

In [ ]:

```
1 *****
2 #
3 #           Assignment : Predicting Survival in Titanic Data Set Using Decision Tree
4 #
5 *****
6 # Problem statement : Predict the Passenger survived /not , for columns
7 # Pclass, Sex, Age, SibSp (Siblings aboard), Parch (Parents/children aboard), and Fare
8
9 # Step 0 : importing the libraries
10 import numpy as np
11 import pandas as pd
12 import seaborn as sb
13 import matplotlib.pyplot as plt
14 import sklearn
15 from pandas import Series, DataFrame
16 from pylab import rcParams
17 from sklearn import tree, metrics, preprocessing
18 from sklearn.model_selection import KFold, cross_val_score
19 from sklearn.metrics import classification_report
20
21 Url= "https://raw.githubusercontent.com/BigDataGal/Python-for-DataScience/master/titanic-train.csv"
22
23 # Step 1 : Loading Titanic Data
24 titanic = pd.read_csv("titanic-train.csv")
25
26 print(" The Titanic Data for Passengers to predict Survival \n", '-'*80)
27 titanic.head()
28
```

```
In [ ]: 1 #Step 2: Analysing the Data
        2
        3 print(" Analysing the Data by seeing functions like shape, info , describe, checking for null values \n",'-'*80)
        4 print(titanic.columns)
        5
        6
        7 # number of rows and features
        8 print('-'*80)
        9 print(titanic.shape)
       10 # there are 891 rows and 12 columns
       11
       12 print('-'*80)
       13 print(titanic.info())
       14
       15 print('-'*80)
       16 print(titanic.describe())
       17
       18 print('-'*80)
       19 print(titanic.isnull().sum())
       20
       21
```

```
In [ ]: 1 #Step 3: Preprocessing the Data
2
3 # Since Sex column is categorical , Factorize Sex Column to numeric value
4
5 print(titanic.Sex.unique())
6 titanic['Sex_num'], _ = pd.factorize(titanic['Sex'])
7
8 print("\n Data Set after changing the Sex Field to numeric in Sex_num column ")
9 titanic.head()
10
11 # we observe that Age and Cabin fields have nulls , as the features to be focussed doesnt contain cabin , we can leave
12 # Filling the null value with mean age for passengers
13
14 titanic = titanic.fillna({'Age':titanic.Age.mean()})
15
16 print('-'*80)
17 print(titanic.isnull().sum())
18
19 print(titanic.Age.mean())
20
21
22
```

```
In [ ]: 1 #Step 4: Visualizing the Data using pairplot to understand the relationships
2
3 #print(titanic.head())
4
5 predictors = ['Pclass', 'Sex_num', 'Age', 'SibSp', 'Parch', 'Fare']
6 target = 'Survived'
7
8 sb.pairplot(titanic, x_vars=predictors, y_vars='Survived', kind="reg")
9
10
```

```
In [ ]: 1 #Step 5: Splitting the Data to train and Test Set with 30 % test data
2
3 X= titanic[predictors]
4 y = titanic[target]
5
6 print(" The number of rows and columns before splitting the Data \n",'-'*80)
7 print(X.shape)
8 print(y.shape)
9
10 y = y.values.reshape(-1,1)
11
12 from sklearn.model_selection import train_test_split
13
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
15
16 print(" The training set and test set rows and columns after splitting the Data \n",'-'*80)
17 print(X_train.shape)
18 print(y_train.shape)
19 print(X_test.shape)
20 print(y_test.shape)
21
22
```

```
In [ ]: 1 # Step 6 : Modeling the Data by using DEcision Tree Classifier and Evaluation by Sklearn metrics
2
3 # measure Entropy for the Tree
4 dtree = tree.DecisionTreeClassifier(criterion='entropy', max_depth=3, random_state=1)
5
6 # # train the tree with training Data Set
7 dtree.fit(X_train,y_train)
8
9 # # predicting the results for the Test data set
10 y_pred = dtree.predict(X_test)
11 #print(y_pred)
12
13 y_pred = y_pred.reshape(-1,1)
14
15 print(y_test.shape)
16 print(y_pred.shape)
17
18 misclassified_count = (y_test != y_pred).sum()
19 print('Misclassified samples: {} / {} and the percentage is {}'.format(misclassified_count,y_test.shape[0],round((mis
20
21
22 accuracy = round(metrics.accuracy_score(y_test, y_pred),2)*100
23
24 print('Accuracy of prediction For Passengers Survival is : {:.2f} percent'.format(accuracy))
25
```

```
In [ ]: 1 # Step 7 : Evaluation by using KFold crossvalidation
2
3 from sklearn.model_selection import KFold
4
5 print(X.shape)
6 print(y.shape)
7
8
9 y = pd.DataFrame(y)
10
11
12 dtree1 = tree.DecisionTreeClassifier(criterion='entropy', max_depth=3, random_state=0)
13
14
15 kf = KFold(n_splits=10, random_state=12)
16
17 print(kf)
18 print(kf.get_n_splits(X))
19
20 fold_accuracy = []
21
22
23 for train_fold, valid_fold in kf.split(X):
24     train = X.loc[train_fold] # Extract train data with cv indices
25     valid = X.loc[valid_fold] # Extract valid data with cv indices
26
27     train_y = y.loc[train_fold]
28     valid_y = y.loc[valid_fold]
29
30
31     train = train.dropna()
32     valid = valid.dropna()
33
34     train_y = train_y.dropna()
35     valid_y = valid_y.dropna()
36
37     model = dtree1.fit(train, train_y)
38     valid_acc = model.score(X = valid, y = valid_y)
39     fold_accuracy.append(valid_acc)
40
41 print("Accuracy per fold: ", fold_accuracy, "\n")
```

```
42  
43 acc = round(sum(fold_accuracy)/len(fold_accuracy),2)*100  
44 print("Average accuracy by using KFold is : ", acc)  
45  
46
```

```
In [ ]: 1 # Model Evaluation by using cross_val_score function  
2 from sklearn.model_selection import cross_val_score  
3  
4 print(X.shape)  
5 print(y.shape)  
6  
7 scores = cross_val_score(dtrees1, X, y, scoring = "accuracy", cv=15)  
8 print("Accuracy per fold: ")  
9 print(scores)  
10 print("Average accuracy using the cross_val_Score function : ", round(scores.mean(),2))
```

```
In [ ]: 1
```