs -cat /user/ec2-user/country_kpi/part-00004-9b289b0c-7588-4db9-a936-396897670239-c000.json`


```
hadoop@ip-172-31-40-239:~$
{"window":{"start":"2022-08-12T11:27:00.000Z","end":"2022-08-12T11:28:00.000Z"},"country":"EIRE","OPM":1,"total_sale_volume":11.850000
381469727,"rate_of_return":0.0}
{"window":{"start":"2022-08-11T11:25:00.000Z","end":"2022-08-11T11:26:00.000Z"},"country":"United Kingdom","OPM":11,"total_sale_volume
":69.05999946594238,"rate_of_return":0.0}
{"window":{"start":"2022-08-12T03:09:00.000Z","end":"2022-08-12T03:10:00.000Z"},"country":"Germany","OPM":1,"total_sale_volume":2.4900
00009536743,"rate_of_return":0.0}
{"window":{"start":"2022-08-12T08:44:00.000Z","end":"2022-08-12T08:45:00.000Z"},"country":"Netherlands","OPM":1,"total_sale_volume":4.
159999847412109,"rate_of_return":0.0}
{"window":{"start":"2022-08-12T00:43:00.000Z","end":"2022-08-12T00:44:00.000Z"},"country":"Germany","OPM":1,"total_sale_volume":10.079
999923706055,"rate_of_return":0.0}
{"window":{"start":"2022-08-11T22:27:00.000Z","end":"2022-08-11T22:28:00.000Z"},"country":"United Kingdom","OPM":9,"total_sale_volume"
":110.23999774456024,"rate_of_return":0.1111111111111111}
{"window":{"start":"2022-08-13T10:02:00.000Z","end":"2022-08-13T10:03:00.000Z"},"country":"EIRE","OPM":1,"total_sale_volume":6.75,"rat
e_of_return":0.0}
{"window":{"start":"2022-08-13T15:40:00.000Z","end":"2022-08-13T15:41:00.000Z"},"country":"United Kingdom","OPM":5,"total_sale_volume"
":83.75000166893005,"rate_of_return":0.0}
{"window":{"start":"2022-08-13T10:11:00.000Z","end":"2022-08-13T10:12:00.000Z"},"country":"United Kingdom","OPM":9,"total_sale_volume"
":73.55000150203705,"rate_of_return":0.1111111111111111}
{"window":{"start":"2022-08-12T02:28:00.000Z","end":"2022-08-12T02:29:00.000Z"},"country":"Hong Kong","OPM":1,"total_sale_volume":5.09
999904632568,"rate_of_return":0.0}
{"window":{"start":"2022-08-13T00:27:00.000Z","end":"2022-08-13T00:28:00.000Z"},"country":"United Kingdom","OPM":6,"total_sale_volume"
":-63.17000770568848,"rate_of_return":0.3333333333333333}
{"window":{"start":"2022-08-12T17:12:00.000Z","end":"2022-08-12T17:13:00.000Z"},"country":"United Kingdom","OPM":5,"total_sale_volume"
":405.4799699783325,"rate_of_return":0.0}
{"window":{"start":"2022-08-12T17:24:00.000Z","end":"2022-08-12T17:25:00.000Z"},"country":"EIRE","OPM":1,"total_sale_volume":19.899999
618530273,"rate_of_return":0.0}
{"window":{"start":"2022-08-12T00:17:00.000Z","end":"2022-08-12T00:18:00.000Z"},"country":"United Kingdom","OPM":13,"total_sale_volume"
":166.170001745224,"rate_of_return":0.0}
{"window":{"start":"2022-08-12T05:55:00.000Z","end":"2022-08-12T05:56:00.000Z"},"country":"France","OPM":1,"total_sale_volume":2.5,"ra
te_of_return":0.0}
{"window":{"start":"2022-08-13T02:43:00.000Z","end":"2022-08-13T02:44:00.000Z"},"country":"Germany","OPM":1,"total_sale_volume":19.559
999465942383,"rate_of_return":0.0}
{"window":{"start":"2022-08-13T10:16:00.000Z","end":"2022-08-13T10:17:00.000Z"},"country":"United Kingdom","OPM":15,"total_sale_volume"
":101.48000064492226,"rate_of_return":0.06666666666666667}
{"window":{"start":"2022-08-13T12:58:00.000Z","end":"2022-08-13T12:59:00.000Z"},"country":"United Kingdom","OPM":9,"total_sale_volume"
":182.19000148773193,"rate_of_return":0.0}
{"window":{"start":"2022-08-11T12:45:00.000Z","end":"2022-08-11T12:46:00.000Z"},"country":"EIRE","OPM":1,"total_sale_volume":4.1599998
47412109,"rate_of_return":0.0}
hadoop@ip-172-31-40-239:~$
```

Transferring files from hadoop using get command and then using WinSCP transferring it to local machine

Command for getting time based KPI

`hadoop fs -get /user/ec2-user/time_kpi`

Command for getting time based KPI

`hadoop fs -get /user/ec2-user/country_kpi`