

In [2]:

```
# Importing the Dependencies

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
```

In [ ]:

```
# Importing the Boston House Price Dataset

house_price_dataset = sklearn.datasets.load_boston()
```

In [7]:

```
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, 1
6.5, 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, 18.4, 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, 16. , 22.2, 25. , 33. , 23.5,
 19.4, 22. , 17.4, 20.9, 24.2, 21.7, 22.8, 23.4, 24.1, 21.4, 20. ,
 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, 20.4, 19.8, 19.4,
 21.7, 22.8, 18.8, 18.7, 18.5, 18.3, 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, 29. , 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. , 33.2, 33.1, 29.1, 35.1, 45.4, 35.4, 46. , 50. , 32.2, 22. ,
 20.1, 23.2, 22.3, 24.8, 28.5, 37.3, 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, 21.1, 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. , 13.8,
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, 19.9, 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, 18.3, 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='<U7'), 'DESCR': ".. _boston_dataset:\n\nBo
ston house prices dataset\n-----\n\n**Data Set Characteristics:**
\n\n :Number of Instances: 506 \n\n :Number of Attributes: 13 numeric/categorical p
redictive. Median Value (attribute 14) is usually the target.\n\n :Attribute Informati
on (in order):\n - CRIM per capita crime rate by town\n - ZN pro
portion of residential land zoned for lots over 25,000 sq.ft.\n - INDUS proport
ion of non-retail business acres per town\n - CHAS Charles River dummy variabl
e (= 1 if tract bounds river; 0 otherwise)\n - NOX nitric oxides concentratio
n (parts per 10 million)\n - RM average number of rooms per dwelling\n
- AGE proportion of owner-occupied units built prior to 1940\n - DIS wei
ghted distances to five Boston employment centres\n - RAD index of accessibil
ity to radial highways\n - TAX full-value property-tax rate per $10,000\n
- PTRATIO pupil-teacher ratio by town\n - B 1000(Bk - 0.63)^2 where Bk is
the proportion of black people by town\n - LSTAT % lower status of the populati
on\n - MEDV Median value of owner-occupied homes in $1000's\n\n :Missing At
tribute Values: None\n\n :Creator: Harrison, D. and Rubinfeld, D.L.\n\nThis is a copy
of UCI ML housing dataset.\nhttps://archive.ics.uci.edu/ml/machine-learning-databases/hou
sing/\n\nThis dataset was taken from the StatLib library which is maintained at Carnegi
e Mellon University.\n\nThe Boston house-price data of Harrison, D. and Rubinfeld, D.L. '
Hedonic\nprices and the demand for clean air', J. Environ. Economics & Management,\nvol.5
, 81-102, 1978. Used in Belsley, Kuh & Welsch, 'Regression diagnostics\n...', Wiley, 19
80. N.B. Various transformations are used in the table on\npages 244-261 of the latter.
\n\nThe Boston house-price data has been used in many machine learning papers that addres
s regression\nproblems. \n\n .. topic:: References\n\n - Belsley, Kuh & Welsch, '
Regression diagnostics: 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, Univ
ersity of Massachusetts, Amherst. Morgan Kaufmann.\n", 'filename': 'boston_house_prices.c
sv', 'data_module': 'sklearn.datasets.data'})

```

In [9]:

```

# Loading the dataset to a Pandas DataFrame
house_price_dataframe = pd.DataFrame(house_price_dataset.data, columns = house_price_dat
aset.feature_names)

```

In [12]:

```

# Print First 5 rows of our DataFrame
house_price_dataframe.head()

```

Out[12]:

	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 [14]:

```

# add the target (price) column to the DataFrame

```

```
house_price_dataframe['price'] = house_price_dataset.target
```

In [15]:

```
house_price_dataframe.head()
```

Out[15]:

	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 [16]:

```
# checking the number of rows and Columns in the data frame
house_price_dataframe.shape
```

Out[16]:

(506, 14)

In [17]:

```
# check for missing values
house_price_dataframe.isnull().sum()
```

Out[17]:

```
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 [18]:

```
# statistical measures of the dataset
house_price_dataframe.describe()
```

Out[18]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	3.613524	11.363636	11.136779	0.069170	0.554695	6.284634	68.574901	3.795043	9.549407	408.2371
std	8.601545	23.322453	6.860353	0.253994	0.115878	0.702617	28.148861	2.105710	8.707259	168.5371
min	0.006320	0.000000	0.460000	0.000000	0.385000	3.561000	2.900000	1.129600	1.000000	187.0000
25%	0.082045	0.000000	5.190000	0.000000	0.449000	5.885500	45.025000	2.100175	4.000000	279.0000
50%	0.256510	0.000000	9.690000	0.000000	0.538000	6.208500	77.500000	3.207450	5.000000	330.0000
75%	3.677083	12.500000	18.100000	0.000000	0.624000	6.623500	94.075000	5.188425	24.000000	666.0000
max	88.976200	100.000000	27.740000	1.000000	0.871000	8.780000	100.000000	12.126500	24.000000	711.0000

## Understanding the correlation between various features in the dataset

### 1.Positive Correlation

### 2.Negative Correlation

In [20]:

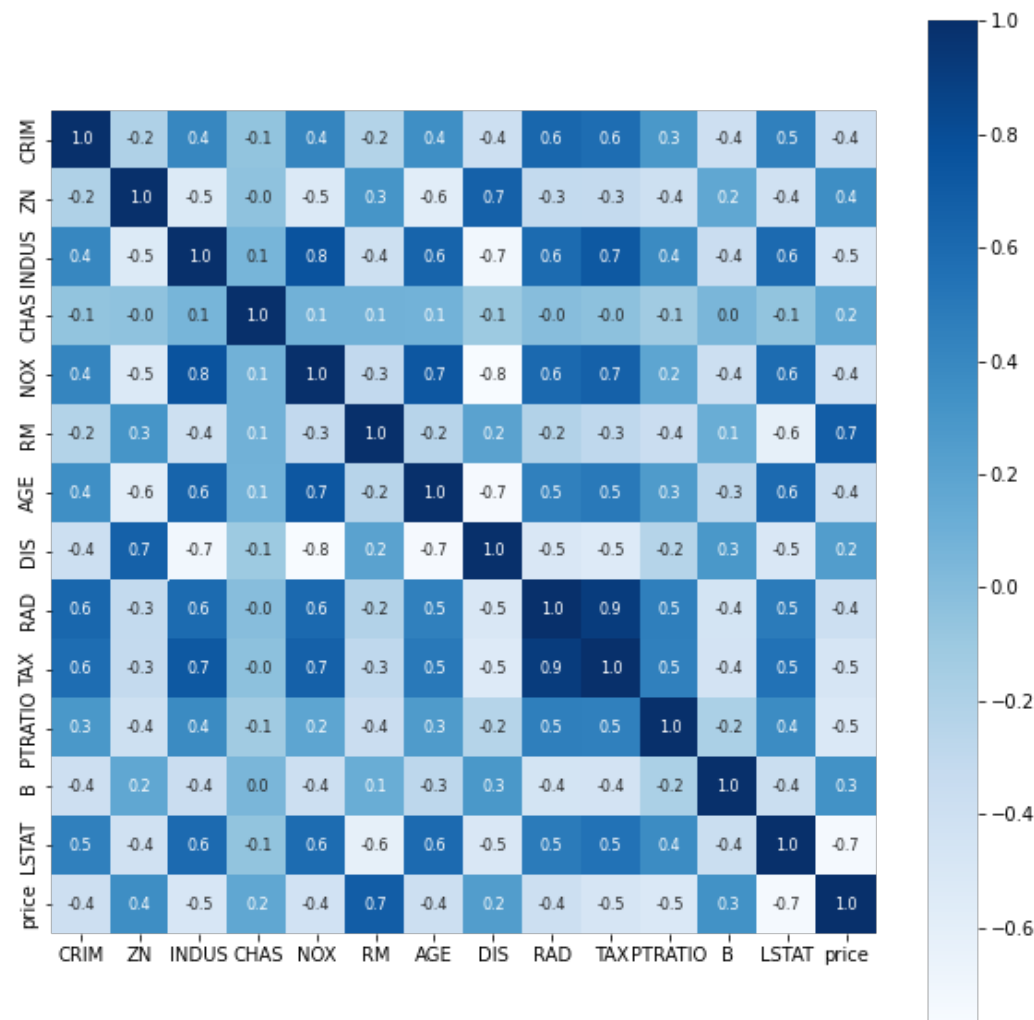
```
correlation = house_price_dataframe.corr()
```

In [21]:

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

Out[21]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f169ed17dd0>



In [22]:

```
# Splitting the data and Target
X = house_price_dataframe.drop(['price'], axis=1)
Y = house_price_dataframe['price']
```

In [23]:

```
print(X)
print(Y)
```

```
   CRIM  ZN  INDUS  CHAS  NOX  ...  RAD  TAX  PTRATIO  B  LSTAT
0  0.00632  18.0   2.31   0.0  0.538  ...   1.0  296.0   15.3  396.90  4.98
1  0.02731   0.0   7.07   0.0  0.469  ...   2.0  242.0   17.8  396.90  9.14
```

2	0.02729	0.0	7.07	0.0	0.469	...	2.0	242.0	17.8	392.83	4.03
3	0.03237	0.0	2.18	0.0	0.458	...	3.0	222.0	18.7	394.63	2.94
4	0.06905	0.0	2.18	0.0	0.458	...	3.0	222.0	18.7	396.90	5.33
...	...	...	...	...	...	...	...	...	...	...	...
501	0.06263	0.0	11.93	0.0	0.573	...	1.0	273.0	21.0	391.99	9.67
502	0.04527	0.0	11.93	0.0	0.573	...	1.0	273.0	21.0	396.90	9.08
503	0.06076	0.0	11.93	0.0	0.573	...	1.0	273.0	21.0	396.90	5.64
504	0.10959	0.0	11.93	0.0	0.573	...	1.0	273.0	21.0	393.45	6.48
505	0.04741	0.0	11.93	0.0	0.573	...	1.0	273.0	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     20.6
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 [24]:

```
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.2, random_state = 2)
```

In [25]:

```
print(X.shape, X_train.shape, X_test.shape)
```

```
(506, 13) (404, 13) (102, 13)
```

In [26]:

```
# Model Training
# XGBoost Regressor

# loading the model
model = XGBRegressor()
```

In [ ]:

```
# training the model with X_train
model.fit(X_train, Y_train)
```

In [28]:

```
# Evaluation
# Prediction on training data

# accuracy for prediction on training data
training_data_prediction = model.predict(X_train)
```

In [29]:

```
print(training_data_prediction)
```

[23.360205	22.462858	20.84804	33.77895	15.333282	13.616525
21.71274	15.175322	11.724756	21.836252	16.08508	7.52517
31.094206	48.56228	32.623158	20.546066	22.177324	20.500404
31.666502	20.551508	25.74269	8.247894	45.200817	22.069397
20.698004	20.100042	19.873472	26.242834	23.39618	31.927258
21.493471	9.280926	18.504272	21.87202	12.504413	10.578829
13.054951	23.541336	19.164755	15.888303	23.768887	28.454714
15.539753	18.049202	16.23671	14.08383	25.33273	17.575668
19.566167	16.990675	21.739977	22.935113	16.125738	22.15393

10.308407	10.308407	21.750777	32.353143	10.125750	22.453593
20.776966	20.042227	22.898897	38.124043	30.607079	32.607468
20.919416	47.348038	14.524615	8.126455	19.581661	9.030508
26.462107	17.69918	20.546162	46.312218	39.689137	34.387108
22.11083	34.568977	24.873934	50.078335	14.5669775	20.525211
20.62971	23.202105	49.514477	23.12061	24.795782	20.319666
43.869396	17.110266	32.165016	34.75202	7.313497	20.309446
18.038298	12.008462	24.216425	47.90671	37.94349	20.759708
40.182804	18.249052	15.611586	26.39461	21.0571	20.421682
18.377089	17.338768	21.223648	22.653662	17.560051	32.635715
16.683764	13.004857	18.488163	20.659714	16.501846	20.648884
48.62411	15.977999	15.97522	18.581459	14.893438	32.871964
14.236945	43.612328	33.881115	19.073408	15.747335	9.4903965
10.153891	14.812717	18.655546	8.596755	22.666656	10.941623
20.534616	49.324417	22.710459	19.99658	31.663935	21.78586
30.9277	30.507492	15.054665	15.854853	48.532074	21.108742
15.687305	12.403721	49.90245	31.557863	11.709707	20.22495
26.214525	32.90807	22.90362	9.542897	24.487959	24.46598
22.509142	14.704502	27.895067	33.619015	14.888735	19.147383
26.40218	32.77208	29.293688	23.638102	10.448805	22.518728
21.47825	35.32415	23.002241	20.470022	18.918747	10.328174
22.244467	17.69918	20.918488	11.913417	42.572548	46.803394
14.652036	20.633188	23.285368	15.295161	20.861048	23.587011
32.94382	21.090906	24.898489	18.465925	31.454802	14.421506
15.421497	21.890705	23.64799	17.40471	26.111868	24.977922
27.56308	22.964123	18.823803	28.856464	14.080684	19.785515
17.007908	42.90537	26.354216	21.719929	23.784258	18.4141
17.923422	20.337881	22.936398	25.297531	17.572325	14.486319
20.739832	21.733093	11.1917715	18.290442	20.70475	20.929468
18.990923	8.7798395	21.141748	21.021317	15.49217	24.455221
31.499088	22.668139	14.862843	19.69585	24.746317	22.913176
48.144817	19.950285	30.148172	49.98047	16.743952	16.218952
9.891141	20.452726	17.06055	14.73646	17.539606	19.555712
30.26191	27.037518	18.43813	20.100842	24.147627	10.21256
25.064299	48.283043	20.977459	23.265625	20.141813	11.87677
17.84212	15.1286955	14.9789295	23.502743	16.092314	21.276255
26.55347	16.940031	23.485325	14.927286	20.90435	19.254526
24.397417	27.566774	23.607512	17.905067	22.675825	25.12203
15.141896	18.460642	23.440636	16.4928	23.372946	30.389936
15.330368	24.69199	17.316717	14.531138	10.496169	24.805672
15.659789	38.916733	20.403166	42.113743	8.544421	22.536352
15.654481	15.709977	17.263374	23.888586	21.690222	46.16276
15.304819	31.137545	25.326769	18.969254	26.29209	11.722559
40.65201	20.52522	17.135836	24.829275	15.565665	23.360205
8.280649	24.018639	19.57025	20.865868	23.611485	22.455328
17.646477	17.687094	14.59732	25.61237	13.333718	22.577513
20.657572	14.8804865	16.539358	23.276703	24.873934	22.52675
23.107155	31.871576	19.262531	19.536154	28.251024	23.817226
12.874959	22.59372	12.234834	10.024989	20.419611	10.369816
45.84478	24.873934	12.357825	16.367088	14.355771	28.338346
18.669233	20.334248	10.546778	21.30952	21.00914	20.669264
23.91886	25.009733	26.945326	13.288843	18.277857	20.95568
18.233625	23.807056	13.400126	23.875198	33.050533	27.785492
25.296518	19.071947	20.950756	11.507434	22.855497	15.573306
22.33747	20.807749	22.41908	17.212593	12.645366	35.121113
18.852188	48.823723	22.462465	24.267456	21.375692	19.38756
8.561088	20.726429	23.400837	21.41578	17.63176	25.232733
21.164701	26.444288	14.49171	49.559753	30.693232	23.20531
22.950115	16.84211	30.982431	16.259336	23.613512	20.93225
20.178421	22.782583	]			

In [30]:

```
# R squared error
score_1 = metrics.r2_score(Y_train, training_data_prediction)

# Mean Absolute Error
score_2 = metrics.mean_absolute_error(Y_train, training_data_prediction)

print("R squared error : ", score_1)
print('Mean Absolute Error : ', score_2)
```

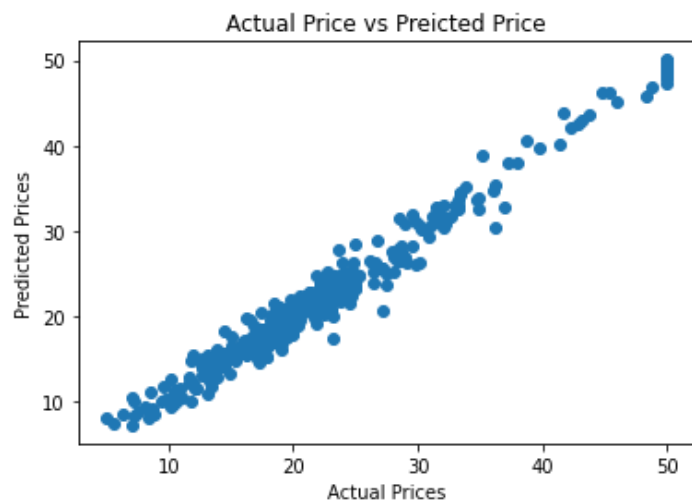
R squared error : 0.9733349094832763

R squared error : 0.975543094032703  
Mean Absolute Error : 1.145314053261634

In [31]:

```
# Visualizing the actual Prices and predicted prices

plt.scatter(Y_train, training_data_prediction)
plt.xlabel("Actual Prices")
plt.ylabel("Predicted Prices")
plt.title("Actual Price vs Preicted Price")
plt.show()
```



## Prediction on Test Data

In [32]:

```
# accuracy for prediction on test data
test_data_prediction = model.predict(X_test)
```

In [34]:

```
# R squared error
score_1 = metrics.r2_score(Y_test, test_data_prediction)
# Mean Absolute Error
score_2 = metrics.mean_absolute_error(Y_test, test_data_prediction)

print("R squared error : ", score_1)
print('Mean Absolute Error : ', score_2)
```

R squared error : 0.9115937697657654  
Mean Absolute Error : 1.9922956859364223

In [ ]: