# DATA PREPROCESSING

**DATAFRAME CREATION:**

#creating data frame from the csv file

import pandas as pd

df=pd.read_csv(r"C:\\Users\Nithesh\Downloads\shopping_trends.csv")
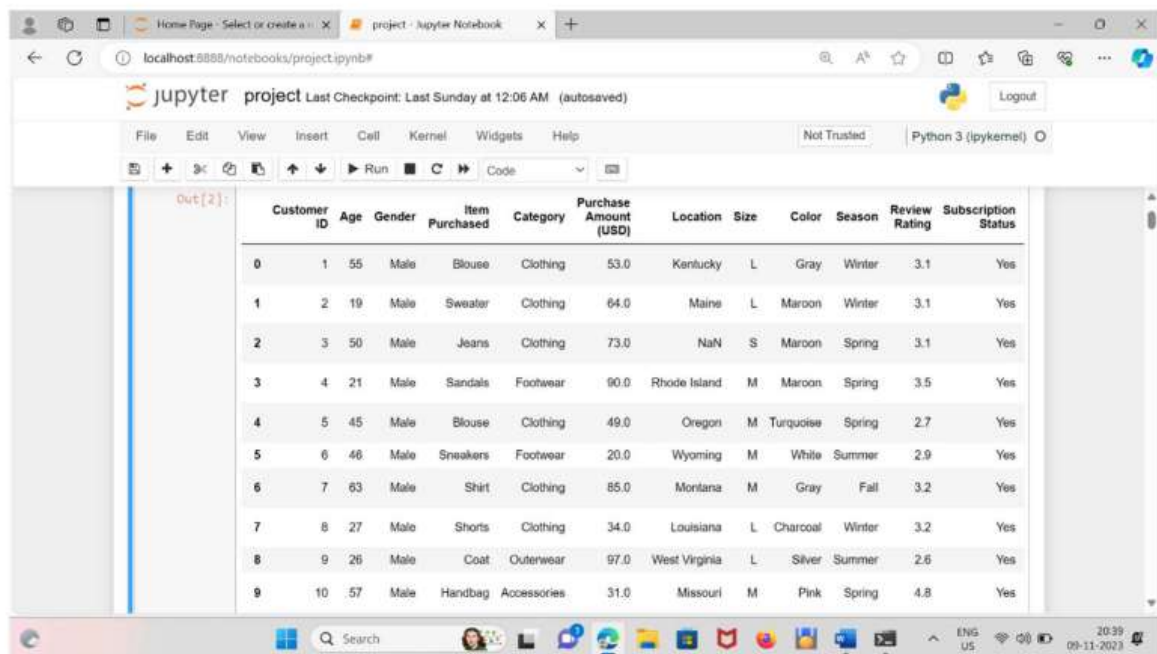
df

**OUTPUT:**



Data Frame of Shopping trends dataset

## Data Frame of Shopping trends dataset



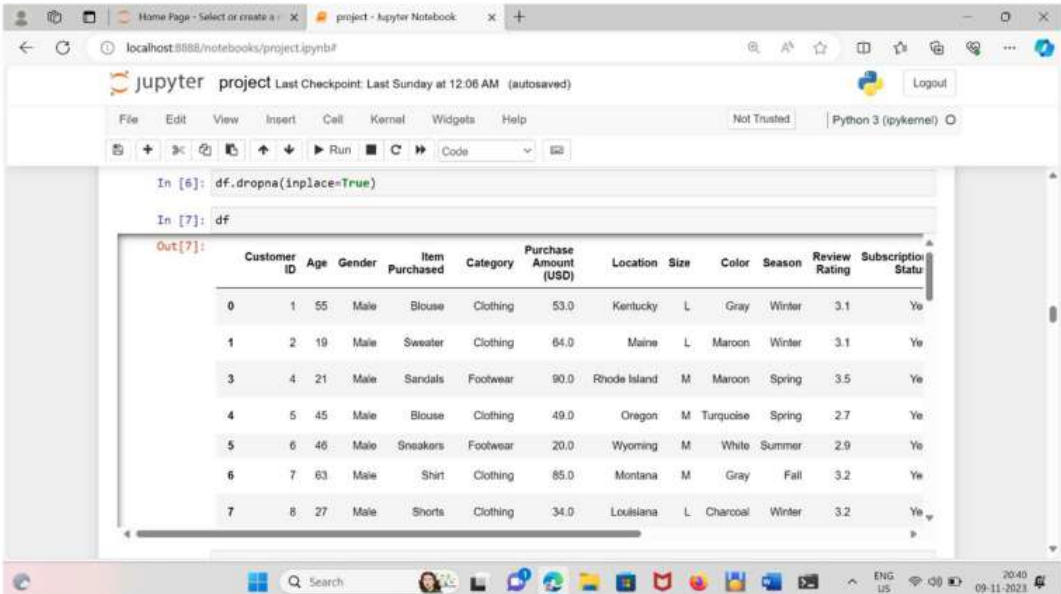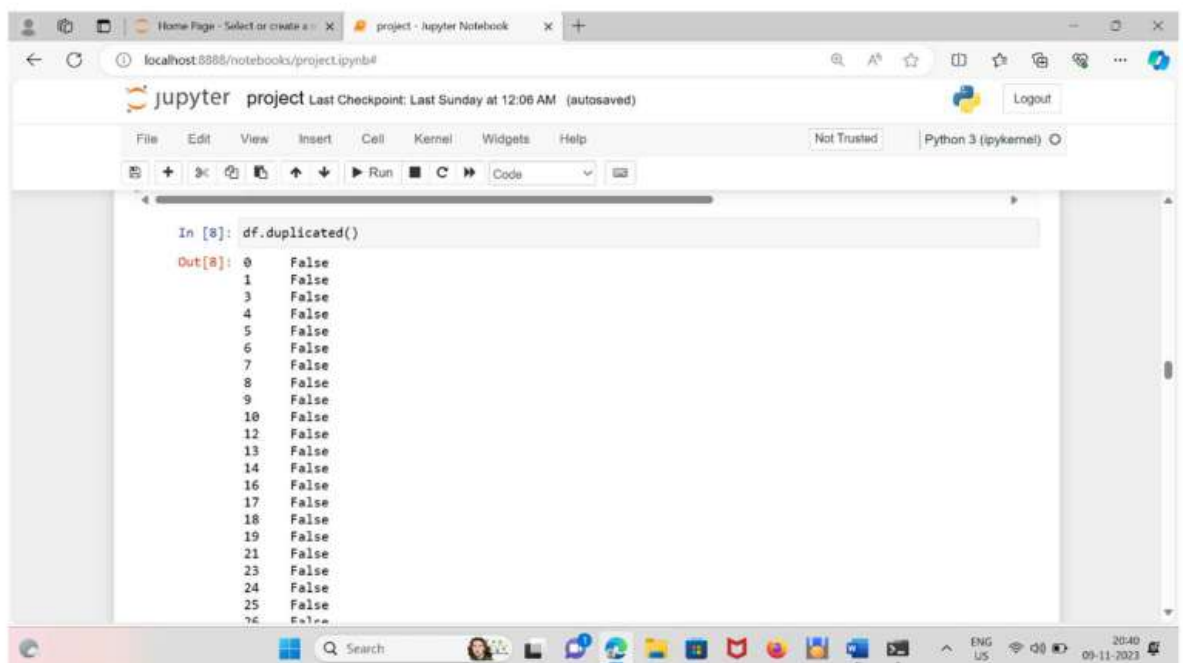| | Customer ID | Age | Gender | Item Purchased | Category | Purchase Amount (USD) | Location | Size | Color | Season | Review Rating | Subscription Status |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 55 | Male | Blouse | Clothing | 53.0 | Kentucky | L | Gray | Winter | 3.1 | Yes |
| 1 | 2 | 19 | Male | Sweater | Clothing | 64.0 | Maine | L | Maroon | Winter | 3.1 | Yes |
| 2 | 3 | 50 | Male | Jeans | Clothing | 73.0 | NaN | S | Maroon | Spring | 3.1 | Yes |
| 3 | 4 | 21 | Male | Sandals | Footwear | 90.0 | Rhode Island | M | Maroon | Spring | 3.5 | Yes |
| 4 | 5 | 45 | Male | Blouse | Clothing | 49.0 | Oregon | M | Turquoise | Spring | 2.7 | Yes |
| 5 | 6 | 46 | Male | Sneakers | Footwear | 20.0 | Wyoming | M | White | Summer | 2.9 | Yes |
| 6 | 7 | 63 | Male | Shirt | Clothing | 85.0 | Montana | M | Gray | Fall | 3.2 | Yes |
| 7 | 8 | 27 | Male | Shorts | Clothing | 34.0 | Louisiana | L | Charcoal | Winter | 3.2 | Yes |
| 8 | 9 | 26 | Male | Coat | Outerwear | 97.0 | West Virginia | L | Silver | Summer | 2.6 | Yes |
| 9 | 10 | 57 | Male | Handbag | Accessories | 31.0 | Missouri | M | Pink | Spring | 4.8 | Yes |

## DATA PREPROCESSING:

#Handling Missing Values

import pandas as pd

df=pd.read_csv(r"C:\\Users\Nithesh\Downloads\shopping_trends.csv")

df.dropna(inplace=True)

## OUTPUT:



Handling Missing values

#Identifying the duplicates in a entire dataframe

df.duplicated()

**OUTPUT:**



Duplicates in a Data Frame

#Identifying the duplicates in a particular column

df.duplicated(subset=['Item Purchased'])

**OUTPUT:**



Duplicates in a particular column

#Removing the duplicates in a particular colum
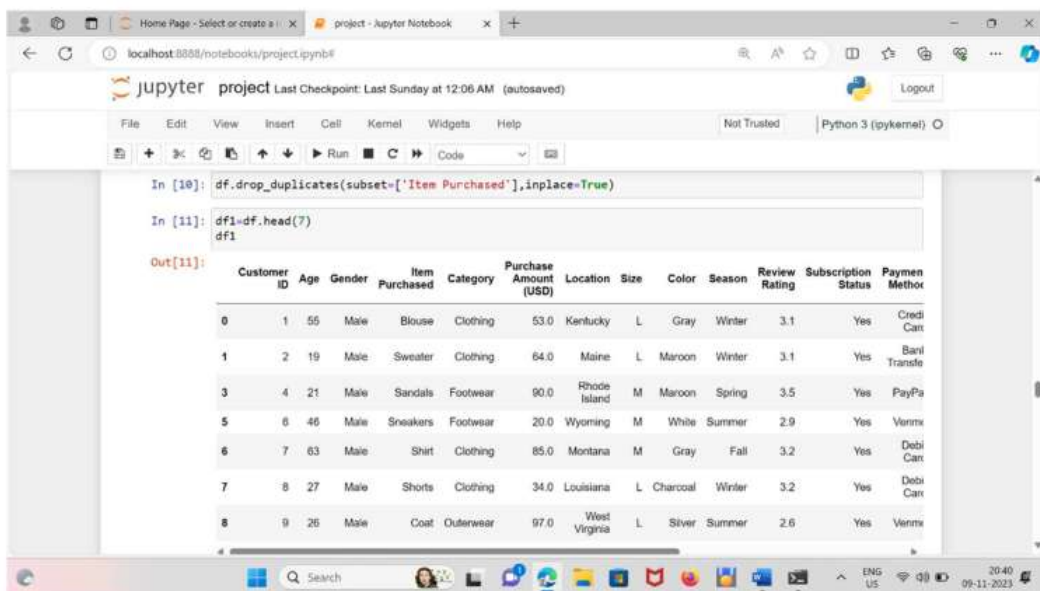
df.drop_duplicates(subset=['Item Purchased'],inplace=True)

#Accesing the top elements in a data frame using head

df1=df.head(7)

df1

**OUTPUT:**



Accessing top elements using head

#Accessing the bottom elements in a data frame using tail

df.tail()

**OUTPUT:**



Accessing bottom elements using tail