```python
import numpy as np
import pandas as pd
from numpy import log2 as log


dataset = [
    ['<21', 'High', 'Male', 'Single', 'No'],
    ['<21', 'High', 'Male', 'Married', 'No'],
    ['21-35', 'High', 'Male', 'Single', 'Yes'],
    ['>35', 'Medium', 'Male', 'Single', 'Yes'],
    ['>35', 'Low', 'Female', 'Single', 'Yes'],
    ['>35', 'Low', 'Female', 'Married', 'No'],
    ['21-35', 'Low', 'Female', 'Married', 'Yes'],
    ['<21', 'Medium', 'Male', 'Single', 'No'],
    ['<21', 'Low', 'Female', 'Married', 'Yes'],
    ['>35', 'Medium', 'Female', 'Single', 'Yes'],
    ['<21', 'Medium', 'Female', 'Married', 'Yes'],
    ['21-35', 'Medium', 'Male', 'Married', 'Yes'],
    ['21-35', 'High', 'Female', 'Single', 'Yes'],
    ['>35', 'Medium', 'Male', 'Married', 'No']
]


columns = ['Age', 'Income', 'Gender', 'Marital Status', 'Buys']
df = pd.DataFrame(dataset,columns=columns)
df
```

|    | Age   | Income | Gender | Marital Status | Buys |
|----|-------|--------|--------|----------------|------|
| 0  | <21   | High   | Male   | Single         | No   |
| 1  | <21   | High   | Male   | Married        | No   |
| 2  | 21-35 | High   | Male   | Single         | Yes  |
| 3  | >35   | Medium | Male   | Single         | Yes  |
| 4  | >35   | Low    | Female | Single         | Yes  |
| 5  | >35   | Low    | Female | Married        | No   |
| 6  | 21-35 | Low    | Female | Married        | Yes  |
| 7  | <21   | Medium | Male   | Single         | No   |
| 8  | <21   | Low    | Female | Married        | Yes  |
| 9  | >35   | Medium | Female | Single         | Yes  |
| 10 | <21   | Medium | Female | Married        | Yes  |
| 11 | 21-35 | Medium | Male   | Married        | Yes  |
| 12 | 21-35 | High   | Female | Single         | Yes  |
| 13 | >35   | Medium | Male   | Married        | No   |

```python
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
for i in range(5):
    df[columns[i]] = le.fit_transform(df[columns[i]])
df
```

| | Age | Income | Gender | Marital Status | Buys |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 | 0 | 0 |
| 2 | 0 | 0 | 1 | 1 | 1 |
| 3 | 2 | 2 | 1 | 1 | 1 |
| 4 | 2 | 1 | 0 | 1 | 1 |
| 5 | 2 | 1 | 0 | 0 | 0 |
| 6 | 0 | 1 | 0 | 0 | 1 |
| 7 | 1 | 2 | 1 | 1 | 0 |
| 8 | 1 | 1 | 0 | 0 | 1 |
| 9 | 2 | 2 | 0 | 1 | 1 |
| 10 | 1 | 2 | 0 | 0 | 1 |
| 11 | 0 | 2 | 1 | 0 | 1 |
| 12 | 0 | 0 | 0 | 1 | 1 |
| 13 | 2 | 2 | 1 | 0 | 0 |

```python
test_data=[[0, 0, 0, 0]]
test = pd.DataFrame(test_data,columns=['Age', 'Income', 'Gender', 'Marital Status'])
test
```

| | Age | Income | Gender | Marital Status |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |

```python
eps = np.finfo(float).eps
```

$$E(S) = \sum_{i=1}^{c} -p_i \log_2 p_i$$

```python
# Calculate the Cost Function that is Entropy
def find_entropy(df):
    Class = df.keys()[-1]
    entropy = 0
    values = df[Class].unique()
    for value in values:
```

```
        fraction = df[Class].value_counts()[value]/len(df[Class])
        entropy += -fraction*np.log2(fraction)
        print("Class: ", Class, " E(S): ", entropy)
    return entropy
```

$$E(T, X) = \sum_{c \in X} P(c)E(c)$$

```
#Find entropy of the attribute (Each Columns)
def find_entropy_attribute(df,attribute):
    Class = df.keys()[-1]
    target_variables = df[Class].unique()
    variables = df[attribute].unique()
    entropy2 = 0
    for variable in variables:
        entropy = 0
        for target_variable in target_variables:
            num = len(df[attribute][df[attribute]==variable][df[Class]==target_variable])
            den = len(df[attribute][df[attribute]==variable])
            fraction = num/(den+eps)
            entropy += -fraction*log(fraction+eps)
        fraction2 = den/len(df)
        entropy2 += -fraction2*entropy
        print("Class: ", Class, " E(T,X): ", entropy2)
    return abs(entropy2)
```

$$Gain(T, X) = Entropy(T) - Entropy(T, X)$$

```
#Find Root Node
def find_winner(df):
    IG = []
    for key in df.keys()[:-1]:
        IG.append(find_entropy(df)-find_entropy_attribute(df,key))
        print(np.argmax(IG))
    return df.keys()[:-1][np.argmax(IG)]


def get_subtable(df, node,value):
    return df[df[node] == value].reset_index(drop=True)


def buildTree(df,tree=None):
    Class = df.keys()[-1]
    #Build Decision Tree

    #Get attribute with maximum information gain
    node = find_winner(df)
    print("node with max info gain: ",node)

    #Get distinct value of that attribute
    attValue = np.unique(df[node])
    print("distinct values found: ", attValue)
```

```python
    #Create an empty dictionary to create tree
    if tree is None:
        tree={}
        tree[node] = {}

    #Check if the subset is pure and stops if it is.
    for value in attValue:
        subtable = get_subtable(df,node,value)
        print("subtable: ", subtable)
        clValue,counts = np.unique(subtable['Buys'],return_counts=True)
        print("clValue: ", clValue)
        print("counts: ", counts)

        if len(counts)==1: #Checking purity of subset
            tree[node][value] = clValue[0]
        else:
            tree[node][value] = buildTree(subtable) #Calling the function recursively

    return tree


dtree = buildTree(df)
dtree
```

```
Class:  Buys  E(T,X):  -0.5509775004326932
2
node with max info gain:  Gender
distinct values found:  [0 1]
subtable:      Age  Income  Gender  Marital Status  Buys
0    1       1       0                0     1
1    1       2       0                0     1
clValue:  [1]
counts:  [2]
subtable:      Age  Income  Gender  Marital Status  Buys
0    1       0       1               1     0
1    1       0       1               0     0
2    1       2       1               1     0
clValue:  [0]
counts:  [3]
subtable:      Age  Income  Gender  Marital Status  Buys
0    2       2       1               1     1
1    2       1       0               1     1
2    2       1       0               0     0
3    2       2       0               1     1
4    2       2       1               0     0
clValue:  [0 1]
counts:  [2 3]
Class:  Buys  E(S):  0.44217935649972373
Class:  Buys  E(S):  0.9709505944546686
Class:  Buys  E(T,X):  -0.970950594454668
0
Class:  Buys  E(S):  0.44217935649972373
Class:  Buys  E(S):  0.9709505944546686
Class:  Buys  E(T,X):  -0.5509775004326933
Class:  Buys  E(T,X):  -0.950977500432693
1
Class:  Buys  E(S):  0.44217935649972373
Class:  Buys  E(S):  0.9709505944546686
Class:  Buys  E(T,X):  -0.3999999999999974
```

```
Class:  Buys  E(T,X):  -0.950977500432693

1
Class:  Buys  E(S):  0.44217935649972373
Class:  Buys  E(S):  0.9709505944546686
Class:  Buys  E(T,X):  1.9220559022889504e-16
Class:  Buys  E(T,X):  3.2034265038149176e-16
3
node with max info gain:  Marital Status
distinct values found:  [0 1]
subtable:      Age  Income  Gender  Marital Status  Buys
0    2       1       0                   0        0
1    2       2       1                   0        0
clValue:  [0]
counts:  [2]
subtable:      Age  Income  Gender  Marital Status  Buys
0    2       2       1                   1        1
1    2       1       0                   1        1
2    2       2       0                   1        1
clValue:  [1]
counts:  [3]
{'Age': {0: 1,
  1: {'Gender': {0: 1, 1: 0}},
  2: {'Marital Status': {0: 0, 1: 1}}}}
```

```python
def predict(inst,tree):
    #Recursively we going through the tree that built earlier
    for nodes in tree.keys():
        value = inst[nodes]
        tree = tree[nodes][value]
        prediction = 0

        if type(tree) is dict:
            prediction = predict(inst, tree)
        else:
            prediction = tree
            break;

    return prediction


tester = test.iloc[0]
Prediction = predict(tester,dtree)


Prediction
```

```
1
```

```python
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.tree import plot_tree
sklearn_dtree=DecisionTreeClassifier(criterion="entropy")


df1 = df.copy()
df1.drop('Buys', axis=1, inplace=True)
```
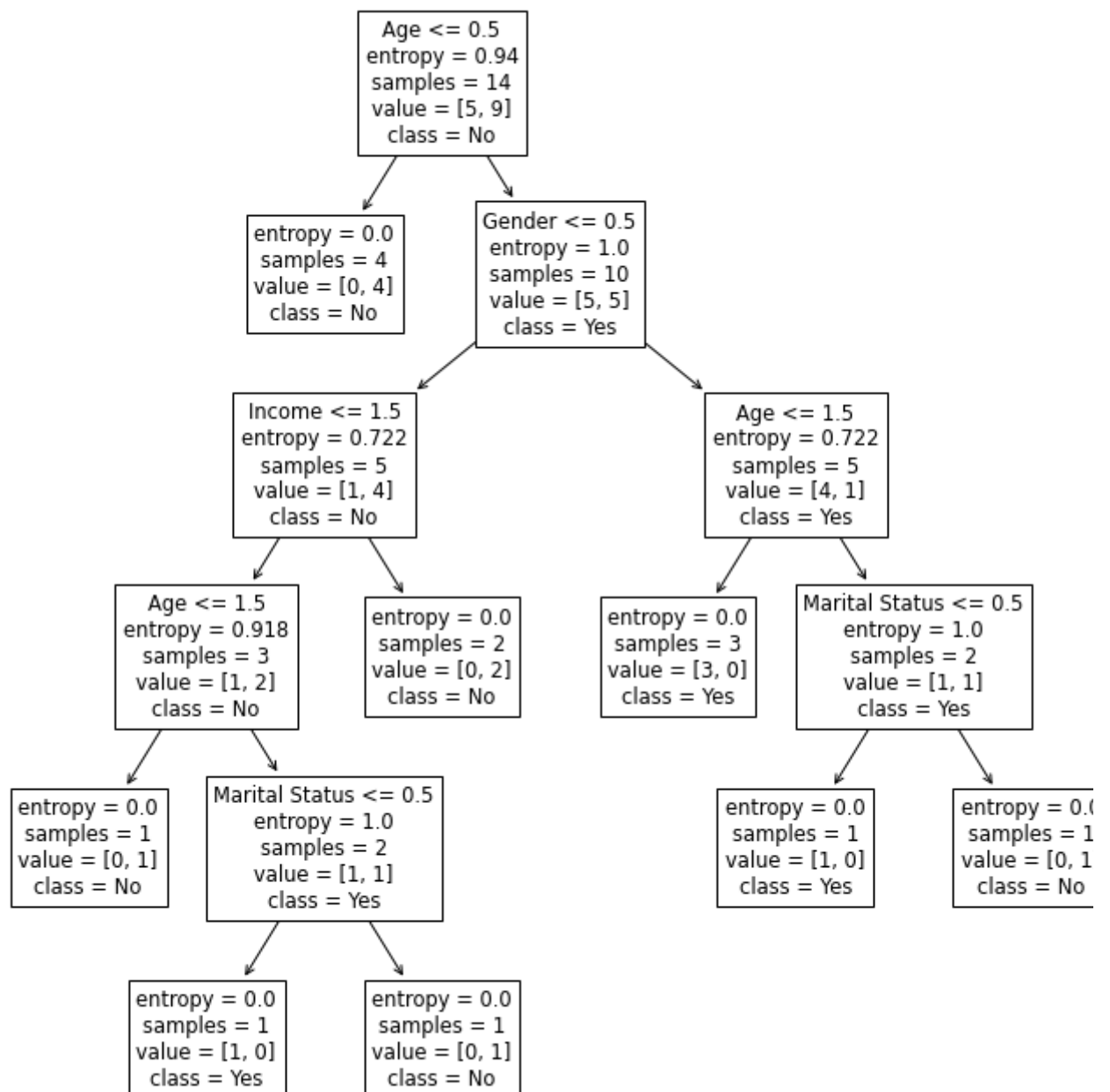
```
X=df1

sklearn_dtree.fit(X, df['Buys'])
sklearn_dtree.predict(test)

    array([1])


import matplotlib.pyplot as plt
plt.figure(figsize=(12,12))
dec_tree = plot_tree(decision_tree=sklearn_dtree, feature_names = df.columns, class_names
plt.show()
```



```
dtree

    {'Age': {0: 1,
      1: {'Gender': {0: 1, 1: 0}},
      2: {'Marital Status': {0: 0, 1: 1}}}}
```

✓ 0s    completed at 6:15 PM