

Decision Tree classifier in Scikit-Learn

June 22, 2021

0.1 22 June 2021

0.2 ML Lab 2

0.2.1 Decision Tree classifier in Scikit-Learn

0.2.2 Dr Neeraj Gupta

```
[3]: import pandas as pd
      from sklearn.tree import DecisionTreeClassifier
      from sklearn.model_selection import train_test_split
      from sklearn import metrics
```

```
[4]: pima = pd.read_csv("image/diabetes.csv")
```

```
[5]: pima.head()
```

```
[5]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	\
0	6	148	72	35	0	33.6	
1	1	85	66	29	0	26.6	
2	8	183	64	0	0	23.3	
3	1	89	66	23	94	28.1	
4	0	137	40	35	168	43.1	

	DiabetesPedigreeFunction	Age	Outcome
0	0.627	50	1
1	0.351	31	0
2	0.672	32	1
3	0.167	21	0
4	2.288	33	1

```
[9]: feature_cols = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'DiabetesPedigreeFunction', 'Age']
      X = pima[feature_cols]
      y = pima.Outcome
      print(X)
      print(y)
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	\
0	6	148	72	35	0	33.6	

1	1	85	66	29	0	26.6
2	8	183	64	0	0	23.3
3	1	89	66	23	94	28.1
4	0	137	40	35	168	43.1
..
763	10	101	76	48	180	32.9
764	2	122	70	27	0	36.8
765	5	121	72	23	112	26.2
766	1	126	60	0	0	30.1
767	1	93	70	31	0	30.4

	DiabetesPedigreeFunction	Age
0	0.627	50
1	0.351	31
2	0.672	32
3	0.167	21
4	2.288	33
..
763	0.171	63
764	0.340	27
765	0.245	30
766	0.349	47
767	0.315	23

[768 rows x 8 columns]

0	1
1	0
2	1
3	0
4	1
..	..
763	0
764	0
765	0
766	1
767	0

Name: Outcome, Length: 768, dtype: int64

```
[10]: #Split dataset into training and test set
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.30,
↳random_state=1 )
```

```
[12]: #Create Descision Tree Classifier Object
clf = DecisionTreeClassifier()

#Train Decision Tree Classifier
clf.fit(X_train, y_train)
```

```
#Predict the response for test dataset
y_pred = clf.predict(X_test)
```

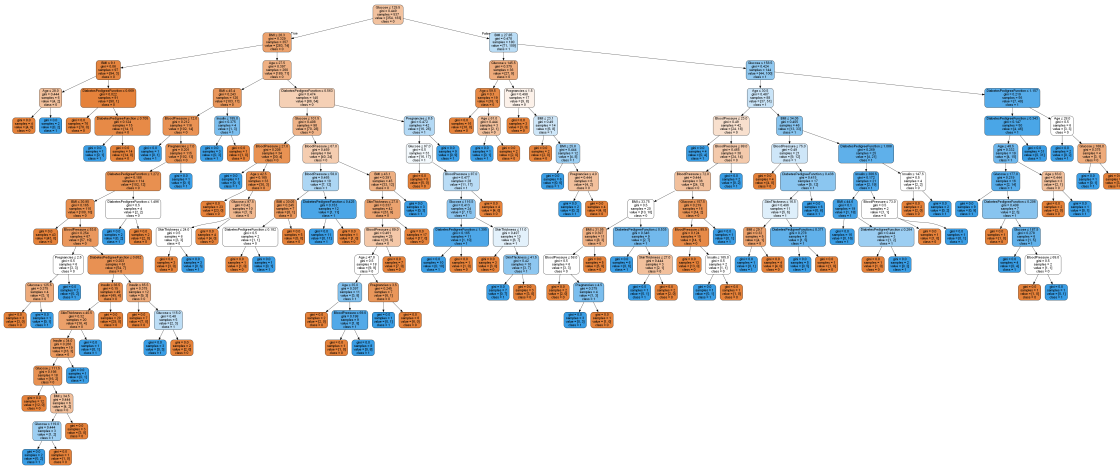
```
[13]: #Model Accuracy
print("Accuracy:", metrics.accuracy_score(y_test, y_pred))
```

Accuracy: 0.696969696969697

```
[14]: from sklearn.tree import export_graphviz
from sklearn.externals.six import StringIO
from IPython.display import Image
import pydotplus
dot_data = StringIO()
export_graphviz(clf, out_file=dot_data,
filled=True, rounded=True,
special_characters=True,feature_names =feature_cols,class_names=['0','1'])
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
graph.write_png('diabetes.png')
Image(graph.create_png())
```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\externals\six.py:31:
DeprecationWarning: The module is deprecated in version 0.21 and will be removed
in version 0.23 since we've dropped support for Python 2.7. Please rely on the
official version of six (<https://pypi.org/project/six/>).
"(https://pypi.org/project/six/).", DeprecationWarning)

[14]:



1 Optimizing Decision Tree Performance

```
[15]: # Create Decision Tree classifier object
      clf = DecisionTreeClassifier(criterion="entropy", max_depth=3)

      # Train Decision Tree Classifier
      clf = clf.fit(X_train,y_train)

      #Predict the response for test dataset
      y_pred = clf.predict(X_test)

      # Model Accuracy, how often is the classifier correct?
      print("Accuracy:",metrics.accuracy_score(y_test, y_pred))
```

Accuracy: 0.7705627705627706

```
[16]: from sklearn.externals.six import StringIO
      from IPython.display import Image
      from sklearn.tree import export_graphviz
      import pydotplus
      dot_data = StringIO()
      export_graphviz(clf, out_file=dot_data,
                      filled=True, rounded=True,
                      special_characters=True, feature_names =feature_cols,class_names=['0','1'])
      graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
      graph.write_png('diabetes.png')
      Image(graph.create_png())
```

[16]:

