

## LAB - 3

ID 3 code.

Code:-

```
import numpy as np
```

```
import pandas as pd
```

```
from graphviz import Digraph
```

```
def entropy(dataset):
```

```
class_counts = dataset.iloc[:, -1].value_counts
```

```
prob = class_counts / len(dataset)
```

```
return -np.sum(prob * np.log2(prob))
```

```
def info_gain(dataset, feature):
```

```
tot_entropy = entropy(dataset)
```

```
feature_value = dataset[feature].value_counts()
```

```
weight_col_entropy = 0
```

```
for value, count in feature_value.items():
```

```
subset = dataset[dataset[feature] == value]
```

```
weight_col_entropy += (count / len(dataset)) *
```

```
entropy(subset)
```

```
return tot_entropy - weight_col_entropy
```

```
def best_feature(dataset):
```

```
feature = dataset.columns[-1]
```

```
best_info_gain = -1
```

```
best_feature = None
```

```
for feature in feature:
```

```
info_gain = information gain / dataset[feature]
```

```
if info_gain > best_info_gain:
```

```
best_info_gain = info_gain
```

but feature - > feature

> return but\_feature

```
def id3(dataset, max_depth=None, depth=0):
    if len(dataset.iloc[:, -1].unique()) == 1:
```

> return dataset.iloc[0, -1]

```
if len(dataset.columns) == 1:
```

> return dataset.iloc[0, -1].node()

best = best\_feature(dataset)

true = & best: 1 0 0

```
for value in dataset[best].unique():
```

subset = dataset[dataset[best] == value]

```
def create_decision_tree(true, slot=None,
```

parent=Root, parent\_value="0"):

if slot is None:

dot = Digraph(format="png", engine="dot")

if isinstance(true, dict):

for feature, branches in true.items():

feature\_name = f'{parent.name}-{feature}'

dot.node(feature\_name, feature)

dot.edge(parent\_name, parent\_value, label="slot")

- parent\_value)

else:

dot.node(parent\_name + " leaf", "leaf")

dot.edge(parent\_name, parent\_value, label="slot")

(label="leaf"))

> return dot

data = { 'outlook': ['Sunny', 'Sunny', 'Overcast', 'Rainy',  
 'Rainy', 'Rainy', 'Overcast'], 'Temp': ['Hot', 'Hot', 'Hot', 'Mild', 'Mild', 'Mild', 'Mild'],  
 'Humidity': ['High', 'High', 'Normal', 'Normal', 'Normal', 'High', 'Normal'],  
 'Wind': ['Strong', 'Weak', 'Weak', 'Weak', 'Strong', 'Strong', 'Weak'],  
 'Play': ['Yes', 'Yes', 'Yes', 'Yes', 'No', 'No', 'Yes'] }

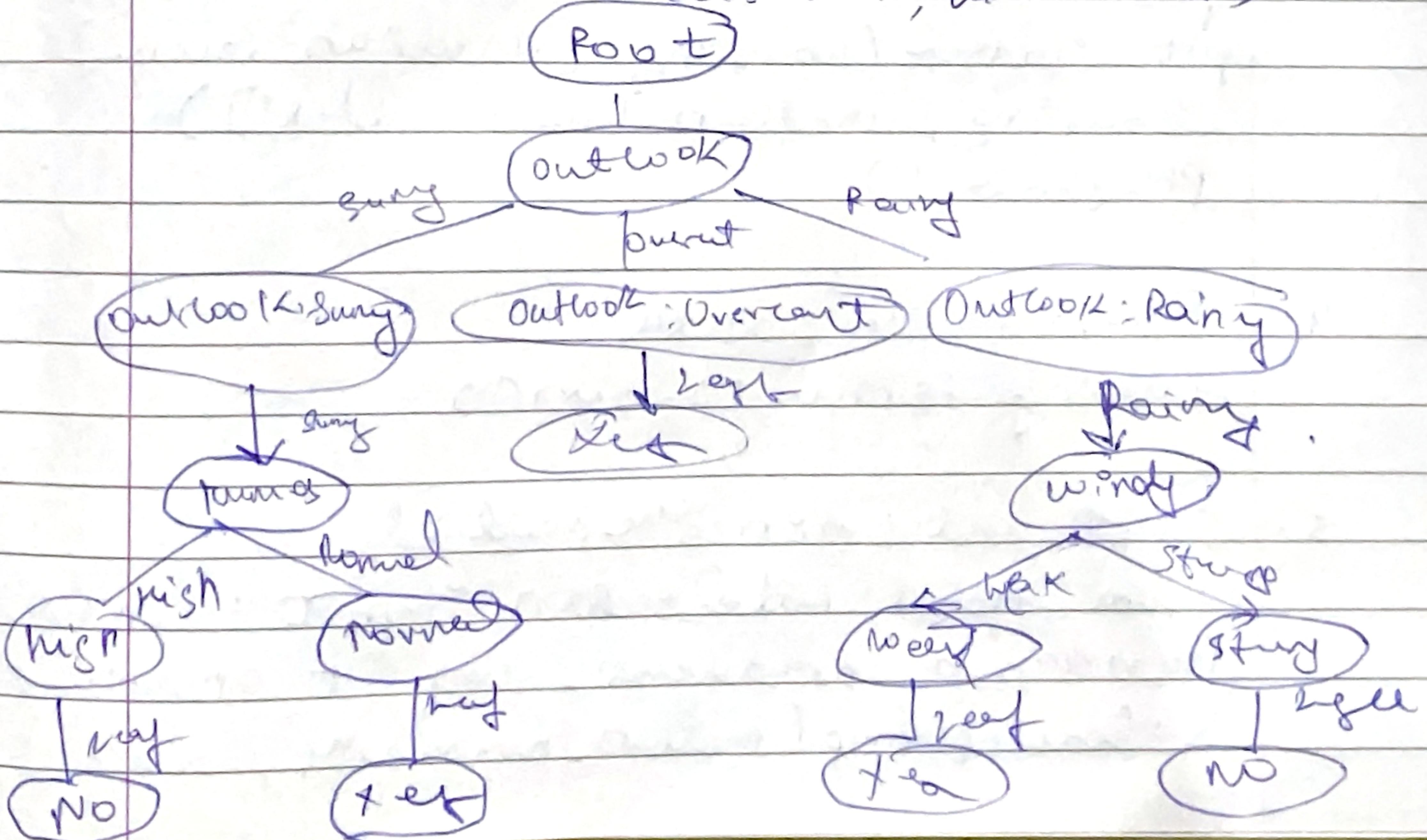
"Temp": ["Hot", "Hot", "Hot", "Mild", "Mild", "Mild", "Mild"],  
 "Humidity": ["High", "High", "Normal", "Normal", "Normal", "High", "Normal"],  
 "Wind": ["Strong", "Weak", "Weak", "Weak", "Strong", "Strong", "Weak"],  
 "Play": ["Yes", "Yes", "Yes", "Yes", "No", "No", "Yes"] }

df = pd.DataFrame(data)

tree = id3(df, max\_depth=3)

dot = create\_decision\_tree(tree)

dot.render("decisionTree", viewer=True).



To Do:- END to END Machine Learning project

1. Get the data:

```
import pandas as pd
```

```
housing = pd.read_csv ("sample_data/  
california_housing_train.csv")
```

2. Discover the data.

```
housing.head()
```

```
housing.info()
```

```
housing.describe()
```

3. Visualize the data:

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns,
```

```
plt.hist(housing['median_income'])
```

```
plt.show()
```

```
plt.scatter(housing['median_income'],
```

```
housing['median_house_value']).
```

```
plt.show()
```

4. Prepare the data

```
housing.isnull().sum()
```

5. Select and train the model

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.preprocessing import OneHotEncoder
```

```
X = housing.drop(['median_house_value'], axis=1)
```

`y = housing['medie_house_value']`

`x_train, y-test, x-train, y-test = train test split (x,y, testsize=0.2, random_state=42)`  
`from sklearn.linear_model import Linear Regression`  
`model = Linear Regression()`  
`model.fit(x_train, y-train)`

### b. Fine Tune your model.

`from sklearn.metrics import root_mean_squared_error`

`import numpy as np`

`y-pred = model.predict(x-test)`

~~`rms = root_mean_squared_error(y-test, y-pred)`~~

~~`printf("RMSE": %f", rms)`~~

ND  
17/3/15