Retail Business Analytics.

Description

Customers can purchase products or services from Amazon for consumption and usage. Amazon usually sells products and services in-store; however, some may be sold online or over the phone and shipped to the customer. Clothing, medicine, supermarket, and convenience stores are examples of their retail operations.

Objective:

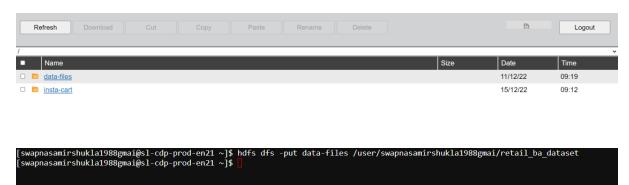
The objective is to analyze the "retail_db" dataset, provide reports on the total completed orders, and perform customer and product analytics.

Step 1: Upload the "data-files" file to the HDFS

- 1.1 Download the relevant dataset from the "Course Resources" section or the project description
- 1.2 Upload the dataset to the "FTP" lab from your local system
- 1.3 To move the dataset to "HDFS" from the "Webconsole" use the put command

Commands:

hdfs dfs -put data-files /user/swapnasamirshukla1988gmai/retail_ba_dataset



Step 2: Perform the below tasks on the uploaded dataset using PySpark:

Login to the Pyspark shell

Task 2.2:

Explore the customer records saved in the "customers-tab-delimited" directory on HDFS REQUIREMENT:

- 2.2.1 Show the client information for those who live in California
- 2.2.2 The final output must be in text format
- 2.2.3 Save the results in the result/scenario1/solution folder
- 2.2.4 Only records with the state value "CA" should be included in the result
- 2.2.5 Only the customer's entire name should be included in the output

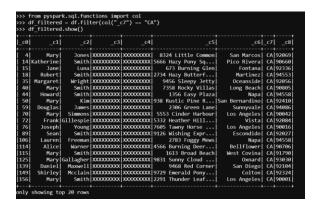
Example: "Robert Hudson"

Read the file, assuming it is tab-separated and does not have a header

```
file_path = "retail_ba_dataset/customers-tab-delimited/part-m-00000" df = spark.read.option("header", "false").option("sep", "\t").csv(file_path) df.printSchema() df.show(10)
```

```
file path = "retail ba dataset/customers-tab-delimited/part-m-00000
>>> df = spark.read.option("header", "false").option("sep",
                                                           "\t").csv(file_path)
>>> df.show(10)
                                                           _c5|
        _c1
                  _c2|
                            _c3|
                                      _c4
                                                                                 _c8|
 c0|
                                                                       _c6|_c7|
  1 | Richard | Hernandez | XXXXXXXXX | XXXXXXXXX
                                            6303 Heather Plaza Brownsville TX 78521
  2
       Mary
              Barrett | XXXXXXXXX | XXXXXXXXX | 9526 Noble Embers...
                                                                 Littleton|
                                                                            CO | 80126
  3
        Ann
                Smith XXXXXXXXX XXXXXXXXX 3422 Blue Pioneer...
                                                                    Caguas |
                                                                            PR | 00725
       Mary
                Jones XXXXXXXXX XXXXXXXX
                                           8324 Little Common
                                                                San Marcos
                                                                            CA 192069
  4
     Robert
               Caguas
                                                                            PR
                                                                               00725
                Smith XXXXXXXXX XXXXXXXXX 3151 Sleepy Quail...
                                                                            NJ 07055
  6
       Mary
                                                                   Passaic
  7|Melissa|
               Wilcox XXXXXXXXX XXXXXXXXX 9453 High Concession
                                                                    Caguas |
                                                                            PR | 00725
  8
                Smith XXXXXXXXX XXXXXXXXX 3047 Foggy Forest...
                                                                  Lawrence|
                                                                            MA | 01841
      Megan
                Perez | XXXXXXXXX | XXXXXXXXX
                                          3616 Quaking Street
                                                                            PR | 00725
  9
       Mary
                                                                    Caguas
                Smith XXXXXXXXX XXXXXXXXX 8598 Harvest Beac...
 10 Melissa
                                                                  Stafford
                                                                            VA 22554
only showing top 10 rows
```

from pyspark.sql.functions import col # Show people residing in California df_filtered = df.filter(col("_c6") == "CA") # Show the filtered DataFrame df_filtered.show()



```
from pyspark.sql.functions import concat_ws, col

# Concat First and second names to get full name
full_name = concat_ws(" ", col("_c1"), col("_c2"))

# Filter for clients in California (assuming _c7 is the state column)
ca_clients = df.filter(col("_c7") == "CA").select(full_name.alias("Full Name"))

# Show the result
ca_clients.show()
```

```
>>> from pyspark.sql.functions import concat_ws, col
>>>
>>> full_name = concat_ws(" ", col("_c1"), col("_c2"))
>>> ca_clients = df.filter(col("_c7") == "CA").select(full_name.alias("Full Name"))
>>> ca_clients.show()
       Full Name
      Mary Jones
Katherine Smith
        Jane Luna
    Robert Smith
Margaret Wright
      Mary Smith
    Howard Smith
        Mary Kim
   Douglas James
    Mary Simmons
Frank Gillespie
    Joseph Young
      Sean Smith
  Lauren Freeman
    Alice Warner
      Mary Smith
  Mary Gallagher
  Daniel Maxwell
 Shirley Mcclain
      Mary Smith
only showing top 20 rows
```

```
# Specify the path where you want to save the results
output_path = "/user/swapnasamirshukla1988gmai/retail_ba/result/scenario1/solution"

# Save the DataFrame as a text file
# coalesce(1) is used to save the output into a single file.
ca_clients.coalesce(1).write.mode("overwrite").text(output_path)
```

```
>>> output_path = "/user/swapnasamirshukla1988gmai/retail_ba/result/scenario1/solution"
>>> ca_clients.coalesce(1).write.mode("overwrite").text(output_path)
>>> |
```

hdfs dfs -ls /user/swapnasamirshukla1988gmai/retail_ba/result/scenario1/solution

The file is stored as part-00000-fbf80159-9c66-4aa4-9c6f-879289198745-c000.txt # Rename file using my command

 $hdfs\ dfs\ -mv\ /user/swapnasamirshukla 1988gmai/retail_ba/result/scenario 1/solution/part-00000-fbf80159-9c66-4aa4-9c6f-879289198745-c000.txt$

/user/swapnasamirshukla1988gmai/retail_ba/result/scenario1/solution/retai_sol1.txt

Task 2.3:

Explore the order records saved in the "orders parquet" directory on HDFS

REQUIREMENT:

- 2.3.1 Show all orders with the order status value "COMPLETE"
- 2.3.2 The output should be in JSON format
- 2.3.3 Save the data in the "result/scenario2/solution" directory on HDFS
- 2.3.4 The "order date" column should be in the "YYYY-MM-DD" format
- 2.3.5 Use GZIP compression to compress the output
- 2.3.6 Only the column names listed below should be included in the output:
 - 2.3.6.1 Order number
 - 2.3.6.2 Order date
 - 2.3.6.3 Current situation

file_path = "retail_ba_dataset/orders_parquet/741ca897-c70e-4633-b352-5dc3414c5680.parquet" df = spark.read.parquet (file_path) df.printSchema() df.show(10)

```
>>> file_path = "retail_ba_dataset/orders_parquet/741ca897-c70e-4633-b352-5dc3414c5680.parquet"
>>> df = spark.read.parquet (file_path)
24/02/04 17:13:38 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
24/02/04 17:13:38 WARN shortcircuit.DomainSocketFactory: The short-circuit local reads feature cannot be used because libhadoop cannot be loaded.
>>> df.printSchema()
root
|-- order_id: integer (nullable = true)
|-- order_date: long (nullable = true)
|-- order_date: long (nullable = true)
|-- order_customer_id: integer (nullable = true)
|-- order_status: string (nullable = true)
|-- order_date| order_customer_id| order_status
|-- order_date| order_customer_id| order_status
|-- | 1 | 1374710400000 | 1159 | CLOSED |
|-- | 2 | 1374710400000 | 12111 | COMPLETE |
|-- | 4 | 1374710400000 | 1318 | COMPLETE |
|-- | 4 | 1374710400000 | 1318 | COMPLETE |
|-- | 6 | 1374710400000 | 250 | PENDING_PAYMENT |
|-- | 7 | 1374710400000 | 2511 | PROCESSING |
|-- | 9 | 1374710400000 | 563 | PENDING_PAYMENT |
|-- | 10 | 1374710400000 | 5648 | PENDING_PAYMENT |
|-- | 10 | 1374710400000 | 5648 | PENDING_PAYMENT |
|-- | 10 | 1374710400000 | 5648 | PENDING_PAYMENT |
|-- | 10 | 1374710400000 | 5648 | PENDING_PAYMENT |
|-- | 10 | 1374710400000 | 5648 | PENDING_PAYMENT |
```

Filter for orders with the status "COMPLETE" df filtered = df.filter(df["order status"] == "COMPLETE")

```
>>> df_filtered = df.filter(df["order_status"] == "COMPLETE")
>>> df_filtered.show(10)
order_id
             order_date|order_customer_id|order_status|
                                       12111
        3 | 1374710400000 |
                                                 COMPLETE
        5 | 1374710400000
                                       11318
                                                 COMPLETE
        6 | 1374710400000 |
                                       7130
                                                 COMPLETE
        7 | 1374710400000 |
                                        4530
                                                 COMPLETE
       15 1374710400000
                                        2568
                                                 COMPLETE
       17 | 1374710400000 |
                                        2667
                                                 COMPLETE
       22 | 1374710400000
                                        333
                                                 COMPLETE
       26 | 1374710400000
                                                 COMPLETE
                                        7562
       28 1374710400000
                                        656
                                                 COMPLETE
       32 | 1374710400000 |
                                        3960
                                                 COMPLETE |
only showing top 10 rows
```

from pyspark.sql.functions import from_unixtime, col

Convert order_date from milliseconds to "YYYY-MM-DD" format

df_with_converted_date = df_filtered.withColumn("order_date", from_unixtime(col("order_date") / 1000, "yyyy-MM-dd"))

```
>>> from pyspark.sql.functions import from_unixtime, col
>>> df_with_converted_date = df_filtered.withColumn("order_date", from_unixtime(col("order_date") / 1000, "yyyy-MM-dd"))
>>> df_with_converted_date.show(10)
|order_id|order_date|order_customer_id|order_status|
         3 2013-07-25
                                      12111
                                                  COMPLETE
        5 2013 - 07 - 25
                                      11318
                                                  COMPLETE
        6 2013-07-25
                                       7130
4530
                                                  COMPLETE
        7 2013-07-25
                                                  COMPLETE
       15 2013-07-25
                                        2568
                                                  COMPLETE
       17 2013-07-25
                                        2667
                                                  COMPLETE
       22 2013-07-25
                                        333
                                                  COMPLETE
       26 2013-07-25
                                                  COMPLETE
       28 2013-07-25
       32 2013-07-25
                                                  COMPLETE
only showing top 10 rows
```

Selecting only the required columns

```
required_columns = df_with_converted_date.select(
    col("order_id").alias("Order number"),
    col("order_date").alias("Order date"),
    col("order_status").alias("Current situation")
```

```
required columns = df with converted date.select(
        col("order id").alias("Order number"),
        col("order_date").alias("Order date"),
        col("order_status").alias("Current situation")
   required columns.show(10)
Order number | Order date | Current situation |
             3 2013-07-25
                                    COMPLETE
            5 2013-07-25
                                    COMPLETE
            6 2013-07-25
                                    COMPLETE
             7 | 2013-07-25 |
                                    COMPLETE
            15 2013-07-25
                                    COMPLETE
            17 2013-07-25
                                    COMPLETE
            22|2013-07-25|
                                    COMPLETE
            26 | 2013-07-25 |
                                    COMPLETE
            28 | 2013-07-25 |
                                    COMPLETE
            32 | 2013-07-25 |
                                    COMPLETE
only showing top 10 rows
```

- # Specify the path where you want to save the results output_path = "/user/swapnasamirshukla1988gmai/retail_ba/result/scenario2/solution"
- # Save the data in JSON format with GZIP compression

required_columns.write.mode("overwrite").option("compression", "gzip").json(output_path)

```
>>> output_path = " /user/swapnasamirshukla1988gmai/retail_ba/result/scenario2/solution"
>>> required_columns.write.mode("overwrite").option("compression", "gzip").json(output_path)
>>> [
```

Check the hdfs location to see the file

hdfs dfs -ls /user/swapnasamirshukla1988gmai/retail_ba/result/scenario2/solution

```
[swapnasamirshukla1988gmai@sl-cdp-prod-en10 ~]$ hdfs dfs -ls /user/swapnasamirshukla1988gmai/retail_ba/result/scenario2/solution
Found 2 Items

0 2024-02-04 17:59 /user/swapnasamirshukla1988gmai/retail_ba/result/scenario2/solution/_SUCCESS

79705 2024-02-04 17:59 /user/swapnasamirshukla1988gmai/retail_ba/result/scenario2/solution/part-00000-785622e8-8

0 2024-02-04 17:59 /user/swapnasamirshukla1988gmai/retail_ba/result/scenario2/solution/SUCCESS

79705 2024-02-04 17:59 /user/swapnasamirshukla1988gmai/retail_ba/result/scenario2/solution/SUCCESS

79705 2024-02-04 17:59 /user/swapnasamirshukla1988gmai/retail_ba/result/scenario2/solution/Fetai_sol.2-json.gz
```

- # The file is stored as part-00000-785622e8-80b0-4ac6-9fe7-67165f99b7de-c000.json.gz
- # Rename file using mv command

 $hdfs\ dfs\ -mv\ /user/swapnasamirshukla 1988 gmai/retail_ba/result/scenario 2/solution/part-00000-785622 e8-80b0-4ac6-9 fe7-67165 f99b7 de-c000. json. gz$

 $/user/swapnasamirshukla 1988 gmai/retail_ba/result/scenario 2/solution/retai_sol 2. json. gz$

Explore the customer records saved in the "customers-tab-delimited" directory on HDFS

REQUIREMENT:

- 2.4.1 Produce a list of all consumers who live in the city of "Caguas"
- 2.4.2 Save the data in the result/scenario3/solution directory on HDFS
- 2.4.3 The result should only contain records with the value "Caguas" for the customer city
- 2.4.4 Use snappy compression to compress the output
- 2.4.5 Save the file in the orc format

Read the file, assuming it is tab-separated and does not have a header

```
file_path = "retail_ba_dataset/customers-tab-delimited/part-m-00000" df = spark.read.option("header", "false").option("sep", "\t").csv(file_path) df.printSchema() df.show(10)
```

```
>>> file_path = "retail_ba_dataset/customers-tab-delimited/part-m-00000
>>> df = spark.read.option("header", "false").option("sep", "\t").csv(file path)
>>> df.printSchema()
root
    _c0: string (nullable = true)
    _c1: string
                (nullable = true)
    _c2: string (nullable = true)
    _c3: string (nullable = true)
     _c4: string
                (nullable = true)
    c5: string (nullable = true)
    _c6: string (nullable = true)
    _c7: string (nullable = true)
    _c8: string (nullable = true)
>>> df.show(10)
 c0
        c1
                                                          c5|
                                                                      c6| c7|
                  c2
                            c3
                                      c4
  6303 Heather Plaza|Brownsville| TX|78521|
                                                                Littleton
                                                                          CO 80126
  зİ
                Smith XXXXXXXXX XXXXXXXXX 3422 Blue Pioneer...
        Ann
                                                                   Caguas |
                                                                          PR | 00725
                Jones XXXXXXXXX XXXXXXXX
  4
       Mary
                                                              San Marcos
                                          8324 Little Common
                                                                          CA 92069
  5
     Robert
               Hudson|XXXXXXXXX|XXXXXXXXXX|10 Crystal River ...
                                                                   Caguas
                                                                          PR | 00725
                Smith XXXXXXXXX XXXXXXXXX 3151 Sleepy Quail...
  6
       Mary
                                                                  Passaic|
                                                                          NJ | 07055
  7 Melissa
               Wilcox XXXXXXXXX XXXXXXXXX 9453 High Concession
                                                                  Caguas
                                                                          PR | 00725
                Smith XXXXXXXXX XXXXXXXXX 3047 Foggy Forest...
  8
      Megan
                                                                 Lawrence
                                                                          MA | 01841
       Mary
                Perez XXXXXXXXX XXXXXXXXX 3616 Quaking Street
                                                                   Caguas |
                                                                          PR | 00725
  10 Melissa
                Smith XXXXXXXXX XXXXXXXXX 8598 Harvest Beac...
                                                                 Stafford|
                                                                          VA 22554
only showing top 10 rows
```

Filter for customers who live in "Caguas" customers_in_caguas = df.filter(col("_c6") == "Caguas")

```
>>> customers_in_caguas = df.filter(col("_c6") == "Caguas")
   customers_in_caguas.show(5)
  c0
                                                                c5|
          c1
                   c2|
                                         c4
                                                                        c6| c7|
                Smith XXXXXXXXX XXXXXXXXX 3422 Blue Pioneer... | Caguas | PR | 00725
   3|
         Ann |
      Robert | Hudson | XXXXXXXXX | XXXXXXXXX | 10 Crystal River ... | Caguas | PR | 00725
   7|Melissa| Wilcox|XXXXXXXXX|XXXXXXXXX|9453 High Concession|Caguas| PR|00725
                Perez | XXXXXXXXX | XXXXXXXXX | 3616 Quaking Street | Caguas |
        Mary | Huffman | XXXXXXXXX | XXXXXXXXX |
                                                 3169 Stony Woods | Caguas | PR | 00725 |
only showing top 5 rows
```

Save the filtered data in ORC format with Snappy compression to the specified HDFS directory output_path = "/user/swapnasamirshukla1988gmai/retail_ba/result/scenario3/solution" customers_in_caguas.write.mode("overwrite").option("compression", "snappy").orc(output_path)

```
>>> output_path = "/user/swapnasamirshukla1988gmai/retail_ba/result/scenario3/solution"
>>> customers_in_caguas.write.mode("overwrite").option("compression", "snappy").orc(output_path)
>>>
```

Check the hdfs location to see the file

hdfs dfs -ls /user/swapnasamirshukla1988gmai/retail_ba/result/scenario3/solution

```
[swapnasamirshukla1988gmai@sl-cdp-prod-en10 ~]$ hdfs dfs -ls /user/swapnasamirshukla1988gmai/retail_ba/result/scenario3/solution
Found 2 items
-nw-rw-r-- 2 swapnasamirshukla1988gmai swapnasamirshukla1988gmai
-nw-rw-r-- 2 swapnasamirshukla1988gmai swap
```

Explore all the category records stored in the "categories" directory on HDFS

REQUIREMENT:

- 2.5.1 Save the result files in CSV format
- 2.5.2 Save the data in the result/scenario4/solution directory on HDFS
- 2.5.3 Use Iz4 compression to compress the output

```
file_path = "retail_ba_dataset/categories/ part-m-00000" df = spark.read.option("header", "false").option("sep", ",").csv(file_path) df.printSchema() df.show(10)
```

```
file_path = "retail_ba_dataset/categories/part-m-00000"
df = spark.read.option("header", "false").option("sep", ",").csv(file_path)
>>> df = spark.read.option("header",
>>> df.printSchema()
       _c0: string (nullable = true)
_c1: string (nullable = true)
_c2: string (nullable = true)
 >>> df.show(10)
 _c0|_c1|
                                       c2
   1|
2|
3|
4|
5|
6|
7|
8|
9|
          2
                               Football
          2 | Soccer
2 | Baseball & Softball
          2
2
2
2
                            Basketball
                              Lacrosse
                  Tennis & Racquet
                                 Hockey
                          More Sports
                  Cardio Equipment
                Strength Training
  nly showing top 10 rows
```

Save the filtered data in text format with lz4 compression to the specified HDFS directory output_path = "/user/swapnasamirshukla1988gmai/retail_ba/result/scenario4/solution" df.write.mode("overwrite").option("compression", "gzip").csv(output_path)

```
>>> output_path = "/user/swapnasamirshukla1988gmai/retail_ba/result/scenario4/solution"
>>> df.write.mode("overwrite").option("compression", "gzip").csv(output_path)
>>>
```

Spark's built-in CSV data source does not directly support LZ4 compression. One might need to save the data uncompressed and then compress it using external tools, or choose a compression codec supported by Spark for CSV, such as bzip2,gzip or deflate.

Check the hdfs location to see the file

hdfs dfs -ls /user/swapnasamirshukla1988gmai/retail_ba/result/scenario4/solution

```
Ssword for swapnasamirshukla1988gmai@CDP-ENV.GNE4-RUTX.CLOUDERA.SITE:
[swapnasamirshukla1988gmai@sl-cdp-prod-en21 ~]$ hdfs dfs -ls /user/swapnasamirshukla1988gmai/retail_ba/result/scenario4/solution
Found 2 Items
Frw.Frw.Fr-- 2 swapnasamirshukla1988gmai swapnasamirshukla1988gmai
Frw.Frw.Fr-- 2 swapnasamirshukla1988gmai swapnasamirshukla1988gmai
621-4922-b6e4-578ebb3834a9-c000.csv.g2
[swapnasamirshukla1988gmai@sl-cdp-prod-en21 ~]$
```

Explore all product records that are saved in the "products" avro" database

REQUIREMENT:

- 2.6.1 Only products with a price of more than 1000.0 should be included in the output
- 2.6.2 Save the output files in parquet format
- 2.6.3 Remove data from the table if the product price is lesser than 1000.0
- 2.6.4 Save the data in the result/scenario5/solution directory on HDFS
- 2.6.5 Use snappy compression to compress the output

```
# Path to the Avro files
file paths = [
    "retail_ba_dataset/products_avro/part-m-00000.avro",
    "retail_ba_dataset/products_avro/part-m-00001.avro",
    "retail_ba_dataset/products_avro/part-m-00002.avro",
    "retail_ba_dataset/products_avro/part-m-00003.avro"
# Read the Avro files into a DataFrame
df = spark.read.format("avro").load(file_paths)
# Remove completely duplicated rows
df_deduplicated = df.dropDuplicates()
     file_paths = [
    "retail_ba_dataset/products_avro/part-m-00000.avro",
          "retail ba_dataset/products_avro/part-m-00001.avro",
"retail_ba_dataset/products_avro/part-m-00002.avro",
"retail_ba_dataset/products_avro/part-m-00003.avro"
...]
>>> df = spark.read.format("avro").load(file_paths)
>>> df_deduplicated = df.dropDuplicates()
>>> df_deduplicated.printSchema()
      product_id: integer (nullable =
    - product_rane: Integer (nullable = true)
- product_name: string (nullable = true)
- product_description: string (nullable = true)
- product_price: float (nullable = true)
- product_image: string (nullable = true)
>>> df_deduplicated.show(5)
|product_id|product_category_id|
                                                     product_name|product_description|product_price|
                                                                                                                              product_image|
                                                                                                           349.99 http://images.acm...
99.0 http://images.acm...
54.99 http://images.acm...
          1033|
                                      46|YETI Tundra 45 Ch...
         1034
1042
                                      47
                                             Nike+ Fuelband SE
                                      47 Under Armour Hust...
50 Nike Men's New Yo...
                                                                                                              34.0 http://images.acm...
         1279
                                      57 PUMA Men's evoPOW...
                                                                                                           189.99 http://images.acm..
only showing top 5 rows
```

Filter for products with price more than 1000.0 filtered_df = df_deduplicated.filter(col("product_price") >= 1000.0)

Save the filtered data in parquet format with snappy compression to the specified HDFS directory output_path = "/user/swapnasamirshukla1988gmai/retail_ba/result/scenario5/solution" filtered_df.write.mode("overwrite").option("compression", "snappy").parquet(output_path)

```
overline in the second product of the s
```

Check the hdfs location to see the file

hdfs dfs -ls /user/swapnasamirshukla1988gmai/retail_ba/result/scenario5/solution

```
Password for swapnasamirshukla1988gmai@CDP-ENV.CM4-RUIX.CLOUDERA.SITE:
[sawapnasamirshukla1988gmai@s]-cdp-prod-en10 ~]$ hdfs dfs -1s /user/swapnasamirshukla1988gmai/retail_ba/result/scenario5/solution
Found 2 items
-rw-ru-r-- 2 swapnasamirshukla1988gmai swapnasamirshukla1988gmai
-rw-ru-r-- 2 swapnasamirshukla1988gmai swapnasamirshukla1988gmai
-rw-ru-r-- 2 swapnasamirshukla1988gmai swapnasamirshukla1988gmai
-rw-ru-r-- 2 swapnasamirshukla1988gmai swapnasamirshukla1988gmai
-rw-ru-r-- 2 swapnasamirshukla1988gmai/retail_ba/result/scenario5/solution/part-00000-725426c9-
-05e-4939-9496-d30227e0ab94-c000.snappy.parquet
[swapnasamirshukla1988gmai@sl-cdp-prod-en10 ~]$ []
```

Explore the "products avro" stored in product records

REQUIREMENT:

- 2.7.1 Only products with a price of more than 1000.0 should be in the output
- 2.7.2 The pattern "Treadmill" appears in the product name
- 2.7.3 Save the output files in parquet format
- 2.7.4Save the data in the result/scenario6/solution directory on HDFS
- 2.7.5 Use GZIP compression to compress the output

```
# Path to the Avro files
file_paths = [
    "retail_ba_dataset/products_avro/part-m-00000.avro",
    "retail_ba_dataset/products_avro/part-m-00001.avro",
    "retail_ba_dataset/products_avro/part-m-00002.avro",
    "retail ba dataset/products avro/part-m-00003.avro"
# Read the Avro files into a DataFrame
df = spark.read.format("avro").load(file paths)
# Remove completely duplicated rows
df deduplicated = df.dropDuplicates()
     "retail_ba_dataset/products_avro/part-m-00000.avro",
   "retail_ba_dataset/products_avro/part-m-00001.avro",
   "retail_ba_dataset/products_avro/part-m-00002.avro",
   "retail_ba_dataset/products_avro/part-m-00003.avro"
product_id: integer (nullable = true)
product_category_id: integer (nullable = true)
product_name: string (nullable = true)
product_description: string (nullable = true)
product_price: float (nullable = true)
product_image: string (nullable = true)
>>> df_deduplicated.show(5)
                                                         product_name|product_description|product_price|
|product_id|product_category_id|
                                                                                                                                       product image
                                                                                                                   349.99 http://images.acm...
99.0 http://images.acm...
54.99 http://images.acm...
34.0 http://images.acm...
189.99 http://images.acm...
          1033
                                         46 YETI Tundra 45 Ch...
                                         47 Nike+ Fuelband SE
47 Under Armour Hust...
50 Nike Men's New Yo...
57 PUMA Men's evoPOW...
          1034
          1042
          1126
          1279
only showing top 5 rows
```

Filter for products with price more than 1000.0 filtered df = df deduplicated.filter(col("product price") >= 1000.0)

Now, filter for product names containing 'Treadmill'

```
treadmill_df = filtered_df.filter(col("product_name").contains("Treadmill"))
>>> treadmill_df = filtered_df.filter(col("product_name").contains("Treadmill"))
>>> treadmill_df.show()
|product_id|product_category_id|
                                                           product_name|product_description|product_price|
                                                                                                                                               product_image|
                                                                                                                        1799.99|http://images.acm...
1799.99|http://images.acm...
1799.99|http://images.acm...
                                             22|SOLE F85 Treadmill|
            496
                                             10 SOLE F85 Treadmill
             66
                                              4 SOLE F85 Treadmill
```

Save the filtered data in parquet format with gzip compression to the specified HDFS directory output_path = "/user/swapnasamirshukla1988gmai/retail_ba/result/scenario6/solution" filtered_df.write.mode("overwrite").option("compression", "gzip").parquet(output_path)

```
>>> output path = "/user/swapnasamirshukla1988gmai/retail ba/result/scenario6/solution"
>>> filtered_df.write.mode("overwrite").option("compression", "gzip").parquet(output_path)
>>>
```

Check the hdfs location to see the file

hdfs dfs -ls /user/swapnasamirshukla1988gmai/retail_ba/result/scenario6/solution

```
swapnasamirshukla1988gmai swapnas
swapnasamirshukla1988gmai swapnas
8402ca31a407-c000.gz.parquet
ukla1988gmai@s1-cdp-prod-en10 ~]$
```

Explore the order records that are saved in the "orders parquet" table on HDFS

REQUIREMENT:

- 2.8.1 Output all PENDING orders in July 2013
- 2.8.2 Output files should be in JSON format
- 2.8.3 Save the data in the result/scenario7/solution directory on HDFS.
- 2.8.4 Only entries with the order status value of "PENDING" should be included in the result
- 2.8.5 Order date should be in the YYY-MM-DD format
- 2.8.6 Use snappy compression to compress the output, which should just contain the order date and order status

file_path = "retail_ba_dataset/orders_parquet/741ca897-c70e-4633-b352-5dc3414c5680.parquet" df = spark.read.parquet (file_path) df.printSchema() df.show(10)

from pyspark.sql.functions import from_unixtime, col

Convert order_date from milliseconds to "YYY-MM-DD" format

df_with_converted_date = df.withColumn("order_date", from_unixtime(col("order_date") / 1000, "yyy-MM-dd"))

```
>>> df_with_converted_date = df.w
>>> df_with_converted_date.show()
                                         df.withColumn("order_date", from_unixtime(col("order_date") / 1000, "yyy-MM-dd");
order_id|order_date|order_customer_id|
                                                             order_status
          2 | 2013-07-25 |
3 | 2013-07-25 |
4 | 2013-07-25 |
5 | 2013-07-25 |
                                                256 PENDING_PAYMENT
12111 COMPLETE
                                                 8827
11318
                                                                      CLOSED
                                                  7130
4530
              2013-07-25
                                                                    COMPLETE
                                                  2911 PROCESSING
5657 PENDING_PAYMENT
5648 PENDING_PAYMENT
           8 2013-07-25
              2013-07-25
         10 | 2013 - 07 - 25
         11 2013-07-25
12 2013-07-25
                                                  918
1837
                                                          PAYMENT_REVIEW
CLOSED
                                                  9149 PENDING_PAYMENT
9842 PROCESSING
          13 | 2013-07-25
                                                  2568 COMPLETE
7276 PENDING_PAYMENT
          15 | 2013-07-25
              2013-07-25
                                                                    COMPLETE
                                                  1205 CLOSED
9488 PENDING_PAYMENT
nly showing top 20 rows
```

```
pending_orders_july_2013 = df_with_converted_date \
    .filter(col("order_status") == "PENDING") \
    .filter(col("order_date").between("2013-07-01", "2013-07-31"))
```

```
... pending_orders_july_2013 = df_with_converted_date \
         .filter(col("order_status") == "PENDING") \
.filter(col("order_date").between("2013-07-01", "2013-07-31"))
>>> pending_orders_july_2013.show(5)
|order_id|order_date|order_customer_id|order_status|
       21 2013-07-25
                                                   PENDING
                                       2711
        36 2013-07-25
                                       5649
                                                   PENDING
        39 2013 - 07 - 25
                                       8214
                                                   PENDING
       42 | 2013 - 07 - 25 |
                                       9776
                                                   PENDING
       44 2013-07-25
                                      10500
                                                   PENDING
only showing top 5 rows
```

```
pending_orders = df_with_converted_date \
    .filter(col("order_status") == "PENDING")
```

```
pending orders = df_with_converted_date \
        .filter(col("order status") == "PENDING")
>>>
>>> pending orders.show()
order id|order_date|order_customer_id|order_status|
       21 2013-07-25
                                     2711
                                                PENDING
       36 2013-07-25
                                     5649
                                                PENDING
       39 2013-07-25
                                     8214
                                                PENDING
       42 | 2013 - 07 - 25 |
                                     9776
                                                PENDING
       44 2013 - 07 - 25
                                    10500
                                                PENDING
       49 2013-07-25
                                     1871
                                                PENDING
       55 2013-07-25
                                     2052
                                                PENDING
       68 2013-07-25
                                     4320
                                                PENDING
       85 | 2013 - 07 - 25 |
                                     1485
                                                PENDING
       96 2013 - 07 - 25
                                     8683
                                                PENDING
       97 2013-07-25
                                    10784
                                                PENDING
      121 2013-07-26
                                     2074
                                                PENDING
      132 | 2013 - 07 - 26 |
                                      289
                                                PENDING
      158 2013 - 07 - 26
                                    12345
                                                PENDING
      167 | 2013 - 07 - 26 |
                                     1347
                                                PENDING
      181 2013-07-26
                                     7473
                                                PENDING
      188 | 2013 - 07 - 26 |
                                     2889
                                                PENDING
      189 2013-07-26
                                    10177
                                                PENDING
      190 2013 - 07 - 26
                                    11115
                                                PENDING
      206 | 2013 - 07 - 26 |
                                     8994
                                                PENDING
only showing top 20 rows
```

```
pending_orders = df_with_converted_date \
    .filter(col("order_status") == "PENDING") \
    .select("order_date", "order_status")
```

```
>>> pending_orders = df_with_converted_date \
        .filter(col("order_status") == "PENDING") \
        .select("order_date", "order_status")
>>> pending_orders.show()
order_date|order_status|
2013-07-25
                 PENDING
2013-07-25
                 PENDING
2013-07-25
                 PENDING
2013-07-25
                 PENDING
2013-07-25
                 PENDING
2013-07-25
                 PENDING
 2013-07-25
                 PENDING
2013-07-25
                 PENDING
2013-07-25
                 PENDING
2013-07-25
                 PENDING
2013-07-25
                 PENDING
2013-07-26
                 PENDING
only showing top 20 rows
```

Save the filtered data in json format with snappy compression to the specified HDFS directory output_path = "/user/swapnasamirshukla1988gmai/retail_ba/result/scenario7/solution"

pending_orders.write.mode("overwrite").option("compression", "gzip").json(output_path)

it's important to note that JSON itself doesn't support Snappy compression directly within the Spark DataFrame API. Snappy compression is typically used with columnar storage formats like Parquet or ORC, which are better suited for compressing tabular data.

Check the hdfs location to see the file

hdfs dfs -ls /user/swapnasamirshukla1988gmai/retail_ba/result/scenario7/solution

```
[swapnasamirshukla1988gmai@sl-cdp-prod-en21 ~]$
[swapnasamirshukla1988gmai@sl-cdp-prod-en21 ~]$ hdfs dfs -ls /user/swapnasamirshukla1988gmai/retail_ba/result/scenario7/solution
Found 2 items
Found 3 items
Found 4 ```