

Code Logic - Retail Data Analysis

Submitted By:

Manish Pramanik (pramanik.manish@gmail.com)
Nandanavanam Dinakar (itsdinakar98@gmail.com)
Abhay Kumar (abhaykumariiit@gmail.com)

Set Up:

Initiate an AWS EMR cluster with Hadoop and spark service.

Environment set up to integrate kafka with spark job,

- export SPARK_KAFKA_VERSION=0.10

Execute Stream Job:

Create the script file **spark-streaming.py** to to run a Spark Streaming job where you are using connecting with a Kafka stream

Command used for run spark streaming job:

- spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.5 spark-streaming.py | tee Console-output

spark-submit to execute spark job, --package to download dependency package required to run job and | tee ConsoleOutput would save printed streaming data on console to ConsoleOutput file

Code Logic of Script **spark-streaming.py** :

1. Import required pyspark packages
2. Initiate spark session
3. Define bootstrap server and kafka-topic provided, from which we would read stream retail data
4. **orderRawStream** : To read raw stream from defined bootstrap server and kafka-topic

5. **orderSchema** : Define Spark schema to parse json format stream from raw stream of data
6. **orderStream** : Parsing of stream of json format of stream using orderRawStream schema.
7. To obtain derived fields, utility functions are required and defined. The functions are:
 - is_order: flag 1 if order type is ORDER else 0
 - is_return: flag 1 if order type is RETURN else 0
 - total_item_count: to count the number of items in each order
 - total_cost: total cost which computed from item quantity * item price for each order. Taken as negative if order type is RETURN.
8. Define User Defined function using UDFs using above created utility functions.
9. **expandedOrderStream**: Define stream with required derived columns. The stream constitutes base columns invoice_no, country, timestamp and derived columns total_cost, total_items, is_order, is_return
10. **orderStreamConsole**: Generate stream and print in console. Stream generated in append mode with 1 minute process time.
11. **aggStreamByTime**: Generate stream on computing aggregation to obtain Key Performance Indicators(KPI) based on window of time.
An event time columns taken for window time. Here timestamp column is event time columns. Aggregation of total_sale_volume, average_transaction_size, rate_of_return calculated for windows of 1 minute and 1 minute sliding interval. To aggregate based on tumbling window, window and sliding interval need to be kept equal which is here 1 minute. Required interval 10 minutes, thus processing time kept 10 minutes.

The stream executed in json mode and stored in hdfs location ../Timebased-KPI

To get KPI stored in json files from hdfs to local storage run

➤ `hadoop fs -get /user/hadoop/Timebased-KPI ./`

12. **aggStreamByTimeCountry**: Country and Time based KPI computed. Here, grouping is done on window of time as well as additionally 'country' column. The stream in json format to be saved on hdfs location ../Country-and-timebased-KPI

To get KPI stored in json files from hdfs to local storage run

➤ `hadoop fs -get /user/hadoop/Country-and-timebased-KPI ./`

Console Output:

```
hadoop@ip-172-31-6-40:~/retail_analysis
22/12/11 14:28:18 INFO JettyUtils: Adding filter org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter to /SQL/execution.
22/12/11 14:28:18 INFO JettyUtils: Adding filter org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter to /SQL/execution/json.
22/12/11 14:28:18 INFO JettyUtils: Adding filter org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter to /static/sql.
22/12/11 14:28:19 INFO StateStoreCoordinatorRef: Registered StateStoreCoordinator endpoint

Batch: 0
-----
|invoice_no|country|timestamp|total_cost|total_items|is_order|is_return|
|-----|-----|-----|-----|-----|-----|-----|
-----

Batch: 1
-----
|invoice_no|country|timestamp|total_cost|total_items|is_order|is_return|
|-----|-----|-----|-----|-----|-----|-----|
|154132552491375|United Kingdom|2022-12-11 14:22:22|19.130001|13|1|0|
|154132552491376|United Kingdom|2022-12-11 14:22:36|30.0|24|1|0|
|154132552491377|United Kingdom|2022-12-11 14:22:37|64.51|23|1|0|
|154132552491378|United Kingdom|2022-12-11 14:22:39|140.23|57|1|0|
|154132552491379|United Kingdom|2022-12-11 14:22:43|334.97998|64|1|0|
|154132552491380|United Kingdom|2022-12-11 14:22:57|5.04|12|1|0|
|154132552491381|United Kingdom|2022-12-11 14:23:06|35.51|21|1|0|
|154132552491382|United Kingdom|2022-12-11 14:23:09|41.300003|10|1|0|
|154132552491383|United Kingdom|2022-12-11 14:23:16|524.32|206|1|0|
|-----|-----|-----|-----|-----|-----|-----|

Batch: 2
-----
|invoice_no|country|timestamp|total_cost|total_items|is_order|is_return|
|-----|-----|-----|-----|-----|-----|-----|
|154132552491384|United Kingdom|2022-12-11 14:23:21|68.1|30|1|0|
|154132552491385|United Kingdom|2022-12-11 14:23:25|16.75|3|1|0|
|154132552491386|United Kingdom|2022-12-11 14:23:27|2.46|1|1|0|
|154132552491387|United Kingdom|2022-12-11 14:23:28|31.82|26|1|0|
|154132552491388|United Kingdom|2022-12-11 14:23:30|618.16|104|1|0|
|154132552491389|United Kingdom|2022-12-11 14:23:38|119.52|46|1|0|
|154132552491390|United Kingdom|2022-12-11 14:23:49|396.21|165|1|0|
|154132552491391|United Kingdom|2022-12-11 14:23:53|166.27|108|1|0|
|154132552491392|United Kingdom|2022-12-11 14:23:57|17.849998|3|1|0|
|154132552491393|United Kingdom|2022-12-11 14:23:59|30.76|4|1|0|
|154132552491394|United Kingdom|2022-12-11 14:24:01|23.13|14|1|0|
|154132552491395|United Kingdom|2022-12-11 14:24:02|70.24|16|1|0|
|154132552491396|United Kingdom|2022-12-11 14:24:06|29.99|33|1|0|
|154132552491397|United Kingdom|2022-12-11 14:24:11|18.82|11|1|0|
|-----|-----|-----|-----|-----|-----|-----|
```

Retrieval of generated data from local EMR instance:

zip -r Timebased-KPI.zip Timebased-KPI

zip -r Country-and-timebased-KPI.zip Country-and-timebased-KPI

Timebased-KPI.zip and Country-and-timebased-KPI.zip obtained using data transfer tool winscp.