

## question1.py

```
1 import numpy as np
2 import pandas as pd
3 import math
4
5 # Define a Pandas dataframe with the contents of the table from the problem statement.
6 data = pd.DataFrame([
7     [1, 1, 0, 0, 1, 1],
8     [1, 1, 1, 0, 1, 1],
9     [0, 0, 1, 0, 0, 0],
10    [0, 1, 1, 0, 1, 0],
11    [0, 1, 1, 0, 0, 1],
12    [0, 0, 1, 1, 1, 1],
13    [1, 0, 0, 0, 1, 0],
14    [0, 1, 0, 1, 1, 1],
15    [0, 0, 1, 0, 1, 1],
16    [1, 0, 0, 0, 0, 0],
17    [1, 1, 1, 0, 0, 1],
18    [0, 1, 1, 1, 1, 0],
19    [0, 0, 0, 0, 1, 0],
20    [1, 0, 0, 1, 0, 1],
21 ], columns=['Early', 'Finished HMK', 'Senior', 'Likes Coffee', 'Liked The Last Jedi', 'A'])
22
23 # Define an entropy function. The input is an attribute (one column from the table).
24 def entropy(attr):
25     total = len(attr)
26     # If there are no values, the entropy is 0
27     if total == 0:
28         return 0
29     # Attributes are binary (0 and 1), so proportions are calculated by summing and dividing.
30     # Proportion of attr = 1
31     p1 = sum(attr) / total
32     # Proportion of attr = 0
33     p0 = 1 - p1
34     # Calculate entropy
35     if p1 == 0 or p0 == 0:
36         # Entropy is 0 if either proportion is 0.
37         e = 0
38     else:
39         # Otherwise, use the formula.
40         e = - (p1 * math.log2(p1) + p0 * math.log2(p0))
41     return e
42
43 # Compute the entropy of the dataset.
44 # This value is global so it can be reference in the information_gain
45 # function defined below.
46 H_S = entropy(data['A'])
47
```

```
48 # Define a function to calculate information gain for an attribute.
49 def information_gain(data, attribute, target):
50     # Get the unique values for the attribute (in this case, 0 and 1)
51     values = data[attribute].unique()
52     # Initialize the weighted entropy to zero.
53     weighted_entropy = 0
54     # Loop through each unique value.
55     for v in values:
56         # Extract the subset that matches this value.
57         subset = data[data[attribute] == v][target]
58         # Calculate the weighted entropy of the subset, and add it to
59         # the overall weighted entropy.
60         weighted_entropy += (len(subset) / len(data)) * entropy(subset)
61     return H_S - weighted_entropy
62
63 # Compute the information gain for each attribute.
64 attributes = ['Early', 'Finished HMK', 'Senior', 'Likes Coffee', 'Liked The Last Jedi']
65 i_g_values = {attr: information_gain(data, attr, 'A') for attr in attributes}
66
67 print(f"Information gain for all attributes at depth 1:\n{i_g_values}")
68 # Select the best attribute (max information gain) to split at depth 1.
69 best_attr = max(i_g_values, key=i_g_values.get)
70
71 # Split data based on that best attribute
72 # The splits are in a hash where the key is the 0/1 value of the split attribute,
73 # and the value is the subset of the original table for that attribute value.
74 depth_1_split = {v: data[data[best_attr] == v] for v in data[best_attr].unique()}
75 print("-----")
76 print(f"depth_1_split: {depth_1_split}")
77
78 # Select the best attributes to split at depth 2.
79 # These splits will also go in a hash as before.
80 depth_2_splits = {}
81 # Loop through the keys/values of the depth 1 split.
82 for v, subset in depth_1_split.items():
83     # If there is only one target value in the split, then no further
84     # splits are needed or possible.
85     if len(subset['A'].unique()) == 1:
86         depth_2_splits[v] = None
87     else:
88         # Generate a list of the remaining attributes.
89         remaining_attrs = [attr for attr in attributes if attr != best_attr]
90         # For each of the remaining attributes, calculate its information gain.
91         # Store those in a hash where the key is the attribute.
92         ig_sub = {attr: information_gain(subset, attr, 'A') for attr in remaining_attrs}
93         print(f"Information gain at depth 2 for {best_attr} = {v}:\n{ig_sub}")
94         # Pick the highest information gain.
95         best_sub_attr = max(ig_sub, key=ig_sub.get)
96         # Store that attribute.
```

```
97         depth_2_splits[v] = best_sub_attr
98
99     print("-----")
100    print(f"depth_2_splits: {depth_2_splits}")
101    print("-----")
102
103    # Construct decision trees for depth 1 and depth 2
104    decision_tree_depth_1 = {best_attr: {v: {'Entropy': entropy(subset['A']), 'Positives':
sum(subset['A']), 'Negatives': len(subset) - sum(subset['A'])} for v, subset in
depth_1_split.items()}}
105    decision_tree_depth_2 = {best_attr: {v: {'Entropy': entropy(subset['A']), 'Positives':
sum(subset['A']), 'Negatives': len(subset) - sum(subset['A']), 'Next Split':
depth_2_splits[v]} for v, subset in depth_1_split.items()}}
106
107    # Extend to depth 3 for part C of this section.
108
109    # The depth_2_split hash that I created above does not have the table values. It only has
the attributes.
110    # So to compute the depth 3 split, I have to start again with the depth 1 split and work down
to it.
111    # This is inefficient and could be improved.
112
113    depth_3_splits = {}
114    for v, subset in depth_1_split.items():
115        if len(subset['A'].unique()) == 1:
116            depth_3_splits[v] = None
117        else:
118            remaining_attrs = [attr for attr in attributes if attr != best_attr]
119            ig_sub = {attr: information_gain(subset, attr, 'A') for attr in remaining_attrs}
120            best_sub_attr = max(ig_sub, key=ig_sub.get)
121
122            split_2 = {vv: subset[subset[best_sub_attr] == vv] for vv in
subset[best_sub_attr].unique()}
123
124            depth_3_splits[v] = {best_sub_attr: {}}
125
126            for vv, subset_2 in split_2.items():
127                if len(subset_2['A'].unique()) == 1:
128                    depth_3_splits[v][best_sub_attr][vv] = None
129                else:
130                    remaining_attrs_2 = [attr for attr in remaining_attrs if attr !=
best_sub_attr]
131                    ig_sub_2 = {attr: information_gain(subset_2, attr, 'A') for attr in
remaining_attrs_2}
132                    print(f"Information gain at depth 3 for {best_attr} = {v}, {best_sub_attr} =
{vv}: \n{ig_sub_2}")
133                    best_sub_attr_2 = max(ig_sub_2, key=ig_sub_2.get)
134                    depth_3_splits[v][best_sub_attr][vv] = best_sub_attr_2
135
136    # Construct decision tree for depth 3
```

```
137 decision_tree_depth_3 = {best_attr: {}}
138 for v, subset in depth_1_split.items():
139     decision_tree_depth_3[best_attr][v] = {
140         'Entropy': entropy(subset['A']),
141         'Positives': sum(subset['A']),
142         'Negatives': len(subset) - sum(subset['A']),
143         'Next Split': depth_2_splits[v],
144         'Depth 3 Splits': depth_3_splits[v] if depth_3_splits[v] is not None else "Leaf Node"
145     }
146
147 print("-----")
148 print(f"depth_3_splits: {depth_3_splits}")
149 print("-----")
150
151 # Display the decision tree results.
152 print("Decision tree, depth = 1:")
153 print(decision_tree_depth_1)
154 print("Decision tree, depth = 2")
155 print(decision_tree_depth_2)
156 print("Decision tree, depth = 3:")
157 print(decision_tree_depth_3)
158
```