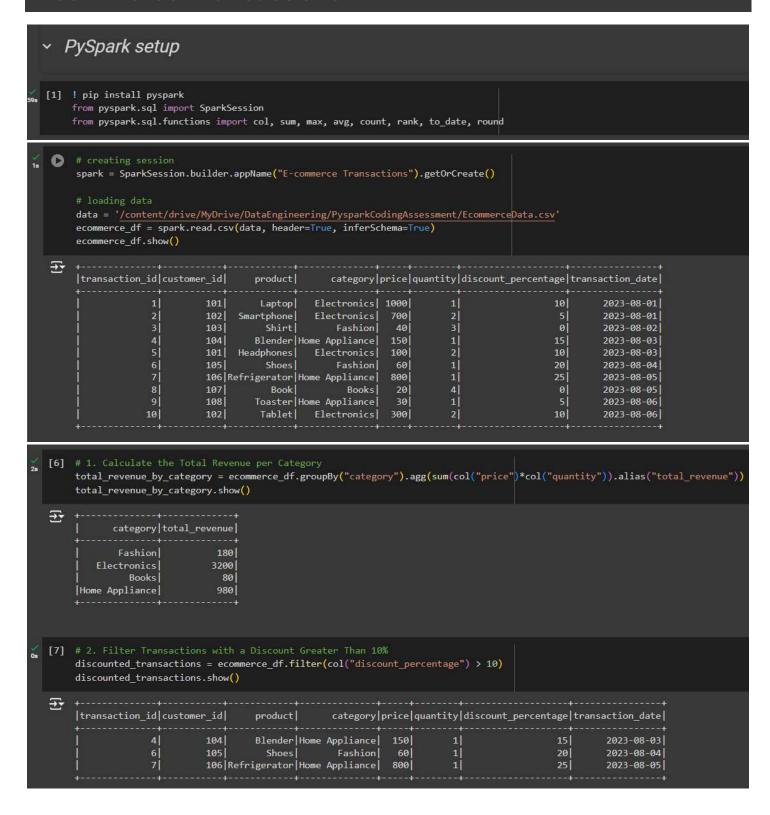
E-commerce Transactions



```
most_expensive_product = ecommerce_df.orderBy(col("price").desc()).limit(1)
      most_expensive_product.show()
 Ŧ
      |transaction_id|customer_id|product| category|price|quantity|discount_percentage|transaction_date|
                                                                      1
                               101 | Laptop|Electronics | 1000 |
                                                                                            101
                                                                                                      2023-08-01
     # 4. Calculate the Average Quantity of Products Sold per Category
 0
      avg_quantity_by_category = ecommerce_df.groupBy("category").agg(avg("quantity").alias("avg_quantity"))
      avg_quantity_by_category.show()
 -
             category|avg_quantity|
                                2.0
          Electronics
                               1.75
                                4.0
               Books
      |Home Appliance|
                                1.0
     customers_with_multiple_products = ecommerce_df.groupBy("customer_id").agg(count("*").alias("total_products")).filter(col("total_products") > 1)
     customers_with_multiple_products.show()
3
     |customer_id|total_products|
 # 6. Find the Top 3 Highest Revenue Transactions
     from pyspark.sql.window import Window
     {\tt ecommerce\_df = ecommerce\_df.withColumn("transaction\_date", to\_date(col("transaction\_date"), "yyyy-MM-dd"))}
    window = Window.orderBy(col("final_price").desc())
    highest_revenue_transactions = ecommerce_df.withColumn("final_price", col("price")*col('quantity') - (col("price") * col("discount_percentage") / 100)) \
        .withColumn("rank", rank().over(window)) \
.filter(col("rank") <= 3)</pre>
     highest_revenue_transactions.show()
₹
                                               category|price|quantity|discount_percentage|transaction_date|final_price|rank|
                                  product
                          102 | Smartphone
                                                                                                            1365.0
                                            Electronics 1000
                          101
                                   Laptop
                                                                                             2023-08-01
                                                                                                             900.0
                                                                                                             600.0
                          106 | Refrigerator | Home Appliance | 800 |
[19] # 7. Calculate the Total Number of Transactions per Day
      total_transactions_per_day = ecommerce_df.groupBy("transaction_date").agg(count("*").alias("total_transactions"))
      total_transactions_per_day.show()
Ŧ
      |transaction_date|total_transactions|
              2023-08-03
                                                2
              2023-08-06
                                                2
              2023-08-01
              2023-08-05
                                                2
              2023-08-04
              2023-08-02
                                                1
```

```
0
      total_spending_by_customer = ecommerce_df.groupBy("customer_id") \
                                .agg(sum(col("price")*col('quantity') - (col("price") * col("discount_percentage") / 100)).alias("total_spending")) \
.orderBy(col("total_spending").desc()).limit(1)
      total_spending_by_customer.show()
  =
       |customer_id|total_spending|
                         1935.0
               102
 [24] # 9. Calculate the Average Discount Given per Product Category
       avg_discount_by_category = ecommerce_df.groupBy("category").agg(avg("discount_percentage").alias("avg_discount"))
      avg_discount_by_category.show()
  ±
             category|avg_discount|
              Fashion
                              10.0
                               0.0
               Books
       Home Appliance
                              15.0
0s 0
        ecommerce_df = ecommerce_df.withColumn("final_price", col("price")*col('quantity') - (col("price") * col("discount_percentage") / 100))
        ecommerce_df.show()
   ∓
        |transaction_id|customer_id|
                                                         category|price|quantity|discount_percentage|transaction_date|final_price|
                                 101
                                           Laptop
                                                                                                              2023-08-01
                                                                                                                                900.0
                       2
3
                                 102
                                                                     700
                                                                                2
                                                                                                              2023-08-01
                                                                                                                               1365.0
                                        Smartphone
                                 103
                                            Shirt
                                                                      40
                                                                                                              2023-08-02
                                                                                                                                120.0
                       4
5
                                 104
                                           Blender | Home Appliance
                                                                     150
                                                                                                              2023-08-03
                                 101
                                       Headphones
                                                                     100|
                                                                                                     10
                                                                                                              2023-08-03
                                                                                                                                190.0
                       6
7
                                 105
                                                         Fashion
                                                                     60
                                                                                                     20
                                                                                                              2023-08-04
                                                                                                                                 48.0
                                 106 Refrigerator Home Appliance
                                                                     800
                                                                                                              2023-08-05
                                                                                                                                600.0
                                 107
                                             Book
                                                            Books
                                                                      20
                                                                                                              2023-08-05
                                                                                                                                80.0
                                 108
                                           Toaster Home Appliance
                                                                      301
                                                                                                              2023-08-06
                                                                                                                                 28.5
                      10
                                 102
                                           Tablet | Electronics |
                                                                     300
                                                                                                     10
                                                                                                              2023-08-06
                                                                                                                                570.0
```