

Import the Dependencies

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import sklearn.datasets
from sklearn.model_selection import train_test_split
from xgboost import XGBRegressor
from sklearn import metrics
import warnings
warnings.filterwarnings('ignore')
```

Importing the Boston House Price Dataset

```
In [2]: house_price_dataset = sklearn.datasets.load_boston()

In [3]: print(house_price_dataset)
{'data': array([[6.3200e-03, 1.8000e+01, 2.3100e+00, ..., 1.5300e+01, 3.9690e+02,
4.9800e+00],
[2.7310e-02, 0.0000e+00, 7.0700e+00, ..., 1.7800e+01, 3.9690e+02,
9.1400e+00],
[2.7290e-02, 0.0000e+00, 7.0700e+00, ..., 1.7800e+01, 3.9283e+02,
4.0300e+00],
...,
[6.0760e-02, 0.0000e+00, 1.1930e+01, ..., 2.1000e+01, 3.9690e+02,
5.6400e+00],
[1.0959e-01, 0.0000e+00, 1.1930e+01, ..., 2.1000e+01, 3.9345e+02,
6.4800e+00],
[4.7410e-02, 0.0000e+00, 1.1930e+01, ..., 2.1000e+01, 3.9690e+02,
7.8800e+00], 'target': array([24., 21.6, 34.7, 33.4, 36.2, 28.7, 22.9, 27.1, 16.5, 18.9, 15.,
18.9, 21.7, 20.4, 18.2, 19.9, 23.1, 17.5, 20.2, 18.2, 13.6, 19.6,
15.2, 14.5, 15.6, 13.9, 16.6, 14.8, 18.7, 21., 12.7, 14.5, 13.2,
13.1, 13.5, 18.9, 20., 21., 24.7, 30.8, 34.9, 26.6, 25.3, 24.7,
21.2, 19.3, 20., 16.6, 14.4, 19.4, 19.7, 20.5, 25., 23.4, 18.9,
35.4, 24.7, 31.6, 23.3, 19.6, 18.7, 18., 22.2, 25., 33., 23.5,
19.4, 22., 17.4, 28.8, 24.2, 21.7, 22.8, 23.4, 24.1, 21.4, 28.,
20.8, 21.2, 20.3, 28., 23.9, 24.8, 22.9, 23.9, 26.6, 22.5, 22.2,
23.6, 28.7, 22.6, 22., 22.9, 25., 20.6, 28.4, 21.4, 38.7, 43.8,
33.2, 27.5, 26.5, 18.6, 19.3, 20.1, 19.5, 19.5, 28.4, 19.8, 19.4,
21.7, 22.8, 18.8, 18.7, 18.5, 18.5, 21.2, 19.2, 20.4, 19.3, 22.,
20.3, 20.5, 17.3, 18.8, 21.4, 15.7, 16.2, 18., 14.3, 19.2, 19.6,
23., 18.4, 15.6, 18.1, 17.4, 17.1, 13.3, 17.8, 14., 14.4, 13.4,
15.6, 11.8, 13.8, 15.6, 14.6, 17.8, 15.4, 21.5, 19.6, 15.3, 19.4,
17., 15.6, 13.1, 41.3, 24.3, 23.3, 27., 50., 50., 50., 22.7,
25., 50., 23.8, 23.8, 22.3, 17.4, 19.1, 23.1, 23.6, 22.6, 29.4,
23.2, 24.6, 29.9, 37.2, 39.8, 36.2, 37.9, 32.5, 26.4, 29.6, 50.,
32., 29.8, 34.9, 37., 30.5, 36.4, 31.1, 29.1, 50., 33.3, 30.3,
34.6, 34.9, 32.9, 24.1, 42.3, 48.5, 50., 22.6, 24.4, 22.5, 24.4,
20., 21.7, 19.3, 22.4, 28.1, 23.7, 25., 23.3, 28.7, 21.5, 23.,
26.7, 21.7, 27.5, 30.1, 44.8, 50., 37.6, 31.6, 46.7, 31.5, 24.3,
31.7, 41.7, 48.3, 28., 24., 25.1, 31.5, 23.7, 23.3, 22., 20.1,
22.2, 23.7, 17.6, 18.5, 24.3, 20.5, 24.5, 26.2, 24.4, 24.8, 29.6,
42.8, 21.9, 20.9, 44., 50., 36., 30.1, 33.8, 43.1, 48.8, 31.,
36.5, 22.8, 30.7, 50., 43.5, 20.7, 21.1, 25.2, 24.4, 35.2, 32.4,
32., 31.2, 33.1, 28.1, 35.1, 45.4, 35.4, 46., 50., 32.2, 22.,
20.1, 23.2, 22.3, 24.8, 28.5, 37.9, 27.9, 23.9, 21.7, 28.6, 27.1,
20.3, 22.5, 29., 24.8, 22., 26.4, 33.1, 36.1, 28.4, 33.4, 28.2,
22.8, 20.3, 16.1, 22.1, 19.4, 21.6, 23.8, 16.2, 17.8, 19.8, 23.1,
21., 23.8, 23.1, 20.4, 18.5, 25., 24.6, 23., 22.2, 19.3, 22.6,
19.8, 17.1, 19.4, 22.2, 20.7, 21., 19.5, 18.5, 20.6, 19., 18.7,
32.7, 16.5, 23.9, 31.2, 17.5, 17.2, 23.1, 24.5, 26.6, 22.9, 24.1,
18.6, 30.1, 18.2, 20.6, 17.8, 21.7, 22.7, 22.6, 25., 19.9, 20.8,
16.8, 21.9, 27.5, 21.9, 23.1, 50., 50., 50., 50., 50., 16.4, 17.7,
13.8, 15., 13.9, 13.3, 13.1, 10.2, 10.4, 10.9, 11.3, 12.3, 8.8,
7.2, 10.5, 7.4, 10.2, 11.5, 15.1, 23.2, 9.7, 13.8, 12.7, 13.1,
12.5, 8.5, 5., 6.3, 5.6, 7.2, 12.1, 8.3, 8.5, 5., 11.9,
27.9, 17.2, 27.5, 15., 17.2, 17.9, 16.3, 7., 7.2, 7.5, 10.4,
8.8, 8.4, 16.7, 14.2, 20.8, 13.4, 11.7, 8.3, 10.2, 10.9, 11.,
9.5, 14.5, 14.1, 16.1, 14.3, 11.7, 13.4, 9.6, 8.7, 8.4, 12.8,
10.5, 17.1, 18.4, 15.4, 10.8, 11.8, 14.9, 12.6, 14.1, 13., 13.4,
15.2, 16.1, 17.8, 14.9, 14.1, 12.7, 13.5, 14.9, 20., 16.4, 17.7,
19.5, 20.2, 21.4, 10.9, 19., 19.1, 19.1, 20.1, 19.9, 19.6, 23.2,
29.8, 13.8, 13.3, 16.7, 12., 14.6, 21.4, 23., 23.7, 25., 21.8,
20.6, 21.2, 19.1, 20.6, 15.2, 7., 8.1, 13.6, 20.1, 21.8, 24.5,
23.1, 19.7, 19.3, 21.2, 17.5, 16.8, 22.4, 20.6, 23.9, 22., 11.9]),
'feature_names': array(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD',
'TAX', 'PTRATIO', 'B', 'LSTAT'], dtype='<U?>'), 'DESCR': '...', '_boston_dataset': '\\\\nBoston house prices dataset\\\\n-----\\\\n\\\\n''Data Set Characteristics:''
\\\\n
\\\\n Number of Instances: 506
\\\\n Number of Attributes: 13 numeric/categorical predictive. Median value (attribute 14) is usually the target.
\\\\n Attribute Information
(in order):
\\\\n - CRIM per capita crime rate by town
\\\\n - ZN proportion of residential land zoned for lots over 25,000 sq.ft.
\\\\n - INDUS proportion of non-retail business acres per town
\\\\n - CHAS Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
\\\\n - NOX nitric oxides concentration (parts per million)
\\\\n - RM average number of rooms per dwelling
\\\\n - AGE proportion of owner-occupied units built prior to 1940
\\\\n - DIS weighted distances to five Boston employment centres
\\\\n - RAD index of accessibility to radial highways
\\\\n - TAX full-value property-tax rate per $10,000
\\\\n - PTRATIO pupil-teacher ratio by town
\\\\n - B lower status of the population
\\\\n - LSTAT % lower status of the population
MEDV Median value of owner-occupied homes in $1000's
\\\\n Missing Attribute Values: None
\\\\n Creator: Harrison, D. and Rubinfeld, D.L.
\\\\n This is a copy of UCI ML housing dataset.
\\\\n https://archive.ics.uci.edu/ml/machine-learning-databases/housing/
\\\\n This dataset was taken from the StatLib library which is maintained at Carnegie Mellon University.
\\\\n The Boston house-price data of Harrison, D. and Rubinfeld, D.L. 'Hedonic prices and the demand for clean air', J. Environ. Economics & Management, vol.5, 81-102, 1978.
Used in Belsley, Kuh & Welsch, 'Regression diagnostics...', Wiley, 1980. N.B. Various transformations are used in the table on pages 244-262 of the latter.
\\\\n The Boston house price data has been used in many machine learning papers that address regression problems
\\\\n Topic: References
\\\\n Belsley, Kuh & Welsch, 'Regression diagnostics
s: Identifying Influential Data and Sources of Collinearity', Wiley, 1980.
244-261.
\\\\n Quinlan, R. (1993). Combining Instance-Based and Model-Based Learning. In Proceedings on the Tenth International Conference of Machine Learning, 236-243, University of Massachusetts, Amherst.
Morgan Kaufmann.
\\\\n', 'filename': 'boston_house_prices.csv', 'data_module': 'sklearn.datasets.data'})
```

```
In [4]: # Loading the dataset to a Pandas DataFrame
house_price_dataframe = pd.DataFrame(house_price_dataset.data, columns = house_price_dataset.feature_names)
```

```
In [5]: # Print First 5 rows of our DataFrame
house_price_dataframe.head()
```

Out[5]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT
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

```
In [6]: # add the target (price) column to the DataFrame
house_price_dataframe['price'] = house_price_dataset.target
```

```
In [7]: house_price_dataframe.head()
```

Out[7]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	price
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	24.0
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	21.6
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	34.7
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	33.4
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	36.2

```
In [8]: # checking the number of rows and columns in the data frame
house_price_dataframe.shape
```

```
Out[8]: (506, 14)
```

```
In [9]: # check for missing values
house_price_dataframe.isnull().sum()
```

Out[9]:

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

```
In [10]: # statistical measures of the dataset
house_price_dataframe.describe()
```

Out[10]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	price
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	3.611324	11.363636	11.36779	0.069170	0.564695	6.286334	68.574901	3.795043	9.549407	408.237154	18.455534	356.674032	12.653063	22.532806
std	8.601545	23.322453	6.860353	0.253994	0.115878	0.702617	28.148861	2.105710	8.707259	168.537116	2.164466	91.294864	7.141062	9.197104
min	0.006320	0.000000	0.460000	0.000000	0.395000	3.561000	2.900000	1.129600	1.000000	187.000000	12.600000	0.320000	1.730000	5.000000
25%	0.082045	0.000000	5.190000	0.000000	0.449000	5.885500	45.025000	2.100175	4.000000	279.000000	17.400000	375.377500	6.950000	17.025000
50%	0.265610	0.000000	9.690000	0.000000	0.538000	6.208500	77.500000	3.207450	5.000000	330.000000	19.050000	381.440000	11.360000	21.200000
75%	3.877083	12.500000	18.100000	0.000000	0.624000	6.623500	94.075000	5.188425	24.000000	666.000000	20.200000	396.225000	16.955000	25.000000
max	88.9575200	100.000000	27.740000	1.000000	0.871000	8.780000	100.000000	12.126500	24.000000	711.000000	22.000000	396.900000	37.970000	50.000000

Understanding the correlation between various features in the dataset

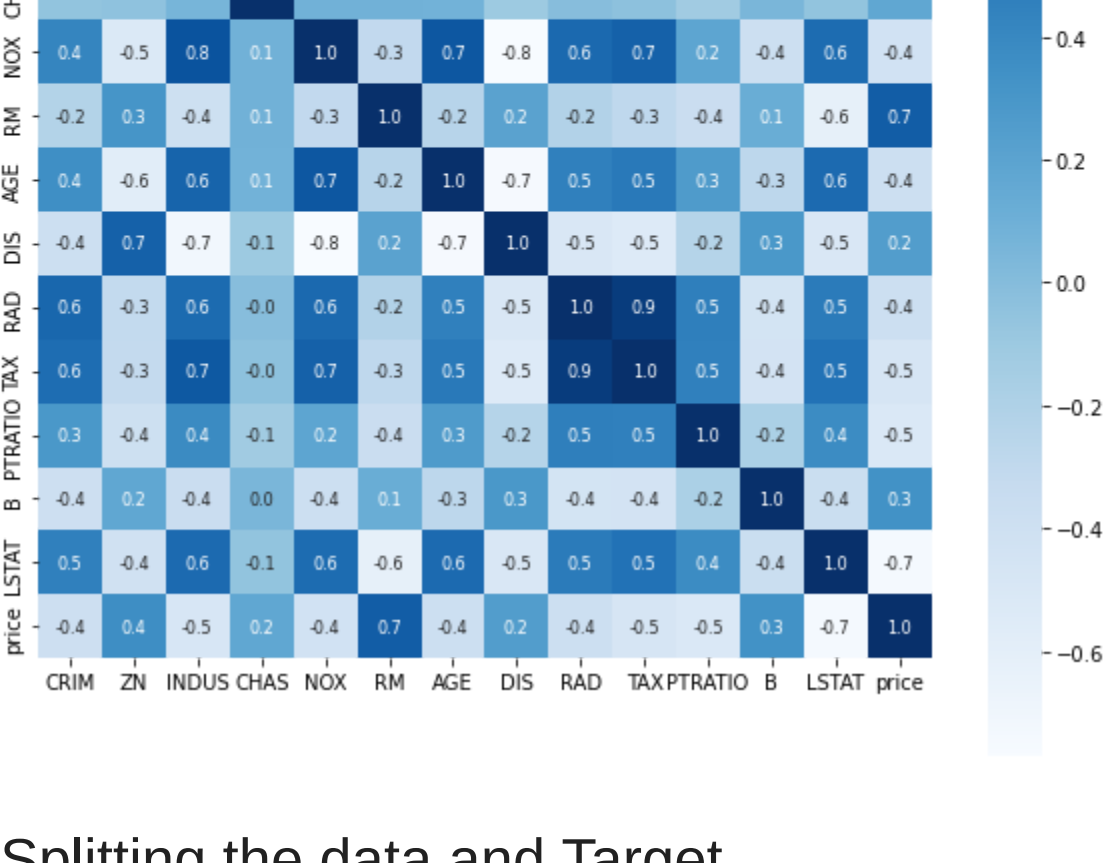
1. Positive Correlation

2. Negative Correlation

```
In [11]: correlation = house_price_dataframe.corr()
```

```
In [12]: # constructing a heatmap to understand the correlation
plt.figure(figsize=(10,10))
sns.heatmap(correlation, char=True, square=True, fmt='.1f', annot=True, annot_kws={'size':8}, cmap='Blues')
```

```
Out[12]: <AxesSubplot>
```



Splitting the data and Target

```
In [13]: X = house_price_dataframe.drop(['price'], axis=1)
Y = house_price_dataframe['price']
```

```
In [14]: print(X)
print(Y)
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	\
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	
...	
501	0.06263	0.0	11.93	0.0	0.573	6.193	69.1	2.4770	1.0	273.0	
502	0.04527	0.0	11.93	0.0	0.573	6.120	76.7	2.2875	1.0	273.0	
503	0.06078	0.0	11.93	0.0	0.573	6.976	91.0	2.1675	1.0	273.0	
504	0.10959	0.0	11.93	0.0	0.573	6.794	89.3	2.3889	1.0	273.0	
505	0.04741	0.0	11.93	0.0	0.573	6.030	80.8	2.5050	1.0	273.0	
...	
	PTRATIO	B	LSTAT								
0	15.3	396.90	4.98								
1	17.8	396.90	9.14								
2	17.8	392.83	4.03								
3	18.7	394.63	2.94								
4	18.7	396.90	5.33								
...								
501	21.0	391.99	9.67								
502	21.0	396.90	9.08								
503	21.0	396.90	5.64								
504	21.0	395.45	6.48								
505	21.0	396.90	7.88								
...								
	[506 rows x 13 columns]										
0	24.0										
1	21.6										
2	34.7										
3	33.4										
4	36.2										
...	...										
501	22.4										
502	28.6										
503	23.9										
504	22.0										
505	11.9										
	Name: price, Length: 506, dtype: float64										

Splitting the data into Training data and Test data

```
In [15]: X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.2, random_state = 2)
```

```
In [16]: print(X.shape, X_train.shape, X_test.shape)
(506, 13) (404, 13) (102, 13)
```

Model Training

XGBoost Regressor

```
In [17]: # Loading the model
model = XGBRegressor()
```

```
In [18]: # training the model with X_train
model.fit(X_train, Y_train)
```

```
Out[18]: XGBRegressor
XGBRegressor(base_score
```