

Classification

K-NN(K-Nearest Neighbors)

```
In [14]: 1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.model_selection import train_test_split
3 import pandas as pd
4 df=pd.read_csv("Datasets/Churn.csv")
5 df=df.dropna()
6 # df.info()
7 x=df.drop(["Unnamed: 0", "customerID", "Churn"],axis=1)
8 y=df['Churn']
9 x=pd.get_dummies(x)
10 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=2)
11 knn=KNeighborsClassifier(n_neighbors=5)
12 knn.fit(x_train,y_train)
13 ypred=knn.predict(x_test)
14 print(ypred)
```

```
['Yes' 'No' 'No' ... 'No' 'Yes' 'No']
```

Overfitting->

- it occurs when a model learns the training data too well including noise and irrelevant details resulting in poor performance on new data
- high accuracy on training data and poor accuracy on new and unseen data

Underfitting

- it happens when a model is too simple to capture the underlying patterns in the training data resulting in poor performance on both training as well as new data

KNN(K-Nearest Neighbors)

- based on its closeness or similarities in a given range(k) of neighbors the algorithm assigns the new data to a class or category in the dataset
- to calculate the closeness __ formula is used
- key point->always use odd number for k

Confusion Matrix (201%)

```
In [4]: 1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.metrics import confusion_matrix
3 from sklearn.model_selection import train_test_split
4 import pandas as pd
5 df=pd.read_csv("Datasets/Churn.csv")
6 df=df.dropna()
7 x=df.drop(["Unnamed: 0", "customerID", "Churn"],axis=1)
8 y=df['Churn']
9 x=pd.get_dummies(x)
10 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=2)
11 knn=KNeighborsClassifier(n_neighbors=5)
12 knn.fit(x_train,y_train)
13 ypred=knn.predict(x_test)
14 conf_mat=confusion_matrix(y_test,ypred)
15 print(conf_mat)
```

```
[[909 124]
 [207 167]]
```

- **Confusion matrix is simple table that show the how well the classification model by comparing its prediction to the actual result**
- its break down prediction into 4 categories
- 1.TP(True positive) ->correctly predicted positive outcome
- 2.TN(true negative) ->correctly predicted negative outcome
- 3.FP(false positive) ->incorrectly predicted a positive outcome(the actual outcome was negative) it is also known as type-1 error
- 4.FN(False negative) ->incorrectly predicted a negative outcome(the actual outcome was positive) it is also known as type-2 error

1.Acuracy ->it measure overall performance of the model

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$$

2.Precision ->it focus is on the quality of the models positive predictions(capacity of a model to correctly predict positive instances)

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

3.Recall ->It show the proportion of true positive detected out of all the actual positive instances

$$\text{recall} = \text{TP} / (\text{TP} + \text{FN})$$

4.specificity -> measures the ability of a model to correctly

identify negative instances

```
In [16]: 1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.metrics import confusion_matrix
3 from sklearn.model_selection import train_test_split
4 import pandas as pd
5 df=pd.read_csv("Datasets/Churn.csv")
6 df=df.dropna()
7 x=df.drop(["Unnamed: 0", "customerID", "Churn"],axis=1)
8 y=df['Churn']
9 x=pd.get_dummies(x)
10 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.002,random_state
11 knn=KNeighborsClassifier(n_neighbors=5)
12 knn.fit(x_train,y_train)
13 ypred=knn.predict(x_test)
14 print(ypred)
15 conf_mat=confusion_matrix(y_test,ypred)
16 print(conf_mat)
17
18 tn=conf_mat[0][0]
19 fp=conf_mat[0][1]
20 fn=conf_mat[1][0]
21 tp=conf_mat[1][1]
22 print("accuracy=", (tp+tn)/(tp+tn+fp+fn))
23 print("Precision=", (tp)/(tp+fp))
24 print("recall=", (tp)/(tp+fn))
25 print("specificity=", (tn)/(tn+fp))
```

['Yes' 'No' 'No' 'No' 'No' 'No' 'No' 'No' 'No' 'No' 'Yes' 'Yes' 'No' 'Yes'
 'No']
 [[9 3]
 [2 1]]
 accuracy= 0.6666666666666666
 Precision= 0.25
 recall= 0.3333333333333333
 specificity= 0.75

```
In [24]: 1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score
3 from sklearn.model_selection import train_test_split
4 import pandas as pd
5 df=pd.read_csv("Datasets/Churn.csv")
6 df=df.dropna()
7 x=df.drop(["Unnamed: 0", "customerID", "Churn"], axis=1)
8 y=df['Churn']
9 x=pd.get_dummies(x)
10 x_train, x_test, y_train, y_test=train_test_split(x, y, test_size=0.002, random_state=42)
11 knn=KNeighborsClassifier(n_neighbors=5)
12 knn.fit(x_train, y_train)
13 ypred=knn.predict(x_test)
14 print(ypred)
15 conf_mat=confusion_matrix(y_test, ypred)
16 print(conf_mat)
17 print("accuracy=", accuracy_score(y_test, ypred))
18 print("Precision=", precision_score(y_test, ypred, pos_label='Yes'))
19 print("recall=", recall_score(y_test, ypred, pos_label="Yes"))
20 print("specificity=", recall_score(y_test, ypred, pos_label="No"))
```

```
['Yes' 'No' 'No' 'No' 'No' 'No' 'No' 'No' 'No' 'No' 'Yes' 'Yes' 'No' 'Yes'
 'No']
[[9 3]
 [2 1]]
accuracy= 0.6666666666666666
Precision= 0.25
recall= 0.3333333333333333
specificity= 0.75
```

```
In [40]: 1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score
3 from sklearn.model_selection import train_test_split
4 import pandas as pd
5 df=pd.read_csv("Datasets\diabetes_unclean.csv")
6 df= df.dropna()
7 x=df.drop(["CLASS", "ID", "No_Pation"],axis=1)
8 # x = pd.get_dummies(df.drop(["CLASS", "ID", "No_Pation"], axis=1))
9 # print(x.dtypes)
10 y=df["CLASS"]
11 x=pd.get_dummies(x)
12 xtrain,xtest,ytrain,ytest = train_test_split(x,y,test_size=0.2,random_state=110)
13 knn=KNeighborsClassifier(n_neighbors=5)
14 knn.fit(xtrain,ytrain)
15 ypred=knn.predict(xtest)
16 # print(ypred)
17 conf_mat=confusion_matrix(ytest,ypred)
18 print(conf_mat)
19 print("accuracy=",accuracy_score(ytest,ypred))
20 print("Precision=",precision_score(ytest,ypred,average='macro'))
21 print("recall=",recall_score(ytest,ypred,average="macro"))
22 print("specificity=",recall_score(ytest,ypred,average="macro"))
```

```
[[ 15   0   8   0]
 [  3   6   2   0]
 [  8   2 154   0]
 [  0   0   1   0]]
```

```
accuracy= 0.8793969849246231
```

```
Precision= 0.5650641025641026
```

```
recall= 0.5341632121854816
```

```
specificity= 0.5341632121854816
```

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\metrics\_classification.py:122
1: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels
with no predicted samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
```

```
In [43]: 1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score
3 from sklearn.model_selection import train_test_split
4 import pandas as pd
5 df=pd.read_csv("Datasets\diabetes_unclean.csv")
6 df= df.dropna()
7 x=df.drop(["CLASS", "ID", "No_Pation"],axis=1)
8 y=df["CLASS"]
9 maxacc=0;
10 x=pd.get_dummies(x)
11 xtrain,xtest,ytrain,ytest = train_test_split(x,y,test_size=0.2,random_state=110)
12 for i in range(1,11):
13     knn=KNeighborsClassifier(n_neighbors=i)
14     knn.fit(xtrain,ytrain)
15     ypred=knn.predict(xtest)
16     conf_mat=confusion_matrix(ytest,ypred)
17     acc=accuracy_score(ytest,ypred)
18     if acc>maxacc:
19         maxacc=acc
20         index=i
21 print(f"max acc {acc} for {index}")
```

max acc 0.8542713567839196 for 1

Decision Tree

```
In [15]: 1 import pandas as pd
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import confusion_matrix, accuracy_score, precision_score,
5 from sklearn.preprocessing import LabelEncoder
6
7 # Load dataset
8 df = pd.read_csv("csv/Churn.csv")
9
10 df.dropna(inplace=True)
11 # Drop unwanted columns
12 df.drop(["Unnamed: 0", "customerID"], axis=1, inplace=True)
13
14 # Convert target variable to numeric
15 df["Churn"] = df["Churn"].replace({"No": 0, "Yes": 1})
16
17 # Handle categorical variables
18 for column in df.columns:
19     if df[column].dtype == 'object':
20         df[column] = LabelEncoder().fit_transform(df[column])
21
22 # Split features and target
23 x = df.drop("Churn", axis=1)
24 y = df["Churn"]
25
26 # Split into train and test
27 xtrain, xtest, ytrain, ytest = train_test_split(x, y, test_size=0.2, random_state=42)
28
29 # Train Decision Tree
30 dt = DecisionTreeClassifier(criterion="entropy", max_depth=2)
31 dt.fit(xtrain, ytrain)
32
33 # Predict and evaluate
34 ypr = dt.predict(xtest)
35 print("Predictions:", ypr)
36 print("Accuracy:", accuracy_score(ytest, ypr))
37 print("Precision:", precision_score(ytest, ypr))
38 print("Recall:", recall_score(ytest, ypr))
39 print("Confusion Matrix:\n", confusion_matrix(ytest, ypr))
40
41 c = classification_report(ytest, ypr)
42 print(c)
```

Predictions: [1 0 0 ... 0 1 1]

Accuracy: 0.7299218194740583

Precision: 0.49427480916030536

Recall: 0.6925133689839572

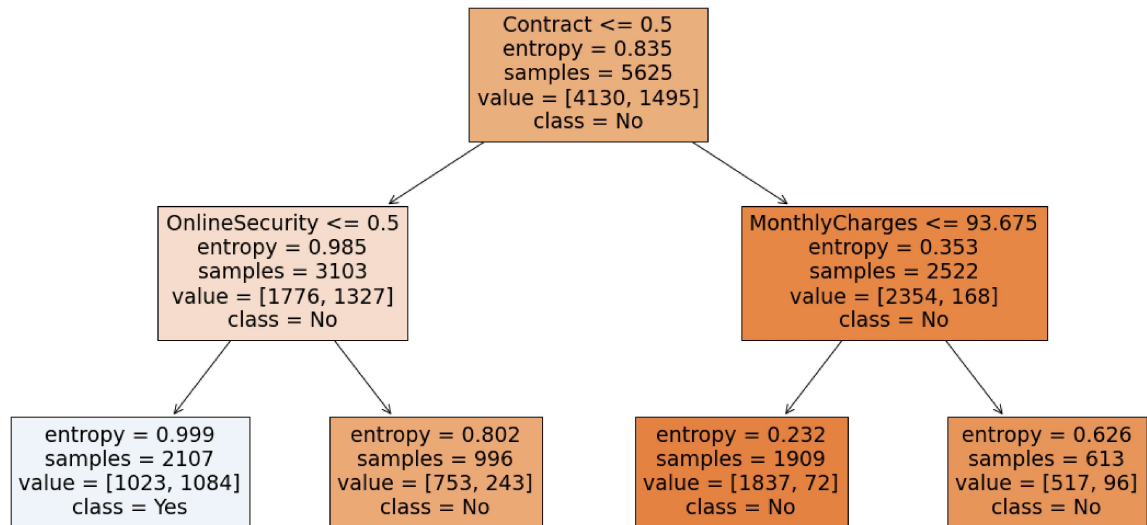
Confusion Matrix:

[[768 265]

[115 259]]

	precision	recall	f1-score	support
0	0.87	0.74	0.80	1033
1	0.49	0.69	0.58	374
accuracy			0.73	1407
macro avg	0.68	0.72	0.69	1407
weighted avg	0.77	0.73	0.74	1407

```
In [13]: 1 from sklearn.tree import plot_tree
2 import matplotlib.pyplot as plt
3
4 plt.figure(figsize=(20,10))
5 plot_tree(dt, feature_names=x.columns, class_names=["No", "Yes"], filled=True)
6 plt.show()
```



P.b. 253 Encode Categorical variables using one hot encoding

- `pd.get_dummies(x)`
- **sensitivity = recall** those are synonyms

In []:

1

In []:

1