

- 3 Write a program to demonstrate the working of the decision tree based ID3 algorithm. Use an appropriate dataset for building the decision tree and apply this knowledge to classify a new sample.

```
import pandas as pd
from pandas import DataFrame
df-tennis = pd.read_csv('C:/Users/hp/Desktop/4MT17CS044-JOVITA/
                        PlayTennis.csv')
```

```
attribute-names = list(df-tennis.columns)
```

```
attribute-names.remove('Play Tennis')
```

```
print(attribute-names)
```

```
def entropy-of-list(lst):
```

```
    from collections import Counter
```

```
    count = Counter(x for x in lst)
```

```
    num-instances = len(lst) * 1
```

```
    probs = [x/num-instances for x in count.values()]
```

```
    return entropy(probs)
```

```
def entropy(probs):
```

```
    import math
```

```
    return sum([-prob * math.log(prob, 2) for prob in probs])
```

```
total-entropy = entropy-of-list(df-tennis['Play Tennis'])
```

```
def information-gain(df, split-attribute-name, target-attribute-name, trace=0):
```

```
    df-split = df.groupby(split-attribute-name)
```

```
    nobs = len(df.index) + 1
```

```
    df-agg-ent = df-split.agg([target-attribute-name : [entropy-of-list,
                                                         lambda x : len(x)/nobs]])
```

```

df_agg_ent.columns = ['Entropy', 'probobservations']
new_entropy = sum(df_agg_ent['Entropy'] * df_agg_ent['probobservations'])
old_entropy = entropy_of_list(df[target_attribute_name])
print('Split attribute-name, 'IGI:', old_entropy - new_entropy)
return old_entropy - new_entropy

```

```

def id3(df, target_attribute_name, attribute_names, default_class = None):
    from collections import Counter
    count = Counter(x for x in df[target_attribute_name])
    if len(count) == 1:
        return next(iter(count))
    elif df.empty or (not attribute_names):
        return default_class
    else:
        default_class = max(count.keys())
        gain = [information_gain(df, attr, target_attribute_name) for
                 attr in attribute_names]
        index_of_max = gain.index(max(gain))
        best_attr = attribute_names[index_of_max]
        tree = {best_attr: {}

```

```

        remaining_attribute_names = [i for i in attribute_names if i != best_attr]

```

```

    for attr_val, data_subset in df.groupby(best_attr):
        subtree = id3(data_subset, target_attribute_name, remaining_attribute_names, default_class)
        tree[best_attr][attr_val] = subtree
    return tree

```

```
from pprint import pprint
tree = id3(df-tennis, 'Play-Tennis', attribute-names)
print("\n The Resultant decision tree is :\n")
pprint(tree)
```


Output

['outlook', 'temperature', 'humidity', 'wind']

outlook IG: 0.2467498197744392

Temperature IG: 0.02992225658954647

Humidity IG: 0.15183550136239136

wind IG: 0.04812703040826927

Temperature IG: 0.01997309402197489

Humidity IG: 0.01997309402197489

wind IG: 0.9709505994546686

Temperature IG: 0.5709505944596686

Humidity IG: 0.9709505944596686

wind IG: 0.01997309402197489

The resultant decision tree is:

{ 'outlook' : { 'overcast' : 'yes',

'Rain' : { 'wind' : { 'strong' : 'No', 'weak' : 'yes' } },

'Sunny' : { 'humidity' : { 'high' : 'No', 'Normal' : 'yes' } } }