

```

import pandas as pd
import numpy as np
import json
import matplotlib.pyplot as plt
import seaborn as sns
import re # For Regular expression

# Load datasets
customer_file_path = "/content/drive/MyDrive/PEI DataSets/Customers.xlsx"
order_file_path = "/content/drive/MyDrive/PEI DataSets/Orders.csv"
shipping_file_path = "/content/drive/MyDrive/PEI DataSets/Shipping.json"

# customer_df = pd.read_excel(customer_file_path, engine="xlrd")
# Engine= xlrd as the file is in xls format which is old one
# customer_file_path = "/content/drive/MyDrive/PEI DataSets/Customers.xls"
#customer_df = pd.read_excel(customer_file_path, engine="xlrd")

# Load Customer Data
customer_df = pd.read_excel(customer_file_path)

# Load Order Data
order_df = pd.read_csv(order_file_path)

# Load Shipping Data
shipping_df = pd.read_json(shipping_file_path)

# Function to Perform EDA + Data Cleaning
# https://emojibd.org/stats-emojis emojis or icons are taken from this website for better look and feel

def perform_eda_and_clean(df, name):
    print(f"\n📊 EDA + Data Cleaning for {name} Dataset:")

    # 📄 1. Columns and Data Types
    print("\n📄 Columns and Data Types:")
    print(df.info())

    # 🔍 2. Printing First 5 Rows
    print("\n🔍 First 5 Rows:")
    print(df.head())

    # 🚩 3. Check for Missing values
    print("\n🚩 Missing Values Count:")
    print(df.isnull().sum())

    # 📊 4. Summary Statistics for Numerical Data
    print("\n📊 Summary Statistics (Numerical Data):")
    print(df.describe())

    # ✅ 5. Unique Values Per Column
    print("\n✅ Unique Values Per Column:")
    print(df.nunique())

    # 🔍 6. Check for Special Characters in String Columns
    print("\n🔍 Special Character Check:")

    # Define regex pattern for special characters (excluding space, a-z, A-Z, 0-9, and basic punctuation)
    special_char_pattern = re.compile(r'[^A-Za-z0-9\s.,]')

    for col in df.select_dtypes(include=["object"]).columns:
        # Find all special characters in the column
        special_chars = df[col].astype(str).apply(lambda x: set(re.findall(special_char_pattern, x)))

        # Get unique special characters found in the column
        unique_special_chars = set().union(*special_chars)

        if unique_special_chars:
            print(f"⚠️ Column `{col}` contains {len(unique_special_chars)} unique special characters: {unique_special_chars}")
        else:
            print(f"✅ Column `{col}` has no special characters.")

    # 🔥 7. Data Cleaning - Remove Special Characters
    df[col] = df[col].apply(lambda x: re.sub(special_char_pattern, '', str(x)))

```

```

# 🚀 8. Handle Missing Values
for col in df.columns:
    if df[col].isnull().sum() > 0: # If missing values exist
        if df[col].dtype == "object":
            df[col].fillna("Unknown", inplace=True) # Fill text columns with "Unknown"
        else:
            df[col].fillna(df[col].median(), inplace=True) # Fill numeric columns with median

# 🚀 9. Remove Duplicate Rows
before = len(df)
df.drop_duplicates(inplace=True)
after = len(df)
print(f"\n✅ Removed {before - after} duplicate rows.")

# 🚀 10. Ensure Correct Data Types
if "Age" in df.columns:
    df["Age"] = df["Age"].astype(int) # Convert Age to integer

if "Amount" in df.columns:
    df["Amount"] = df["Amount"].astype(float) # Convert Amount to float

print("\n✅ Data Cleaning Completed! Dataset is Ready for Analysis 🚀")
return df # Return cleaned DataFrame

# Perform EDA on each dataset
customer_df = perform_eda_and_clean(customer_df, "Customer")

# Perform EDA on Order dataset
order_df = perform_eda_and_clean(order_df, "Order")

```



EDA + Data Cleaning for Order Dataset:

Columns and Data Types:

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 250 entries, 0 to 249
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Order_ID    250 non-null   int64
1   Item        250 non-null   object
2   Amount      250 non-null   int64
3   Customer_ID 250 non-null   int64
dtypes: int64(3), object(1)
memory usage: 7.9+ KB
None

```

First 5 Rows:

|   | Order_ID | Item     | Amount | Customer_ID |
|---|----------|----------|--------|-------------|
| 0 | 1        | Keyboard | 400    | 139         |
| 1 | 2        | Mouse    | 300    | 250         |
| 2 | 3        | Monitor  | 12000  | 239         |
| 3 | 4        | Keyboard | 400    | 153         |
| 4 | 5        | Mousepad | 250    | 153         |

Missing Values Count:

```

Order_ID    0
Item        0
Amount      0
Customer_ID 0
dtype: int64

```

Summary Statistics (Numerical Data):

|       | Order_ID   | Amount       | Customer_ID |
|-------|------------|--------------|-------------|
| count | 250.000000 | 250.000000   | 250.000000  |
| mean  | 125.500000 | 2130.000000  | 130.404000  |
| std   | 72.312977  | 3575.43493   | 69.192711   |
| min   | 1.000000   | 200.000000   | 4.000000    |
| 25%   | 63.250000  | 300.000000   | 71.500000   |
| 50%   | 125.500000 | 400.000000   | 125.500000  |
| 75%   | 187.750000 | 1500.000000  | 190.750000  |
| max   | 250.000000 | 12000.000000 | 250.000000  |

Unique Values Per Column:

```

Order_ID    250
Item        8
Amount      9
Customer_ID 160
dtype: int64

```

- 🔍 Special Character Check:
- ✅ Column `Item` has no special characters.
- ✅ Removed 0 duplicate rows.
- ✅ Data Cleaning Completed! Dataset is Ready for Analysis 🚀

```
shipping_df = perform_eda_and_clean(shipping_df, "Shipping")
```

```
...
```

```
#Check for Duplicates in each dataset
```

```
df = order_df
```

```
duplicates = df[df.duplicated(keep=False)] # Get all duplicate rows
```

```
total_duplicates = df.duplicated().sum() # Count duplicate rows
```

```
print(f"\n🇮🇹 Checking Duplicates in {df} Dataset:")
```

```
print(f"📄 Total Duplicate Rows: {total_duplicates}") '''
```