# ID3 Implementation

```python
import pandas as pd
import numpy as np
import math
from graphviz import Digraph

def calculate_entropy(data, target_column):
    total_samples = len(data)
    class_counts = data[target_column].value_counts()
    entropy = 0
    for count in class_counts:
        probability = count / total_samples
        entropy -= probability * math.log2(probability)
    return entropy

def calculate_information_gain(data, feature, target_column):
    total_samples = len(data)
    weighted_entropy = 0
    for value in data[feature].unique():
        subset = data[data[feature] == value]
        subset_entropy = calculate_entropy(subset, target_column)
        weighted_entropy += (len(subset)/total_samples)
                        * subset_entropy
    return
    return calculate_entropy(data, target_column) - weighted_entropy

def build_tree(data, target_column, features, parent_node_
                            class = None):
    if len(data[target_column].unique()) <= 1:
        return data[target_column].unique()[0]
    if len(features) == 0:
        return parent_node_class
    parent_node_class = data[target_column].mode()[0]
```

```
best feature = max (features, key = lambda feature:
        Calculate information-gain (data, feature,
                        target column.))

tree = { best-feature : {}}
features. remove (best feature)

for value in data (best-feature). unique ():
    subset = data [data [best-feature] == value]
    subtree = build tree (subset, target-column,
    features.copy(), parent-node-class
    tree [best-feature][value]: subtree
return tree


def visualize tree (tree, dot = None, node-name = 'Root'):
    if not dot:
        dot = Digraph (comment = 'Decision Tree')
    if isinstance (tree, dict):
        feature = list (tree. keys ()) [0]
        dot.node (node-name, label = feature)

        for value, subtree in tree [feature].items
            ():

            child-node-name = node-name + "_" & sts(value)
            dot.edge (node-name, child-node-name, labl = sts(
                                    value))

            visualize tree (subtree, dot, child-node
                                    name)

    else:
        dot.node (node-name, label= sts (tree ))
        return dot
data = { 'Outlook' : [.........
            ]}
```

```
df = pd. DataFrame (data)

target_colum = 'Play Tennis'
features = list (df. columns)
feature. remove (target_column)
tree = build-tree (df, target_column, features)
dot = visualize_tree (tree)
dot. render ('decision_tree')
```
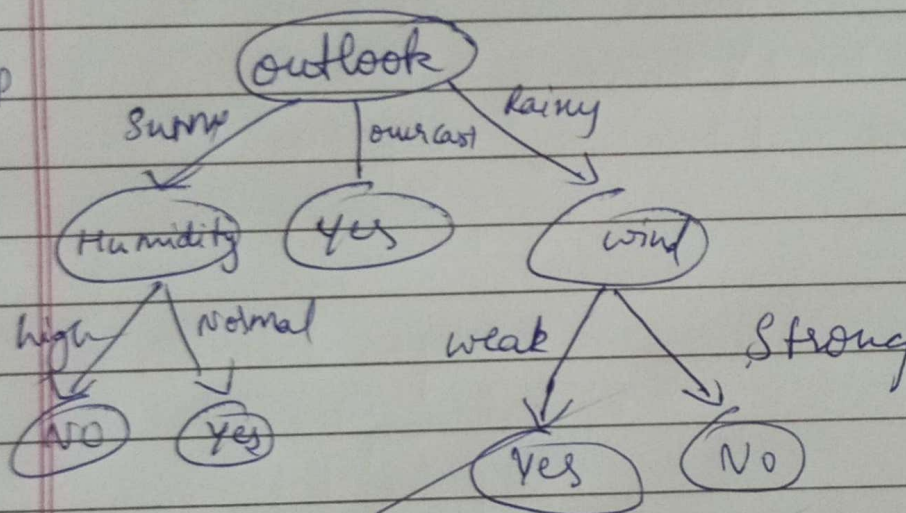
o/p

Outlook

Sunny — Humidity

overcast — Yes

Rainy — Wind

Humidity: high → No, Normal → Yes

Wind: weak → Yes, Strong → No