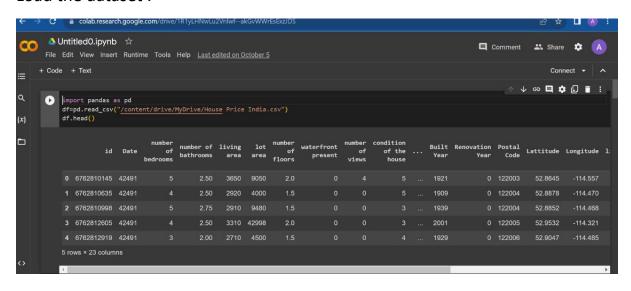# DA ASSIGNMENT – 3

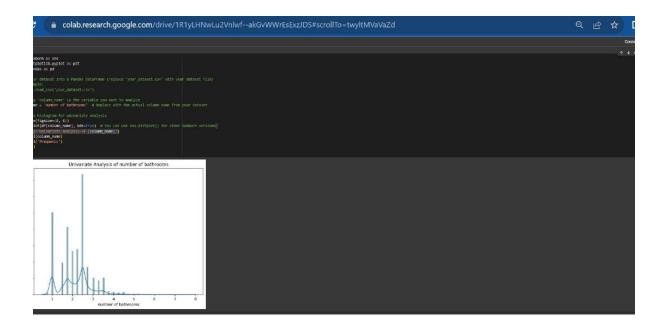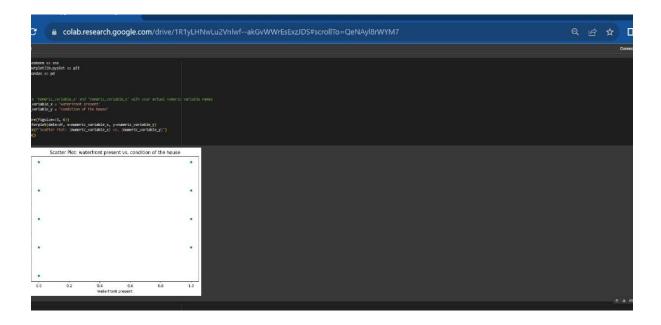**NAME**           :   G.POOJITH

**REGISTER NO**      : 723920243013

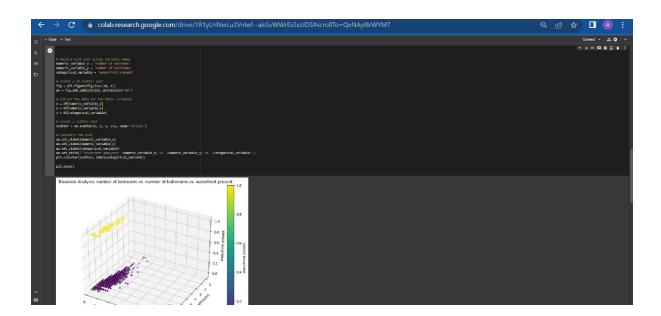Load the dataset :



Univariate Analysis :

## Bi - Variate Analysis :



## Multi-Variate Analysis :

Descriptive statistics on the dataset :



Handle the Missing values :



```python
# Define a function to handle missing values
def handle_missing_values(df):
    # Impute missing values in numeric columns with the mean
    numeric_columns = df.select_dtypes(include=['number']).columns
    df[numeric_columns] = df[numeric_columns].fillna(df[numeric_columns].mean())

    # Impute missing values in categorical columns with a specified value (e.g., 'Unknown')
    categorical_columns = df.select_dtypes(include=['object']).columns
    df[categorical_columns] = df[categorical_columns].fillna('Unknown')

    return df

# Apply the function to handle missing values
df = handle_missing_values(df)

# Check if there are any remaining missing values
missing_values_count = df.isnull().sum().sum()
if missing_values_count == 0:
    print("All missing values have been handled.")
else:
    print(f"There are still {missing_values_count} missing values in the dataset.")

# Save the cleaned dataset to a new file (optional)
df.to_csv('cleaned_dataset.csv', index=False)  # Replace with your desired file name
```

```
All missing values have been handled.
```