## Importing Libraries

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import re
from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error
from sklearn.datasets import load_boston
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
```

## Problem Statement

- You have been given a dataset that describes the houses in Boston. Now, based on the given features, you have to predict the house price.

## Creating a DataFrame

```python
boston = load_boston()
df = pd.DataFrame(boston.data)
```

## EDA - Exploratory Data Analysis

```python
df.head()
```

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|----|----|----|
| 0 | 0.00632 | 18.0 | 2.31 | 0.0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1.0 | 296.0 | 15.3 | 396.90 | 4.98 |
| 1 | 0.02731 | 0.0 | 7.07 | 0.0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2.0 | 242.0 | 17.8 | 396.90 | 9.14 |
| 2 | 0.02729 | 0.0 | 7.07 | 0.0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2.0 | 242.0 | 17.8 | 392.83 | 4.03 |
| 3 | 0.03237 | 0.0 | 2.18 | 0.0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3.0 | 222.0 | 18.7 | 394.63 | 2.94 |
| 4 | 0.06905 | 0.0 | 2.18 | 0.0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3.0 | 222.0 | 18.7 | 396.90 | 5.33 |

- Adding columns

```python
df.columns = boston.feature_names
```

```python
df.head()
```

|   | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LS |
|---|------|----|-------|------|-----|----|----|-----|-----|-----|---------|---|----|
| 0 | 0.00632 | 18.0 | 2.31 | 0.0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1.0 | 296.0 | 15.3 | 396.90 | ⋮ |
| 1 | 0.02731 | 0.0 | 7.07 | 0.0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2.0 | 242.0 | 17.8 | 396.90 | ⋮ |
| 2 | 0.02729 | 0.0 | 7.07 | 0.0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2.0 | 242.0 | 17.8 | 392.83 | ⋮ |
| 3 | 0.03237 | 0.0 | 2.18 | 0.0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3.0 | 222.0 | 18.7 | 394.63 | ⋮ |
| 4 | 0.06905 | 0.0 | 2.18 | 0.0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3.0 | 222.0 | 18.7 | 396.90 | ⋮ |

## Columns Informations

- CRIM per capita crime rate by town
- ZN proportion of residential land zoned for lots over 25,000 sq.ft.
- INDUS proportion of non-retail business acres per town
- CHAS Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
- NOX nitric oxides concentration (parts per 10 million)
- RM average number of rooms per dwelling
- AGE proportion of owner-occupied units built prior to 1940
- DIS weighted distances to five Boston employment centres

- RAD index of accessibility to radial highways
- TAX full-value property-tax rate per 10,000usd
- PTRATIO pupil-teacher ratio by town
- B 1000(Bk - 0.63)^2 where Bk is the proportion of blacks by town
- LSTAT % lower status of the population

## Adding the target column into the DataFrame

```
df['PRICE'] = boston.target
```

```
df.head()
```

|   | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LS |
|---|------|----|-------|------|-----|----|-----|-----|-----|-----|---------|---|----|
| 0 | 0.00632 | 18.0 | 2.31 | 0.0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1.0 | 296.0 | 15.3 | 396.90 | |
| 1 | 0.02731 | 0.0 | 7.07 | 0.0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2.0 | 242.0 | 17.8 | 396.90 | |
| 2 | 0.02729 | 0.0 | 7.07 | 0.0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2.0 | 242.0 | 17.8 | 392.83 | |
| 3 | 0.03237 | 0.0 | 2.18 | 0.0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3.0 | 222.0 | 18.7 | 394.63 | |
| 4 | 0.06905 | 0.0 | 2.18 | 0.0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3.0 | 222.0 | 18.7 | 396.90 | |

```
df.tail()
```

|   | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | L |
|---|------|----|-------|------|-----|----|-----|-----|-----|-----|---------|---|---|
| 501 | 0.06263 | 0.0 | 11.93 | 0.0 | 0.573 | 6.593 | 69.1 | 2.4786 | 1.0 | 273.0 | 21.0 | 391.99 | |
| 502 | 0.04527 | 0.0 | 11.93 | 0.0 | 0.573 | 6.120 | 76.7 | 2.2875 | 1.0 | 273.0 | 21.0 | 396.90 | |
| 503 | 0.06076 | 0.0 | 11.93 | 0.0 | 0.573 | 6.976 | 91.0 | 2.1675 | 1.0 | 273.0 | 21.0 | 396.90 | |
| 504 | 0.10959 | 0.0 | 11.93 | 0.0 | 0.573 | 6.794 | 89.3 | 2.3889 | 1.0 | 273.0 | 21.0 | 393.45 | |
| 505 | 0.04741 | 0.0 | 11.93 | 0.0 | 0.573 | 6.030 | 80.8 | 2.5050 | 1.0 | 273.0 | 21.0 | 396.90 | |

```
df.shape
```

```
(506, 14)
```

```
df.columns
```

```
Index(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX',
       'PTRATIO', 'B', 'LSTAT', 'PRICE'],
      dtype='object')
```

```
df.dtypes
```

```
CRIM       float64
ZN         float64
INDUS      float64
CHAS       float64
NOX        float64
RM         float64
AGE        float64
DIS        float64
RAD        float64
TAX        float64
PTRATIO    float64
B          float64
LSTAT      float64
PRICE      float64
dtype: object
```

```
df.nunique()
```

```
CRIM       504
ZN          26
INDUS       76
CHAS         2
NOX         81
RM         446
AGE        356
DIS        412
RAD          9
TAX         66
PTRATIO     46
```

```
B          357
LSTAT      455
PRICE      229
dtype: int64
```

df.isnull()

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | False | False | False | False | False | False | False | False | False | False | False | False | |
| **1** | False | False | False | False | False | False | False | False | False | False | False | False | |
| **2** | False | False | False | False | False | False | False | False | False | False | False | False | |
| **3** | False | False | False | False | False | False | False | False | False | False | False | False | |
| **4** | False | False | False | False | False | False | False | False | False | False | False | False | |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| **501** | False | False | False | False | False | False | False | False | False | False | False | False | |
| **502** | False | False | False | False | False | False | False | False | False | False | False | False | |
| **503** | False | False | False | False | False | False | False | False | False | False | False | False | |
| **504** | False | False | False | False | False | False | False | False | False | False | False | False | |
| **505** | False | False | False | False | False | False | False | False | False | False | False | False | |

506 rows × 14 columns

df.isnull().sum()

```
CRIM       0
ZN         0
INDUS      0
CHAS       0
NOX        0
RM         0
AGE        0
DIS        0
RAD        0
TAX        0
PTRATIO    0
B          0
LSTAT      0
PRICE      0
dtype: int64
```

df.describe()

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AG |
|---|---|---|---|---|---|---|---|
| **count** | 506.000000 | 506.000000 | 506.000000 | 506.000000 | 506.000000 | 506.000000 | 506.00000 |
| **mean** | 3.613524 | 11.363636 | 11.136779 | 0.069170 | 0.554695 | 6.284634 | 68.57490 |
| **std** | 8.601545 | 23.322453 | 6.860353 | 0.253994 | 0.115878 | 0.702617 | 28.14886 |
| **min** | 0.006320 | 0.000000 | 0.460000 | 0.000000 | 0.385000 | 3.561000 | 2.90000 |
| **25%** | 0.082045 | 0.000000 | 5.190000 | 0.000000 | 0.449000 | 5.885500 | 45.02500 |
| **50%** | 0.256510 | 0.000000 | 9.690000 | 0.000000 | 0.538000 | 6.208500 | 77.50000 |
| **75%** | 3.677083 | 12.500000 | 18.100000 | 0.000000 | 0.624000 | 6.623500 | 94.07500 |
| **max** | 88.976200 | 100.000000 | 27.740000 | 1.000000 | 0.871000 | 8.780000 | 100.00000 |

df.corr

```
<bound method DataFrame.corr of          CRIM    ZN  INDUS  CHAS    NOX  ...    TAX  PTRATIO       B  LSTAT  PRICE
0     0.00632  18.0   2.31   0.0  0.538  ...  296.0     15.3  396.90   4.98   24.0
1     0.02731   0.0   7.07   0.0  0.469  ...  242.0     17.8  396.90   9.14   21.6
2     0.02729   0.0   7.07   0.0  0.469  ...  242.0     17.8  392.83   4.03   34.7
3     0.03237   0.0   2.18   0.0  0.458  ...  222.0     18.7  394.63   2.94   33.4
4     0.06905   0.0   2.18   0.0  0.458  ...  222.0     18.7  396.90   5.33   36.2
..        ...   ...    ...   ...    ...  ...    ...      ...     ...    ...    ...
501   0.06263   0.0  11.93   0.0  0.573  ...  273.0     21.0  391.99   9.67   22.4
502   0.04527   0.0  11.93   0.0  0.573  ...  273.0     21.0  396.90   9.08   20.6
503   0.06076   0.0  11.93   0.0  0.573  ...  273.0     21.0  396.90   5.64   23.9
504   0.10959   0.0  11.93   0.0  0.573  ...  273.0     21.0  393.45   6.48   22.0
505   0.04741   0.0  11.93   0.0  0.573  ...  273.0     21.0  396.90   7.88   11.9

[506 rows x 14 columns]>
```
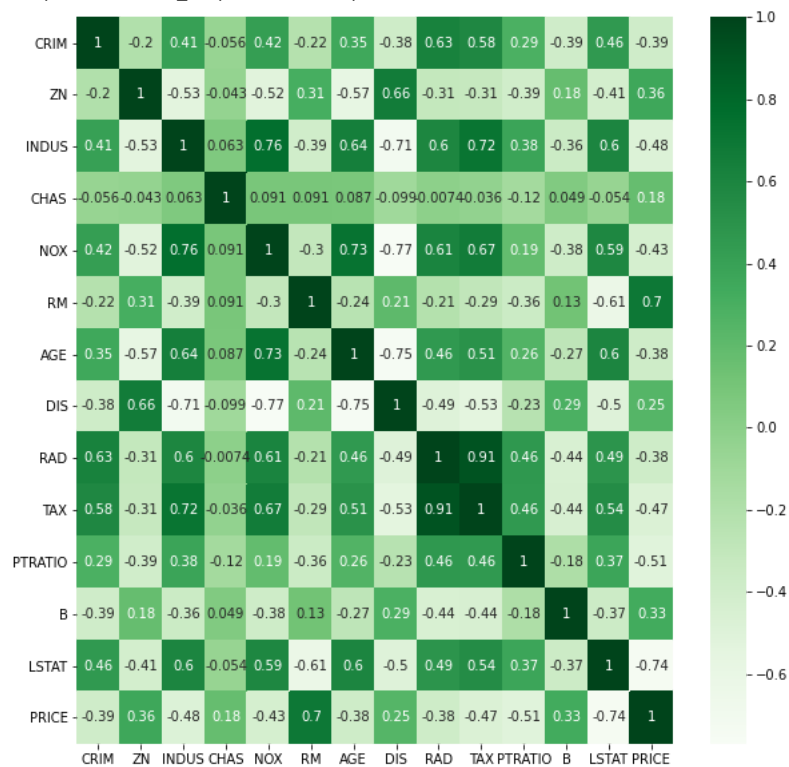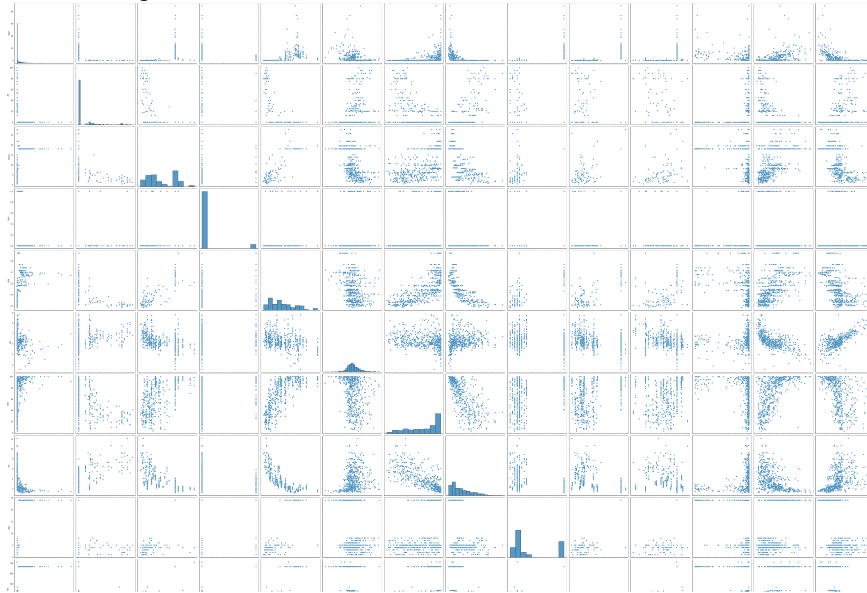
```python
plt.figure(figsize=(10,10))
sns.heatmap(data=df.corr(), annot=True, cmap='Greens')
```

<matplotlib.axes._subplots.AxesSubplot at 0x7fe05ed6ff50>



```python
sns.pairplot(df, size=5)
```
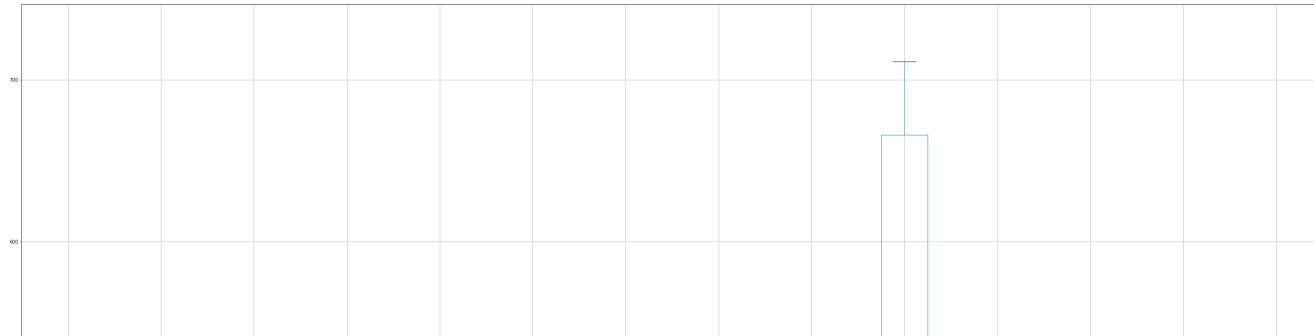
```
/usr/local/lib/python3.7/dist-packages/seaborn/axisgrid.py:1969: UserWarning: The `si
  warnings.warn(msg, UserWarning)
<seaborn.axisgrid.PairGrid at 0x7fe057603890>
```



```
# Plot a Boxplot
plt.figure(figsize=(50,50))
df.boxplot()
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fe04b646950>
```



```
# Minimum Price
df.PRICE.min()
```

```
    5.0
```

```
# Masximum Price
df.PRICE.max()
```

```
    50.0
```

```
# Standard Deviation
df.PRICE.std()
```

```
    9.19710408737982
```

## ⌄ Export the dataset

```
df.to_csv('boston_datset.csv',)
```

## ⌄ Machine Learning - Linear Regression

```
df.head()
```

|   | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LS |
|---|------|-----|-------|------|-------|-------|------|--------|-----|-------|---------|--------|----|
| 0 | 0.00632 | 18.0 | 2.31 | 0.0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1.0 | 296.0 | 15.3 | 396.90 | |
| 1 | 0.02731 | 0.0 | 7.07 | 0.0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2.0 | 242.0 | 17.8 | 396.90 | |
| 2 | 0.02729 | 0.0 | 7.07 | 0.0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2.0 | 242.0 | 17.8 | 392.83 | |
| 3 | 0.03237 | 0.0 | 2.18 | 0.0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3.0 | 222.0 | 18.7 | 394.63 | |
| 4 | 0.06905 | 0.0 | 2.18 | 0.0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3.0 | 222.0 | 18.7 | 396.90 | |

```
X = np.array(df.drop('PRICE', axis=1))
y = np.array(df.PRICE)
```

```
# X = boston.data
# y = boston.target
```

## ⌄ Splitting the data

```
X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2,
                                        random_state=42)
```

```
len(X_train)
```

```
    404
```

```
len(y_train)
```

```
    404
```

```
len(X_test)
```

```
    102
```

```
len(y_test)
```

```
    102
```

## ⌄ Choosing the model

```
model = LinearRegression()
```

## ⌄ Fitting/Train the model

```
model.fit(X_train,y_train)
```

```
    LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
# Intercept Value
model.intercept_
```

```
    30.24675099392392
```

```
# Coefficient Value
model.coef_
```

```
    array([-1.13055924e-01,  3.01104641e-02,  4.03807204e-02,  2.78443820e+00,
           -1.72026334e+01,  4.43883520e+00, -6.29636221e-03, -1.44786537e+00,
            2.62429736e-01, -1.06467863e-02, -9.15456240e-01,  1.23513347e-02,
           -5.08571424e-01])
```

## ⌄ Prediction

```
y_test
```

```
    array([23.6, 32.4, 13.6, 22.8, 16.1, 20. , 17.8, 14. , 19.6, 16.8, 21.5,
           18.9,  7. , 21.2, 18.5, 29.8, 18.8, 10.2, 50. , 14.1, 25.2, 29.1,
           12.7, 22.4, 14.2, 13.8, 20.3, 14.9, 21.7, 18.3, 23.1, 23.8, 15. ,
           20.8, 19.1, 19.4, 34.7, 19.5, 24.4, 23.4, 19.7, 28.2, 50. , 17.4,
           22.6, 15.1, 13.1, 24.2, 19.9, 24. , 18.9, 35.4, 15.2, 26.5, 43.5,
           21.2, 18.4, 28.5, 23.9, 18.5, 25. , 35.4, 31.5, 20.2, 24.1, 20. ,
           13.1, 24.8, 30.8, 12.7, 20. , 23.7, 10.8, 20.6, 20.8,  5. , 20.1,
           48.5, 10.9,  7. , 20.9, 17.2, 20.9,  9.7, 19.4, 29. , 16.4, 25. ,
           25. , 17.1, 23.2, 10.4, 19.6, 17.2, 27.5, 23. , 50. , 17.9,  9.6,
           17.2, 22.5, 21.4])
```

```
y_pred = model.predict(X_test)
```

```
y_pred
```

```
    array([28.99672362, 36.02556534, 14.81694405, 25.03197915, 18.76987992,
           23.25442929, 17.66253818, 14.34119   , 23.01320703, 20.63245597,
           24.90850512, 18.63883645, -6.08842184, 21.75834668, 19.23922576,
           26.19319733, 20.64773313,  5.79472718, 40.50033966, 17.61289074,
           27.24909479, 30.06625441, 11.34179277, 24.16077616, 17.86058499,
           15.83609765, 22.78148106, 14.57704449, 22.43626052, 19.19631835,
           22.43383455, 25.21979081, 25.93909562, 17.70162434, 16.76911711,
           16.95125411, 31.23340153, 20.13246729, 23.76579011, 24.6322925 ,
           13.94204955, 32.25576301, 42.67251161, 17.32745046, 27.27618614,
           16.99310991, 14.07009109, 25.90341861, 20.29485982, 29.95339638,
           21.28860173, 34.34451856, 16.04739105, 26.22562412, 39.53939798,
           22.57950697, 18.84531367, 32.72531661, 25.0673037 , 12.88628956,
           22.68221908, 30.48287757, 31.52626806, 15.90148607, 20.22094826,
           16.71089812, 20.52384893, 25.96356264, 30.61607978, 11.59783023,
           20.51232627, 27.48111878, 11.01962332, 15.68096344, 23.79316251,
            6.19929359, 21.6039073 , 41.41377225, 18.76548695,  8.87931901,
           20.83076916, 13.25620627, 20.73963699,  9.36482222, 23.22444271,
           31.9155003 , 19.10228271, 25.51579303, 29.04256769, 20.14358566,
           25.5859787 ,  5.70159447, 20.09474756, 14.95069156, 12.50395648,
           20.72635294, 24.73957161, -0.164237  , 13.68486682, 16.18359697,
           22.27621999, 24.47902364])
```

## ∨ Testing the model performance

```
model.score(X_test,y_test)
```

```
     0.6687594935356307
```

```
# R squared
r2_score(y_test,y_pred)
```

```
     0.6687594935356307
```

```
# Adjusted R squared
```

```
# MSE
mean_squared_error(y_test,y_pred)
```

```
     24.291119474973616
```

```
# MAE
mean_absolute_error(y_test,y_pred)
```

```
     3.1890919658878523
```

```
# RMSE
np.sqrt(mean_squared_error(y_test,y_pred))
```

```
     4.9286021826653466
```

```
plt.scatter(y_test,y_pred)
plt.xlabel("Actual Price")
plt.ylabel("Predicted Price")
plt.grid()
plt.plot([min(y_test),max(y_test)],[min(y_pred),max(y_pred)], color='red')
plt.title('Actual Price V/s Predicted Price')
```

```
     Text(0.5, 1.0, 'Actual Price V/s Predicted Price')
```