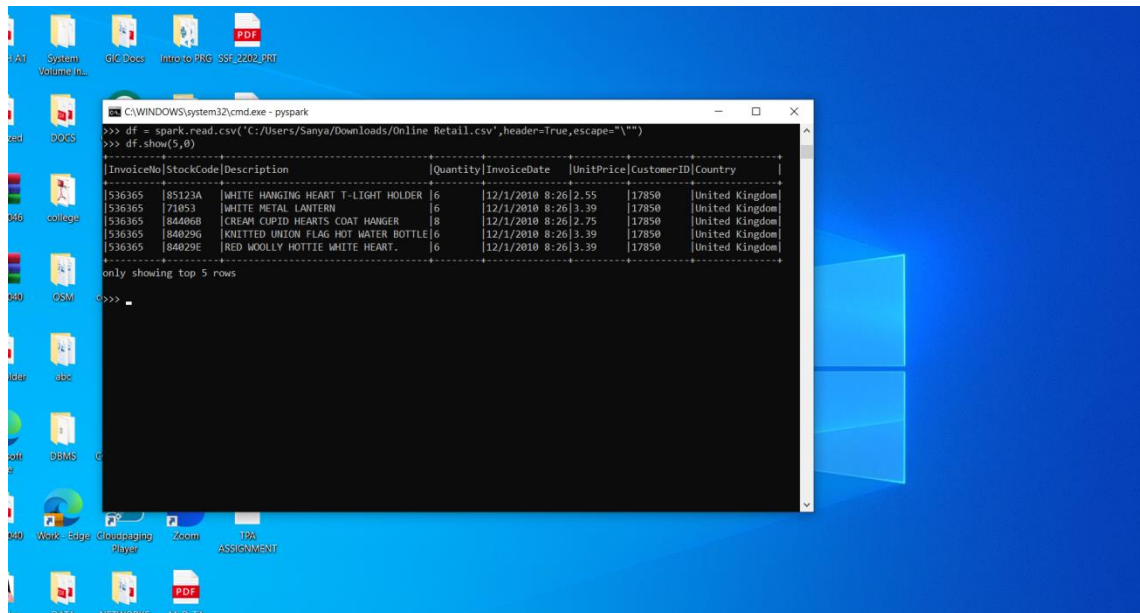# Machine Learning with PySpark: Customer Segmentation

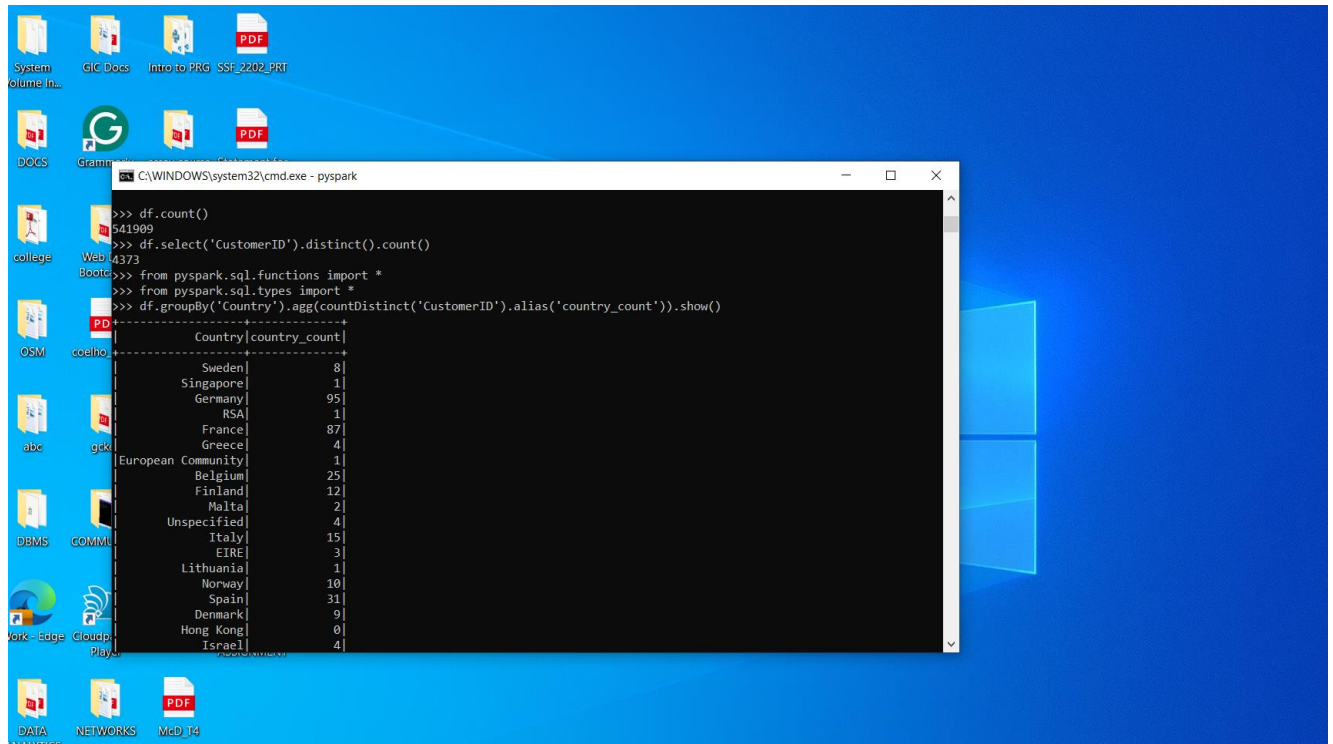**SCREENSHOTS:**

# TASK 1: Creating a SparkSession

# TASK 2: CREATING A DATA FRAME



# Task 3: Exploratory Data Analysis

## 1. Counting the number of rows
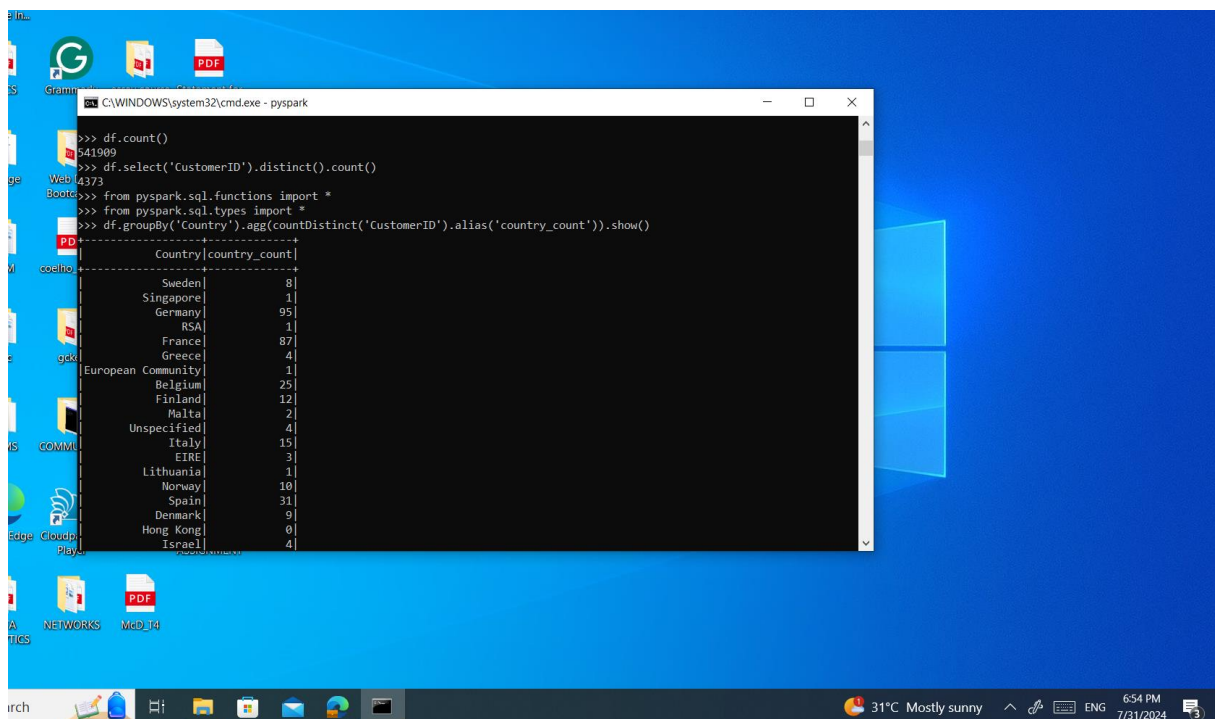
## 2. Counting unique customers

# 3. Country which purchases the maximum quantity



```
C:\WINDOWS\system32\cmd.exe - pyspark

>>> df.count()
541909
>>> df.select('CustomerID').distinct().count()
4373
>>> from pyspark.sql.functions import *
>>> from pyspark.sql.types import *
>>> df.groupBy('Country').agg(countDistinct('CustomerID').alias('country_count')).show()
+------------------+-------------+
|           Country|country_count|
+------------------+-------------+
|            Sweden|            8|
|         Singapore|            1|
|           Germany|           95|
|               RSA|            1|
|            France|           87|
|            Greece|            4|
|European Community|            1|
|           Belgium|           25|
|           Finland|           12|
|             Malta|            2|
|       Unspecified|            4|
|             Italy|           15|
|              EIRE|            3|
|         Lithuania|            1|
|            Norway|           10|
|             Spain|           31|
|           Denmark|            9|
|         Hong Kong|            0|
|            Israel|            4|
+------------------+-------------+
```
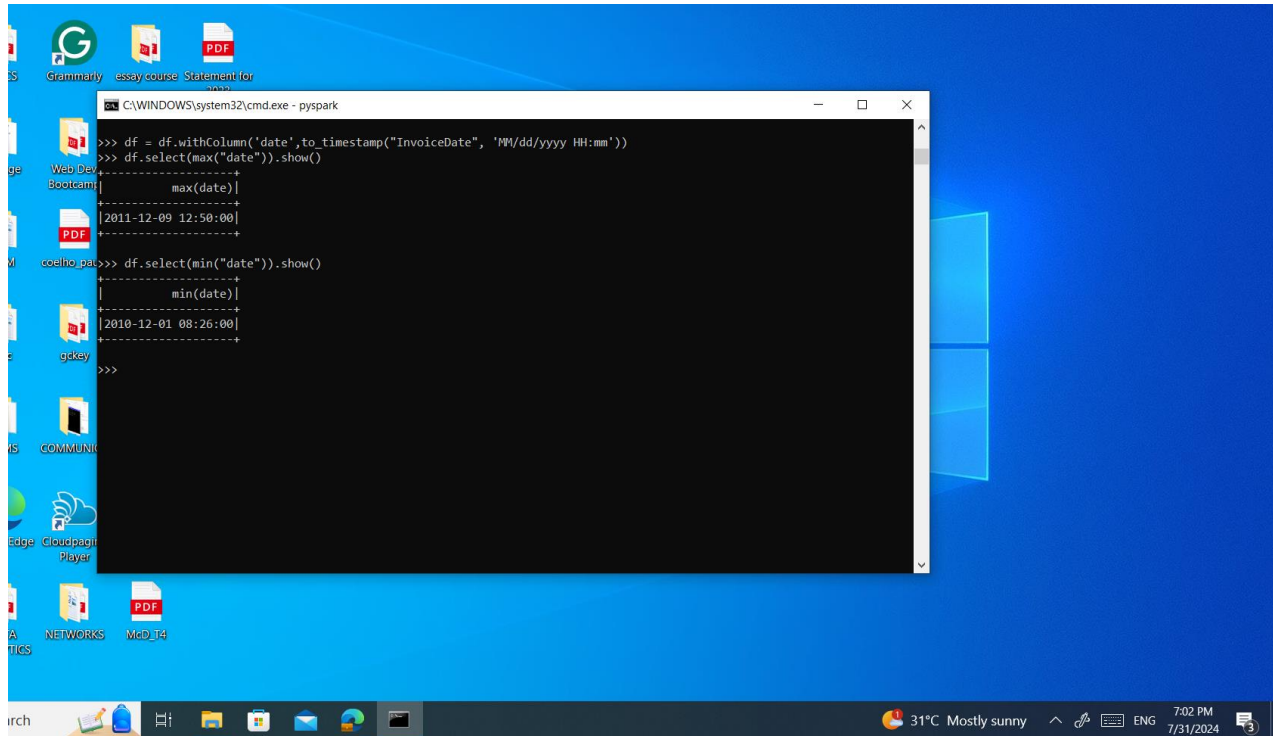
# MOST RECENT PURCHASE AND MINIMUM DATE



# TASK 4: DATA PRE-PROCESSING – RFM

# RECENCY

# PRINTING SCHEMA



# FREQUENCY

```
>>> m_val = df3.withColumn('TotalAmount',col("Quantity") * col("UnitPrice"))
>>> m_val = m_val.groupBy('CustomerID').agg(sum('TotalAmount').alias('monetary_value'))
>>> finaldf = m_val.join(df3,on='CustomerID',how='inner')
>>> finaldf = finaldf.select(['recency','frequency','monetary_value','CustomerID']).distinct()
>>> finaldf.show(5,0)
+---------+---------+---------------------+----------+
|recency  |frequency|monetary_value       |CustomerID|
+---------+---------+---------------------+----------+
|-30405240|10       |153.0                |17714     |
|-25294740|10       |163.3                |16250     |
|-33761640|43       |306.84               |17551     |
|-31843740|37       |236.01999999999995   |13187     |
|-30299640|30       |215.78               |15052     |
+---------+---------+---------------------+----------+
only showing top 5 rows

>>>
```

# MONETARY VALUE



```
>>> m_val = m_val.groupBy('CustomerID').agg(sum('TotalAmount').alias('monetary_value'))
>>> finaldf = m_val.join(df3,on='CustomerID',how='inner')
>>> finaldf = finaldf.select(['recency','frequency','monetary_value','CustomerID']).distinct()
>>> finaldf.show(5,0)
+---------+---------+---------------------+----------+
|recency  |frequency|monetary_value       |CustomerID|
+---------+---------+---------------------+----------+
|4586760  |10       |153.0                |17714     |
|9697260  |10       |163.3                |16250     |
|1230360  |43       |306.84               |17551     |
|3148260  |37       |236.01999999999995   |13187     |
|4692360  |30       |215.78               |15052     |
+---------+---------+---------------------+----------+
only showing top 5 rows

>>>
```

# STANDARDIZATION: TO ENSURE VARIABLES ARE AROUND SAME SCALE

# TASK 5: BUILILDNG THE MACHINE LEARNING MODEL USING K-MEANS CLUSTERING

## 1. FINDING THE NUMBER OF CLUSTERS TO USE





**FROM THE ABOVE GRAPH, WE INFER THE NUMBER OF CLUSTERS TO BE 4.**

# BUILDING THE K-MEANS CLUSTERING MODEL AND MAKING PREDICTIONS