

Code Logic - Retail Data Analysis

Setup:

1. Setting up spark environment before running the spark job

```
export SPARK_KAFKA_VERSION=0.10
```

2. Run the spark job to execute spark-streaming.py script
spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.5 spark-streaming.py

org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.5 - Spark Kafka Sql Jar version spark-streaming.py - Script to be executed

Code Logic:

spark-streaming.py script executes sequentially with following steps:

→ Import the necessary modules

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

→ Initialising Spark Session

```
spark = SparkSession \
    .builder \
    .appName("Retail_Project") \
    .getOrCreate()
```

→ Setting the log level to be ERROR

```
spark.sparkContext.setLogLevel('ERROR')
```

→ Reading data from Kafka

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

Format - Kafka

Bootstrap Server - 18.211.252.152:9092

Subscribe to Topic - real-time-project**Starting Offsets – Earliest**

→ Construct schema in a format provided as a input:

Input

```
{
  "invoice_no": 154132541653705,
  "country": "United Kingdom",
  "timestamp": "2020-09-18 10:55:23",
  "type": "ORDER",
  "items": [
    { "SKU": "21485",
      "title": "RETROSPOT HEART HOT WATER BOTTLE",
      "unit_price": 4.95,
      "quantity": 6
    },
    { "SKU": "23499",
      "title": "SET 12 VINTAGE DOILY CHALK",
      "unit_price": 0.42,
      "quantity": 2
    }
  ]
}
```

Schema

```
schema = 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())
    ])))
```

→ Read the input data with respect to schema provided

```
order_stream = raw_data.select(from_json(col("value").cast("string"),
    schema).alias("data")).select("data.*")
```

→ Add the following utility functions

is_order - A user defined function to return whether it's type is ORDER or not

```
def is_order(type):  
    if type=="ORDER":  
        return 1  
    else:  
        return 0
```

is_return - A user defined function to return whether it's type is RETURN or not

```
def is_return(type):  
    if type=="RETURN":  
        return 1  
    else:  
        return 0
```

total_items_count - A user defined function to return the total count of items ordered.

```
def total_items_count(items):  
    total_count = 0  
    for item in items:  
        total_count = total_count + item['quantity']  
    return total_count
```

calculate_total_cost - A user defined function to return the total cost of items ordered.

```
def calculate_total_cost(items,type):  
    total_price = 0  
    for item in items:  
        total_price = total_price + item['unit_price'] * item['quantity']  
    if type=="RETURN":  
        return total_price * -1  
    else:  
        return total_price
```

Define the UDFs with the utility functions

```
flag_order = udf(is_order, IntegerType())  
flag_return = udf(is_return, IntegerType())  
calculate_item = udf(total_items_count, IntegerType())  
calculate_cost = udf(calculate_total_cost, FloatType())
```

→ Append additional columns to the order_stream with the user defined functions

```
order_stream_extended = order_stream \  
.withColumn("total_items", calculate_item(order_stream.items)) \  
.withColumn("total_cost",calculate_cost(order_stream.items,order_stream.type))\  
.withColumn("is_order", flag_order(order_stream.type)) \  
.withColumn("is_return", flag_return(order_stream.type))
```

→ Select columns required to be logged as console. Every 1 minute the new data gets processed in a stream.

```
order_table_console = order_stream_extended \
    .select("invoice_no", "country",
"timestamp", "type", "total_items", "total_cost", "is_order", "is_return") \
    .writeStream \
    .outputMode("append") \
    .format("console") \
    .option("truncate", "false") \
    .trigger(processingTime="1 minute") \
    .start()
```

→ Calculate Time and Country based KPIs

Calculate time based KPIs

```
agg_time = order_stream_extended \
    .withWatermark("timestamp", "1 minute") \
    .groupby(window("timestamp", "1 minute", "1 minute")) \
    .agg(sum("total_cost").alias("total_volume_of_sales"),
    avg("total_cost").alias("average_transaction_size"),
    avg("is_return").alias("rate_of_return")) \
    .select("window.start", "window.end", "total_volume_of_sales", "average_transactio
n_size", "rate_of_return")
```

Calculate time and country based KPIs

```
agg_time_country = order_stream_extended \
    .withWatermark("timestamp", "1 minute") \
    .groupBy(window("timestamp", "1 minute", "1 minute"), "country") \
    .agg(sum("total_cost").alias("total_volume_of_sales"),
    count("invoice_no").alias("OPM"), avg("is_return").alias("rate_of_return")) \
    .select("window.start", "window.end", "country",
"OPM", "total_volume_of_sales", "rate_of_return")
```

→ Logging Time and Country based KPIs

Write time based KPI values

```
console_by_time = agg_time.writeStream \
    .format("json") \
    .outputMode("append") \
    .option("truncate", "false") \
    .option("path", "time_based_kpi/") \
    .option("checkpointLocation", "time_based_kpi/cp/") \
    .trigger(processingTime="1 minute") \
```

```
.start()
```

Write time and country based KPI values

```
console_by_country = agg_time_country.writeStream \
```

```
.format("json") \
```

```
.outputMode("append") \
```

```
.option("truncate", "false") \
```

```
.option("path", "time_country_based_kpi/") \
```

```
.option("checkpointLocation", "time_country_based_kpi/cp/") \
```

```
.trigger(processingTime="1 minute") \
```

```
.start()
```

→ Await Termination execution for writestream queries

```
order_table_console.awaitTermination()
```

```
console_by_time.awaitTermination()
```

```
console_by_country.awaitTermination()
```

```
hadoop@ip-172-31-38-124:~$ spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.5 spark-streaming.py
[bin]hadoop@ip-172-31-38-124:~$ spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.5 spark-streaming.py
Ivy Default Cache set to: /home/hadoop/.ivy2/cache
The jars for the packages stored in: /home/hadoop/.ivy2/jars
:: loading settings :: url = jar:file:/usr/lib/spark/jars/ivy-2.4.0.jar!/org/apache/ivy/core/settings/ivysettings.xml
org.apache.spark:spark-sql-kafka-0-10_2.11 added as a dependency
:: resolving dependencies :: org.apache.spark:spark-submit-parent-77d61317-f9e3-49ea-931e-e85137fcdcf3;1.0
  confs: [default]
    found org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.5 in central
    found org.apache.kafka:kafka-clients:2.0.0 in central
    found org.scala-lang:scala-library:2.10.4 in central
    found org.xerial.snappy:snappy-java:1.1.7.3 in central
    found org.slf4j:slf4j-api:1.7.16 in central
    found org.spark-project.spark:unused:1.0.0 in central
:: resolution report :: resolve 560ms :: artifacts dl 28ms
  :: modules in use:
    org.apache.kafka:kafka-clients:2.0.0 from central in [default]
    org.apache.spark:spark-sql-kafka-0-10_2.11:2.4.5 from central in [default]
    org.scala-lang:scala-library:2.10.4 from central in [default]
    org.xerial.snappy:snappy-java:1.1.7.3 from central in [default]
    org.spark-project.spark:unused:1.0.0 from central in [default]
    org.xerial.snappy:snappy-java:1.1.7.3 from central in [default]
  ::-----+-----+
  |      conf      |      | modules |      | artifacts |
  |      |      |      | number| search|dwnlded|evicted|| number|dwnlded|
  ::-----+-----+
  | default        | 6    | 0       | 0     | 0         | 0      | 6       | 0     |
  :: retrieving :: org.apache.spark:spark-submit-parent-77d61317-f9e3-49ea-931e-e85137fcdcf3
  confs: [default]
  0 artifacts copied, 6 already retrieved (0KB/15MB)
23/02/12 15:22:28 INFO SparkContext: Running Spark version 2.4.5-amzn-0
23/02/12 15:22:29 INFO SparkContext: Submitted application: Retail Project
23/02/12 15:22:29 INFO SecurityManager: Changing view acls to: hadoop
23/02/12 15:22:29 INFO SecurityManager: Changing modify acls to: hadoop
23/02/12 15:22:29 INFO SecurityManager: Changing view acls groups to:
23/02/12 15:22:29 INFO SecurityManager: Changing modify acls groups to:
23/02/12 15:22:29 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: Set(hadoop); groups with view permissions: Set(); users with modify permissions: Set(hadoop); groups with modify permissions: Set()
23/02/12 15:22:29 INFO Utils: Successfully started service 'sparkDriver' on port 41633.
23/02/12 15:22:29 INFO SparkEnv: Registering MapOutputTracker
23/02/12 15:22:29 INFO SparkEnv: Registering BlockManagerMaster
23/02/12 15:22:29 INFO BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
23/02/12 15:22:29 INFO BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
23/02/12 15:22:29 INFO DiskBlockManager: Created local directory at /mnt/tmp/blockmgr-f0b246ab-bff5-4753-98b6-405f8049559e
23/02/12 15:22:29 INFO MemoryStore: MemoryStore started with capacity 1038.8 MB
23/02/12 15:22:29 INFO SparkEnv: Registering OutputCommitCoordinator
23/02/12 15:22:30 WARN Utils: Service 'SparkUI' could not bind on port 4040. Attempting port 4041.
23/02/12 15:22:30 INFO Utils: Successfully started service 'SparkUI' on port 4041.
23/02/12 15:22:30 INFO SparkUI: Bound SparkUI to 0.0.0.0, and started at http://ip-172-31-38-124.ec2.internal:4041
23/02/12 15:22:30 INFO Utils: Using initial executors = 40, max of spark.dynamicAllocation.initialExecutors, spark.dynamicAllocation.minExecutors and spark.executor.instances
23/02/12 15:22:31 INFO RMPProxy: Connecting to ResourceManager at ip-172-31-38-124.ec2.internal/172.31.38.124:8032
23/02/12 15:22:31 INFO Client: Requesting a new application from cluster with 2 NodeManagers
23/02/12 15:22:31 INFO Client: Verifying our application has not requested more than the maximum memory capability of the cluster (6144 MB per container)
23/02/12 15:22:31 INFO Client: Will allocate AM container, with 896 MB memory including 384 MB overhead
23/02/12 15:22:31 INFO Client: Setting up container launch context for our AM
23/02/12 15:22:31 INFO Client: Setting up the launch environment for our AM container
23/02/12 15:22:31 INFO Client: Preparing resources for our AM container
23/02/12 15:22:31 WARN Client: Hadoop spark.yarn.jar not spark.yarn.archive is set, falling back to uploading libraries under SPARK HOME.
23/02/12 15:22:34 INFO Client: Uploading resource file:/mnt/cp/spark-865e74-4850-44e3-8410-e9a5d0ddbf0/_spark_libs_8808450727500007382.zip -> hdfs://ip-172-31-38-124.ec2.internal:8020/user/hadoop/.sparkStaging/application/1676213549055_0002/_spark_libs_8808450727500007382.zip
23/02/12 15:22:36 INFO Client: Uploading resource file:/home/hadoop/.ivy2/jars/org.apache.spark:spark-sql-kafka-0-10_2.11-2.4.5.jar -> hdfs://ip-172-31-38-124.ec2.internal:8020/user/hadoop/.sparkStaging/application/1676213549055_0002/org.apache.spark:spark-sql-kafka-0-10_2.11-2.4.5.jar
23/02/12 15:22:36 INFO Client: Uploading resource file:/home/hadoop/.ivy2/jars/org.apache.kafka:kafka-clients-2.0.0.jar -> hdfs://ip-172-31-38-124.ec2.internal:8020/user/hadoop/.sparkStaging/
```

 `hadoop@ip-172-31-38-124:~`

```
-----
Batch: 2
```

hadoop@ip-172-31-38-124:~

```
[154132553226654|United Kingdom|2023-02-12 15:23:25|ORDER|336|544.8|1|0|
[154132553226655|United Kingdom|2023-02-12 15:22:58|ORDER|25|43.75|1|0|
[154132553226656|United Kingdom|2023-02-12 15:23:01|ORDER|4|24.38|1|0|
[154132553226657|Belgium|2023-02-12 15:23:02|ORDER|15|35.71|1|0|
[154132553226658|United Kingdom|2023-02-12 15:23:03|ORDER|38|56.74|1|0|
[154132553226659|United Kingdom|2023-02-12 15:23:04|ORDER|19|33.79|1|0|
+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows
```

Batch: 2

invoice_no	country	timestamp	type	total_items	total_cost	is_order	is_return
[154132553226669 United Kingdom 2023-02-12 15:23:37 ORDER 1 10.42 1 0							
[154132553226670 United Kingdom 2023-02-12 15:23:37 ORDER 12 19.700001 1 0							
[154132553226671 United Kingdom 2023-02-12 15:23:41 ORDER 12 41.74 1 0							
[154132553226672 United Kingdom 2023-02-12 15:23:50 ORDER 15 68.01 1 0							
[154132553226673 United Kingdom 2023-02-12 15:23:57 ORDER 98 266.74 1 0							
[154132553226674 United Kingdom 2023-02-12 15:23:59 ORDER 1 4.25 1 0							
[154132553226675 United Kingdom 2023-02-12 15:24:02 ORDER 10 30.8 1 0							
[154132553226676 United Kingdom 2023-02-12 15:24:05 ORDER 26 54.6 1 0							
[154132553226677 United Kingdom 2023-02-12 15:24:10 ORDER 29 207.05 1 0							
[154132553226678 United Kingdom 2023-02-12 15:24:18 ORDER 1 2.1 1 0							
[154132553226679 United Kingdom 2023-02-12 15:24:27 ORDER 22 57.260002 1 0							
[154132553226680 United Kingdom 2023-02-12 15:24:28 RETURN 7 -5.18 0 1							
[154132553226681 United Kingdom 2023-02-12 15:24:06 ORDER 44 49.53 1 0							
[154132553226682 United Kingdom 2023-02-12 15:24:09 ORDER 3 25.800001 1 0							
[154132553226683 United Kingdom 2023-02-12 15:24:09 ORDER 68 87.6 1 0							
[154132553226684 Netherlands 2023-02-12 15:24:10 ORDER 24 19.92 1 0							
[154132553226685 United Kingdom 2023-02-12 15:24:12 ORDER 13 18.49 1 0							
[154132553226686 United Kingdom 2023-02-12 15:24:15 ORDER 3 12.39 1 0							
[154132553226687 United Kingdom 2023-02-12 15:24:21 ORDER 433 600.24 1 0							
[154132553226688 United Kingdom 2023-02-12 15:24:22 ORDER 16 16.68 1 0							

only showing top 20 rows

Batch: 3

invoice_no	country	timestamp	type	total_items	total_cost	is_order	is_return
[154132553226692 United Kingdom 2023-02-12 15:25:02 ORDER 31 68.47 1 0							
[154132553226693 United Kingdom 2023-02-12 15:25:07 ORDER 49 118.64 1 0							
[154132553226694 Germany 2023-02-12 15:25:15 ORDER 21 22.05 1 0							
[154132553226695 United Kingdom 2023-02-12 15:25:28 ORDER 1 3.75 1 0							
[154132553226696 United Kingdom 2023-02-12 15:25:28 RETURN 141 -166.73 0 1							
[154132553226697 United Kingdom 2023-02-12 15:25:29 ORDER 54 97.6 1 0							
[154132553226698 United Kingdom 2023-02-12 15:24:57 ORDER 53 214.5 1 0							
[154132553226699 United Kingdom 2023-02-12 15:24:58 ORDER 7 11.75 1 0							
[154132553226700 United Kingdom 2023-02-12 15:25:07 ORDER 20 42.2 1 0							
[154132553226701 United Kingdom 2023-02-12 15:25:17 ORDER 10 21.22 1 0							
[154132553226702 United Kingdom 2023-02-12 15:25:27 ORDER 4 2.5 1 0							
[154132553226703 United Kingdom 2023-02-12 15:25:29 ORDER 6 5.1000004 1 0							
[154132553226704 United Kingdom 2023-02-12 15:25:29 ORDER 14 18.71 1 0							
[154132553226705 United Kingdom 2023-02-12 15:25:39 ORDER 6 25.5 1 0							
[154132553226706 United Kingdom 2023-02-12 15:25:40 ORDER 22 39.51 1 0							
[154132553226707 United Kingdom 2023-02-12 15:25:46 ORDER 85 138.83 1 0							

█


```
hadoop@ip-172-31-38-124:~
|154132553226706|United Kingdom|2023-02-12 15:25:40|ORDER|22|39.51|1|0|
|154132553226707|United Kingdom|2023-02-12 15:25:46|ORDER|85|138.83|1|0|
+-----+-----+-----+-----+-----+-----+-----+-----+
Batch: 4
+-----+-----+-----+-----+-----+-----+-----+-----+
|invoice_no|country|timestamp|type|total_items|total_cost|is_order|is_return|
+-----+-----+-----+-----+-----+-----+-----+-----+
|154132553226708|United Kingdom|2023-02-12 15:25:29|ORDER|30|89.44|1|0|
|154132553226709|United Kingdom|2023-02-12 15:25:34|ORDER|12|38.5|1|0|
|154132553226710|United Kingdom|2023-02-12 15:25:39|ORDER|2|1.7|1|0|
|154132553226711|United Kingdom|2023-02-12 15:25:46|ORDER|2|5.9|1|0|
|154132553226712|Channel Islands|2023-02-12 15:25:50|RETURN|31|-55.71|0|1|
|154132553226713|United Kingdom|2023-02-12 15:25:52|ORDER|2|1.48|1|0|
|154132553226714|United Kingdom|2023-02-12 15:25:55|ORDER|40|124.8|1|0|
|154132553226715|United Kingdom|2023-02-12 15:26:04|ORDER|126|118.98|1|0|
|154132553226716|United Kingdom|2023-02-12 15:26:07|ORDER|200|416.0|1|0|
|154132553226717|United Kingdom|2023-02-12 15:26:09|RETURN|5|-55.65|0|1|
|154132553226718|United Kingdom|2023-02-12 15:26:10|ORDER|15|20.85|1|0|
|154132553226719|United Kingdom|2023-02-12 15:26:14|ORDER|57|30.75|1|0|
|154132553226720|United Kingdom|2023-02-12 15:26:18|ORDER|5|24.75|1|0|
|154132553226721|United Kingdom|2023-02-12 15:26:26|ORDER|42|97.82|1|0|
|154132553226722|United Kingdom|2023-02-12 15:26:27|ORDER|1|0.19|1|0|
|154132553226723|United Kingdom|2023-02-12 15:25:54|ORDER|3|4.15|1|0|
|154132553226724|Belgium|2023-02-12 15:26:04|ORDER|10|79.259995|1|0|
|154132553226725|United Kingdom|2023-02-12 15:26:05|ORDER|4|18.599998|1|0|
|154132553226726|United Kingdom|2023-02-12 15:26:09|ORDER|86|107.32|1|0|
|154132553226727|United Kingdom|2023-02-12 15:26:13|ORDER|21|46.879997|1|0|
+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows

Batch: 5
+-----+-----+-----+-----+-----+-----+-----+-----+
|invoice_no|country|timestamp|type|total_items|total_cost|is_order|is_return|
+-----+-----+-----+-----+-----+-----+-----+-----+
|154132553226738|United Kingdom|2023-02-12 15:26:28|ORDER|1|1.65|1|0|
|154132553226739|United Kingdom|2023-02-12 15:26:29|ORDER|34|66.46|1|0|
|154132553226740|United Kingdom|2023-02-12 15:26:42|ORDER|24|106.299995|1|0|
|154132553226741|United Kingdom|2023-02-12 15:26:51|ORDER|30|28.29|1|0|
|154132553226742|United Kingdom|2023-02-12 15:26:51|ORDER|12|15.0|1|0|
|154132553226743|United Kingdom|2023-02-12 15:26:52|ORDER|11|41.379997|1|0|
|154132553226744|France|2023-02-12 15:27:06|ORDER|25|42.87|1|0|
|154132553226745|United Kingdom|2023-02-12 15:27:06|ORDER|34|62.66|1|0|
|154132553226746|United Kingdom|2023-02-12 15:27:09|ORDER|29|100.2|1|0|
|154132553226747|United Kingdom|2023-02-12 15:27:13|ORDER|11|36.83|1|0|
|154132553226748|United Kingdom|2023-02-12 15:27:21|RETURN|153|-249.53|0|1|
|154132553226749|United Kingdom|2023-02-12 15:26:59|ORDER|1|2.1|1|0|
|154132553226750|United Kingdom|2023-02-12 15:27:00|RETURN|57|-134.24|0|1|
|154132553226751|United Kingdom|2023-02-12 15:27:05|ORDER|31|31.47|1|0|
|154132553226752|United Kingdom|2023-02-12 15:27:07|ORDER|2|16.94|1|0|
|154132553226753|United Kingdom|2023-02-12 15:27:10|ORDER|6|12.48|1|0|
|154132553226754|United Kingdom|2023-02-12 15:27:13|ORDER|10|29.5|1|0|
|154132553226755|United Kingdom|2023-02-12 15:27:29|ORDER|2|4.2|1|0|
|154132553226756|United Kingdom|2023-02-12 15:27:29|ORDER|64|62.34|1|0|
|154132553226757|United Kingdom|2023-02-12 15:27:30|ORDER|3|8.75|1|0|
+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows
```



```
hadoop@ip-172-31-38-124:~
```

```
Batch: 6
```

invoice_no	country	timestamp	type	total_items	total_cost	is_order	is_return
154132553226761	United Kingdom	2023-02-12 15:27:31	ORDER	2	9.9	1	0
154132553226762	United Kingdom	2023-02-12 15:27:32	ORDER	34	105.54	1	0
154132553226763	United Kingdom	2023-02-12 15:27:41	ORDER	13	51.95	1	0
154132553226764	United Kingdom	2023-02-12 15:27:48	ORDER	103	221.73999	1	0
154132553226765	United Kingdom	2023-02-12 15:27:54	ORDER	19	86.58	1	0
154132553226766	United Kingdom	2023-02-12 15:28:05	ORDER	76	154.70999	1	0
154132553226767	United Kingdom	2023-02-12 15:28:09	ORDER	27	44.5	1	0
154132553226768	United Kingdom	2023-02-12 15:28:12	ORDER	3	4.3500004	1	0
154132553226769	United Kingdom	2023-02-12 15:28:24	ORDER	131	126.49	1	0
154132553226770	United Kingdom	2023-02-12 15:27:55	ORDER	38	58.1	1	0
154132553226771	United Kingdom	2023-02-12 15:27:57	ORDER	12	45.0	1	0
154132553226772	United Kingdom	2023-02-12 15:28:00	ORDER	9	75.74	1	0
154132553226773	United Kingdom	2023-02-12 15:28:05	ORDER	7	22.21	1	0
154132553226774	United Kingdom	2023-02-12 15:28:07	ORDER	12	22.68	1	0
154132553226775	United Kingdom	2023-02-12 15:28:10	ORDER	6	9.1	1	0
154132553226776	United Kingdom	2023-02-12 15:28:17	ORDER	1	2.95	1	0
154132553226777	United Kingdom	2023-02-12 15:28:18	ORDER	32	36.78	1	0
154132553226778	Iceland	2023-02-12 15:28:23	ORDER	28	36.16	1	0
154132553226779	United Kingdom	2023-02-12 15:28:29	ORDER	8	24.119999	1	0
154132553226780	United Kingdom	2023-02-12 15:28:30	ORDER	9	23.8	1	0

```
only showing top 20 rows
```

```
Batch: 7
```

invoice_no	country	timestamp	type	total_items	total_cost	is_order	is_return
154132553226784	United Kingdom	2023-02-12 15:28:26	ORDER	17	58.07	1	0
154132553226785	United Kingdom	2023-02-12 15:28:31	ORDER	12	25.9	1	0
154132553226786	United Kingdom	2023-02-12 15:28:38	ORDER	11	43.71	1	0
154132553226787	United Kingdom	2023-02-12 15:28:39	ORDER	15	25.449999	1	0
154132553226788	United Kingdom	2023-02-12 15:28:48	ORDER	29	55.92	1	0
154132553226789	United Kingdom	2023-02-12 15:28:53	ORDER	13	21.1	1	0
154132553226790	United Kingdom	2023-02-12 15:28:56	ORDER	12	19.8	1	0
154132553226791	United Kingdom	2023-02-12 15:28:57	ORDER	22	52.23	1	0
154132553226792	United Kingdom	2023-02-12 15:28:57	ORDER	6	12.48	1	0
154132553226793	United Kingdom	2023-02-12 15:28:58	ORDER	60	75.0	1	0
154132553226794	United Kingdom	2023-02-12 15:29:01	ORDER	2	6.58	1	0
154132553226795	United Kingdom	2023-02-12 15:29:01	ORDER	2	6.58	1	0
154132553226796	United Kingdom	2023-02-12 15:29:01	ORDER	14	35.54	1	0
154132553226797	United Kingdom	2023-02-12 15:29:07	ORDER	12	24.48	1	0
154132553226798	United Kingdom	2023-02-12 15:29:10	ORDER	15	40.32	1	0
154132553226799	United Kingdom	2023-02-12 15:29:24	ORDER	10	22.74	1	0
154132553226800	United Kingdom	2023-02-12 15:29:02	ORDER	4	10.33	1	0
154132553226801	United Kingdom	2023-02-12 15:29:14	ORDER	33	172.228996	1	0
154132553226802	United Kingdom	2023-02-12 15:29:19	ORDER	6	10.700001	1	0
154132553226803	United Kingdom	2023-02-12 15:29:22	ORDER	36	23.4	1	0

```
only showing top 20 rows
```

```
Batch: 8
```

invoice_no	country	timestamp	type	total_items	total_cost	is_order	is_return
------------	---------	-----------	------	-------------	------------	----------	-----------

I checked HDFS to make sure the KPI files were present.

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

```
hadoop@ip-172-31-38-124 ~]$ hadoop fs -ls user/ec2-user/
Found 4 items
drwxr-xr-x - hadoop hadoop          0 2023-02-12 15:39 user/ec2-user/country_kpi
drwxr-xr-x - hadoop hadoop          0 2023-02-12 15:23 user/ec2-user/country_kpi_checkpoints
drwxr-xr-x - hadoop hadoop          0 2023-02-12 15:39 user/ec2-user/time_kpi
drwxr-xr-x - hadoop hadoop          0 2023-02-12 15:23 user/ec2-user/time_kpi_checkpoints
hadoop@ip-172-31-38-124 ~]$
```

I also checked the folders to see the JSON files:

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

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

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

And used 'cat' command to take a look at the data:

hadoop fs -cat user/ec2-user/time_kpi/part*

```
[hadoop@ip-172-31-38-124 ~]$ hadoop fs -cat user/ec2-user/time_kpi/part*
{"start":2023-02-12T15:35:00.000Z,"end":2023-02-12T15:36:00.000Z,"total_volume_of_sales":940.980011651421,"average_transaction_size":39.20750046521425,"rate_of_return":0.125}
{"start":2023-02-12T15:34:00.000Z,"end":2023-02-12T15:35:00.000Z,"total_volume_of_sales":1557.2100145816003,"average_transaction_size":62.2584005326721,"rate_of_return":0.04}
{"start":2023-02-12T15:23:00.000Z,"end":2023-02-12T15:24:00.000Z,"total_volume_of_sales":1776.429989510684,"average_transaction_size":65.025924413534,"rate_of_return":0.03703703703703}
{"start":2023-02-12T15:43:00.000Z,"end":2023-02-12T15:44:00.000Z,"total_volume_of_sales":1545.900016794668,"average_transaction_size":67.2130420802905,"rate_of_return":0.0}
{"start":2023-02-12T15:36:00.000Z,"end":2023-02-12T15:37:00.000Z,"total_volume_of_sales":1332.8799983263016,"average_transaction_size":63.47047611077627,"rate_of_return":0.17647619047619047}
{"start":2023-02-12T15:24:00.000Z,"end":2023-02-12T15:25:00.000Z,"total_volume_of_sales":1552.3499898910522,"average_transaction_size":81.70263104689749,"rate_of_return":0.05263157894736842}
{"start":2023-02-12T15:27:00.000Z,"end":2023-02-12T15:28:00.000Z,"total_volume_of_sales":1197.549947249103,"average_transaction_size":81.41521509833957,"rate_of_return":0.00695652173913043}
{"start":2023-02-12T15:33:00.000Z,"end":2023-02-12T15:34:00.000Z,"total_volume_of_sales":1527.0099925994973,"average_transaction_size":32.93812453746796,"rate_of_return":0.1875}
{"start":2023-02-12T15:29:00.000Z,"end":2023-02-12T15:30:00.000Z,"total_volume_of_sales":1468.869951171875,"average_transaction_size":127.58059794806985,"rate_of_return":0.0}
{"start":2023-02-12T15:40:00.000Z,"end":2023-02-12T15:41:00.000Z,"total_volume_of_sales":12765.3599593639374,"average_transaction_size":81.331164518805,"rate_of_return":0.029411764705882353}
{"start":2023-02-12T15:31:00.000Z,"end":2023-02-12T15:32:00.000Z,"total_volume_of_sales":1361.83999729156484,"average_transaction_size":127.83345945504997,"rate_of_return":0.0}
{"start":2023-02-12T15:31:00.000Z,"end":2023-02-12T15:32:00.000Z,"total_volume_of_sales":11759.2798622404636,"average_transaction_size":117.28533082803051,"rate_of_return":0.0}
{"start":2023-02-12T15:22:00.000Z,"end":2023-02-12T15:23:00.000Z,"total_volume_of_sales":386.07998752593994,"average_transaction_size":48.25999844074249,"rate_of_return":0.0}
{"start":2023-02-12T15:37:00.000Z,"end":2023-02-12T15:38:00.000Z,"total_volume_of_sales":1161.7598833011627,"average_transaction_size":68.33882254712722,"rate_of_return":0.17647058823529413}
{"start":2023-02-12T15:41:00.000Z,"end":2023-02-12T15:42:00.000Z,"total_volume_of_sales":1477.820042848597,"average_transaction_size":82.0733957138104,"rate_of_return":0.05535353535353535}
{"start":2023-02-12T15:35:00.000Z,"end":2023-02-12T15:36:00.000Z,"total_volume_of_sales":147.6100089558018,"average_transaction_size":129.4361858884372,"rate_of_return":0.09090909090909091}
{"start":2023-02-12T15:42:00.000Z,"end":2023-02-12T15:43:00.000Z,"total_volume_of_sales":12038.139984846115,"average_transaction_size":175.4866610541167,"rate_of_return":0.11111111111111111}
{"start":2023-02-12T15:28:00.000Z,"end":2023-02-12T15:29:00.000Z,"total_volume_of_sales":1139.4199714660645,"average_transaction_size":43.82384505638709,"rate_of_return":0.0}
{"start":2023-02-12T15:39:00.000Z,"end":2023-02-12T15:40:00.000Z,"total_volume_of_sales":1486.8099828985811,"average_transaction_size":19.5306239309907,"rate_of_return":0.0}
{"start":2023-02-12T15:26:00.000Z,"end":2023-02-12T15:27:00.000Z,"total_volume_of_sales":1514.49990210426,"average_transaction_size":152.2241372507194,"rate_of_return":0.034482758620489655}
{"start":2023-02-12T15:32:00.000Z,"end":2023-02-12T15:33:00.000Z,"total_volume_of_sales":1784.989966716766,"average_transaction_size":55.780937395898965,"rate_of_return":0.03125}
{"start":2023-02-12T15:38:00.000Z,"end":2023-02-12T15:39:00.000Z,"total_volume_of_sales":873.7299938201904,"average_transaction_size":145.98579148431074,"rate_of_return":0.0}
[hadoop@ip-172-31-38-124 ~]$
```

hadoop fs -cat user/ec2-user/country_kpi/part*

```
[hadoop@ip-172-31-38-124 ~]$ hadoop fs -cat user/ec2-user/country_kpi/part*
{"start":2023-02-12T15:41:00.000Z,"end":2023-02-12T15:42:00.000Z,"total_volu
-bash: ${ume_of_sales}161.7599833011627,average_transaction_size:68.3382254712722,rate_of_return:0.17647058823529413}{start:2023-02-12T15:41:00.000Z,end:2023-02-12T15:42:00.000Z,total_volu": 0
[hadoop@ip-172-31-38-124 ~]$
[hadoop@ip-172-31-38-124 ~]$ hadoop fs -cat user/ec2-user/country_kpi/part*
{"start":2023-02-12T15:31:00.000Z,"end":2023-02-12T15:32:00.000Z,"country":"United Kingdom","ORM":15,"total volume of sales":1759.2799624204636,"rate_of_return":0.0}
{"start":2023-02-12T15:41:00.000Z,"end":2023-02-12T15:41:00.000Z,"country":"Netherlands","ORM":11,"total volume of sales":145.82000732421975,"rate_of_return":0.0}
{"start":2023-02-12T15:31:00.000Z,"end":2023-02-12T15:32:00.000Z,"country":"France","ORM":1,"total volume of sales":6.4600003146973,"rate_of_return":0.0}
{"start":2023-02-12T15:29:00.000Z,"end":2023-02-12T15:30:00.000Z,"country":"United Kingdom","ORM":17,"total volume of sales":1468.869951171875,"rate_of_return":0.0}
{"start":2023-02-12T15:40:00.000Z,"end":2023-02-12T15:41:00.000Z,"country":"Germany","ORM":1,"total volume of sales":3.75,"rate_of_return":0.0}
{"start":2023-02-12T15:40:00.000Z,"end":2023-02-12T15:41:00.000Z,"country":"France","ORM":1,"total volume of sales":2.549999952316284,"rate_of_return":0.0}
{"start":2023-02-12T15:31:00.000Z,"end":2023-02-12T15:32:00.000Z,"country":"United Kingdom","ORM":24,"total volume of sales":1532.4180088656,"rate_of_return":0.04166666666666664}
{"start":2023-02-12T15:28:00.000Z,"end":2023-02-12T15:29:00.000Z,"country":"United Kingdom","ORM":124,"total volume of sales":1097.8299717903137,"rate_of_return":0.0}
{"start":2023-02-12T15:37:00.000Z,"end":2023-02-12T15:38:00.000Z,"country":"IRE","ORM":1,"total volume of sales":30.0,"rate_of_return":0.0}
{"start":2023-02-12T15:36:00.000Z,"end":2023-02-12T15:37:00.000Z,"country":"Germany","ORM":1,"total volume of sales":21.329999923706055,"rate_of_return":0.0}
{"start":2023-02-12T15:24:00.000Z,"end":2023-02-12T15:25:00.000Z,"country":"Germany","ORM":1,"total volume of sales":19.79999937060547,"rate_of_return":0.0}
{"start":2023-02-12T15:26:00.000Z,"end":2023-02-12T15:27:00.000Z,"country":"France","ORM":1,"total volume of sales":15.4289998218633,"rate_of_return":0.2}
{"start":2023-02-12T15:27:00.000Z,"end":2023-02-12T15:28:00.000Z,"country":"United Kingdom","ORM":21,"total volume of sales":1709.9799513816933,"rate_of_return":0.09523809523809523}
{"start":2023-02-12T15:27:00.000Z,"end":2023-02-12T15:28:00.000Z,"country":"France","ORM":1,"total volume of sales":42.869998931884766,"rate_of_return":0.0}
{"start":2023-02-12T15:33:00.000Z,"end":2023-02-12T15:34:00.000Z,"country":"United Kingdom","ORM":16,"total volume of sales":1827.0099925994973,"rate_of_return":0.1875}
{"start":2023-02-12T15:37:00.000Z,"end":2023-02-12T15:38:00.000Z,"country":"United Kingdom","ORM":15,"total volume of sales":1096.489992238223,"rate_of_return":0.0}
{"start":2023-02-12T15:26:00.000Z,"end":2023-02-12T15:27:00.000Z,"country":"United Kingdom","ORM":128,"total volume of sales":1435.2399857640266,"rate_of_return":0.03571428571428571}
{"start":2023-02-12T15:27:00.000Z,"end":2023-02-12T15:28:00.000Z,"country":"IRE","ORM":1,"total volume of sales":119.69999694824219,"rate_of_return":0.0}
{"start":2023-02-12T15:39:00.000Z,"end":2023-02-12T15:40:00.000Z,"country":"United Kingdom","ORM":15,"total volume of sales":1445.8598021329117,"rate_of_return":0.0}
{"start":2023-02-12T15:41:00.000Z,"end":2023-02-12T15:42:00.000Z,"country":"United Kingdom","ORM":17,"total volume of sales":1457.6400424948113,"rate_of_return":0.05882352941764705}
{"start":2023-02-12T15:39:00.000Z,"end":2023-02-12T15:40:00.000Z,"country":"Germany","ORM":1,"total volume of sales":50.55000076293945,"rate_of_return":0.0}
{"start":2023-02-12T15:25:00.000Z,"end":2023-02-12T15:26:00.000Z,"country":"United Kingdom","ORM":20,"total volume of sales":681.2700088024139,"rate_of_return":0.05}
{"start":2023-02-12T15:36:00.000Z,"end":2023-02-12T15:36:00.000Z,"country":"France","ORM":12,"total volume of sales":72.59000033140182,"rate_of_return":0.5}
{"start":2023-02-12T15:24:00.000Z,"end":2023-02-12T15:25:00.000Z,"country":"United Kingdom","ORM":16,"total volume of sales":1532.4180088656,"rate_of_return":0.05555555555555555}
{"start":2023-02-12T15:42:00.000Z,"end":2023-02-12T15:43:00.000Z,"country":"France","ORM":2,"total volume of sales":60.67000158364288,"rate_of_return":0.0}
{"start":2023-02-12T15:36:00.000Z,"end":2023-02-12T15:36:00.000Z,"country":"United Kingdom","ORM":120,"total volume of sales":807.480010963281,"rate_of_return":0.1}
{"start":2023-02-12T15:24:00.000Z,"end":2023-02-12T15:25:00.000Z,"country":"Netherlands","ORM":1,"total volume of sales":19.920000076293945,"rate_of_return":0.0}
{"start":2023-02-12T15:38:00.000Z,"end":2023-02-12T15:39:00.000Z,"country":"United Kingdom","ORM":17,"total volume of sales":1805.689991950988,"rate_of_return":0.0}
{"start":2023-02-12T15:44:00.000Z,"end":2023-02-12T15:45:00.000Z,"country":"France","ORM":1,"total volume of sales":72.22999991823629,"rate_of_return":0.0}
{"start":2023-02-12T15:42:00.000Z,"end":2023-02-12T15:43:00.000Z,"country":"IRE","ORM":1,"total volume of sales":1461.639941308594,"rate_of_return":0.0}
{"start":2023-02-12T15:38:00.000Z,"end":2023-02-12T15:39:00.000Z,"country":"Switzerland","ORM":1,"total volume of sales":157.08000183105465,"rate_of_return":0.0}
{"start":2023-02-12T15:28:00.000Z,"end":2023-02-12T15:29:00.000Z,"country":"Iceland","ORM":1,"total volume of sales":36.15999984741211,"rate_of_return":0.0}
{"start":2023-02-12T15:43:00.000Z,"end":2023-02-12T15:44:00.000Z,"country":"United Kingdom","ORM":1,"total volume of sales":1505.889990545265,"rate_of_return":0.0}
{"start":2023-02-12T15:34:00.000Z,"end":2023-02-12T15:35:00.000Z,"country":"Germany","ORM":1,"total volume of sales":35.0699996482422,"rate_of_return":0.0}
{"start":2023-02-12T15:40:00.000Z,"end":2023-02-12T15:41:00.000Z,"country":"Hong Kong","ORM":1,"total volume of sales":13.979999542236328,"rate_of_return":0.0}
{"start":2023-02-12T15:36:00.000Z,"end":2023-02-12T15:37:00.000Z,"country":"United Kingdom","ORM":19,"total volume of sales":1276.039996266365,"rate_of_return":0.05263157894736842}
{"start":2023-02-12T15:41:00.000Z,"end":2023-02-12T15:41:00.000Z,"country":"United Kingdom","ORM":16,"total volume of sales":139.4180088656,"rate_of_return":0.0}
{"start":2023-02-12T15:43:00.000Z,"end":2023-02-12T15:44:00.000Z,"country":"United Kingdom","ORM":120,"total volume of sales":1443.1800020323364,"rate_of_return":0.0}
{"start":2023-02-12T15:44:00.000Z,"end":2023-02-12T15:45:00.000Z,"country":"IRE","ORM":1,"total volume of sales":109.08000183105465,"rate_of_return":0.0}
{"start":2023-02-12T15:43:00.000Z,"end":2023-02-12T15:44:00.000Z,"country":"France","ORM":2,"total volume of sales":127.44999809265137,"rate_of_return":0.0}
{"start":2023-02-12T15:22:00.000Z,"end":2023-02-12T15:23:00.000Z,"country":"United Kingdom","ORM":16,"total volume of sales":339.4180088656,"rate_of_return":0.0}
{"start":2023-02-12T15:22:00.000Z,"end":2023-02-12T15:23:00.000Z,"country":"Portugal","ORM":1,"total volume of sales":26.85999970300293,"rate_of_return":0.0}
{"start":2023-02-12T15:30:00.000Z,"end":2023-02-12T15:31:00.000Z,"country":"United Kingdom","ORM":113,"total volume of sales":361.83999729156494,"rate_of_return":0.0}
{"start":2023-02-12T15:42:00.000Z,"end":2023-02-12T15:43:00.000Z,"country":"Sweden","ORM":1,"total volume of sales":39.34000015258789,"rate_of_return":0.0}
{"start":2023-02-12T15:35:00.000Z,"end":2023-02-12T15:36:00.000Z,"country":"IRE","ORM":12,"total volume of sales":60.90999984741211,"rate_of_return":0.0}
{"start":2023-02-12T15:41:00.000Z,"end":2023-02-12T15:42:00.000Z,"country":"Portugal","ORM":1,"total volume of sales":75.26000033140182,"rate_of_return":0.0}
{"start":2023-02-12T15:23:00.000Z,"end":2023-02-12T15:24:00.000Z,"country":"United Kingdom","ORM":126,"total volume of sales":1720.7199868261814,"rate_of_return":0.039461538461538464}
{"start":2023-02-12T15:36:00.000Z,"end":2023-02-12T15:37:00.000Z,"country":"Poland","ORM":1,"total volume of sales":35.51000213623047,"rate_of_return":0.0}
{"start":2023-02-12T15:23:00.000Z,"end":2023-02-12T15:24:00.000Z,"country":"Belgium","ORM":1,"total volume of sales":35.709999804472656,"rate_of_return":0.0}
{"start":2023-02-12T15:42:00.000Z,"end":2023-02-12T15:43:00.000Z,"country":"Belgium","ORM":1,"total volume of sales":9.1579999237060547,"rate_of_return":1.0}
{"start":2023-02-12T15:37:00.000Z,"end":2023-02-12T15:38:00.000Z,"country":"France","ORM":1,"total volume of sales":34.55000076293945,"rate_of_return":0.0}
{"start":2023-02-12T15:25:00.000Z,"end":2023-02-12T15:26:00.000Z,"country":"Channel Islands","ORM":1,"total volume of sales":1.55.709999804472656,"rate_of_return":1.0}
{"start":2023-02-12T15:40:00.000Z,"end":2023-02-12T15:41:00.000Z,"country":"United Kingdom","ORM":30,"total volume of sales":2599.259952545166,"rate_of_return":0.03333333333333333}
{"start":2023-02-12T15:42:00.000Z,"end":2023-02-12T15:43:00.000Z,"country":"United Kingdom","ORM":122,"total volume of sales":1486.2899978160888,"rate_of_return":0.09090909090909091}
{"start":2023-02-12T15:28:00.000Z,"end":2023-02-12T15:29:00.000Z,"country":"Germany","ORM":1,"total volume of sales":29.04999937060547,"rate_of_return":0.0}
{"start":2023-02-12T15:26:00.000Z,"end":2023-02-12T15:27:00.000Z,"country":"Belgium","ORM":1,"total volume of sales":78.25999450683594,"rate_of_return":0.0}
{"start":2023-02-12T15:32:00.000Z,"end":2023-02-12T15:33:00.000Z,"country":"United Kingdom","ORM":32,"total volume of sales":1784.989966716766,"rate_of_return":0.03125}
[hadoop@ip-172-31-38-124 ~]$
```

Transfer of files from EMR Instance on AWS to my system, using WinSCP :

First, I needed to transfer the JSON files from HDFS into the the EC2 system I created directories for time-based and then time-and-country-based KPIs as ec2-user.

Using the 'get' command I copied the contents of the output folders into the EC2 system.

mkdir timebased-KPI

hadoop fs -get user/ec2-user/time_kpi /home/hadoop/timebased-KPI

```
[hadoop@ip-172-31-38-124 ~]$ mkdir timebased-KPI
[hadoop@ip-172-31-38-124 ~]$ ls
read_from_kafka.py  spark-streaming.py  timebased-KPI
[hadoop@ip-172-31-38-124 ~]$ hadoop fs -get user/ec2-user/time_kpi /home/ec2-user/timebased-KPI
get: mkdir '/home/ec2-user/timebased-KPI': Input/output error
[hadoop@ip-172-31-38-124 ~]$ pwd
/home/hadoop
[hadoop@ip-172-31-38-124 ~]$ hadoop fs -get user/ec2-user/time_kpi /home/hadoop/timebased-KPI
[hadoop@ip-172-31-38-124 ~]$
```

mkdir country-and-timebased-KPI

hadoop fs -get user/ec2-user/country_kpi /home/hadoop/country-and-timebased-KPI

```
[hadoop@ip-172-31-38-124 ~]$ mkdir country-and-timebased-KPI
[hadoop@ip-172-31-38-124 ~]$ hadoop fs -get user/ec2-user/country_kpi /home/hadoop/country-and-timebased-KPI
[hadoop@ip-172-31-38-124 ~]$
```

Thereafter I used WinSCP to establish a connection between the EC2 instance and my local file system to transfer all the required files into my system.