

Code Logic - Retail Data Analysis

Spark Streaming Code Implementation:

1. Import the necessary **PySpark** modules such as **SparkSession** class, all the functions and data types from **pyspark.sql.functions** and **pyspark.sql.types** modules respectively.

```
# Import necessary libraries
from pyspark.sql import SparkSession
from pyspark.sql.functions import *
from pyspark.sql.types import *
```

2. Create a Spark Session which is an entry point to interact with Spark application.

```
# Create a spark session
spark = SparkSession \
    .builder \
    .appName("Retail Data Analysis - Spark Streaming") \
    .getOrCreate()

# Set log level to ERROR
spark.sparkContext.setLogLevel('ERROR')
```

3. Connect to Kafka server and read the retail data from it (raw data is in JSON format), define proper schema or data structure to convert it to structured data frame.

```
# Read Retail Dataset from Kafka topic : real-time-project
retail_data_source = spark \
    .readStream \
    .format("kafka") \
    .option("kafka.bootstrap.servers", "18.211.252.152:9092") \
    .option("subscribe", "real-time-project") \
    .load()

# Defining the Data Structure for source retail data (json format).
retail_datastructure = StructType([
    StructField("invoice_no", LongType()),
    StructField("country", StringType()),
    StructField("timestamp", TimestampType()),
    StructField("type", StringType()),
    StructField("items", ArrayType(StructType([
        StructField("SKU", StringType()),
        StructField("title", StringType()),
        StructField("unit_price", FloatType()),
        StructField("quantity", IntegerType())
    ])))
])

# Converting source data from json format to dataframe (table).
retail_data_formatted = retail_data_source \
    .select(from_json(col("value").cast("string"), retail_datastructure).alias("retail_data")) \
    .select("retail_data.*")
```

4. Define the utility functions to calculate the following metrics:

- a. **Total Cost:**

Calculate the total cost of all items present in each invoice. Cost for each items list can be calculated by multiplying quantity with unit price. If the invoice type is

“RETURN” then we convert the total cost to negative value which indicates loss of amount.

```
def calculate_total_cost(items, type):
    """
    Calculate total cost associated for each invoice.
    total_cost = sum of (unit_price * quantity) for all items.
    If the type is Return then the total cost will be negative which represents amount loss.
    """
    total_cost = 0
    for item in items:
        total_cost = total_cost + (item[2] * item[3])
    if type == 'RETURN':
        return total_cost * -1
    return total_cost
```

b. Total Items:

Calculate the total quantity for all the items list in an invoice be it an “ORDER” or “RETURN” type.

```
def calculate_total_items(items):
    """
    Calculate total no. of items present in each invoice (quantity).
    """
    total_items = 0
    for item in items:
        total_items = total_items + item[3]
    return total_items
```

c. Check “ORDER” type:

This function is to check whether the type is ORDER or not. If it is ORDER then we assign the flag as 1, else 0.

```
def check_order_type(type):
    """
    To check whether the type is Order or not
    1 - ORDER Type
    0 - Not an Order Type
    """
    if type == 'ORDER':
        return 1
    return 0
```

d. Check “RETURN” type:

This function is to check whether the type is RETURN or not. If it is RETURN then we assign the flag as 1, else 0.

```
def check_return_type(type):
    '''
        To check whether the type is Return or not
        1 - Return Type
        0 - Not a Return Type
    '''
    if type == 'RETURN':
        return 1
    return 0
```

5. Convert the utility functions to UDFs which can then be used to derive the metrics in further transformations.

```
# Converting utility functions to UDFs
get_total_cost = udf(calculate_total_cost, DoubleType())
get_total_items = udf(calculate_total_items, IntegerType())
get_order_flag = udf(check_order_type, IntegerType())
get_return_flag = udf(check_return_type, IntegerType())
```

6. Derive the attributes/metrics such as total_cost, total_items, is_order and is_return based on the source data using the UDFs defined and select the following columns as final data frame and write the output to the console:
 - a. invoice_no
 - b. country
 - c. timestamp
 - d. total_cost
 - e. total_items
 - f. is_order
 - g. is_return

```
# Derive the columns/metrics: total_cost, total_items, is_order, is_return from the source dataframe and select the required columns
retail_data_transformed = retail_data_formatted \
    .withColumn("total_cost", get_total_cost(retail_data_formatted.items, retail_data_formatted.type)) \
    .withColumn("total_items", get_total_items(retail_data_formatted.items)) \
    .withColumn("is_order", get_order_flag(retail_data_formatted.type)) \
    .withColumn("is_return", get_return_flag(retail_data_formatted.type)) \
    .select("invoice_no", "country", "timestamp", "total_cost", "total_items", "is_order", "is_return")

# Write the retail transformed data to console
retail_data_sink = retail_data_transformed \
    .writeStream \
    .outputMode("append") \
    .format("console") \
    .option("truncate", "false") \
    .trigger(processingTime="1 minute") \
    .start()
```

7. Calculate the following KPIs based on timestamp with one minute window range and one minute sliding interval:
 - a. **Orders Per Minute (OPM)**
 - b. **total_sale_volume:** Total amount for each minute (window range) – overall sales
 - c. **average_transaction_size:** Average amount per minute – overall sales
 - d. **rate_of_return:** Average number of returned items per minute.

Write the Time-based KPI output in HDFS (Path: /user/hadoop/Timebased-KPI) in the form of JSON files.

```
# Calculate time based KPIs: OPM (Orders Per Minute), total_sale_volume, average_transaction_size, rate_of_return.
retail_time_kpi = retail_data_transformed \
    .withWatermark("timestamp", "1 minute") \
    .groupBy(window("timestamp", "1 minute", "1 minute")) \
    .agg(count("invoice_no").alias("OPM"),
         sum("total_cost").alias("total_sale_volume"),
         avg("total_cost").alias("average_transaction_size"),
         avg("is_return").alias("rate_of_return")) \
    .select("window",
            "OPM",
            "total_sale_volume",
            "average_transaction_size",
            "rate_of_return")

# Write retail data time based KPI in HDFS (JSON format)
retail_time_kpi_sink = retail_time_kpi \
    .writeStream \
    .format("json") \
    .outputMode("append") \
    .option("truncate", "false") \
    .option("path", "Timebased-KPI") \
    .option("checkpointLocation", "Timebased-Checkpoint") \
    .trigger(processingTime="1 minute") \
    .start()
```

8. Calculate the same above KPIs except average_transaction_size but based on time and country wise with one minute window range and one minute sliding interval. Write Time and Country based KPI output in HDFS (Path: /user/hadoop/Country-and-timebased-KPI) in the form of JSON files.

```
# Calculate time-country based KPIs: OPM (Orders Per Minute), total_sale_volume, rate_of_return.
retail_time_country_kpi = retail_data_transformed \
    .withWatermark("timestamp", "1 minute") \
    .groupBy(window("timestamp", "1 minute", "1 minute"), "country") \
    .agg(count("invoice_no").alias("OPM"),
         sum("total_cost").alias("total_sale_volume"),
         avg("is_return").alias("rate_of_return")) \
    .select("window",
            "country",
            "OPM",
            "total_sale_volume",
            "rate_of_return")

# Write retail data time and country based KPI in HDFS (JSON format)
retail_time_country_kpi_sink = retail_time_country_kpi \
    .writeStream \
    .format("json") \
    .outputMode("append") \
    .option("truncate", "false") \
    .option("path", "Country-and-timebased-KPI") \
    .option("checkpointLocation", "Country-and-timebased-Checkpoint") \
    .trigger(processingTime="1 minute") \
    .start()
```

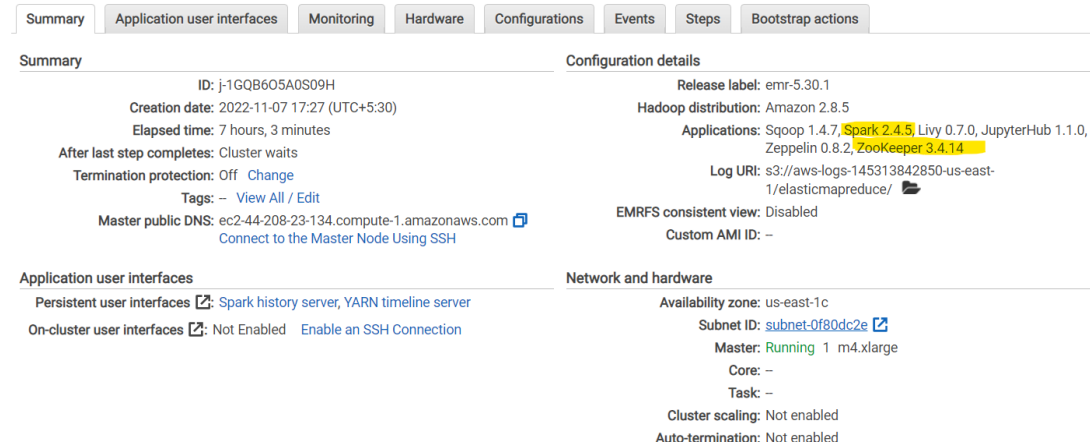
9. Keep all the streams alive until there is an interruption.

```
# Keep all the streams alive until it is terminated
retail_data_sink.awaitTermination()
retail_time_kpi_sink.awaitTermination()
retail_time_country_kpi_sink.awaitTermination()
```

Code Deployment and Execution Steps:

1. Create an EMR Instance with **Spark** and **ZooKeeper** applications installed in it.

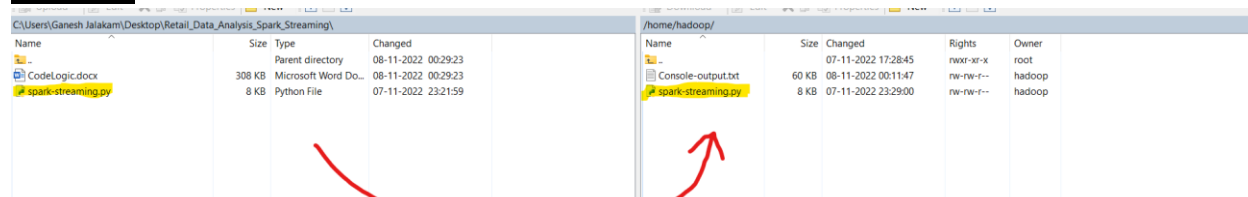
Cluster: Retail-Data-Analysis-SparkStreaming Waiting Cluster ready after last step completed.



The screenshot shows the AWS EMR console for a cluster named 'Retail-Data-Analysis-SparkStreaming'. The cluster is in a 'Waiting' state, indicating it is ready after the last step completed. The console displays various tabs: Summary, Application user interfaces, Monitoring, Hardware, Configurations, Events, Steps, and Bootstrap actions. The 'Summary' tab is active, showing details such as the cluster ID (j-1GQB605A0S09H), creation date (2022-11-07 17:27 UTC+5:30), and elapsed time (7 hours, 3 minutes). It also lists the applications installed: Sqoop 1.4.7, Spark 2.4.5, Livy 0.7.0, JupyterHub 1.1.0, Zeppelin 0.8.2, and ZooKeeper 3.4.14. The 'Configuration details' tab shows the release label (emr-5.30.1), Hadoop distribution (Amazon 2.8.5), and the EMRFS consistent view (Disabled). The 'Application user interfaces' tab shows persistent user interfaces for Spark history server and YARN timeline server. The 'Network and hardware' tab shows the availability zone (us-east-1c), subnet ID (subnet-0f80dc2e), and the master node (Running 1 m4.xlarge).

2. Copy the python file (spark-streaming.py) under the path: /home/hadoop using WinSCP files transfer application (Windows OS).

WinSCP:



The screenshot shows the WinSCP application interface. On the left, the local file system is displayed, showing a file named 'spark-streaming.py' in the directory 'C:\Users\Ganesh Jalakam\Desktop\Retail_Data_Analysis_Spark_Streaming\'. On the right, the remote file system is displayed, showing the directory '/home/hadoop/'. A red arrow indicates the transfer of the 'spark-streaming.py' file from the local desktop to the remote directory.

EMR Instance:

```
[hadoop@ip-172-31-82-226 ~]$ ls -ltr *.py
-rw-rw-r-- 1 hadoop hadoop 7374 Nov  7 17:59 spark-streaming.py
[hadoop@ip-172-31-82-226 ~]$ pwd
/home/hadoop
[hadoop@ip-172-31-82-226 ~]$
```

3. Run the following command to enable Kafka Integration with Apache Spark.
export SPARK_KAFKA_VERSION=0.10
4. Execute the python file using spark-submit command providing Kafka jar package as an argument. Save the console output in a text file (Console-output.txt).
spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.5 spark-streaming.py > Console-output.txt
5. Read the console output file using the command: cat Console-output.txt and check whether we can see the transformed data as per our requirement.

```

-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
Batch: 1
-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|invoice_no|country|timestamp|total_cost|total_items|is_order|is_return|
|-----|-----|-----|-----|-----|-----|-----|
|154132551996985|United Kingdom|2022-11-07 19:19:03|0.8500000238418579|1|1|0|
|154132551996986|United Kingdom|2022-11-07 19:19:05|35.40000057220459|12|1|0|
|154132551996987|United Kingdom|2022-11-07 19:19:15|106.2000002861023|60|1|0|
|154132551996988|United Kingdom|2022-11-07 19:19:31|4.199999809265137|2|1|0|
|154132551996989|United Kingdom|2022-11-07 19:19:31|44.780001163482666|26|1|0|
|154132551996990|United Kingdom|2022-11-07 19:19:31|4.130000114440918|1|1|0|
|154132551996991|United Kingdom|2022-11-07 19:19:40|43.52000045776367|16|1|0|
|154132551996992|United Kingdom|2022-11-07 19:19:42|13.000000178813934|24|1|0|
|154132551996993|United Kingdom|2022-11-07 19:19:46|256.83999959206581|65|1|0|
|154132551996994|United Kingdom|2022-11-07 19:19:48|37.899999380111694|22|1|0|
|154132551996995|United Kingdom|2022-11-07 19:19:51|25.0|2|1|0|
|154132551996996|United Kingdom|2022-11-07 19:19:57|40.63999992609024|38|1|0|
-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+

Batch: 2
-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|invoice_no|country|timestamp|total_cost|total_items|is_order|is_return|
|-----|-----|-----|-----|-----|-----|-----|
|154132551996997|United Kingdom|2022-11-07 19:19:59|60.79999828338623|14|1|0|
|154132551996998|United Kingdom|2022-11-07 19:20:00|121.48000180721283|52|1|0|
|154132551996999|United Kingdom|2022-11-07 19:20:02|5.879999935626984|8|1|0|
|154132551997000|United Kingdom|2022-11-07 19:20:06|66.0999984741211|21|1|0|
|154132551997001|United Kingdom|2022-11-07 19:20:15|34.55000019073486|11|1|0|
|154132551997002|France|2022-11-07 19:20:16|222.59999465942383|28|1|0|
|154132551997003|United Kingdom|2022-11-07 19:20:29|48.910000920295715|9|1|0|
|154132551997004|United Kingdom|2022-11-07 19:20:30|1.649999976158142|1|1|0|
|154132551997005|Belgium|2022-11-07 19:20:33|1.659999966621399|1|1|0|
|154132551997006|United Kingdom|2022-11-07 19:20:37|9.990000009536743|7|1|0|
|154132551997007|United Kingdom|2022-11-07 19:20:44|10.809999823570251|13|1|0|
|154132551997008|France|2022-11-07 19:20:44|73.49999916553497|30|1|0|
|154132551997009|United Kingdom|2022-11-07 19:20:48|5099.489902853966|533|1|0|
|154132551997010|Sweden|2022-11-07 19:20:54|82.00000143051147|124|1|0|
-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+

Batch: 3
-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+

```

6. Now check whether all the JSON files are created for Time based and Time-and-Country based KPIs in HDFS (Path: /user/hadoop/).

```

[hadoop@ip-172-31-82-226 ~]$ hadoop fs -ls /user/hadoop/
Found 5 items
drwxr-xr-x - hadoop hadoop 0 2022-11-07 19:35 /user/hadoop/.sparkStaging
drwxr-xr-x - hadoop hadoop 0 2022-11-07 18:05 /user/hadoop/Country-and-timebased-Checkpoint
drwxr-xr-x - hadoop hadoop 0 2022-11-07 19:35 /user/hadoop/Country-and-timebased-KPI
drwxr-xr-x - hadoop hadoop 0 2022-11-07 18:05 /user/hadoop/Timebased-Checkpoint
drwxr-xr-x - hadoop hadoop 0 2022-11-07 19:35 /user/hadoop/Timebased-KPI
[hadoop@ip-172-31-82-226 ~]$

```

Read Timebased-KPI JSON files:

hadoop fs -ls /user/hadoop/Timebased-KPI


```

-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 18:26 /user/hadoop/Timebased-KPI/part-00114-e8eaf430-20d1-407d-b343-e3aa3a30e002-c000.json
-rw-r--r-- 1 hadoop hadoop 209 2022-11-07 18:15 /user/hadoop/Timebased-KPI/part-00115-ad85fec7-df7e-4321-b399-51b2d2af6f51-c000.json
-rw-r--r-- 1 hadoop hadoop 193 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00117-25bfaf90-8042-4d97-1afa-5fa3e36e0252-c000.json
-rw-r--r-- 1 hadoop hadoop 193 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00120-3a9d438a-30a3-450c-bf89-f7a366275c2a-c000.json
-rw-r--r-- 1 hadoop hadoop 193 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00122-0a79d6a6-a1f9-4fb7-9835-3fcfcf7fa721-c000.json
-rw-r--r-- 1 hadoop hadoop 210 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00124-e02b901f-e7b4-479d-b677-3121222f8493-c000.json
-rw-r--r-- 1 hadoop hadoop 193 2022-11-07 18:37 /user/hadoop/Timebased-KPI/part-00126-8edad221-a391-4250-8544-ee66f25e604-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00127-4df11b7b-6071-4e44-92b8-12d0044b6080-c000.json
-rw-r--r-- 1 hadoop hadoop 196 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00128-bbf0f0ffcc-06a3-4b3c-bf23-0b8b8a08c224-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:33 /user/hadoop/Timebased-KPI/part-00129-dbd1d51a8-ae77-443f-99b3-1590b0a531f2-c000.json
-rw-r--r-- 1 hadoop hadoop 209 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00135-bdfe4333-70d9-4576-8862-adc63ef6a75a-c000.json
-rw-r--r-- 1 hadoop hadoop 212 2022-11-07 18:18 /user/hadoop/Timebased-KPI/part-00139-220fd44f-22fb-4716-af21-437bb88ce6a4-c000.json
-rw-r--r-- 1 hadoop hadoop 210 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00139-6e7ee3ce-186e-4820-add6-65c8c12647a3-c000.json
-rw-r--r-- 1 hadoop hadoop 196 2022-11-07 18:27 /user/hadoop/Timebased-KPI/part-00140-7cebe36c-c4a6-4207-9849-d4651b021751-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00142-57e8a1e8-ec97-4188-9bdf-02020a0c78f7-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00146-b801c336-62dd-42ac-aab6-fa45f6aa2a9b-c000.json
-rw-r--r-- 1 hadoop hadoop 210 2022-11-07 18:22 /user/hadoop/Timebased-KPI/part-00148-4cd0368a-4329-48bf-80db-54aa8f42629-c000.json
-rw-r--r-- 1 hadoop hadoop 210 2022-11-07 19:27 /user/hadoop/Timebased-KPI/part-00149-dcf598fa-7b96-4bb5-a98a-f54945f16b47-c000.json
-rw-r--r-- 1 hadoop hadoop 195 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00150-5a019ba1-93a4-4ada-9379-aeb31363f9be-c000.json
-rw-r--r-- 1 hadoop hadoop 197 2022-11-07 18:39 /user/hadoop/Timebased-KPI/part-00151-8c072788-ed17-4249-a888-0aa74ea569a9-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00152-df974d6e-66fb-4116-9284-c025e90b94b4-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00153-c4335a99-80cc-4995-b331-29d7a8b65984-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 18:34 /user/hadoop/Timebased-KPI/part-00155-59bf4475-b821-48ed-9e7f-e2c3b2f12f1-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00156-3c658fcc-9425-4ec6-9820-8cba0d86f473-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 18:40 /user/hadoop/Timebased-KPI/part-00160-e97f4991-a574-4d74-819c-fd0be799bccc-c000.json
-rw-r--r-- 1 hadoop hadoop 193 2022-11-07 18:41 /user/hadoop/Timebased-KPI/part-00163-f37bc758-a00c-482a-be53-1fc49bd7d8ee-c000.json
-rw-r--r-- 1 hadoop hadoop 211 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00164-5da7baef-84da-480f-b85e-acca5807f29f-c000.json
-rw-r--r-- 1 hadoop hadoop 405 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00165-b009cae9-4633-4033-9b02-1b3bf7c1a2fb-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:30 /user/hadoop/Timebased-KPI/part-00166-871a26c9-976e-43d4-alb2-d5da0ff7bac-c000.json
-rw-r--r-- 1 hadoop hadoop 211 2022-11-07 18:14 /user/hadoop/Timebased-KPI/part-00168-ee26c369-7fb5-4dff-aac6-d3d0766745d4-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00169-3c3b7a4f-1ff5-476a-be45-a03506b4056-c000.json
-rw-r--r-- 1 hadoop hadoop 212 2022-11-07 18:08 /user/hadoop/Timebased-KPI/part-00170-44b18051-14f1-4b10-bc01-f0f16dd9d960-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00170-526ce1bd-7818-4078-af2d-9ef57546cd35-c000.json
-rw-r--r-- 1 hadoop hadoop 210 2022-11-07 18:09 /user/hadoop/Timebased-KPI/part-00173-c83b299e-e569-45c3-8c56-b49721c5bc75-c000.json
-rw-r--r-- 1 hadoop hadoop 211 2022-11-07 18:11 /user/hadoop/Timebased-KPI/part-00174-73d0fcfc-d9b2-4cbd-a07a-f9afe74a2919-c000.json
-rw-r--r-- 1 hadoop hadoop 210 2022-11-07 18:11 /user/hadoop/Timebased-KPI/part-00175-8d91dd6d-69eb-4b46-82c1-15908323145c-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00178-6dc3c044-ff63-4796-aal5-04c9b9e5f329-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00180-acc9f1d1-8f79-404d-9176-f38155699c06-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 18:38 /user/hadoop/Timebased-KPI/part-00185-fb7f7367-f57b-4730-87d6-f9cc84f27bb2-c000.json
-rw-r--r-- 1 hadoop hadoop 210 2022-11-07 18:18 /user/hadoop/Timebased-KPI/part-00187-12501d46-81cb-49df-abd9-84e20f2179e5-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 18:30 /user/hadoop/Timebased-KPI/part-00189-c2b20298-215f-4daa-a10d-6e3bf52e4f6c-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00189-ebf49140-4f97-477c-84e8-016076e29bfc-c000.json
-rw-r--r-- 1 hadoop hadoop 194 2022-11-07 19:26 /user/hadoop/Timebased-KPI/part-00190-ac94212c-5d26-423f-9ffe-ae3342daab50-c000.json
-rw-r--r-- 1 hadoop hadoop 196 2022-11-07 18:29 /user/hadoop/Timebased-KPI/part-00191-b55b4e63-3d79-4650-ab24-603fc38c6724-c000.json
-rw-r--r-- 1 hadoop hadoop 193 2022-11-07 18:35 /user/hadoop/Timebased-KPI/part-00199-74d6d0d3-bd84-40df-9133-claa39a30a8d-c000.json

```

hadoop fs -cat /user/hadoop/Timebased-KPI/part*

```

{"window": {"start": "2022-11-07T18:24:00.000Z", "end": "2022-11-07T18:25:00.000Z", "opkm": 6, "total_sale_volume": 169.0399973988533, "average_transaction_size": 28.17333289880884, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T19:00:00.000Z", "end": "2022-11-07T19:01:00.000Z", "opkm": 5, "total_sale_volume": 978.3600056171417, "average_transaction_size": 195.67200112342834, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T18:30:00.000Z", "end": "2022-11-07T18:31:00.000Z", "opkm": 13, "total_sale_volume": 859.5000116229057, "average_transaction_size": 66.1153850945428, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T18:31:00.000Z", "end": "2022-11-07T18:32:00.000Z", "opkm": 5, "total_sale_volume": 729.5399905741215, "average_transaction_size": 145.9079981140243, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T18:39:00.000Z", "end": "2022-11-07T18:40:00.000Z", "opkm": 13, "total_sale_volume": 42.43999069929123, "average_transaction_size": 3.2646146691762485, "rate_of_return": 0.07692307692307693}, {"window": {"start": "2022-11-07T18:52:00.000Z", "end": "2022-11-07T18:53:00.000Z", "opkm": 9, "total_sale_volume": 1015.339966326952, "average_transaction_size": 112.81555181410577, "rate_of_return": 0.1111111111111111}, {"window": {"start": "2022-11-07T18:55:00.000Z", "end": "2022-11-07T18:56:00.000Z", "opkm": 10, "total_sale_volume": 168.06999823451042, "average_transaction_size": 16.806999823451044, "rate_of_return": 0.1}, {"window": {"start": "2022-11-07T19:20:00.000Z", "end": "2022-11-07T19:21:00.000Z", "opkm": 13, "total_sale_volume": 5778.619899213314, "average_transaction_size": 444.5092230164088, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T18:04:00.000Z", "end": "2022-11-07T18:05:00.000Z", "opkm": 12, "total_sale_volume": 187.05999130010605, "average_transaction_size": 15.58833260834217, "rate_of_return": 0.16666666666666666}, {"window": {"start": "2022-11-07T19:02:00.000Z", "end": "2022-11-07T19:03:00.000Z", "opkm": 11, "total_sale_volume": 638.1799947321415, "average_transaction_size": 58.01636315746741, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T17:58:00.000Z", "end": "2022-11-07T17:59:00.000Z", "opkm": 12, "total_sale_volume": 1286.32999593019485, "average_transaction_size": 23.860832994182903, "rate_of_return": 0.08333333333333333}, {"window": {"start": "2022-11-07T19:08:00.000Z", "end": "2022-11-07T19:09:00.000Z", "opkm": 5, "total_sale_volume": 338.95000195503235, "average_transaction_size": 67.79000039100647, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T17:59:00.000Z", "end": "2022-11-07T18:00:00.000Z", "opkm": 15, "total_sale_volume": 923.0799910724163, "average_transaction_size": 61.53866607149442, "rate_of_return": 0.06666666666666667}, {"window": {"start": "2022-11-07T18:02:00.000Z", "end": "2022-11-07T18:03:00.000Z", "opkm": 14, "total_sale_volume": 1095.5799964368343, "average_transaction_size": 78.25571403120246, "rate_of_return": 0.07142857142857142}, {"window": {"start": "2022-11-07T18:01:00.000Z", "end": "2022-11-07T18:02:00.000Z", "opkm": 12, "total_sale_volume": 889.1900013685226, "average_transaction_size": 74.09916678071022, "rate_of_return": 0.08333333333333333}, {"window": {"start": "2022-11-07T18:46:00.000Z", "end": "2022-11-07T18:47:00.000Z", "opkm": 10, "total_sale_volume": 532.0199964940548, "average_transaction_size": 53.20199964940548, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T18:53:00.000Z", "end": "2022-11-07T18:54:00.000Z", "opkm": 10, "total_sale_volume": 507.8799990415573, "average_transaction_size": 50.78799990415573, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T18:28:00.000Z", "end": "2022-11-07T18:29:00.000Z", "opkm": 9, "total_sale_volume": 521.0200008153915, "average_transaction_size": 57.891111201710174, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T18:08:00.000Z", "end": "2022-11-07T18:09:00.000Z", "opkm": 12, "total_sale_volume": 458.82999137043953, "average_transaction_size": 38.2358326142033, "rate_of_return": 0.08333333333333333}, {"window": {"start": "2022-11-07T18:20:00.000Z", "end": "2022-11-07T18:21:00.000Z", "opkm": 8, "total_sale_volume": 976.8099947422743, "average_transaction_size": 122.10124934278429, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T18:45:00.000Z", "end": "2022-11-07T18:46:00.000Z", "opkm": 9, "total_sale_volume": 1014.2899831533432, "average_transaction_size": 126.7862478941679, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T19:01:00.000Z", "end": "2022-11-07T19:02:00.000Z", "opkm": 13, "total_sale_volume": 1980.4299919307232, "average_transaction_size": 75.41769168697871, "rate_of_return": 0.0}, {"window": {"start": "2022-11-07T18:19:00.000Z", "end": "2022-11-07T18:20:00.000Z", "opkm": 8, "total_sale_volume": 340.1900028884411, "average_transaction_size": 42.523750361055136, "rate_of_return": 0.125}, {"window": {"start": "2022-11-07T18:25:00.000Z", "end": "2022-11-07T18:26:00.000Z", "opkm": 9, "total_sale_volume": 1730.93994140625, "average_transaction_size": 216.36749267578125, "rate_of_return": 0.0}

```

Read Country-and-timebased-KPI JSON files:

hadoop fs -ls /user/hadoop/Country-and-timebased-KPI

```
hadoop fs -cat /user/hadoop/Country-and-timebased-KPI/part*
```

© Copyright 2020. upGrad Education Pvt. Ltd. All rights reserved

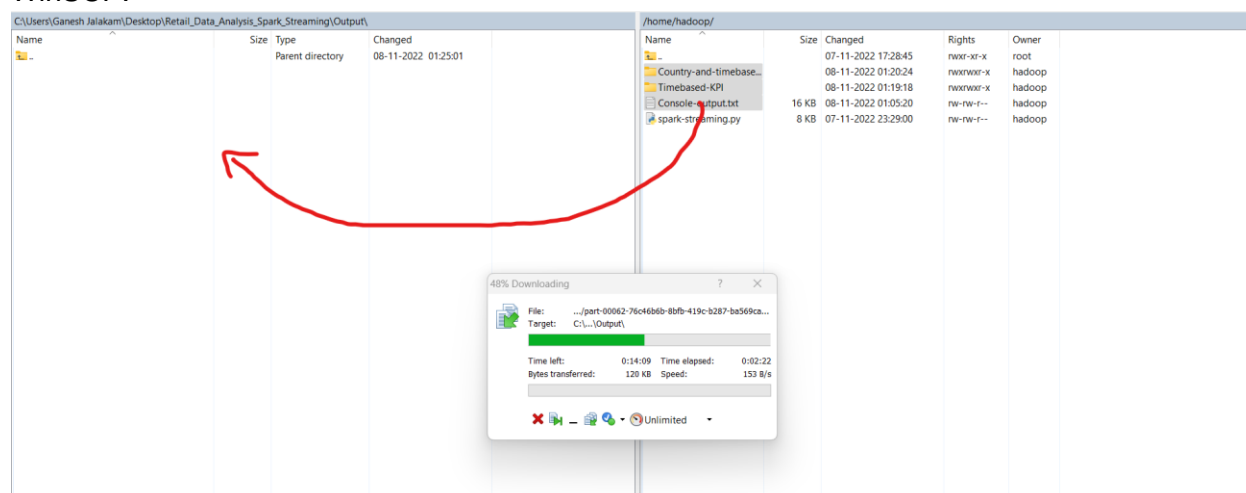
- Copy Time based and Country-and-Time based KPIs directories from HDFS to local path (/home/hadoop) in EMR instance.

```
hadoop fs -get /user/hadoop/Timebased-KPI /home/hadoop/
```

```
hadoop fs -get /user/hadoop/Country-and-timebased-KPI /home/hadoop/
```

```
[hadoop@ip-172-31-82-226 ~]$ ls -ltr
total 60
-rw-rw-r-- 1 hadoop hadoop 7374 Nov 7 17:59 spark-streaming.py
-rw-rw-r-- 1 hadoop hadoop 16349 Nov 7 19:35 Console-output.txt
drwxrwxr-x 3 hadoop hadoop 12288 Nov 7 19:49 Timebased-KPI
drwxrwxr-x 3 hadoop hadoop 16384 Nov 7 19:50 Country-and-timebased-KPI
[hadoop@ip-172-31-82-226 ~]$
```

- Copy all the console and JSON output files from EMR to local machine using WinSCP.



- Post completion of all the above steps, terminate the EMR instance from AWS console.

