

Assignment 1

1) Decision Tree Basics

- a) Given $A = 5$ and $D = 3$, we calculate permutations $P(5, 3)$, meaning that 3 are chosen from a set of 5 where order matters:

$$P(5,3) = \frac{5!}{(5-3)!} = \frac{120}{2} = 60$$

For the general case $A \gg D$, the formula is:

$$P(A,D) = \frac{A!}{(A-D)!}$$

- b) Here is the output for the tree of depth 1:

```
{ 'Finished HMK':  
  {  
    1: { 'Entropy': 0.863120568566631, 'Positives': 5, 'Negatives': 2 },  
    0: { 'Entropy': 0.9852281360342516, 'Positives': 3, 'Negatives': 4 }  
  }  
}
```

Here is the information gain for all attributes at depth 1:

```
{  
  'Early': 0.02024420715375619,  
  'Finished HMK': 0.06105378373381032,  
  'Senior': 0.011265848648557286,  
  'Likes Coffee': 0.03914867190307081,  
  'Liked The Last Jedi': 0.0013397424044413464  
}
```

Here is the output for the tree of depth 2:

```
{ 'Finished HMK':  
  {  
    1: {  
      'Entropy': 0.863120568566631,  
      'Positives': 5,  
      'Negatives': 2,  
      'Next Split': 'Early',  
    },  
    0: {  
      'Entropy': 0.9852281360342516,  
      'Positives': 3,  
      'Negatives': 4,  
      'Next Split': 'Likes Coffee',  
    }  
  }  
}
```

(See next page for information gain)

Information gain at depth 2 for Finished HMK = 1:

```
{  
  'Early': 0.41379956460568024,  
  'Senior': 0.2916919971380598,  
  'Likes Coffee': 0.18385092540042136,  
  'Liked The Last Jedi': 0.2916919971380598  
}
```

Information gain at depth 2 for Finished HMK = 0:

```
{  
  'Early': 0.02024420715375619,  
  'Senior': 0.12808527889139454,  
  'Likes Coffee': 0.46956521111470706,  
  'Liked The Last Jedi': 0.02024420715375619  
}
```

The code for this is in the appendix at the end of the file.

c) Here is the output when the tree is extended to a depth of 3:

```
{'Finished HMK':  
  {  
    1: {  
      'Entropy': 0.863120568566631,  
      'Positives': 5,  
      'Negatives': 2,  
      'Next Split': 'Early',  
      'Depth 3 Splits': {  
        'Early': {  
          1: None,  
          0: 'Senior'}  
        }  
      },  
    0: {  
      'Entropy': 0.9852281360342516,  
      'Positives': 3,  
      'Negatives': 4,  
      'Next Split': 'Likes Coffee',  
      'Depth 3 Splits': {  
        'Likes Coffee': {  
          0: 'Senior',  
          1: None}  
        }  
      }  
    }  
  }  
}
```

Here is the information gain at depth 3:

Information gain at depth 3 for Finished HMK = 1, Early = 0:

```
{  
  'Senior': 0.29650626049338447,  
  'Likes Coffee': -0.014771863965748366,  
  'Liked The Last Jedi': 0.29650626049338447  
}
```

Information gain at depth 3 for Finished HMK = 0, Likes Coffee = 0:

```
{  
  'Early': 0.43425063560155797,  
  'Senior': 0.5852281360342516,  
  'Liked The Last Jedi': 0.43425063560155797  
}
```

The code for this is in the appendix at the end of the file.

Overall, I would choose the depth=2 tree. A depth of 1 likely does not capture enough of the complexity in the dataset, given the number of features, and a depth of 3 may cause overfitting.

- d) A decision tree is realizable if there is some tree that perfectly classifies the given data set. A tree of “no fixed depth” means that there is no set limit to the depth of the tree, so the tree can grow to arbitrary depth.
- A dataset where two instances (table rows) have the same attribute (input) values but a different label (output) value will not be realizable, because no decision tree can be constructed that will label these two instances correctly. This can be proved by contradiction:
- Assume two instances, x_1 and x_2 , that have the same attributes, but x_1 has label y , and x_2 has label z , where $y \neq z$.
- Because x_1 and x_2 have the same attributes, their feature vectors are identical. A deterministic decision tree must assign the same label to identical feature vectors. Since we know that $y \neq z$, we know that the tree does not classify one of x_1 or x_2 correctly. And since a realizable decision tree must classify all instances correctly, no such tree can exist.

2) Application of Decision Tree on Real-World Data Set

a) Here is the training accuracy for cut-off depths between 2 and 10:

```
'Depth': 2, 'Training Accuracy': 0.9444575312119404  
'Depth': 3, 'Training Accuracy': 0.9448083679575788}  
'Depth': 4, 'Training Accuracy': 0.945615292472547  
'Depth': 5, 'Training Accuracy': 0.9492238989991129  
'Depth': 6, 'Training Accuracy': 0.9497551660710796  
'Depth': 7, 'Training Accuracy': 0.9520305929642197  
'Depth': 8, 'Training Accuracy': 0.9526671110598778}  
'Depth': 9, 'Training Accuracy': 0.9536695017617016  
'Depth': 10, 'Training Accuracy': 0.9550177172556548}]
```

Based on these results, a depth of $k = 10$ gives the best training accuracy.

b) Here is the test accuracy for a depth of 10:

Test accuracy for depth=10: 0.9506525530763217

This equates to 95.07% accuracy in classifying the test data.

c) Based on the results from the test data, I do not see any overfitting issues. The accuracy of 95.07% with the test data was very close to the 95.50% accuracy with the training data.

The code for this is in the appendix at the end of this file.

3) Independent Events and Bayes Theorem

a) Bayes theorem states:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Event B can occur in two ways: with event A occurring, and without event A occurring. This means that $P(B)$ is equivalent to the following:

$$P(B) = P(B|A)P(A) + P(B|\neg A)P(\neg A)$$

This gives the expression shown in the denominator of the question. So substituting this expanded expression for $P(B)$ in the definition of Bayes theorem gives the proof.

b)

i) Is X independent of Y ? Why or why not?

X and Y are independent if $P(X, Y) = P(X)P(Y)$

First, compute $P(X)$:

$$P(X = 0) = P(0,0,0) + P(0,1,0) + P(0,0,1) + P(0,1,1)$$

$$P(X = 0) = 0.1 + 0.2 + 0.1 + 0.175 = 0.575$$

$$P(X = 1) = P(1,0,0) + P(1,1,0) + P(1,0,1) + P(1,1,1)$$

$$P(X = 1) = 0.05 + 0.1 + 0.1 + 0.175 = 0.425$$

Next, compute $P(Y)$:

$$P(Y = 0) = P(0,0,0) + P(1,0,0) + P(0,0,1) + P(1,0,1)$$

$$P(Y=0) = 0.1 + 0.05 + 0.1 + 0.1 = 0.35$$

$$P(Y = 1) = P(0,1,0) + P(1,1,0) + P(0,1,1) + P(1,1,1)$$

$$P(Y = 1) = 0.2 + 0.1 + 0.175 + 0.175 = 0.65$$

Next, compute $P(X,Y)$:

$$P(X = 0, Y = 0) = P(0,0,0) + P(0,0,1) = 0.1 + 0.1 = 0.2$$

$$P(X = 0, Y = 1) = P(0,1,0) + P(0,1,1) = 0.2 + 0.175 = 0.375$$

$$P(X = 1, Y = 0) = P(1,0,0) + P(1,0,1) = 0.05 + 0.1 = 0.15$$

$$P(X = 1, Y = 1) = P(1,1,0) + P(1,1,1) = 0.1 + 0.175 = 0.275$$

(continued on next page)

Next, compare $P(X)P(Y)$ with $P(X,Y)$:

X	Y	$P(X, Y)$	$P(X)P(Y)$
0	0	0.2	0.20125
0	1	0.375	0.37375
1	0	0.15	0.14875
1	1	0.275	0.27625

Since $P(X,Y) \neq P(X)P(Y)$, X and Y are not independent.

- ii) Is X conditionally independent of Y , given Z ? Why or why not?
 X and Y are conditionally independent, given Z , if:

$$P(X, Y|Z) = P(X|Z)P(Y|Z)$$

This must be true for all values of X , Y , and Z .

Step 1: Calculate $P(Z)$:

$$\begin{aligned} P(Z=0) &= P(0,0,0) + P(0,1,0) + P(1,0,0) + P(1,1,0) \\ &= 0.1 + 0.2 + 0.05 + 0.1 = 0.45 \end{aligned}$$

$$\begin{aligned} P(Z=1) &= P(0,0,1) + P(0,1,1) + P(1,0,1) + P(1,1,1) \\ &= 0.1 + 0.175 + 0.1 + 0.175 = 0.55 \end{aligned}$$

Step 2: Calculate $P(X|Z)$:

For $X = 0$:

$$P(X=0 | Z=0) = (P(0,0,0) + P(0,1,0)) / P(Z=0) = (0.1 + 0.2) / 0.45 = 0.6667$$

$$P(X=0 | Z=1) = (P(0,0,1) + P(0,1,1)) / P(Z=1) = (0.1 + 0.175) / 0.55 = 0.5000$$

For $X = 1$:

$$P(X=1 | Z=0) = (P(1,0,0) + P(1,1,0)) / P(Z=0) = (0.05 + 0.1) / 0.45 = 0.3333$$

$$P(X=1 | Z=1) = (P(1,0,1) + P(1,1,1)) / P(Z=1) = (0.1 + 0.175) / 0.55 = 0.5000$$

Step 3: Calculate $P(Y|Z)$

For $Y = 0$:

$$P(Y=0 | Z=0) = (P(0,0,0) + P(1,0,0)) / P(Z=0) = (0.1 + 0.05) / 0.45 = 0.3333$$

$$P(Y=0 | Z=1) = (P(0,0,1) + P(1,0,1)) / P(Z=1) = (0.1 + 0.1) / 0.55 = 0.3636$$

For $Y = 1$:

$$P(Y=1 | Z=0) = (P(0,1,0) + P(1,1,0)) / P(Z=0) = (0.2 + 0.1) / 0.45 = 0.6667$$

$$P(Y=1 | Z=1) = (P(0,1,1) + P(1,1,1)) / P(Z=1) = (0.175 + 0.175) / 0.55 = 0.6364$$

Step 4: Calculate $P(X,Y | Z)$:

Using the formula:

$$P(X, Y | Z) = P(X, Y, Z) / P(Z)$$

For $Z = 0$:

$$P(0,0 | 0) = P(0,0,0) / P(Z=0) = 0.1 / 0.45 = 0.2222$$

$$P(0,1 | 0) = P(0,1,0) / P(Z=0) = 0.2 / 0.45 = 0.4444$$

$$P(1,0 | 0) = P(1,0,0) / P(Z=0) = 0.05 / 0.45 = 0.1111$$

$$P(1,1 | 0) = P(1,1,0) / P(Z=0) = 0.1 / 0.45 = 0.2222$$

For $Z = 1$:

$$P(0,0 | 1) = P(0,0,1) / P(Z=1) = 0.1 / 0.55 = 0.1818$$

$$P(0,1 | 1) = P(0,1,1) / P(Z=1) = 0.175 / 0.55 = 0.3182$$

$$P(1,0 | 1) = P(1,0,1) / P(Z=1) = 0.1 / 0.55 = 0.1818$$

$$P(1,1 | 1) = P(1,1,1) / P(Z=1) = 0.175 / 0.55 = 0.3182$$

Step 5: Compute $P(X | Z)P(Y | Z)$:

$$P(0 | 0) P(0 | 0) = (0.6667)(0.3333) = 0.2222$$

$$P(0 | 0) P(1 | 0) = (0.6667)(0.6667) = 0.4444$$

$$P(1 | 0) P(0 | 0) = (0.3333)(0.3333) = 0.1111$$

$$P(1 | 0) P(1 | 0) = (0.3333)(0.6667) = 0.2222$$

$$P(0 | 1) P(0 | 1) = (0.5000)(0.3636) = 0.1818$$

$$P(0 | 1) P(1 | 1) = (0.5000)(0.6364) = 0.3182$$

$$P(1 | 1) P(0 | 1) = (0.5000)(0.3636) = 0.1818$$

$$P(1 | 1) P(1 | 1) = (0.5000)(0.6364) = 0.3182$$

Given that the values obtained in Step 4 and Step 5 all match, we have proven that

$$P(X, Y | Z) = P(X | Z)P(Y | Z)$$

Therefore, X and Y are conditionally independent, given Z .

iii) Calculate $P(X \neq Y | Z = 0)$

$$P(X \neq Y | Z = 0) = \frac{P(X \neq Y, Z = 0)}{P(Z = 0)}$$

$$P(X \neq Y, Z = 0) = P(0,1,0) + P(1,0,0) = 0.2 + 0.05 = 0.25$$

$$P(Z = 0) = P(0,0,0) + P(1,0,0) + P(0,1,0) + P(1,1,0)$$

$$P(Z = 0) = 0.1 + 0.05 + 0.2 + 0.1 = 0.45$$

$$P(X \neq Y | Z) = \frac{0.25}{0.45}$$

$$P(X \neq Y | Z) = 0.556$$

Implementing Naïve Bayes

The following results were obtained from the Naïve Bayes classifier:

Training Accuracy: 0.9693
Testing Accuracy: 0.9823
Training Time: 0.1093 seconds

The code for this implementation is in the appendix.

Appendices

The following pages contain the code for sections 1, 2, and 4.

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
```

question2.py

```
1 import numpy as np
2 import pandas as pd
3 import math
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.preprocessing import OneHotEncoder
6 from sklearn.metrics import accuracy_score
7
8 # The CSV data has no headers. These headers will be added after import.
9 column_headers = [
10     "age",
11     "class_of_worker",
12     "detailed_industry_recode",
13     "detailed_occupation_recode",
14     "education",
15     "wage_per_hour",
16     "enroll_in_edu_inst_last_wk",
17     "marital_stat",
18     "major_industry_code",
19     "major_occupation_code",
20     "race",
21     "hispanic_origin",
22     "sex",
23     "member_of_labor_union",
24     "reason_for_unemployment",
25     "full_or_part_time_employment_stat",
26     "capital_gains",
27     "capital_losses",
28     "dividends_from_stocks",
29     "tax_filer_stat",
30     "region_of_previous_residence",
31     "state_of_previous_residence",
32     "detailed_household_and_family_stat",
33     "detailed_household_summary_in_household",
34     "unknown_value",
35     "migration_code_change_in_msa",
36     "migration_code_change_in_reg",
37     "migration_code_move_within_reg",
38     "live_in_this_house_1_year_ago",
39     "migration_prev_res_in_sunbelt",
40     "num_persons_worked_for_employer",
41     "family_members_under_18",
42     "country_of_birth_father",
43     "country_of_birth_mother",
44     "country_of_birth_self",
45     "citizenship",
46     "own_business_or_self_employed",
47     "fill_inc_questionnaire_for_veterans_admin",
```

```
48     "veterans_benefits",
49     "weeks_worked",
50     "year",
51     "income_50k_plus"
52 ]
53
54 # The names of the files to import.
55 datafiles = ["data/census-income.data.csv", "data/census-income.test.csv"]
56
57 # Initialize an array to hold dataframes for the data and test sets.
58 dataframes = []
59
60 # Load the CSV files.
61 for f in datafiles:
62     # Read the file into a dataframe. As noted, there are no headers.
63     df = pd.read_csv(f, header=None)
64     # Add the column headers.
65     df.columns = column_headers
66
67     # Data cleaning operations:
68
69     # Handle missing values using forward fill and backward fill.
70     df = df.ffill().bfill()
71
72     # Upon analyzing the data, the two values in the "income_50k_plus" column, after
73     # import, look like this: [' - 50000.', ' 50000+.'. I am replacing these values using
74     # one-hot encoding so that less than 50k income is a 0, and 50k+ is a 1.
75
76     # Strip whitespace
77     df["income_50k_plus"] = df["income_50k_plus"].str.strip()
78     # Map the existing values to 0 or 1.
79     df["income_50k_plus"] = df["income_50k_plus"].map({" - 50000.": 0, " 50000+": 1})
80
81     # Upon analyzing the data, the the values in the "race" column contain leading
82     # and/or trailing spaces. These will be removed.
83
84     # Strip whitespace
85     df["race"] = df["race"].str.strip()
86
87     # Add this dataframe to the array
88     dataframes.append(df)
89
90 # Assign the data frames as training or test data.
91 training_data = dataframes[0]
92 test_data = dataframes[1]
93
94 # Define features and target variable
95 target_col = "income_50k_plus"
96 feature_cols = [col for col in training_data.columns if col != target_col]
```



```
97
98 # Identify categorical and numerical features
99 categorical_cols = training_data.select_dtypes(include=["object"]).columns.tolist()
100 numerical_cols = training_data.select_dtypes(include=["int64", "float64"]).columns.tolist()
101
102 # Create a one-hot encoder for categorical features.
103 # Note: The income_50k_plus column was encoded earlier, because that code was carried
104 # over from Assignment 0.
105 encoder = OneHotEncoder(handle_unknown="ignore", sparse_output=False)
106
107 # Encode the categorical columns from the test data.
108 encoded_train = encoder.fit_transform(training_data[categorical_cols])
109
110 # Encode the categorical columns from the test data.
111 encoded_test = encoder.transform(test_data[categorical_cols])
112
113 # Convert the encoded columns to dataframes.
114 encoded_train_df = pd.DataFrame(encoded_train, columns=encoder.get_feature_names_out(categorical_cols))
115 encoded_test_df = pd.DataFrame(encoded_test, columns=encoder.get_feature_names_out(categorical_cols))
116
117 # Assemble the prepared training and test data from the numerical features and the one-hot
118 # encoded
119 # categorical features. Remove the target column from numerical features before during this
120 # process.
121 X_train =
122 pd.concat([training_data[numerical_cols].drop(columns=[target_col]).reset_index(drop=True),
123 encoded_train_df], axis=1)
124 X_test =
125 pd.concat([test_data[numerical_cols].drop(columns=[target_col]).reset_index(drop=True),
126 encoded_test_df], axis=1)
127
128 # Error checking: The code will exit if the target column is in the training or test data.
129 assert target_col not in X_train.columns, "Target column should not be in X_train!"
130 assert target_col not in X_test.columns, "Target column should not be in X_test!"
131
132 # Get the target column from the original training and test datasets.
133 y_train = training_data[target_col]
134 y_test = test_data[target_col]
135
136 # Part A: Train decision trees for depths of 2 - 10 and store accuracies.
137 depth_accuracies = []
138 for depth in range(2, 11):
139     # The random_state value is arbitrary, as long as the same value is used
140     # each time to allow reproducibility.
141     model = DecisionTreeClassifier(max_depth=depth, random_state=42)
142     model.fit(X_train, y_train)
143     train_acc = accuracy_score(y_train, model.predict(X_train))
144     depth_accuracies.append({"Depth": depth, "Training Accuracy": train_acc})
```

```
139
140 print(depth_accuracies)
141 # Find the maximum accuracy and its associated depth.
142 max_accuracy = max([entry["Training Accuracy"] for entry in depth_accuracies])
143 # The optimal depth is the depth associated with max_accuracy.
144 optimal_depth = next(item["Depth"] for item in depth_accuracies if item["Training Accuracy"]
    == max_accuracy)
145 print(f"max_accuracy = {max_accuracy}")
146 print(f"optimal_depth = {optimal_depth}")
147
148 # Part B: use the optimal depth found above to classify the test data.
149 # Re-train the model using the optimal depth found above.
150 model = DecisionTreeClassifier(max_depth=optimal_depth, random_state=42)
151 model.fit(X_train, y_train)
152 # Get the accuracy score for the training data (should be the same as before).
153 train_acc = accuracy_score(y_train, model.predict(X_train))
154 # Now get the accuracy score for the test data.
155 test_acc = accuracy_score(y_test, model.predict(X_test))
156 print(f"Training accuracy for depth=10: {train_acc}")
157 print(f"Test accuracy for depth=10: {test_acc}")
158
159
```

question4.py

```
1 import numpy as np
2 import pandas as pd
3 import time
4 from sklearn.naive_bayes import MultinomialNB
5 from sklearn.metrics import accuracy_score
6
7 # Load training data. The shape is (4527, 5180).
8 X_train = pd.read_csv("nbdata/train.csv", header=None)
9 # Load training target data. The shape is (4527).
10 y_train = pd.read_csv("nbdata/train_labels.txt", header=None).values.ravel()
11
12 # Load test data. The shape is (1806, 5180).
13 X_test = pd.read_csv("nbdata/test.csv", header=None) # Shape: (1806, 5180)
14 # Load test target data. The shape is (1806).
15 y_test = pd.read_csv("nbdata/test_labels.txt", header=None).values.ravel() # Shape: (1806,)
16
17 # Instantiate a Multinomial Naive Bayes Classifier
18 nbc = MultinomialNB()
19
20 # Record the starting time.
21 start_time = time.time()
22 # Train the classifier.
23 nbc.fit(X_train, y_train)
24 # When complete, calculate the training time.
25 training_time = time.time() - start_time
26
27 # Get the predictions for the training and test data.
28 y_train_pred = nbc.predict(X_train)
29 y_test_pred = nbc.predict(X_test)
30
31 # Calculate accuracies.
32 train_accuracy = accuracy_score(y_train, y_train_pred)
33 test_accuracy = accuracy_score(y_test, y_test_pred)
34
35 # Print results
36 print(f"Training Accuracy: {train_accuracy:.4f}")
37 print(f"Testing Accuracy: {test_accuracy:.4f}")
38 print(f"Training Time: {training_time:.4f} seconds")
39
40
```