

```
1 !pip install anytree
```

```
Collecting anytree
```

```
  Downloading anytree-2.12.1-py3-none-any.whl (44 kB)
```

```
44.9/44.9 kB 1.3 MB/s eta 0:00:00
```

```
Requirement already satisfied: six in /usr/local/lib/python3.10/dist-packages (from anytree) (1.16.0)
```

```
Installing collected packages: anytree
```

```
Successfully installed anytree-2.12.1
```

```
1 !pip install pandas numpy typing
```

```
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (1.5.3)
```

```
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (1.23.5)
```

```
Requirement already satisfied: typing in /usr/local/lib/python3.10/dist-packages (3.7.4.3)
```

```
Requirement already satisfied: python-dateutil>=2.8.1 in /usr/local/lib/python3.10/dist-packages (from pandas) (2.8.2)
```

```
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas) (2023.3.post1)
```

```
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.1->pandas)
```

```
1 import pandas as pd
```

```
2 import math
```

```
3 from collections import Counter
```

```
4 from anytree import NodeMixin, RenderTree
```

```
5 from typing import Dict, List
```

```
6 from anytree import RenderTree
```

```
7 from graphviz import Digraph
```

```
8 from graphviz import Source
```

```

1 class DecTree():
2
3     def __init__(self, data: pd.DataFrame, target_attribute: str):
4         self.data = data
5         self.target_attr = target_attribute
6         self.target_attr_vals = data[target_attribute].unique()
7         self.root_node = None
8
9     def pmf_target(self, df: pd.DataFrame) -> Dict[str, float]:
10         target_counts = df[self.target_attr].value_counts(normalize=True)
11         return dict(target_counts)
12
13     def entropy(self, pmf: Dict[str, float]) -> float:
14         entropy_val = -sum(p * math.log2(p) for p in pmf.values())
15         return entropy_val
16
17     def cal_entropy_df(self, df: pd.DataFrame) -> float:
18         pmf = self.pmf_target(df)
19         entropy_val = self.entropy(pmf)
20         return entropy_val
21
22     def info_gain_attribute(self, df: pd.DataFrame, attribute: str) -> float:
23         total_entropy = self.cal_entropy_df(df)
24         attribute_groups = df.groupby(attribute)
25         weighted_entropy = sum(len(subgroup) / len(df) * self.cal_entropy_df(subgroup) for _, subgroup in attribute_groups)
26         information_gain = total_entropy - weighted_entropy
27         return information_gain
28
29     def max_info_gain_attribute(self, df: pd.DataFrame, attributes: List[int]) -> str:
30         info_gain_dict = {attr: self.info_gain_attribute(df, attr) for attr in attributes}
31         max_attr = max(info_gain_dict, key=info_gain_dict.get)
32         print("max_info_gain_attribute = ", max_attr, "INFG-Value = ", info_gain_dict[max_attr])
33         return max_attr
34
35     def build_tree_infgain(self, df: pd.DataFrame, attr_list: List[str], start_node: 'DecTreeNode'):
36
37         if len(df[self.target_attr].unique()) == 1:
38             start_node.name = df[self.target_attr].iloc[0]
39             return
40
41         if not attr_list:
42             start_node.name = df[self.target_attr].value_counts().idxmax()
43             return
44
45         max_attr = self.max_info_gain_attribute(df, attr_list)
46         start_node.attribute = max_attr
47
48         remaining_attrs = [attr for attr in attr_list if attr != max_attr] # Exclude the selected attribute
49         for value in df[max_attr].unique():
50             value_subset = df[df[max_attr] == value]
51             if value_subset.empty:
52                 leaf_node = DecTreeNode(name=df[self.target_attr].mode().values[0], attribute=max_attr, parent=start_node)
53             else:
54                 child_node = DecTreeNode(name=value, attribute=max_attr, parent=start_node)
55                 self.build_tree_infgain(value_subset, attr_list[:], child_node)
56
57     def generate_tree(self):
58         attributes = self.data.columns.to_list()
59         attributes.remove(self.target_attr)
60
61         start_node = DecTreeNode("start", "start")
62         self.build_tree_infgain(self.data, attributes, start_node)
63
64         self.root_node = start_node
65
66     def print_tree(self):
67         for pre, _, node in RenderTree(self.root_node):
68             print(f"{pre} {node.attribute}={node.name}")
69
70
71     def predict(self, X: pd.DataFrame) -> List[str]:
72         predictions = []
73         for _, row in X.iterrows():
74             node = self.root_node
75             while node.children:
76                 attr_value = row[node.attribute.split('=')[0]]
77                 next_child = next((child for child in node.children if child.name == attr_value), None)
78                 if next_child:
79                     node = next_child
80             else:
81                 break
82         predictions.append(node.name)

```

```

83         return predictions
84
85
86 class DecTreeNode(NodeMixin):
87     def __init__(self, name: str, attribute: str, parent: 'DecTreeNode' = None):
88         super(DecTreeNode, self).__init__()
89         self.name = name
90         self.attribute = attribute
91         self.parent = parent
92         self.attr_value = f"{attribute}={name}"

```



```
1 !pwd
```

```
    /content
```

```

1 data = pd.read_csv("data.csv")
2 data

```

	outlook	temp	humidity	windy	play	
0	sunny	hot	high	False	no	
1	sunny	hot	high	True	no	
2	overcast	hot	high	False	yes	
3	rainy	mild	high	False	yes	
4	rainy	cool	normal	False	yes	
5	rainy	cool	normal	True	no	
6	overcast	cool	normal	True	yes	
7	sunny	mild	high	False	no	
8	sunny	cool	normal	False	yes	
9	rainy	mild	normal	False	yes	
10	sunny	mild	normal	True	yes	
11	overcast	mild	high	True	yes	
12	overcast	hot	normal	False	yes	
13	rainy	mild	high	True	no	

```

1 tree = DecTree(data, 'play')
2 tree.generate_tree()
3 tree.print_tree()
4
5
6 dot = Digraph()
7 for pre, _, node in RenderTree(tree.root_node):
8     dot.node(node.attr_value)
9     if node.parent:
10         dot.edge(node.parent.attr_value, node.attr_value)
11 dot.render("decision_tree_visualization", format="png", view=True)
12 Source(dot.source)

```

```
max_info_gain_attribute = outlook TMEG Value = 0.24674091077442032
```

```
1 dt_train = pd.read_csv("cars_train.csv")
```

```
2 dt_test = pd.read_csv("cars_test.csv")
```

```
3
```

```
4 dt_test.head()
```

	buying_price	maintenance_cost	num_doors	num_persons	size_luggage	safet
0	vhigh	med	2	2	small	low
1	high	high	3	4	big	med
2	low	med	4	2	big	low
3	vhigh	high	5more	4	small	low
4	med	med	4	2	med	high

(outlook=sunny) (outlook=overcast) (outlook=rainy)

```
1 # Train
```

```
2 dec_tree = DecTree(dt_train, "decision")
```

```
3 dec_tree.generate_tree()
```

```
4 dec_tree.print_tree()
```

```

├── size_luggage=acc
│   ├── size_luggage=unacc
│   │   ├── buying_price=acc
│   │   │   ├── num_doors=med
│   │   │   │   ├── num_doors=unacc
│   │   │   │   └── num_doors=acc
│   │   └── num_doors=med
│   │       ├── buying_price=4
│   │       │   ├── buying_price=acc
│   │       │   ├── buying_price=acc
│   │       │   └── size_luggage=high
│   │       │       ├── size_luggage=acc
│   │       │       └── size_luggage=unacc
│   │       ├── num_doors=acc
│   │       ├── num_doors=acc
│   │       └── num_doors=unacc
│   └── num_persons=unacc
└── buying_price=4
    ├── maintenance_cost=low
    │   ├── maintenance_cost=acc
    │   ├── maintenance_cost=acc
    │   ├── size_luggage=vhigh
    │   │   ├── size_luggage=unacc
    │   │   ├── size_luggage=unacc
    │   │   └── size_luggage=acc
    │   └── maintenance_cost=acc
    ├── maintenance_cost=vhigh
    │   ├── size_luggage=low
    │   │   ├── size_luggage=acc
    │   │   └── size_luggage=unacc
    │   ├── maintenance_cost=unacc
    │   │   ├── size_luggage=med
    │   │   │   ├── size_luggage=unacc
    │   │   │   ├── size_luggage=acc
    │   │   │   └── size_luggage=unacc
    │   │   └── maintenance_cost=unacc
    │   └── maintenance_cost=med
    │       ├── size_luggage=high
    │       │   ├── size_luggage=acc
    │       │   └── size_luggage=unacc
    │       ├── maintenance_cost=acc
    │       │   ├── num_doors=vhigh
    │       │   │   ├── num_doors=acc
    │       │   │   ├── size_luggage=5more
    │       │   │   │   ├── size_luggage=unacc
    │       │   │   │   └── size_luggage=acc
    │       │   │   └── num_doors=unacc
    │       │   └── maintenance_cost=acc
    │       └── maintenance_cost=high
    │           ├── num_doors=low
    │           │   ├── num_doors=unacc
    │           │   └── num_doors=acc
    │           ├── maintenance_cost=unacc
    │           │   ├── size_luggage=med
    │           │   │   ├── size_luggage=acc
    │           │   │   ├── size_luggage=acc
    │           │   │   └── size_luggage=unacc
    │           └── maintenance_cost=unacc

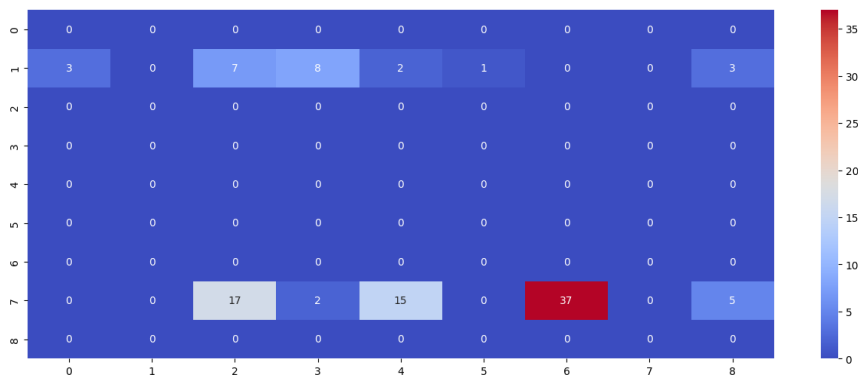
```

```

1 # Test
2 dt_test_x = dt_test.drop(columns="decision")
3 dt_test_y = dt_test["decision"].to_list()
4
5 preds = dec_tree .predict(dt_test_x)

1 # Evaluate results
2 from sklearn.metrics import accuracy_score, confusion_matrix
3 import matplotlib.pyplot as plt
4
5 import seaborn as sns
6
7 acc = accuracy_score(dt_test_y, preds)
8 cm = confusion_matrix(dt_test_y, preds)
9 # print(f"Accuracy: {acc}")
10 # print(cm)
11
12 plt.figure(figsize=(16, 6))
13 fig8= sns.heatmap(cm, annot=True, cmap='coolwarm')
14 fig8.figure.savefig('Confusion_Metrix.png', bbox_inches='tight')

```



```

1 # Create a graphviz Digraph object
2 dot = Digraph()
3
4 # Add nodes and edges to the graphviz object based on your tree structure
5 for pre, _, node in RenderTree(dec_tree.root_node):
6     dot.node(node.attr_value)
7
8     if node.parent:
9         dot.edge(node.parent.attr_value, node.attr_value)
10
11 # Visualize the decision tree
12 dot.render("decision_tree_visualization", format="png", view=True)
13 Source(dot.source)

```