

```
In [55]: import numpy as np
import pandas as pd
import seaborn as sns
import sklearn.datasets
from sklearn.model_selection import train_test_split
from xgboost import XGBRegressor
from sklearn import metrics
import matplotlib.pyplot as plt

In [2]: #Loading the dataset from sklearn
house_price_dataset = sklearn.datasets.fetch_california_housing()
```

```
In [7]: print(house_price_dataset)

{'data': array([[ 8.3252, ..., 41., ..., 6.98412698, ..., 2.55555556,
[ 37.88, ..., -122.23 ], ...,
[ 8.3014, ..., 21., ..., 6.23813708, ..., 2.10984183,
[ 37.86, ..., -122.22 ], ...,
[ 7.2574, ..., 52., ..., 8.28813559, ..., 2.80225989,
[ 37.85, ..., -122.24 ], ...,
...,
[ 1.7, ..., 17., ..., 5.20554273, ..., 2.3256351,
[ 39.43, ..., -121.22 ], ...,
[ 1.8672, ..., 18., ..., 5.32951289, ..., 2.12320917,
[ 39.43, ..., -121.32 ], ...,
[ 2.3886, ..., 16., ..., 5.25471698, ..., 2.61698113,
[ 39.37, ..., -121.24 ]]), 'target': array([4.526, 3.585, 3.521, ..., 0.923, 0.847, 0.894]), 'frame': None, 'target_names': ['MedHouseVal'], 'feature_names':
['MedInc', 'HouseAge', 'AveRooms', 'AveBedrms', 'Population', 'AveOccup', 'Latitude', 'Longitude'], 'DESCR': '.._california_housing_dataset:\n\nCalifornia Housing dataset\n-----\n\nData Set Characteristics:**\n\n :Number of Instances: 20640\n\n :Number of Attributes: 8 numeric, predictive attributes and the target\n\n :Attribute Information:\n - MedInc median income in block group\n - HouseAge median house age in block group\n - AveRooms average number of rooms per household\n - AveBedrms average number of bedrooms per household\n - Population block group population\n - AveOccup average number of household members\n - Latitude block group latitude\n - Longitude block group longitude\n\n :Missing Attribute Values: None\n\nThis dataset was obtained from the Statlib repository.\nhttps://www.dcc.fc.up.pt/~ltorgo/Regression/cal_housing.html\n\nThe target variable is the median house value for California districts, expressed in hundreds of thousands of dollars ($100,000).\n\nThis dataset was derived from the 1990 U.S. census, using one row per census block group. A block group is the smallest geographical unit for which the U.S. Census Bureau publishes sample data (a block group typically has a population of 600 to 3,000 people).\n\nA household is a group of people residing within a home. Since the average number of rooms and bedrooms in this dataset are provided per household, these columns may take surprising large values for block groups with few households and many empty houses, such as vacation resorts.\n\nIt can be downloaded/loaded using the\nfunc: 'sklearn.datasets.fetch_california_housing' function.\n\n.. topic:: References\n\n - Pace, R. Kelley and Ronald Barry, Sparse Spatial Autoregressions,\nStatistics and Probability Letters, 33 (1997) 291-297\n'}
```

```
In [16]: #converting it to dataframe
house_price_dataframe=pd.DataFrame(house_price_dataset.data,columns=house_price_dataset.feature_names)
```

```
In [17]: house_price_dataframe.head()
```

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude
0	8.3252	41.0	6.984127	1.023810	322.0	2.555556	37.88	-122.23
1	8.3014	21.0	6.238137	0.971880	2401.0	2.109842	37.86	-122.22
2	7.2574	52.0	8.288136	1.073446	496.0	2.802260	37.85	-122.24
3	5.6431	52.0	5.817352	1.073059	558.0	2.547945	37.85	-122.25
4	3.8462	52.0	6.281853	1.081081	565.0	2.181467	37.85	-122.25

```
In [18]: #adding the target (price column to DF)
house_price_dataframe['Price']=house_price_dataset.target
```

```
In [19]: house_price_dataframe.head()
```

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude	Price
0	8.3252	41.0	6.984127	1.023810	322.0	2.555556	37.88	-122.23	4.526
1	8.3014	21.0	6.238137	0.971880	2401.0	2.109842	37.86	-122.22	3.585
2	7.2574	52.0	8.288136	1.073446	496.0	2.802260	37.85	-122.24	3.521
3	5.6431	52.0	5.817352	1.073059	558.0	2.547945	37.85	-122.25	3.413
4	3.8462	52.0	6.281853	1.081081	565.0	2.181467	37.85	-122.25	3.422

```
In [21]: #No of row and columns
house_price_dataframe.shape
```

```
Out[21]: (20640, 9)
```

```
In [22]: #checking missing value
house_price_dataframe.isnull().sum()
```

```
Out[22]: MedInc      0
HouseAge    0
AveRooms    0
AveBedrms   0
Population  0
AveOccup    0
Latitude     0
Longitude    0
Price       0
dtype: int64
```

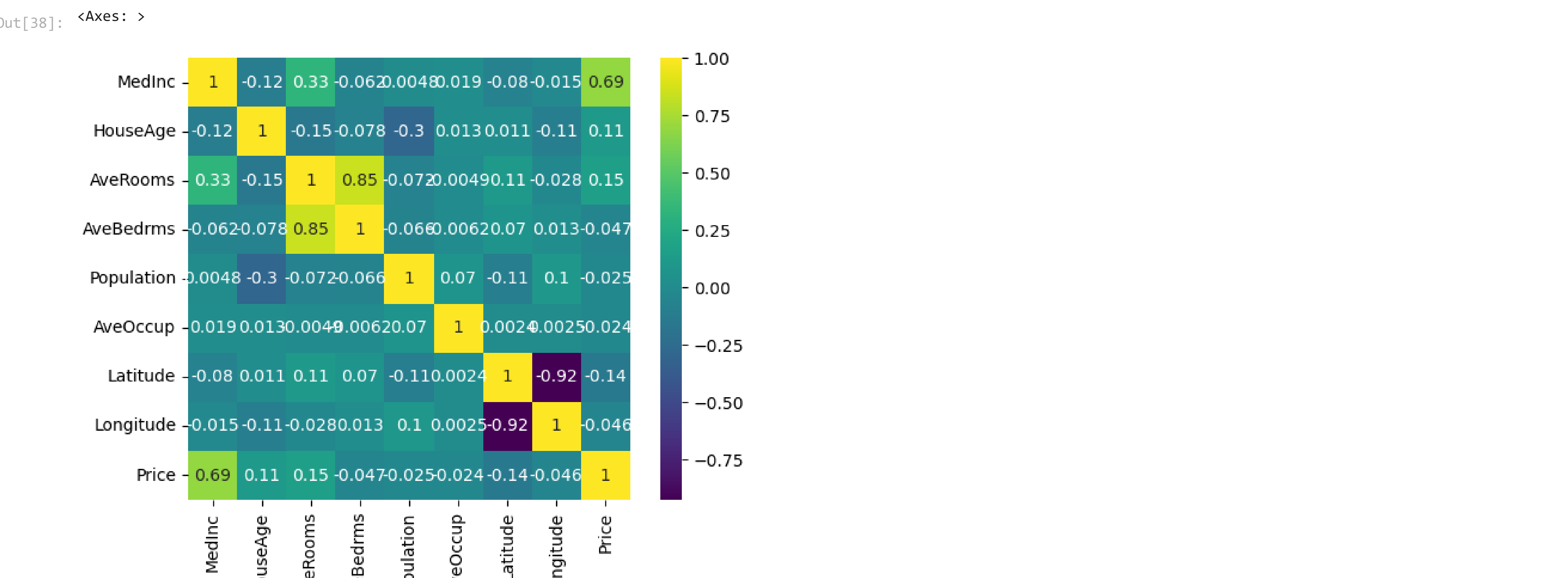
```
In [24]: #checking Statistics of DF
house_price_dataframe.describe()
```

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude	Price
count	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000
mean	3.870671	28.639486	5.429000	1.096675	1425.476744	3.070655	35.631861	-119.569704	2.068558
std	1.899822	12.585558	2.474173	0.473911	1132.462122	10.386050	2.135952	2.003532	1.153956
min	0.499900	1.000000	0.846154	0.333333	3.000000	0.692308	32.540000	-124.350000	0.149990
25%	2.563400	18.000000	4.440716	1.006079	787.000000	2.429741	33.930000	-121.800000	1.196000
50%	3.534800	29.000000	5.229129	1.048780	1166.000000	2.818116	34.260000	-118.490000	1.797000
75%	4.743250	37.000000	6.052381	1.099526	1725.000000	3.282261	37.710000	-118.010000	2.647250
max	15.000100	52.000000	141.909091	34.066667	35682.000000	1243.333333	41.950000	-114.310000	5.000010

```
In [28]: #finding correlation btw various feature in data
house_price_dataframe.corr()
```

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude	Price
MedInc	1.000000	-0.119034	0.326895	-0.062040	0.004834	0.018766	-0.079809	-0.015176	0.688075
HouseAge	-0.119034	1.000000	-0.153277	-0.077747	-0.296244	0.013191	0.011173	-0.108197	0.105623
AveRooms	0.326895	-0.153277	1.000000	0.847621	-0.072213	-0.004852	0.106389	-0.027540	0.151948
AveBedrms	-0.062040	-0.077747	0.847621	1.000000	-0.066197	-0.006181	0.069721	0.013344	-0.046701
Population	0.004834	-0.296244	-0.072213	-0.066197	1.000000	0.069863	-0.108785	0.099773	-0.024650
AveOccup	0.018766	0.013191	-0.004852	-0.006181	0.069863	1.000000	0.002366	0.002476	-0.023737
Latitude	-0.079809	0.011173	0.106389	0.069721	-0.108785	0.002366	1.000000	-0.924664	-0.144160
Longitude	-0.015176	-0.108197	-0.027540	0.013344	0.099773	0.002476	-0.924664	1.000000	-0.045967
Price	0.688075	0.105623	0.151948	-0.046701	-0.024650	-0.023737	-0.144160	-0.045967	1.000000

```
In [38]: sns.heatmap(house_price_dataframe.corr(),cmap="viridis",cbar=True,square=True,annot=True)
```



```
In [39]: #Split the dataset in data & Label/target(price)
X=house_price_dataframe.drop(['Price'],axis=1)
Y=house_price_dataframe['Price']
```

```
In [40]: print(X)
print(Y)
```

```
0      MedInc  HouseAge  AveRooms  AveBedrms  Population  AveOccup  Latitude  Longitude
0      8.3252     41.0    6.984127    1.023810     322.0    2.555556     37.88     -122.23
1      8.3014     21.0    6.238137    0.971880    2401.0    2.109842     37.86     -122.22
2      7.2574     52.0    8.288136    1.073446     496.0    2.802260     37.85     -122.24
3      5.6431     52.0    5.817352    1.073059     558.0    2.547945     37.85     -122.25
4      3.8462     52.0    6.281853    1.081081     565.0    2.181467     37.85     -122.25
...      ...      ...      ...      ...      ...      ...      ...      ...
20635   1.5603     25.0    5.045455    1.133333     845.0    2.560606     39.48     -121.09
20636   2.5568     18.0    6.114835    1.315789     356.0    3.122807     39.49     -121.21
20637   1.7000     17.0    5.205543    1.120092    1007.0    2.325635     39.43     -121.22
20638   1.8672     18.0    5.329513    1.171920     741.0    2.123209     39.43     -121.32
20639   2.3886     16.0    5.254717    1.162264    1387.0    2.616981     39.37     -121.24

[20640 rows x 8 columns]
0      4.526
1      3.585
2      3.521
3      3.413
4      3.422
...      ...
20635   0.781
20636   0.771
20637   0.923
20638   0.847
20639   0.894
Name: Price, Length: 20640, dtype: float64
```

```
In [41]: #Split Into Training And Test Data
X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.2,random_state=2)
```

```
In [44]: print(X.shape,X_test.shape,X_train.shape)

(20640, 8) (4128, 8) (16512, 8)
```

```
In [45]: #MODEL TRAINING_XGBOOST REGRESSOR( Decision tree model)
model=XGBRegressor()
```

```
In [46]: model.fit(X_train,Y_train)
```

```
Out[46]: XGBRegressor(base_score=None, booster=None, callbacks=None,
colsample_bylevel=None, colsample_bynode=None,
colsample_bytree=None, early_stopping_rounds=None,
enable_categorical=False, eval_metric=None, feature_types=None,
gamma=None, gpu_id=None, grow_policy=None, importance_type=None,
interaction_constraints=None, learning_rate=None, max_bin=None,
max_cat_threshold=None, max_cat_to_onehot=None,
max_delta_step=None, max_depth=None, max_leaves=None,
min_child_weight=None, missing=nan, monotone_constraints=None,
n_estimators=100, n_jobs=None, num_parallel_tree=None,
predictor=None, random_state=None, ...)
```

```
In [47]: #EVALUATION & PREDICTION
training_data_pred=model.predict(X_train)
```

```
In [48]: print(training_data_pred)

[0.6893792  2.986824   0.48874274 ... 1.8632544  1.7800125  0.7565893 ]
```

```
In [50]: #R square & mean absolute error metric for regression
score_1=metrics.r2_score(Y_train,training_data_pred)
#Mena absolute error
score_2=metrics.mean_absolute_error(Y_train,training_data_pred)

print('R square error:',score_1)

print('Mean absolute error:',score_2)

R square error: 0.9451221492760822
Mean absolute error: 0.1919170860794262
```

```
In [ ]: #Less size data thatswhy used XGBR
```

```
In [51]: #Prediction on Test Data

test_data_pred=model.predict(X_test)
```

```
In [52]: print(test_data_pred)

[2.787383  1.9628428 0.782536 ... 1.5060123 0.8763797 1.9317917]
```

```
In [53]: #R square & mean absolute error metric for regression
score_3=metrics.r2_score(Y_test,test_data_pred)
#Mena absolute error
score_4=metrics.mean_absolute_error(Y_test,test_data_pred)

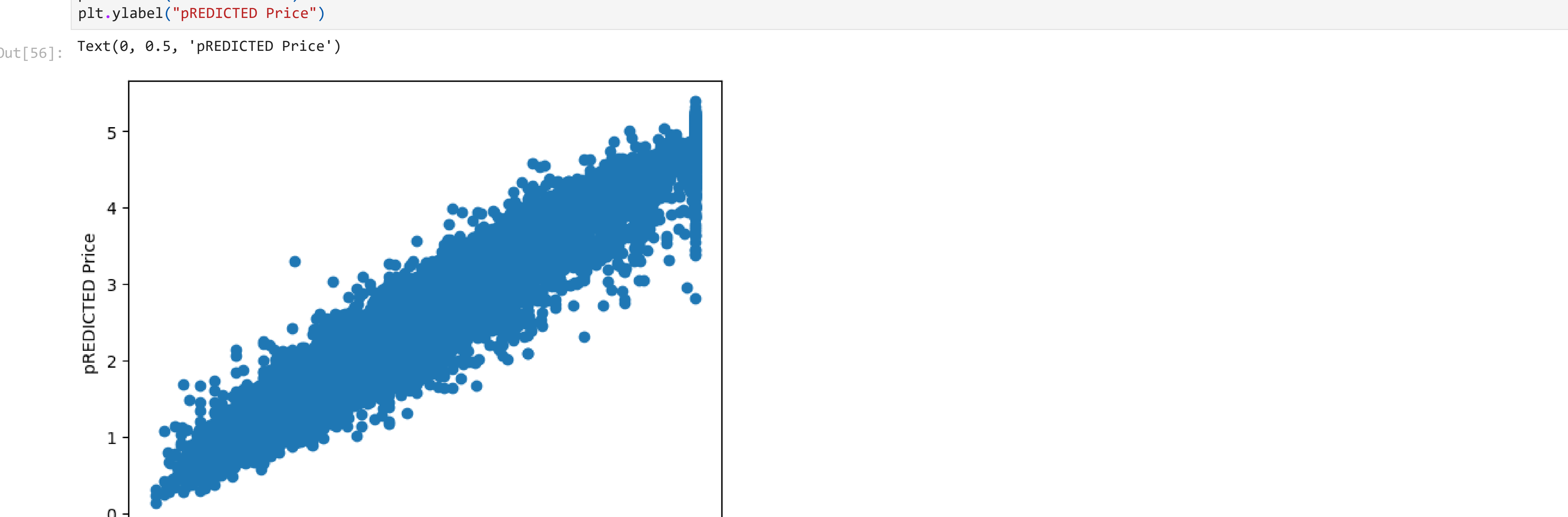
print('R square error:',score_3)

print('Mean absolute error:',score_4)

R square error: 0.8412904408180302
Mean absolute error: 0.30753655785801337
```

```
In [56]: #Visualising the actual and predicted prices
plt.scatter(Y_train,training_data_pred)
plt.xlabel("actual Price")
plt.ylabel("pPREDICTED Price")
```

```
Out[56]: Text(0, 0.5, 'pPREDICTED Price')
```



```
In [ ]: #value predicted are close to acutal prices hence our model works fine.
```