# Assignment No-4

**CODE:-**

```python
import pandas as pd
import numpy as np
from google.colab import files

# Upload file
uploaded = files.upload()
filename = next(iter(uploaded))

# Load dataset
df = pd.read_csv(filename)

# Clean columns
df.columns = df.columns.str.strip().str.lower().str.replace(' ',
'_').str.replace(r'[^\w\s]', '', regex=True)
print("Columns after cleaning:", df.columns.tolist())

# Find sale price column
sale_price_col = None
for col in df.columns:
    if 'price' in col:
        sale_price_col = col
        print(f"Using '{col}' as sale price column")
        break

if sale_price_col is None:
    raise ValueError("No column related to 'price' found in the dataset!")

# Handle missing values
for col in df.columns:
    if df[col].dtype in ['float64', 'int64']:
        df[col].fillna(df[col].median(), inplace=True)
    else:
        df[col].fillna(df[col].mode()[0], inplace=True)
```

```python
# Filter example
if 'sale_date' in df.columns:
    df['sale_date'] = pd.to_datetime(df['sale_date'], errors='coerce')
    df = df[(df['sale_date'].dt.year >= 2018) & (df['sale_date'].dt.year <= 2023)]

if 'property_type' in df.columns:
    df = df[df['property_type'].str.lower() == 'single family']

print(f"Data shape after filtering: {df.shape}")

# Encode categoricals
cat_cols = df.select_dtypes(include=['object']).columns.tolist()
df_encoded = pd.get_dummies(df, columns=cat_cols, drop_first=True)

print(f"Shape after encoding: {df_encoded.shape}")

# Aggregate avg sale price by neighborhood
if 'neighborhood' in df.columns:
    avg_price_by_neighborhood = df.groupby('neighborhood')[sale_price_col].mean().reset_index().rename(columns={sale_price_col: 'avg_sale_price'})
    print(avg_price_by_neighborhood.head())

# Aggregate avg sale price by property type
if 'property_type' in df.columns:
    avg_price_by_property_type = df.groupby('property_type')[sale_price_col].mean().reset_index().rename(columns={sale_price_col: 'avg_sale_price'})
    print(avg_price_by_property_type.head())

# Outlier handling
Q1 = df[sale_price_col].quantile(0.25)
Q3 = df[sale_price_col].quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
```

```python
df_no_outliers = df[(df[sale_price_col] >= lower_bound) &
(df[sale_price_col] <= upper_bound)]
print(f"Shape after removing outliers: {df_no_outliers.shape}")

# Save cleaned data
df_no_outliers.to_csv('Cleaned_RealEstate_Prices.csv', index=False)
files.download('Cleaned_RealEstate_Prices.csv')
```

**OUTPUT :-**