

ASD 1049

File Edit View Insert Runtime Tools Help

Comment Share

Code Test

import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

Decision Tree model on fimi dataset

```
# Load the dataset  
data = pd.read_csv('social_network_data.csv')  
data.head()
```

User ID	Gender	Age	EstimatedSalary	Purchased	
0	1984519	Male	19	19000	0
1	1981064	Male	35	20000	0
2	1988075	Female	26	43000	0
3	1980545	Female	27	57000	0
4	1980402	Male	19	70000	0

```
# Use only Age and EstimatedSalary as our independent variables x  
# because all other features like gender and user ID are irrelevant and have no effect on the purchasing capacity of a person.  
feature_cols = ['age', 'estimatedSalary']  
x = data.loc[:,feature_cols].values  
y = data.loc[:,4].values # purchased is our dependent variable y
```

```
# split the dataset into training and test  
from sklearn.model_selection import train_test_split  
x_train, x_test, y_train, y_test = train_test_split(x,y,test_size = 0.25, random_state=0)
```

```
# feature scaling  
from sklearn.preprocessing import StandardScaler  
sc_X = StandardScaler()  
x_train = sc_X.fit_transform(x_train)  
x_test = sc_X.transform(x_test)
```

```
# Fit the model to the decision tree classifier.  
from sklearn.tree import DecisionTreeClassifier  
classifier = DecisionTreeClassifier()  
classifier = classifier.fit(x_train,y_train)
```

```
# prediction  
y_pred = classifier.predict(x_test)accuracy  
from sklearn import metrics  
print('accuracy score:', metrics.accuracy_score(y_test,y_pred))
```

accuracy score: 0.88

```
# Confusion Matrix  
from sklearn.metrics import confusion_matrix  
cm = confusion_matrix(y_test, y_pred)
```

```
# output:  
# array([[10,  0]  
#       [ 0, 17]])  
# This means 10 observations have been falsely classified.
```

```
print(cm)
```

```
[[10  0]  
 [ 0 17]]
```

Let first visualize the model prediction results.

```
from matplotlib.colors import ListedColormap  
x_set, y_set = x_test, y_test  
X1, X2 = np.meshgrid(np.arange(start = x_set[:,0].min()-1, stop = x_set[:,0].max()+1, step = 0.01),np.arange(start = x_set[:,1].min()-1, stop = x_set[:,1].max()+1, step = 0.01))  
plt.contourf(X1,X2, classifier.predict(np.array([X1.ravel(), X2.ravel()]).ravel()).reshape(X1.shape), alpha=0.75, cmap = ListedColormap(['red','green','blue']))  
plt.xlim(x_set[:,0].min(), x_set[:,0].max())  
plt.ylim(x_set[:,1].min(), x_set[:,1].max())  
for i, j in enumerate(np.unique(y_set)):  
    plt.scatter(x_set[:,0],x_set[:,1],c = ListedColormap(Cmap['green'])(i),label = j)
```

plt.title("Decision Tree Test set")
plt.xlabel("Age")
plt.ylabel("Estimated Salary")
plt.legend()

plt.show()

caption(input is:shorthand for: warning: %s argument looks like a single name; use %s or %s as sequence, which should be avoided as name mapping will have precedence in case its not: starting with %s and %s and %s and %s). s = ListedColormap(['red','green'])(i),label = j)

Let us also visualize the tree:

```
from sklearn.tree import export_graphviz  
from six import StringIO  
from IPython.display import Image  
import pydotplus  
dot_data = StringIO()  
export_graphviz(classifier, 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())  
Image(graph.create_png())
```

In the decision tree chart, each internal node has a decision rule that splits the data. One referred to as the Gini ratio, which measures the impurity of the node. We can say a node is pure when all of its records belong to the same class, such nodes known as the leaf node. Here, the resultant tree is unpruned. This unpruned tree is unexplainable and not easy to understand. In the next section, let's optimize it by pruning.

In sklearn, optimization of decision tree classifier performed by only pre-pruning. The maximum depth of the tree can be used as a control variable for pre-pruning.

```
# Create decision tree classifier object  
classifier = DecisionTreeClassifier(criterion="entropy", max_depth=10) # Create decision tree classifier  
classifier = classifier.fit(x_train,y_train) # Fit the model to the decision tree classifier  
y_pred = classifier.predict(x_test) # Predict accuracy, how well the classifier works?  
print('accuracy:', metrics.accuracy_score(y_test, y_pred))
```

accuracy: 0.84

```
# Now let us again visualize the pruned decision tree after optimization.  
dot_data = StringIO()  
export_graphviz(classifier, 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())  
Image(graph.create_png())
```

The pruned model is less complex, explainable, and easy to understand than the previous decision tree model plot.

Decision Tree model on the dataset discussed in the class

```
# Load the dataset  
data = pd.read_csv('class_data.csv')  
data.head()
```

age	income	student	credit_rating	buys_computer
0	1	5	0	0
1	1	3	0	1
2	2	5	0	0
3	3	3	0	0
4	3	1	1	0

```
# feature_cols = ['age', 'income']  
x = data.loc[:,feature_cols].values  
y = data.loc[:,4].values
```

```
# split the dataset into training and test  
from sklearn.model_selection import train_test_split  
x_train, x_test, y_train, y_test = train_test_split(x,y,test_size = 0.25, random_state=0)
```

```
# feature scaling  
from sklearn.preprocessing import StandardScaler  
sc_X = StandardScaler()  
x_train = sc_X.fit_transform(x_train)  
x_test = sc_X.transform(x_test)
```

```
# Fit the model to the decision tree classifier.  
from sklearn.tree import DecisionTreeClassifier  
classifier = DecisionTreeClassifier()  
classifier = classifier.fit(x_train,y_train)
```

```
# prediction  
y_pred = classifier.predict(x_test)accuracy  
from sklearn import metrics  
print('accuracy score:', metrics.accuracy_score(y_test,y_pred))
```

accuracy score: 1.0

```
# Confusion Matrix  
from sklearn.metrics import confusion_matrix  
cm = confusion_matrix(y_test, y_pred)
```

```
# output:  
# array([[4,  0]  
#       [ 0, 1]])  
# This means 0 observations have been falsely classified.
```

```
print(cm)
```

```
[[4  0]  
 [ 0 1]]
```

```
from matplotlib.colors import ListedColormap  
x_set, y_set = x_test, y_test  
X1, X2 = np.meshgrid(np.arange(start = x_set[:,0].min()-1, stop = x_set[:,0].max()+1, step = 0.01),np.arange(start = x_set[:,1].min()-1, stop = x_set[:,1].max()+1, step = 0.01))  
plt.contourf(X1,X2, classifier.predict(np.array([X1.ravel(), X2.ravel()]).ravel()).reshape(X1.shape), alpha=0.75, cmap = ListedColormap(['red','green','blue']))  
plt.xlim(x_set[:,0].min(), x_set[:,0].max())  
plt.ylim(x_set[:,1].min(), x_set[:,1].max())  
for i, j in enumerate(np.unique(y_set)):  
    plt.scatter(x_set[:,0],x_set[:,1],c = ListedColormap(Cmap['green'])(i),label = j)
```

plt.title("Decision Tree Test set")
plt.xlabel("Age")
plt.ylabel("Income")
plt.legend()

plt.show()

caption(input is:shorthand for: warning: %s argument looks like a single name; use %s or %s as sequence, which should be avoided as name mapping will have precedence in case its not: starting with %s and %s and %s and %s). s = ListedColormap(['red','green'])(i),label = j)

```
from sklearn.tree import export_graphviz  
from six import StringIO  
from IPython.display import Image  
import pydotplus  
dot_data = StringIO()  
export_graphviz(classifier, 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())  
Image(graph.create_png())
```

Codecademy products - Cancel contracts here