

In []:

```

1  #####
2  #
3      Program : To predict House Price based on Random Forest Model
4  #
5  #####
6
7  # Step0 : Loading of Packages and Loading the Data Set
8
9  import numpy as np
10 import pandas as pd
11 import matplotlib.pyplot as plt
12 import seaborn as sns
13 from sklearn.model_selection import train_test_split, cross_val_score
14 from sklearn.preprocessing import StandardScaler
15 from sklearn.ensemble import RandomForestClassifier
16 from sklearn import datasets
17 from sklearn.tree import DecisionTreeClassifier
18 from sklearn.metrics import roc_auc_score, accuracy_score
19
20 # Loading the standard Boston House Data
21 boston = datasets.load_boston()
22
23 # Columns /Features in the Boston Data
24 print(boston.feature_names)
25 print('-'*80)
26
27 # understanding the meaning of the various features/Attributes of Boston Data
28 print(boston.DESCR)
29 print('-'*80)
30 # Creating a DataFrame of the Boston Data
31
32 HouseData = pd.DataFrame(boston.data, columns=boston.feature_names)
33 print(" The Boston Data for Predicting the House price \n", '-'*80)
34 print(HouseData.head())
35 print('-'*80)
36
37 # Housing price
38 targets = boston.target
39

```

```
In [ ]: 1 # Step1 : Analysing the Data
2
3 # including the target column in the main dataframe ,to analyse the data together
4 HouseData['target'] = targets
5
6 print(" Analysing The Data , by using info, describe functions , checking for nulls and size")
7 print(HouseData.info())
8 print('-'*80)
9
10 print(HouseData.shape)
11 print('-'*80)
12
13 print(HouseData.describe())
14 print('-'*80)
15
16 print(HouseData.isnull().sum())
17 print('-'*80)
18
19 print(" Observing the Data , We see no null values , and 506 entries and 14 features \n")
20
```

```
In [ ]: 1 # Step2 : Visualising the Data
2
3 # visualise the data
4
5 features = boston.feature_names
6 y = 'target'
7
8 print(" Visualising the Data to see the relationship with House Price")
9 column_cnt = len(features)
10 i = 0
11 while i < column_cnt:
12     sns.pairplot(HouseData,x_vars=features[i:i+4],y_vars=y,kind="reg")
13     i = i+4
14
15
```

```
In [ ]: 1 # Step3 : Preprocessing the Data
2 # since the target is continuous Data , we will make it discrete by using the median value
3 # converting the target field to binary
4
5 # the median value of the housing price
6 med = np.median(targets)
7
8 HouseData['target_num'] = HouseData.target.map(lambda x: 1 if x>med else 0)
9 print(HouseData.head(2))
10
11
```

```
In [ ]: 1 # Step4 : Creating the Train and Test sets by splitting the Data
2
3 y = HouseData['target_num']
4
5 print(y.shape)
6 y = y.values.reshape(-1,1)
7
8 X = HouseData.drop(['target', 'target_num'],axis=1)
9 features = ['RM', 'AGE', 'DIS', 'B', 'CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RAD', 'TAX', 'PTRATIO', 'LSTAT']
10
11 print(HouseData.head())
12
13 # training set size is 70% and test size is 30%
14 X_train,X_test,y_train,y_test = train_test_split(X[features],y,test_size=0.3, random_state=2)
15
16 print("\n The size after splitting of the Data into Training and Test Set ",'-'*80)
17 print(X_train.shape)
18 print(X_test.shape)
19 print(y_train.shape)
20 print(y_test.shape)
21
22
```

```
In [ ]: 1 # Step5 : Model Selection
2
3 #implemneting the random forest manually , by 3 decision trees and then combining their predictions
4 #
5
6 clf = DecisionTreeClassifier(random_state=1, min_samples_leaf=3)
7 clf.fit(X_train, y_train)
8
9 clf2 = DecisionTreeClassifier(random_state=1, max_depth=4)
10 clf2.fit(X_train, y_train)
11
12 clf3 = DecisionTreeClassifier(random_state=1, max_features=5)
13 clf3.fit(X_train, y_train)
14
15 print(" The accuracy of each Decision Tree for Random Forest build manually is :\n", '-'*80)
16
17 predictions1 = clf.predict(X_test)
18 print(round(roc_auc_score(y_test, predictions1),3))
19
20 predictions2 = clf2.predict(X_test)
21 print(round(roc_auc_score(y_test, predictions2),3))
22
23 predictions3 = clf3.predict(X_test)
24 print(round(roc_auc_score(y_test, predictions3),3))
25
26 # now finding the result when both the trees are combined
27
28 predictions = clf.predict_proba(X_test)[: ,1]
29 predictions2 = clf2.predict_proba(X_test)[: ,1]
30 predictions3 = clf3.predict_proba(X_test)[: ,1]
31 combined = (predictions + predictions2 + predictions3) / 3
32 rounded = np.round(combined,3)
33
34 print(" The accuracy of Manually Build Random Forest is :\n", '-'*80)
35 print(round(roc_auc_score(y_test, rounded),3))
36
37
```

```
In [ ]: 1 # Step6 : Evaluation of the Model with different Random Forest Models
2
3 rfc1 = RandomForestClassifier(max_features=13, random_state=1)
4 rfc1.fit(X_train, y_train)
5 pred1 = rfc1.predict(X_test)
6
7 #y_test = y_test.ravel()
8 print(y_test.shape)
9 print(pred1.shape)
10
11
12 print(" The accuracy of Random Forest with max_features=13 :\n", '-'*80)
13 print(round(roc_auc_score(y_test, pred1),3))
14
15 # play around with the setting for max_features
16 rfc2 = RandomForestClassifier(max_features=5, random_state=1)
17 rfc2.fit(X_train, y_train)
18 pred2 = rfc2.predict(X_test)
19
20 print(" The accuracy of Random Forest with max_features=5 :\n", '-'*80)
21 print(round(roc_auc_score(y_test, pred2),3))
22
23 # play around with the setting for max_features
24 print(" The accuracy of Random Forest with max_features=auto :\n", '-'*80)
25 rfc2 = RandomForestClassifier(max_features="auto", bootstrap=True, random_state=1)
26 rfc2.fit(X_train, y_train)
27 pred2 = rfc2.predict(X_test)
28 print(round(roc_auc_score(y_test, pred2),3))
```