```python
import pandas as pd
import numpy as np
import json
import matplotlib.pyplot as plt
import seaborn as sns
import re


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


# 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
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")
```

```
📊 EDA + Data Cleaning for Customer Dataset:

📝 Columns and Data Types:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 250 entries, 0 to 249
Data columns (total 5 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   Customer_ID  250 non-null    int64
 1   First        250 non-null    object
 2   Last         250 non-null    object
 3   Age          250 non-null    int64
 4   Country      250 non-null    object
dtypes: int64(2), object(3)
memory usage: 9.9+ KB
None

🔍 First 5 Rows:
   Customer_ID    First      Last  Age Country
0            1   Joseph      Rice   43     USA
1            2     Gary     Moore   71     USA
2            3     John    Walker   44      UK
3            4     Eric    Carter   38      UK
4            5  William   Jackson   58     UAE

 Missing Values Count:
Customer_ID    0
First          0
Last           0
Age            0
Country        0
dtype: int64

📈 Summary Statistics (Numerical Data):
       Customer_ID         Age
count   250.000000  250.000000
mean    125.500000   47.576000
std      72.312977   18.978011
min       1.000000   18.000000
25%      63.250000   29.000000
50%     125.500000   47.000000
75%     187.750000   63.000000
max     250.000000   80.000000

✅ Unique Values Per Column:
Customer_ID    250
First          171
Last           189
Age             62
Country          3
dtype: int64

🔍 Special Character Check:
⚠️ Column `First` contains 2 unique special characters: {'@', '!'}
✅ Column `Last` has no special characters.
✅ Column `Country` has no special characters.
```