

**NAME: MAVIA ALAM KHAN (2303.KHI.DEG.017)**

**PAIRING WITH : Mohammad Hussam (2303.KHI.DEG.020)**

**& AQSA TAUHEED(2303.KHI.DEG.011)**

## **ASSIGNMENT NO : 5.1**

### Import libraries

```
In [1]: 1 from pyspark.sql import SparkSession
2 from pyspark.sql.types import IntegerType, DateType
3 from pyspark.sql.functions import sum, col, avg, desc
4

In [2]: 1 scSpark = SparkSession.builder.appName("Spark Example").getOrCreate()
```

Two join conditions are defined:

join\_condition1 specifies that the "CustomerId" column in the transaction DataFrame should match the "CustomerId" column in the customer DataFrame.

join\_condition2 specifies that the "ProductId" column in the transaction DataFrame should match the "ProductId" column in the products DataFrame.

```
In [3]: 1 transaction = scSpark.read.csv("transactions_*.csv", header=True, inferSchema=True)
2 customer = scSpark.read.csv("customers.csv", header=True, inferSchema=True)
3 products = scSpark.read.csv("products.csv", header=True, inferSchema=True)
4
5 # Define join conditions
6 join_condition1 = transaction["CustomerId"] == customer["CustomerId"]
7
8 join_condition2 = transaction["ProductId"] == products["ProductId"]
9
10 # Perform joins
11 joined_df = transaction.join(customer, join_condition1, "left") \
12     .join(products, join_condition2, "left")
13
14 # Show the resulting joined data frame
15 joined_df.show()
```

|    | StoreId       | TransactionId | CustomerId | ProductId | Quantity | TransactionTime     | CustomerId | Name           | Email                |
|----|---------------|---------------|------------|-----------|----------|---------------------|------------|----------------|----------------------|
|    | roductId      | Name          | Category   | UnitPrice |          |                     |            |                |                      |
|    | 3             |               | 454        | 35        | 3        | 2022-12-23 17:36:11 | 35         | Dwayne Johnson | dwayne.johnson@gm... |
| 3  | Blue Shorts   | Shorts        | 118.88     |           |          |                     |            |                |                      |
|    | 3             |               | 524        | 37        | 9        | 2022-12-23 22:02:51 | 37         | Brittany Holt  | brittany.holt@exa... |
| 9  | Green Sandals | Shoes         | 137.53     |           |          |                     |            |                |                      |
|    | 3             |               | 562        | 4         | 3        | 2022-12-23 02:51:50 | 4          | Alevtin Paska  | alevtin.paska@exa... |
| 3  | Blue Shorts   | Shorts        | 118.88     |           |          |                     |            |                |                      |
|    | 3             |               | 581        | 35        | 14       | 2022-12-23 17:05:54 | 35         | Dwayne Johnson | dwayne.johnson@gm... |
| 14 | Red t-shirt   | T-Shirts      | 121.58     |           |          |                     |            |                |                      |
|    | 3             |               | 200        | 34        | 15       | 2022-12-23 07:15:01 | 34         | Avi Shet       | avi.shet@example.com |
| 15 | White t-shirt | T-Shirts      | 131.13     |           |          |                     |            |                |                      |
|    | 3             |               | 506        | 41        | 24       | 2022-12-23 21:26:29 | 41         | Alice Morin    | alice.morin@examp... |
| 24 | Blue Jeans    | Pants         | 173.1      |           |          |                     |            |                |                      |
|    | 3             |               | 278        | 5         | 1        | 2022-12-23 16:41:42 | 5          | Charlotte Wong | charlotte.wong@ex... |
| 1  | Red Shorts    | Shorts        | 89.75      |           |          |                     |            |                |                      |

Now we add total price column in which we multiply quantity & unit price column

```
In [4]: 1 join_df = joined_df.withColumn("total price", col("Quantity") * col("UnitPrice"))
        2 join_df.show()
```

| StoreId   | TransactionId | CustomerId | ProductId | Quantity    | TransactionTime     | CustomerId | Name            | Email                |
|-----------|---------------|------------|-----------|-------------|---------------------|------------|-----------------|----------------------|
| ProductId | Name          | Category   | UnitPrice | total price |                     |            |                 |                      |
| 3         | Blue Shorts   | Shorts     | 118.88    | 356.64      | 2022-12-23 17:36:11 | 35         | Dwayne Johnson  | dwayne.johnson@gm... |
| 3         | Blue Shorts   | Shorts     | 118.88    | 356.64      | 2022-12-23 22:02:51 | 37         | Brittany Holt   | brittany.holt@exa... |
| 9         | Green Sandals | Shoes      | 137.53    | 1512.83     | 2022-12-23 02:51:50 | 4          | Alevtin Paska   | alevtin.paska@exa... |
| 3         | Blue Shorts   | Shorts     | 118.88    | 475.52      | 2022-12-23 17:05:54 | 35         | Dwayne Johnson  | dwayne.johnson@gm... |
| 14        | Red t-shirt   | T-Shirts   | 121.58    | 6808.48     | 2022-12-23 07:15:01 | 34         | Avi Shet        | avi.shet@example.com |
| 15        | White t-shirt | T-Shirts   | 131.13    | 3147.12     | 2022-12-23 21:26:29 | 41         | Alice Morin     | alice.morin@examp... |
| 24        | Blue Jeans    | Pants      | 173.1     | 3288.9      | 2022-12-23 16:41:42 | 5          | Charlotte Wong  | charlotte.wong@ex... |
| 1         | Red Shorts    | Shorts     | 89.75     | 448.75      | 2022-12-23 13:22:55 | 36         | William Nielsen | william.nielsen@e... |
| 23        | Green Chinios | Pants      | 150.93    | 1962.09     | 2022-12-23 16:47:14 | 34         | Avi Shet        | avi.shet@example.com |
| 7         | White Sandals | Shoes      | 160.96    | 482.88      | 2022-12-23 22:36:48 | 19         | Alexia Renaud   | alexia.renaud@exa... |
| 7         | White Sandals | Shoes      | 160.96    | 2092.48     | 2022-12-23 10:11:29 | 48         | Amoli Shenoy    | amoli.shenoy@exam... |
| 18        | Black t-shirt | T-Shirts   | 102.41    | 1228.92     |                     |            |                 |                      |

The daily total sales for the store with id 1

and the withColumn method is used to create a new column in the DataFrame with the casted values, and the original column is replaced by the new one.

```
In [5]: 1 join_df = join_df.withColumn("total price", col("total price").cast("float"))
        2 join_df = join_df.withColumn("TransactionTime", join_df["TransactionTime"].cast(DateType()))
        3 store_df = join_df.filter(join_df["StoreId"] == 1)
        4 daily_total_sales = store_df.groupBy("TransactionTime").agg(sum("total price").alias("TotalSales"))
        5 daily_total_sales.show()
```

| TransactionTime | TotalSales        |
|-----------------|-------------------|
| 2022-12-23      | 41264.00012207031 |

the mean sales for the store with id 2

```
In [6]: 1 store_sales_df = join_df.filter(join_df["StoreId"] == 2)
        2 mean_sales_store2 = store_sales_df.select(avg("total price").alias("MeanSales"))
        3 mean_sales_store2.show()
```

| MeanSales         |
|-------------------|
| 513.4598035625382 |

the email of the client that spent the most when summing up purchases from all of the stores

the sort method is used on total\_sales\_per\_customer to sort the DataFrame in descending order based on the "TotalSales" column. The desc function specifies the sorting order as descending.

```
In [8]: 1 total_sales_per_customer = join_df.groupBy("Email").agg(sum("total price").alias("TotalSales"))
2 highest_show = total_sales_per_customer.sort(desc("TotalSales"))
3 highest_show.show(1)
```

```
+-----+-----+
|      Email      | TotalSales |
+-----+-----+
|dwayne.johnson@gm...|10653.080017089844|
+-----+-----+
only showing top 1 row
```

And 5 products are most frequently bought across all stores

```
In [10]: 1 productCount_df = join_df.groupBy("ProductId").agg(sum("Quantity").alias("totalQuantity"))
2 orderedProductCount_df = productCount_df.orderBy("totalQuantity", ascending=False)
3 top5products = orderedProductCount_df.limit(5)
4 top5products.show()
```

```
+-----+-----+
|ProductId|totalQuantity|
+-----+-----+
|      14 |          82 |
|      24 |          77 |
|      15 |          76 |
|       5 |          75 |
|      19 |          74 |
+-----+-----+
```

Here we print name of the product that purchase most so we join the table

```
In [11]: 1 joinCondition = top5products["ProductId"] == products["ProductId"]
2 joined_df = top5products.join(products,joinCondition)
3 joined_df.orderBy("totalQuantity", ascending=False).show(5)
```

```
+-----+-----+-----+-----+-----+-----+
|ProductId|totalQuantity|ProductId|      Name      |Category|UnitPrice|
+-----+-----+-----+-----+-----+-----+
|      14 |          82 |      14 | Red t-shirt |T-Shirts|   121.58|
|      24 |          77 |      24 | Blue Jeans |Pants   |   173.1 |
|      15 |          76 |      15 | White t-shirt|T-Shirts|   131.13|
|       5 |          75 |       5 | Black Shorts|Shorts  |    74.58|
|      19 |          74 |      19 | Green jacket|Jackets |   223.69|
+-----+-----+-----+-----+-----+-----+
```

---