# Cloud Analytics

# Assignment 2

**Student name:** Ahmed Ibraheem Ahmed Badwy      **Student ID:** 300389393

**Student name:** Anas Ibrahin Ali Elbattra      **Student ID:** 300389368

**Student name:** Esmael Alahmady Ebrahim Ezz      **Student ID:** 300389885

**Student name:** Yousef Abdelhaleem Ahmed Shindy      **Student ID:** 300389365

**University of Ottawa, Ca**

---

By typing in my name and student ID on this form and submitting it electronically, I am attesting to the fact that I have reviewed not only my work but the work of my team member, in its entirety.

I attest to the fact that my work in this project adheres to the fraud policies as outlined in the Academic Regulations in the University's Graduate Studies Calendar. I further attest that I have knowledge of and have respected the "Beware of Plagiarism" brochure for the university. To the best of my knowledge, I also believe that each of my group colleagues has also met the aforementioned requirements and regulations. I understand that if my group assignment is submitted without a completed copy of this Personal Work Statement from each group member, it will be interpreted by the school that the missing student(s) name is confirmation of the nonparticipation of the aforementioned student(s) in the required work.

We, by typing in our names and student IDs on this form and submitting it electronically,

• warrant that the work submitted herein is our own group members' work and not the work of others

• acknowledge that we have read and understood the University Regulations on Academic Misconduct

• acknowledge that it is a breach of University Regulations to give or receive unauthorized and/or unacknowledged assistance on a graded piece of work

# Part 1

**1- GlusterFS:** A distributed file system that provides a global view of a set of files that are spread across multiple servers. scale-out file system that can be scaled up or down by adding or removing servers. fault-tolerant file system which can continue to operate even if one or more servers fail.

- GlusterFS is implemented using a client-server architecture. The client is a FUSE module that allows GlusterFS to be mounted as a regular file system. The server is a daemon that runs on each server in the cluster. The GlusterFS servers communicate with each other to maintain a consistent view of the file system.
- When a client accesses a file on a GlusterFS volume, the client communicates with the servers to determine which server has the most recent copy of the file. The GlusterFS client then reads or writes the file to that server.
- GlusterFS supports a variety of replication and striping strategies. Replication allows GlusterFS to create multiple copies of each file on different servers. This improves the fault tolerance of the file system. Striping allows GlusterFS to divide a file into multiple chunks and store each chunk on a different server. This improves the performance of the file system.

- **Use Cases:**
  - High-performance file sharing
  - Storing and managing large and highly scalable file systems
  - Applications that require high scalability, fault tolerance, and performance

**NFS**: A client-server protocol distributed file system that allows users to access files stored on remote servers as if they were local.

- NFS is a client-server protocol. The client is a program that runs on the client machine and allows the user to access the NFS file system. The server is a program that runs on the server machine and exports the NFS file system to clients.
- When a client wants to access a file on an NFS server, the client makes a request to the server. The server then sends the file to the client. The NFS client then saves the file to the client machine.
- NFS supports a variety of authentication mechanisms, such as Kerberos and NIS. NFS also supports a variety of file locking mechanisms, which ensure that only one client can write to a file at a time.

- **Use Cases**:
  - File sharing between servers and clients
  - Storing and managing small to medium-sized file systems
  - Applications that require simplicity and performance

## 2- Apache Hadoop Map-Reduce:

- **Scalability:** MapReduce is highly scalable, that it can be used to process very large datasets across multiple servers. This is because MapReduce breaks down large datasets into smaller tasks that can be processed in parallel.
- **Fault tolerance:** MapReduce is fault-tolerant, meaning that it can continue to operate even if one or more servers fail. This is because MapReduce replicates data across multiple servers and can reschedule tasks to other servers if a server fails.
- **Simplicity:** MapReduce is a relatively simple programming model, making it easy for developers to write parallel applications. This is because MapReduce takes care of the details of distributing data across servers and managing task failures.

## The Limitations:

- **Performance:** MapReduce is slower than Spark, especially for iterative workloads. This is because MapReduce has to write intermediate results to disk after each map and reduce phase, while Spark can keep intermediate results in memory.
- **Ease of use:** MapReduce is more difficult to program than Spark due to its two-stage programming model. Spark has a more unified programming model that makes it easier to write and debug parallel applications.
- **Versatility:** MapReduce is primarily designed for batch processing, while Spark can be used for both batch and stream processing. Spark also offers a wider range of libraries for machine learning and other data science tasks.

Spark is a more modern and versatile distributed computing framework than MapReduce. However, MapReduce is still a widely used and powerful tool for batch processing large datasets

## 3- Apache Spark:
The main difference between the low-level and high-level APIs is the level of control they offer over the underlying execution of Spark jobs. The low-level API gives users complete control over how their data is processed, while the high-level APIs hide many of the implementation details.
- **Low-level API**: Resilient Distributed Datasets (RDDs), provides a direct way to access and manipulate data in Spark.
  - o **Advantages:**
    - ▪ **Flexibility:** The low-level API provides a more flexible way to represent and manipulate data. For example, RDDs can be used to represent data in a variety of formats, such as lists, tuples, and matrices.
    - ▪ **Performance:** The low-level API can be more efficient for certain workloads, such as iterative workloads. This is because the low-level API gives users more control over how their data is processed.
  - o **Disadvantages:**
    - ▪ **Complexity:** The low-level API can be more difficult to use than the high-level APIs. This is because the low-level API requires users to explicitly manage the partitioning and distribution of their data.
    - ▪ **Error-prone:** The low-level API is more error-prone than the high-level APIs. This is because users are responsible for managing many of the implementation details of Spark jobs.

- o **Examples:**
  - **Machine learning:** The low-level API can be used to implement custom machine learning algorithms. For example, researchers at the University of California, Berkeley used the low-level API to implement a new algorithm for natural language processing.
  - **Real-time stream processing:** The low-level API can be used to implement real-time stream processing applications. For example, the company Foursquare uses the low-level API to process over 1 billion check-ins per day.
  - **Graph processing:** The low-level API can be used to process large graphs. For example, the company Facebook uses the low-level API to process its social graph, which has over 2 billion nodes and 100 trillion edges.

- **High-level API:** DataFrames and Datasets, provide a more abstract view of data and offer a variety of operations for filtering, sorting, aggregating, and joining data.

  - o **Advantages:**
    - **Ease of use:** The high-level APIs are easier to use than the low-level API. This is because the high-level APIs hide many of the implementation details of Spark jobs.
    - **Performance:** The high-level APIs are optimized for performance for most workloads. This is because the high-level APIs use a variety of techniques to improve performance, such as caching and code generation.
    - **Safety:** The high-level APIs are safer to use than the low-level API. This is because the high-level APIs handle many of the potential pitfalls of Spark programming, such as data skew and memory management.

  - o **Disadvantages:**
    - **Flexibility:** The high-level APIs are not as flexible as the low-level API. For example, DataFrames and Datasets can only be used to represent data in a structured format.
    - **Performance:** The high-level APIs may not be as efficient for certain workloads, such as iterative workloads. This is because the high-level APIs hide many of the implementation details of Spark jobs.

  - o **Examples:**
    - **Data warehousing:** The high-level API can be used to build data warehouses. For example, the company Airbnb uses the high-level API to build its data warehouse, which stores over 100 petabytes of data.
    - **Data analytics:** The high-level API can be used to perform data analytics on large datasets. For example, the company Netflix uses the high-level API to analyze its customer data to improve its recommendation system.
    - **Internet of Things (IoT) processing:** The high-level API can be used to process data from IoT devices. For example, the company General Electric uses the high-level API to process data from its jet engines to predict maintenance needs.

**4-**

**A- Immutability:** is a property of an object that means it cannot be changed once it has been created. This is in contrast to mutability, which allows objects to be changed after they have been created.

- o **Advantages:**
  - ▪ **Predictability:** Immutable objects are more predictable because they cannot be changed after they have been created. This makes it easier to reason about the behavior of code that uses immutable objects.
  - ▪ **Thread safety:** Immutable objects are thread safe by default, because they cannot be changed by multiple threads at the same time. This makes them ideal for use in multithreaded applications.
  - ▪ **Performance:** Immutable objects can be more performant than mutable objects, because they do not need to be synchronized or copied when they are shared between threads.

- o **Examples:**
  - ▪ **The hash of a file:** is immutable This means that the hash of a file will always be the same, regardless of when the file was created or modified. This makes hashes useful for verifying the integrity of files.
  - ▪ **The blockchain:** is immutable This means that once a transaction is added to the blockchain, it cannot be changed. This makes the blockchain secure and tamper-proof.

**B- Lazy Evaluation:** is a technique where Spark does not execute transformations on a dataset until an action is called. This means that Spark can build a logical execution plan, called the Directed Acyclic Graph (DAG), incrementally as transformations are applied. Spark can then optimize the DAG before executing it, which can lead to significant performance improvements.

- o **Positive Impacts:**

  - ▪ **Reduced memory usage:** Spark only needs to keep the data in memory that is needed for the current transformation. This can significantly reduce memory usage, especially for large datasets.
  - ▪ **Improved performance:** Spark can optimize the DAG before executing it, which can lead to significant performance improvements. For example, Spark can fuse multiple transformations together or reorder transformations to reduce data shuffling.
  - ▪ **Increased flexibility:** Lazy evaluation allows you to experiment with different transformations without having to commit to a specific execution plan. This can be helpful for debugging and performance tuning your Spark applications.

- o **Negative Impacts:**
  - ▪ **Debugging:** Lazy evaluation can make debugging more difficult because it can be difficult to track down the source of an error, especially if the error occurs in a transformation that has not yet been executed.
  - ▪ **Performance:** Lazy evaluation can sometimes lead to performance problems, especially if the DAG is not optimized correctly. For example, if Spark chooses to shuffle the data between two transformations, this can significantly reduce performance.

- o **Examples:**
  - ▪ **Filtering**
  - ▪ **Sorting**
  - ▪ **Aggregating**
  - ▪ **Joining**
  - ▪ **Mapping**
  - ▪ **Reducing**

**C-**
**SparkContext:** Is a central entry point for Spark functionality. It is responsible for coordinating the execution of tasks across a cluster, managing the distribution of data, and scheduling jobs. The SparkContext is the foundation upon which all other Spark functionality is built, and it is the first thing that needs to be created when starting a Spark application.
- o The SparkContext is mainly used for lower-level operations, such as creating RDDs, accumulators, and broadcasting variables.
- o The SparkContext is a singleton and can only be created once in a Spark application.
- o The SparkContext is created using the SparkConf, which allows you to set various Spark configurations.
- o The SparkContext provides methods for creating RDDs, accumulators, and broadcasting variables, as well as methods for starting tasks on the executors.
- o The SparkContext is used to access the Spark environment and perform operations on it.

**SparkSession:** is the entry point to the Spark functionality. It is the main starting point for interacting with data and creating a DataFrame. SparkSession is the new entry point of Spark that replaces the old SQLContext and HiveContext. It is a unified entry point to interact with structured and semi-structured data and it also provides support for reading from and writing to a variety of data sources such as Avro, Parquet, JSON, and JDBC.
- o The SparkSession is the preferred way to work with Spark data structures such as DataFrames and Datasets, as it provides a more consistent and simpler interface.
- o The SparkSession, can be created multiple times within an application.
- o The SparkSession, does not have a corresponding configuration object, but you can set configurations using the **.config** method of the SparkSession.
- o The SparkSession provide methods for creating DataFrames and Datasets, as well as methods for reading and writing data.
- o The SparkSession, is used to access the data stored in Spark and perform operations on it.

**D- Spark MLlib:** is a machine learning library for Apache Spark. It provides a variety of algorithms for classification, regression, clustering, and collaborative filtering. Spark MLlib algorithms are designed to be scalable and fault-tolerant, making them ideal for processing large datasets.
- o **Transformers: A**re algorithms that transform one DataFrame into another DataFrame. For example, a Transformer could be used to tokenize a text column, scale a numerical column, or encode a categorical column.
  - ▪ **Tokenizer**
  - ▪ **HashingTF**
  - ▪ **IDF**
  - ▪ **StringIndexer**
  - ▪ **OneHotEncoder**
  - ▪ **StandardScaler**
  - ▪ **MinMaxScaler**

- o **Estimators:** are algorithms that fit a model to a DataFrame. The model can then be used to make predictions on new data. For example, an Estimator could be used to train a logistic regression model to predict whether a customer is likely to churn, or a random forest model to predict the price of a house.
    - LogisticRegression
    - RandomForestClassifier
    - RandomForestRegressor
    - KMeans
    - LinearRegression
    - DecisionTreeClassifier
    - DecisionTreeRegressor
- o **Evaluators:** are algorithms that evaluate the performance of a machine learning model on a held-out test set. For example, an Evaluator could be used to calculate the accuracy, precision, recall, and F1 score of a classification model.
    - MulticlassClassificationEvaluator
    - BinaryClassificationEvaluator
    - RegressionEvaluator

**5- Spark Streaming:** Is a batch processing framework that is designed to process streaming data as if it were a batch of data. Spark Streaming divides the stream of data into batches and then processes each batch using Spark's batch processing engine.
**For example:** Spark Streaming can be used to process web traffic data, social media data, and sensor data.
**Spark Structured Streaming:** is a real-time stream processing framework that is designed to process streaming data as it arrives. Spark Structured Streaming does not divide the stream of data into batches. Instead, Spark Structured Streaming keeps the entire stream of data in memory and processes it as it arrives.
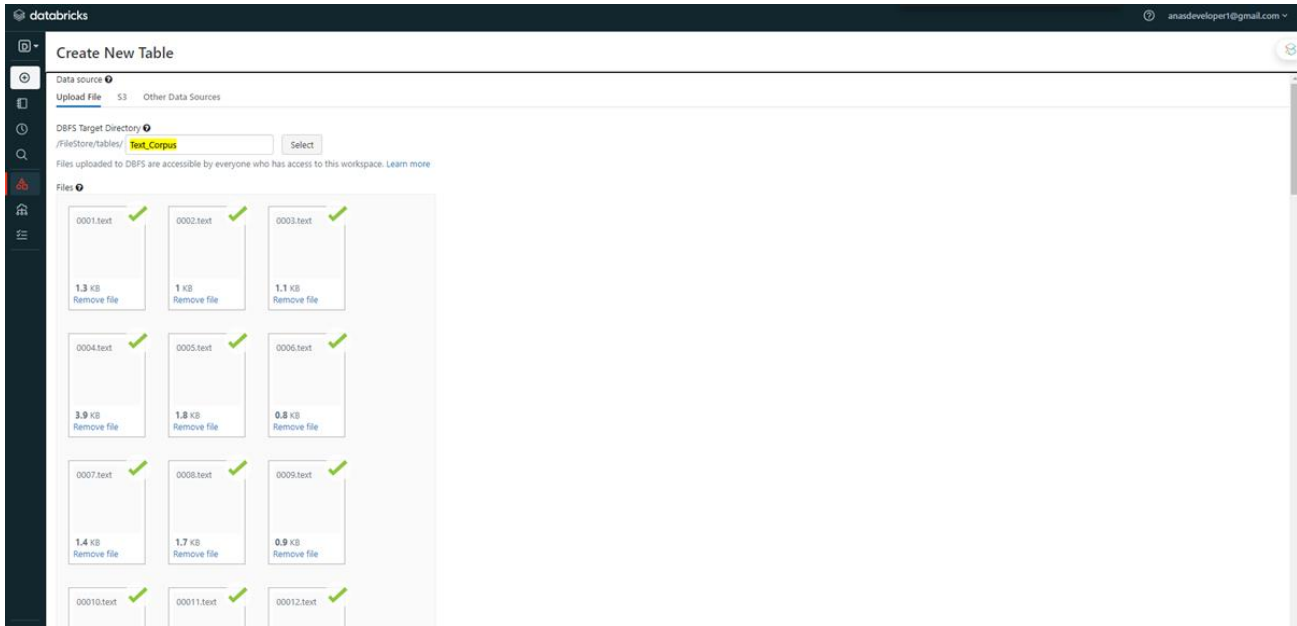**S**park Structured Streaming is a good choice for processing streaming data that is highly stateful.
**For example**, Spark Structured Streaming can be used to process streaming data from financial markets, IoT devices, and gaming applications.

# Part 2

## 1) Data Transformation Pipelines

### a) Uploaded the files to your DBFS table space.



### b) Use Spark Scala to load your data into an RDD.

```
Cmd 1

1    val dbfsDirectoryPath = "/FileStore/tables/Text_Corpus/*.text"
2    val dataRDD = spark.sparkContext.textFile(dbfsDirectoryPath)

dbfsDirectoryPath: String = /FileStore/tables/Text_Corpus/*.text
dataRDD: org.apache.spark.rdd.RDD[String] = /FileStore/tables/Text_Corpus/*.text MapPartitionsRDD[11] at textFile at command-3383544309757279:2
Command took 0.44 seconds -- by anasdeveloper1@gmail.com at 10/31/2023, 4:47:37 PM on Text Corpus cluster
```

### c) Count the number of lines across all the files.

```
Cmd 2

1    dataRDD.count()

▶ (1) Spark Jobs

res4: Long = 1645

Command took 3.62 seconds -- by anasdeveloper1@gmail.com at 10/31/2023, 4:47:40 PM on Text Corpus cluster
```

**d) Find the number of occurrences of the word "antibiotics".**

```scala
1   %scala
2   val keyword = "antibiotics"
3   |
4   val countOfAntibiotics = dataRDD
5     .flatMap(line => line.split(" "))
6     .filter(word => word == keyword)
7     .count()
8
9   println(s"The word '$keyword' appears $countOfAntibiotics times.")
```

▶ (1) Spark Jobs

```
The word 'antibiotics' appears 2 times.
keyword: String = antibiotics
countOfAntibiotics: Long = 2
```

Command took 4.68 seconds -- by anasdeveloper1@gmail.com at 10/31/2023, 4:47:47 PM on Text Corpus cluster

**e) Count the occurrence of the word "patient" and "admitted" on the same line of text. Please ensure that your code contains at least 2 transformation functions in a pipeline.**

```scala
1    %scala
2    val keyword1 = "patient"
3    val keyword2 = "admitted"
4
5    val countOfBothWords = dataRDD
6      .map(line => (line, line.split(" ")))
7      .filter { case (line, words) => words.contains(keyword1) && words.contains(keyword2) }
8      .count()
9
10   println(s"The words '$keyword1' and '$keyword2' appear together on the same line $countOfBothWords times.")
11
```

▶ (1) Spark Jobs

```
The words 'patient' and 'admitted' appear together on the same line 7 times.
keyword1: String = patient
keyword2: String = admitted
countOfBothWords: Long = 7
```
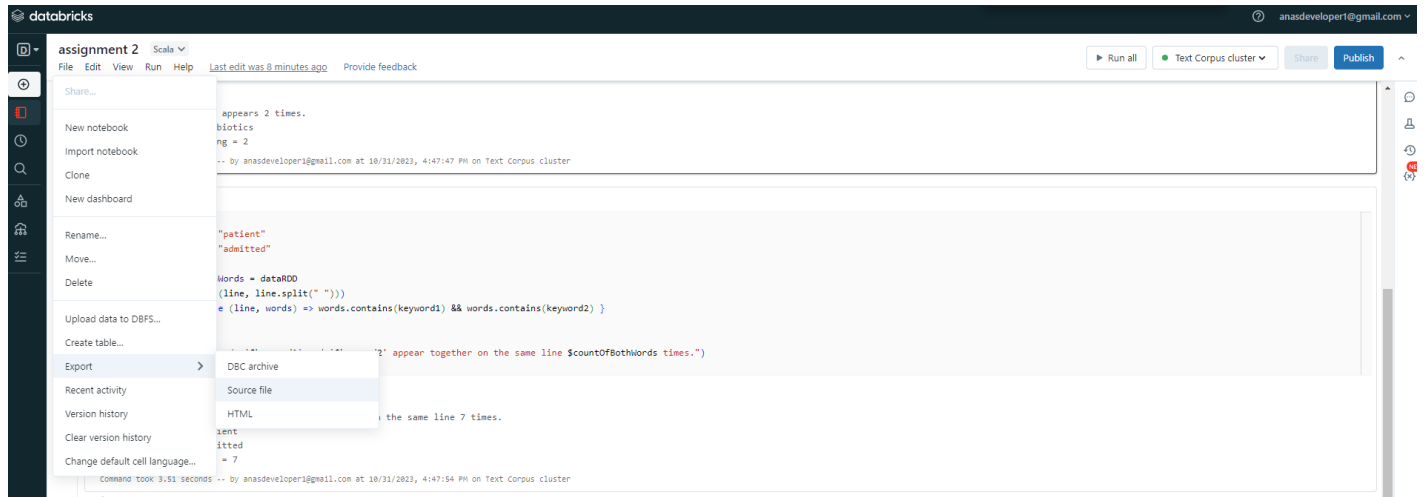
Command took 3.51 seconds -- by anasdeveloper1@gmail.com at 10/31/2023, 4:47:54 PM on Text Corpus cluster
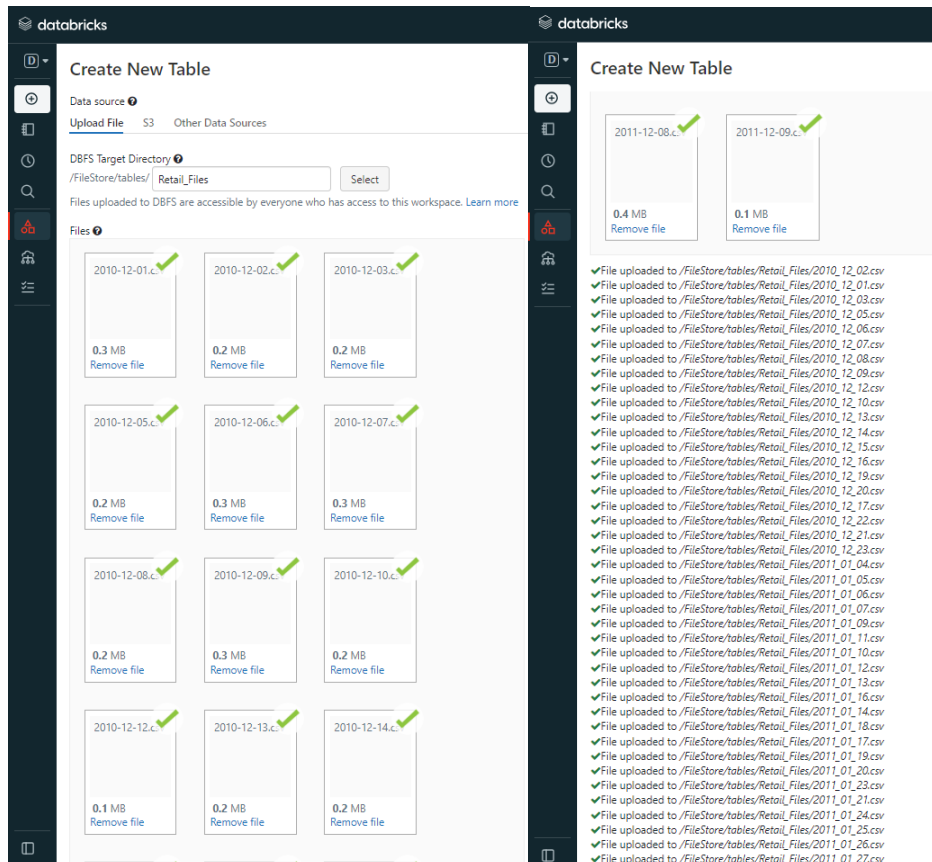
Cmd 5

f) **Upload your exported (.scala format) notebook as part of your submission. Please include the data in your Brightspace submission.**



## 2) Retail Data Analysis.

a) **Uploaded the files to your DBFS table space.**

**b) Output the total number of transactions across all the files and the total value of the transactions.**

```
1   df_count=spark.sql('select COUNT(DISTINCT InvoiceNo)  AS TotalTransactions from retail_data_tabel')
```

▸ ▤ df_count: pyspark.sql.dataframe.DataFrame = [TotalTransactions: long]

Command took 0.08 seconds -- by ahmedeali696@gmail.com at 10/31/2023, 8:41:20 PM on cloud_task

```
1   df_count.show(truncate=False)
```

▸ (3) Spark Jobs

```
+-----------------+
|TotalTransactions|
+-----------------+
|25900            |
+-----------------+
```

Command took 25.76 seconds -- by ahmedeali696@gmail.com at 10/31/2023, 8:41:20 PM on cloud_task

```
1   df_total_value=spark.sql('select FORMAT_NUMBER(SUM(Quantity * UnitPrice), 2) AS TotalValue from retail_data_tabel')
2   |
3
4
```

▸ ▤ df_total_value: pyspark.sql.dataframe.DataFrame = [TotalValue: string]

Command took 0.18 seconds -- by ahmedeali696@gmail.com at 10/31/2023, 8:41:21 PM on cloud_task

```
1   df_total_value.show(truncate=False)
2
```

▸ (2) Spark Jobs

```
+------------+
|TotalValue  |
+------------+
|9,747,747.93|
+------------+
```

Command took 15.55 seconds -- by ahmedeali696@gmail.com at 10/31/2023, 8:41:21 PM on cloud_task

## c) Output the 5 top-selling products.

```
1   df_5_top_selling = spark.sql('SELECT StockCode, SUM(Quantity) AS TotalSold ' +
2                                 'FROM retail_data_tabel ' +
3                                 'GROUP BY StockCode ' +
4                                 'ORDER BY TotalSold DESC ' +
5                                 'LIMIT 5')
```

▶ ▦ df_5_top_selling: pyspark.sql.dataframe.DataFrame = [StockCode: string, TotalSold: double]

Command took 0.19 seconds -- by ahmedeali696@gmail.com at 10/31/2023, 8:41:22 PM on cloud_task

```
1   df_5_top_selling.show(truncate=False)
```

▶ (2) Spark Jobs

```
+---------+---------+
|StockCode|TotalSold|
+---------+---------+
|22197    |56450.0  |
|84077    |53847.0  |
|85099B   |47363.0  |
|85123A   |38830.0  |
|84879    |36221.0  |
+---------+---------+
```

Command took 16.60 seconds -- by ahmedeali696@gmail.com at 10/31/2023, 8:41:22 PM on cloud_task

## d) Output the 5 topmost valuable products.

```
1   df_5_topmost_valuable = spark.sql('SELECT StockCode, Description,CAST(SUM(Quantity * UnitPrice)AS DECIMAL(10, 2))  AS TotalValue ' +
2                                 'FROM retail_data_tabel ' +
3                                 'GROUP BY StockCode, Description ' +
4                                 'ORDER BY TotalValue DESC ' +
5                                 'LIMIT 5')
```

▶ ▦ df_5_topmost_valuable: pyspark.sql.dataframe.DataFrame = [StockCode: string, Description: string ... 1 more field]

Command took 0.13 seconds -- by ahmedeali696@gmail.com at 10/31/2023, 8:51:50 PM on cloud_task

```
1    df_5_topmost_valuable.show(truncate=False)
```

▶ (2) Spark Jobs

```
+---------+------------------------------+----------+
|StockCode|Description                   |TotalValue|
+---------+------------------------------+----------+
|DOT      |DOTCOM POSTAGE                |206245.48 |
|22423    |REGENCY CAKESTAND 3 TIER      |164762.19 |
|47566    |PARTY BUNTING                 |98302.98  |
|85123A   |WHITE HANGING HEART T-LIGHT HOLDER|97715.99  |
|85099B   |JUMBO BAG RED RETROSPOT       |92356.03  |
+---------+------------------------------+----------+
```

Command took 17.69 seconds -- by ahmedeali696@gmail.com at 10/31/2023, 8:52:04 PM on cloud_task

## e)  Output each country and the total value of their purchases.

```
1    df_country_and_totalValue = spark.sql('SELECT Country, CAST(SUM(Quantity * UnitPrice) AS DECIMAL(10, 2)) AS TotalValue ' +
2                                          'FROM retail_data_tabel ' +
3                                          'GROUP BY Country ' )
4
```

▶ 🔲 df_country_and_totalValue: pyspark.sql.dataframe.DataFrame = [Country: string, TotalValue: decimal(10,2)]

Command took 0.09 seconds -- by ahmedeali696@gmail.com at 10/31/2023, 9:18:54 PM on cloud_task

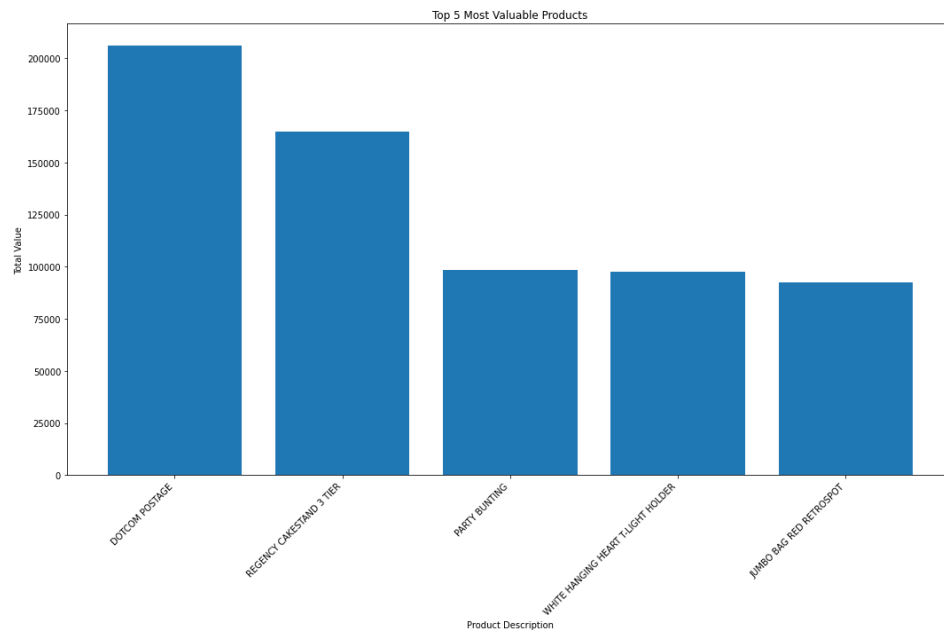```
1    df_country_and_totalValue.show(truncate=False)
```

▶ (2) Spark Jobs

```
+---------------+----------+
|Country        |TotalValue|
+---------------+----------+
|Sweden         |36595.91  |
|Germany        |221698.21 |
|France         |197403.90 |
|Greece         |4710.52   |
|Belgium        |40910.96  |
|Finland        |22326.74  |
|Malta          |2505.47   |
|Unspecified    |4749.79   |
|Italy          |16890.51  |
|EIRE           |263276.82 |
|Norway         |35163.46  |
|Spain          |54774.58  |
|Denmark        |18768.14  |
|Hong Kong      |10117.04  |
|Iceland        |4310.00   |
|Channel Islands|20086.29  |
|USA            |1730.92   |
|Switzerland    |56385.35  |
```

Command took 25.40 seconds -- by ahmedeali696@gmail.com at 10/31/2023, 9:18:56 PM on cloud_task

**f) Use a graphical representation to describe the result from step (d).**
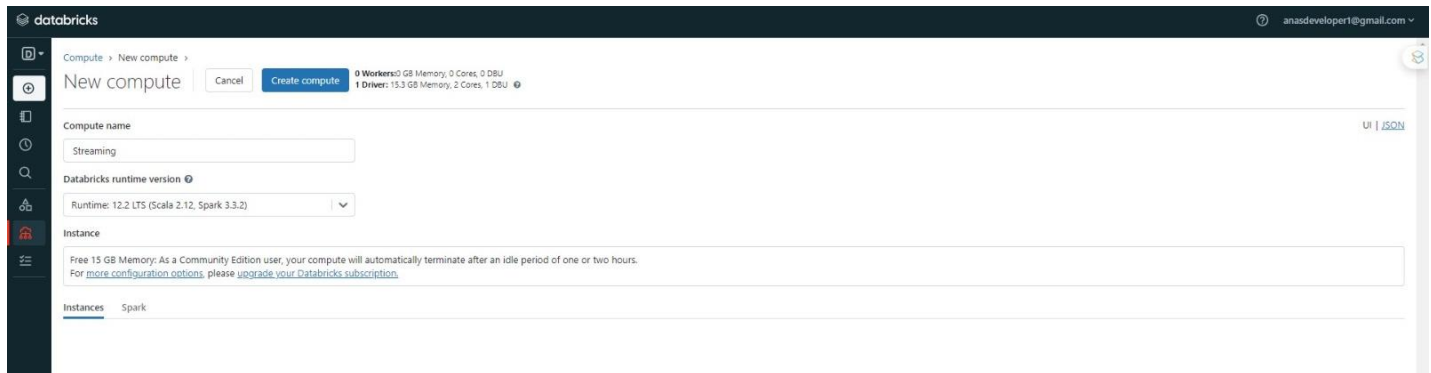
```
1
2    import matplotlib.pyplot as plt
3
4    df_pandas = df_5_topmost_valuable.toPandas()
5
6
7    plt.figure(figsize=(15, 10))
8    plt.bar(df_pandas['Description'], df_pandas['TotalValue'])
9    plt.xlabel('Product Description')
10   plt.ylabel('Total Value')
11   plt.title('Top 5 Most Valuable Products')
12   plt.xticks(rotation=45, ha="right")
13   plt.tight_layout()
14   plt.show()
15
```



Top 5 Most Valuable Products

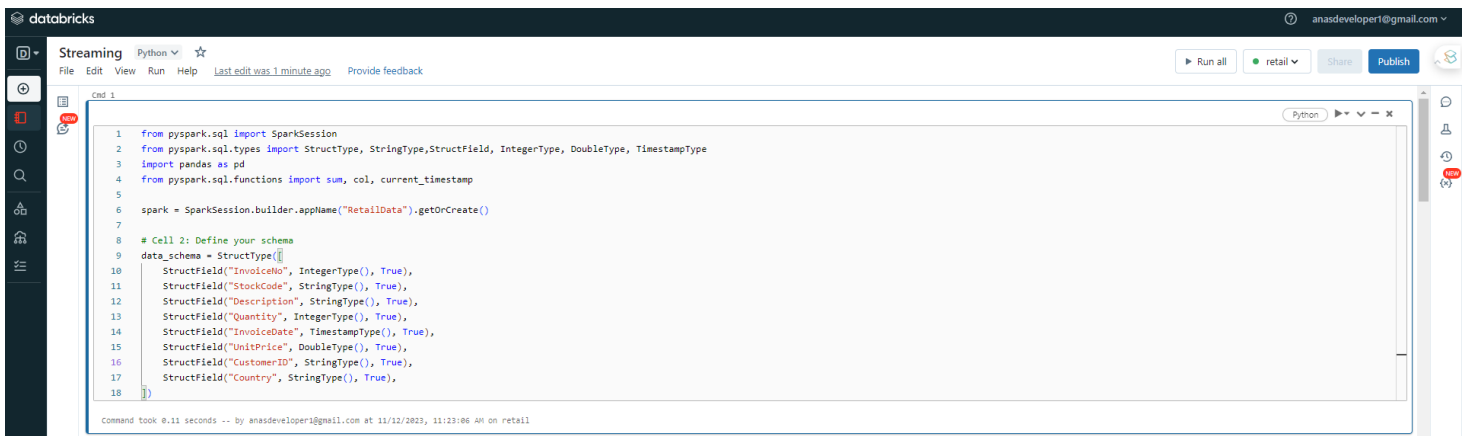Command took 19.77 seconds -- by ahmedeali696@gmail.com at 10/31/2023, 9:29:51 PM on cloud_task

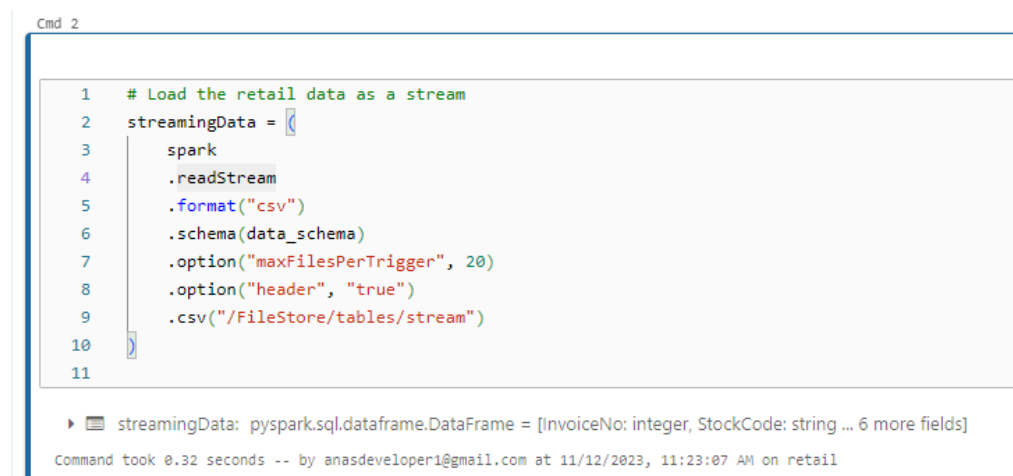## 3) Structured Streaming.

### a) Create a new notebook.



### b) Load the retail data as a stream, at 20 files per trigger. For each batch pulled, capture the customer stock aggregates – total stocks, total value.

- **Define your schema.**



- **Load the retail data as a stream.**

- **Prepare the data.**

```
2    transformed_data = (
3        streamingData
4        .selectExpr(
5            "InvoiceNo",
6            "StockCode",
7            "Description",
8            "cast(Quantity as int) as Quantity",
9            "cast(UnitPrice as double) as UnitPrice",
10           "cast(CustomerID as int) as CustomerID",
11           "Country"
12       )
13   )
14
```

▶ 🖾 transformed_data: pyspark.sql.dataframe.DataFrame = [InvoiceNo: integer, StockCode: string ... 5 more fields]

Command took 0.14 seconds -- by anasdeveloper1@gmail.com at 11/12/2023, 11:23:08 AM on retail

- **Capture customer stock aggregates.**

```
3    aggregated_data = (
4        transformed_data
5        .groupBy("CustomerID")
6        .agg(
7            sum("Quantity").alias("TotalStocks"),
8            sum((col("Quantity") * col("UnitPrice"))).alias("TotalValue")
9        )
10   )
```

▶ 🖾 aggregated_data: pyspark.sql.dataframe.DataFrame = [CustomerID: integer, TotalStocks: long ... 1 more field]

Command took 0.16 seconds -- by anasdeveloper1@gmail.com at 11/12/2023, 11:23:09 AM on retail

- **Start the streaming query and get the results.**

```
2    df = (
3        aggregated_data
4        .writeStream
5        .outputMode("complete")
6        .format("memory")
7        .queryName("aggregated_data")
8        .start()
9    )
10
```

▶ (1) Spark Jobs

▼ ⊘ aggregated_data (id: 8701ec74-9b26-4615-9673-14097a8f7a67)    Last updated: 14 minutes ago

Dashboard    Raw Data



Command complete

- **Display the aggregated data.**

```
2    df = spark.sql("SELECT * FROM aggregated_data")
3    df.show(truncate=False)
```

▶ (1) Spark Jobs

▶ ▦ df: pyspark.sql.dataframe.DataFrame = [CustomerID: integer, TotalStocks: long ... 1 more field]

```
+----------+-----------+--------------------+
|CustomerID|TotalStocks|TotalValue          |
+----------+-----------+--------------------+
|17420     |53         |130.85              |
|14570     |95         |218.05999999999992  |
|17389     |2214       |9447.24             |
|15727     |1203       |2876.0500000000006  |
|16861     |-7         |-22.110000000000003 |
|16503     |443        |1082.8899999999999  |
|13285     |914        |1173.2699999999998  |
|15619     |136        |336.40000000000003  |
|13623     |148        |532.0200000000001   |
|14450     |152        |269.45              |
|15790     |114        |220.84999999999997  |
|15447     |85         |155.17              |
|16339     |23         |109.95000000000002  |
|16386     |26         |61.10000000000001   |
|15100     |58         |635.0999999999999   |
|18161     |315        |540.4200000000001   |
|18282     |28         |77.84               |
|16500     |19         |217.74              |
```

Command took 6.69 seconds -- by anasdeveloper1@gmail.com at 11/12/2023, 11:24:37 AM on retail

c) **For each batch of the input stream, create a new stream that populates another dataframe or dataset with progress for each loaded set of data. This data set should have the columns – TriggerTime (Date/Time), Records Imported, Sale value (Total value of transactions).**

- **Start the Stream**

```python
1   # Initialize DataFrame to store progress information
2   progress_data = pd.DataFrame(columns=["TriggerTime", "RecordsImported", "SaleValue"])
```

Command took 0.12 seconds -- by anasdeveloper1@gmail.com at 11/12/2023, 11:27:10 AM on retail

Cmd 8

```python
1   # Define the function process_batch
2   def process_batch(batch_df, batch_id):
3       total_value = batch_df.selectExpr("sum(Quantity * UnitPrice) as TotalValue").collect()[0]["TotalValue"]
4       current_time_value = batch_df.select(current_timestamp().cast(TimestampType())).collect()[0][0]
5       current_time_python = current_time_value.replace(microsecond=0)
6       progress_df = spark.createDataFrame([(current_time_python, batch_df.count(), total_value)],
7                                            ["TriggerTime", "RecordsImported", "SaleValue"])
8       progress_df.show(truncate=False)
9
10      # Convert the progress_df to pandas and concatenate
11      global progress_data
12      progress_data = pd.concat([progress_data, progress_df.toPandas()])
13
14
15
```

Command took 0.13 seconds -- by anasdeveloper1@gmail.com at 11/12/2023, 11:27:15 AM on retail

```
2   query = (
3        streamingData
4        .writeStream
5        .outputMode("append")
6        .foreachBatch(process_batch)
7        .start()
8   )
9   query.awaitTermination()
10
```

```
+-------------------+---------------+-----------------+
|TriggerTime        |RecordsImported|SaleValue        |
+-------------------+---------------+-----------------+
|2023-11-12 09:27:18|61688          |1018540.1600000046|
+-------------------+---------------+-----------------+


+-------------------+---------------+-----------------+
|TriggerTime        |RecordsImported|SaleValue        |
+-------------------+---------------+-----------------+
|2023-11-12 09:27:24|32867          |581157.0500000012|
+-------------------+---------------+-----------------+


+-------------------+---------------+-----------------+
|TriggerTime        |RecordsImported|SaleValue        |
+-------------------+---------------+-----------------+
|2023-11-12 09:27:28|27444          |448132.1999999975|
+-------------------+---------------+-----------------+


+-------------------+---------------+-----------------+
|TriggerTime        |RecordsImported|SaleValue        |
+-------------------+---------------+-----------------+
  Cancelled
Command took 1.09 minutes -- by anasdeveloper1@gmail.com at 11/12/2023, 11:27:15 AM on retail
```
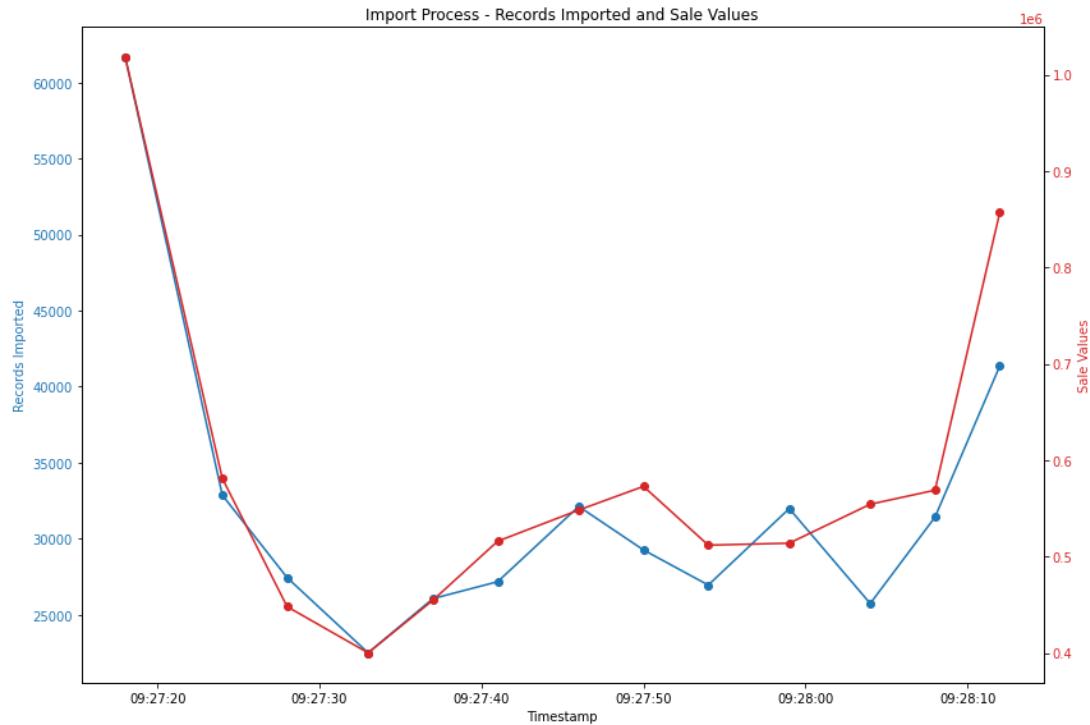
**d) Use the dataset from step (c) to plot a line graph of the import process – showing two timelines – records imported and sale values.**

```python
1    import matplotlib.pyplot as plt
2    fig, ax1 = plt.subplots(figsize=(12, 8))
3    color = 'tab:blue'
4    ax1.set_xlabel('Timestamp')
5    ax1.set_ylabel('Records Imported', color=color)
6    ax1.plot(progress_data['TriggerTime'], progress_data['RecordsImported'], color=color, marker='o')
7    ax1.tick_params(axis='y', labelcolor=color)
8    ax2 = ax1.twinx()
9    color = 'tab:red'
10   ax2.set_ylabel('Sale Values', color=color)
11   ax2.plot(progress_data['TriggerTime'], progress_data['SaleValue'], color=color, marker='o')
12   ax2.tick_params(axis='y', labelcolor=color)
13   fig.tight_layout()
14   plt.title('Import Process - Records Imported and Sale Values')
15   plt.show()
```

Import Process - Records Imported and Sale Values

Command took 0.53 seconds -- by anasdeveloper1@gmail.com at 11/12/2023, 11:43:17 AM on retail

## References:

**[1]** https://hadoop.apache.org

**[2]** https://hadoop.apache.org/docs/stable/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html

**[3]** https://medium.com/analytics-vidhya/spark-rdd-low-level-api-basics-using-pyspark-a9a322b58f6

**[4]** https://medium.com/analytics-vidhya/sparksql-and-dataframe-high-level-api-basics-using-pyspark-eaba6acf944b

**[5]** https://spark.apache.org/docs/latest/api/python/index.html