

University of Waterloo

CS 486, Winter 2024

Assignment 1

Problem 1

a)

```
2 from math import log2
3 from queue import PriorityQueue
4 import matplotlib.pyplot as plt
5 from multiprocessing import Pool

7 ##### Variables #####
8 NUM_FEATURES = 0 # Define later
9 NUM_SAMPLE = 1500
10 ATHEISM_ID = 1
11 BOOKS_ID = 2
12 LABEL_STR = {ATHEISM_ID: "Atheism", BOOKS_ID: "Books" }

14 # We define  $P(x)$  as # belongs to atheism over total.

16 ##### READING FILE #####

18 # Function to read label data, returning a dictionary mapping document ID to label
19 def read_label_data(file_name):
20     label_data = {}
21     with open(file_name, 'r') as file:
22         lineNum = 1
23         for line in file:
24             label_data[lineNum] = int(line.strip())
25             lineNum += 1
26             # docId = line number
27     file.close()
28     return label_data

30 # Function to read word data, returning a dictionary mapping word ID to word
31 def read_word_data(file_name):
32     word_data = {}
33     with open(file_name, 'r') as file:
34         lineNum = 1
35         for line in file:
36             word_data[int(lineNum)] = line.strip()
37             lineNum += 1
38             # wordId = line number
39     # Define the Number of features
40     global NUM_FEATURES
41     NUM_FEATURES = lineNum
42     file.close()
43     return word_data

45 # Function to read train or test data
46 # document_data[n] = array of word_id that n+1'th doc have.
47 def read_document_data(file_name):
```

```

48     document_data = {i: [] for i in range(1, NUM_SAMPLE + 1)} # 1 to 1500
49     with open(file_name, 'r') as file:
50         for line in file:
51             doc_id, word_id = line.strip().split()
52             document_data[int(doc_id)].append(int(word_id))
53     file.close()
54     return document_data

56 ##### COMPUTING #####

58 # functions that does log_2(x/x+y).
59 # the # of bits to encode x.
60 def entropy(x, y):
61     #print(f"x = {x}, y = {y}")

63     # corner case, for our purpose log_2(0) = 0.
64     if x == 0:
65         return 0

67     val = x / (x + y)
68     #print(f"val = {val}")
69     #print(f"log2(val) = {log2(val)}")
70     return - log2(val)

72 def information_content(dataset, doc_subreddit_dict):
73     # calculate the information content of the dataset
74     # dataset needs to be a dictionary

76     atheism_count = 0
77     books_count = 0

79     doc_ids = [key for key, _ in dataset.items()]

81     # Loop over all elements and calculate # belongs to atheism or books respectively
82     for doc_id in doc_ids:
83         if doc_subreddit_dict[doc_id] == ATHEISM_ID:
84             atheism_count += 1
85         elif doc_subreddit_dict[doc_id] == BOOKS_ID:
86             books_count += 1

89     etp1 = entropy(atheism_count, books_count)
90     etp2 = entropy(books_count, atheism_count)

92     # Corner case
93     if (atheism_count + books_count) == 0:
94         return 0

96     return ((atheism_count) / (atheism_count + books_count) * etp1) + ((books_count) / (
        atheism_count + books_count) * etp2)

98 def information_gain(num1, num2):
99     etp1 = entropy(num1, num2)
100    etp2 = entropy(num2, num1)

102    # Corner case
103    if (num1 + num2) == 0:

```

```

104         return 0

106     return ((num1) / (num1 + num2) * etp1) + ((num2) / (num1 + num2) * etp2)

108 # Function to calculate the delta information gain for a given split
109 def delta_information_gain(elements, doc_subreddit_dict, word_to_split, method, debug=False):
110     :

111     atheism_count_E = 0
112     books_count_E = 0

114     doc_ids = [key for key, _ in elements.items()]

116     # Loop over all elements and calculate # belongs to atheism or books respectively
117     for doc_id in doc_ids:
118         if doc_subreddit_dict[doc_id] == ATHEISM_ID:
119             atheism_count_E += 1
120         elif doc_subreddit_dict[doc_id] == BOOKS_ID:
121             books_count_E += 1

123     # TODO: not sure if for IE we going to use different methods.
124     IE = information_gain(atheism_count_E, books_count_E)
125     if debug:
126         print(f"IE={IE}, atheism_count_E={atheism_count_E}, books_count_E={books_count_E}")
127     # We now proceed to split the elements by the word_to_split.

129     # E1
130     has_word_to_split = [key for key, value in elements.items()
131                          if word_to_split in value]
132     # E2
133     not_have_word_to_split = [key for key, value in elements.items()
134                              if word_to_split not in value]

136     #if debug:
137     # print(f"len has_word_to_split = {len(has_word_to_split)}, len not_have_word_to_split = {len(not_have_word_to_split)}")

139     # note the values in the arrays are doc_id has or has not the word to split.
140     atheism_count_E1 = 0
141     books_count_E1 = 0

143     # Loop over all elements and calculate # belongs to atheism or books respectively
144     for doc_id in has_word_to_split:
145         if doc_subreddit_dict[doc_id] == ATHEISM_ID:
146             atheism_count_E1 += 1
147         elif doc_subreddit_dict[doc_id] == BOOKS_ID:
148             books_count_E1 += 1

150     # TODO: not sure if for IE we going to use different methods.
151     IE1 = information_gain(atheism_count_E1, books_count_E1)
152     if debug:
153         print(f"IE1={IE1}, atheism_count_E1={atheism_count_E1}, books_count_E1={books_count_E1}")

155     atheism_count_E2 = 0

```

```

156     books_count_E2 = 0

158     # Loop over all elements and calculate # belongs to atheism or books respectively
159     for doc_id in not_have_word_to_split:
160         if doc_subreddit_dict[doc_id] == ATHEISM_ID:
161             atheism_count_E2 += 1
162         elif doc_subreddit_dict[doc_id] == BOOKS_ID:
163             books_count_E2 += 1

165     # TODO: not sure if for IE we going to use different methods.
166     IE2 = information_gain(atheism_count_E2, books_count_E2)
167     if debug:
168         print(f"IE2={IE2}, atheism_count_E2={atheism_count_E2}, books_count_E2={books_count_E2}")

170     # Finally the calculation
171     if method == 1:
172         # Method 1: Average information gain across the leaves
173         return IE - ((IE1 / 2) + (IE2 / 2))
174     elif method == 2:
175         # Method 2: The one discussed in Class

177         sum_E1 = atheism_count_E1 + books_count_E1
178         sum_E2 = atheism_count_E2 + books_count_E2

180         # Corner case
181         if sum_E1 == 0 and sum_E2 == 0:
182             return IE

184         return IE - ((sum_E1 / (sum_E1 + sum_E2)) * IE1) - ((sum_E2 / (sum_E1 + sum_E2)) * IE2)

186     ##### Decision Tree #####

188     class Node:
189         def __init__(self, dataset, point_estimate, feature_to_split=None, info_gain=0, splitted_feature=[]):
190             self.dataset = dataset # The subset 'E' of the dataset at this node
191             self.feature_to_split = feature_to_split # 'X_prime': The feature to split on at this node
192             self.point_estimate = point_estimate # Point estimation of current node.
193             self.splitted_feature = splitted_feature[:] # array of already splitted feature
194             self.info_gain = info_gain # 'delta_I': Information gain of the split
195             self.left = None # Left child (with feature)
196             self.right = None # Right child (without feature)

198         # This is needed for the PriorityQueue to compare Nodes based on information gain
199         def __lt__(self, other):
200             # since priority queue is a min_heap, so we define lt as gt
201             return self.info_gain > other.info_gain

204     def calculate_point_estimate(dataset, label_dict):
205         # Calculate the dominant label (subreddit) in the dataset
206         label_counts = {ATHEISM_ID: 0, BOOKS_ID: 0}
207         for doc_id in dataset:
208             label = label_dict[doc_id]

```

```

209         label_counts[label] += 1
210         # The point estimate is the label with the highest count
211         point_estimate = max(label_counts, key=label_counts.get)
212         return point_estimate

214 def find_best_to_split(dataset, label_dict, method, words=None, splitted_feature=[]):
215     # Find the best feature to split next
216     # by try to compute the delta_information_gain on all features
217     # appeared in the dataset but not the ones in splitted_feature, and return a tuple containing
218     # the feature to split that will give the biggest delta_information_gain
219     # and the delta_information_gain.
220     best_feature = None
221     best_info_gain = 0 # Start with 0 to ensure any gain is better, but not no split

223     # Iterate through each feature in the dataset to find the best one to split on
224     for feature in range(1, NUM_FEATURES + 1):

226         # skips the feature we already split on
227         if feature in splitted_feature:
228             continue

230         # Compute the information gain for splitting on the current feature
231         current_info_gain = delta_information_gain(dataset, label_dict, feature, method)

233         # If the information gain of the current feature is better than the best one so far, update the_
234         best feature and gain
235         if current_info_gain > best_info_gain:
236             #print(f"better word = {words[feature]}, info gain = {current_info_gain} ")
237             #print(f"delta I = {delta_information_gain(dataset, label_dict, feature, method, True)}")
238             #print("-----")
239             best_feature = feature
240             best_info_gain = current_info_gain

241     #print(f"delta I = {delta_information_gain(dataset, label_dict, best_feature, method, True)}")
242     return best_feature, best_info_gain

244 def split_dataset(dataset, feature_to_split):
245     # Datasets to hold the split
246     dataset_with_feature = {}
247     dataset_without_feature = {}
248     #print(f"Splitting on {feature_to_split}")
249     # Iterate over each entry in the dataset
250     for doc_id, word_ids in dataset.items():
251         # Check if the feature to split on is in the document's word IDs
252         if feature_to_split in word_ids:
253             # Add this document to the dataset with the feature
254             #print(f"with_feature: doc_id = {doc_id}, features = {word_ids}")
255             dataset_with_feature[doc_id] = word_ids[:] # [:] to make a shallow copy

257             #print(f"remove feature: {feature_to_split}")
258             # remove the splitting feature
259             #dataset_with_feature[doc_id].remove(feature_to_split)
260         else:
261             # Add this document to the dataset without the feature
262             #print(f"without_feature: doc_id = {doc_id}, features = {word_ids}")
263             dataset_without_feature[doc_id] = word_ids

```

```

265     #print(f"L size: {str(len(dataset_with_feature))}, R size: {str(len(dataset_without_feature))}")
266     return dataset_with_feature, dataset_without_feature

268 # Function to build the decision tree
269 def build_decision_tree(train_data, train_labels, method, subreddit_dict, max_nodes=100, ↵
    words=None):
270     MAX_NODES = max_nodes
271     pq = PriorityQueue()
272     root = Node(dataset=train_data, point_estimate=calculate_point_estimate(train_data, ↵
        subreddit_dict))

274     best_feature, best_info_gain = find_best_to_split(root.dataset, train_labels, method, ↵
        words=words)

276     # Update root node with split info
277     root.feature_to_split = best_feature
278     root.info_gain = best_info_gain

280     pq.put(root)

282     while not pq.empty() and max_nodes > 0:

284         current_node = pq.get()

286         # We only counts the number of internal nodes with max_nodes
287         max_nodes -= 1
288         #print(f"\ndoing the {MAX_NODES - max_nodes}'th node")
289         #print(f"info gained = {info_gained}, split word = {words[current_node.feature_to_split]}")

291         # Split the dataset based on the best feature
292         current_node.splitted_feature.append(current_node.feature_to_split) # add splitted ↵
        feature

294         #print(f"splitted feature list is {current_node.splitted_feature}")
295         left_dataset, right_dataset = split_dataset(current_node.dataset, current_node.↵
        feature_to_split)

297         best_feature_L, best_info_gain_L = find_best_to_split(left_dataset, train_labels, ↵
        method, words=words, splitted_feature=current_node.splitted_feature)
298         best_feature_R, best_info_gain_R = find_best_to_split(right_dataset, train_labels, ↵
        method, words=words, splitted_feature=current_node.splitted_feature)

301         left_node = Node(dataset=left_dataset, point_estimate=calculate_point_estimate(↵
        left_dataset, subreddit_dict), info_gain=best_info_gain_L, feature_to_split=↵
        best_feature_L, splitted_feature=current_node.splitted_feature)
302         right_node = Node(dataset=right_dataset, point_estimate=calculate_point_estimate(↵
        right_dataset, subreddit_dict), info_gain=best_info_gain_R, feature_to_split=↵
        best_feature_R, splitted_feature=current_node.splitted_feature)

304         # Update current node with split info
305         current_node.left = left_node
306         current_node.right = right_node

308         # Add child nodes to the priority queue
309         if best_feature_L is not None:

```

```

310         pq.put(left_node)
311         #print(f"Adde L feature = {words[best_feature_L]}, info gain = {best_info_gain_L}")
312     if best_feature_R is not None:
313         pq.put(right_node)
314         #print(f"Added R feature = {words[best_feature_R]}, info gain = {best_info_gain_R}")

317     return root

319 def print_tree(node, depth=0, feature_names=None):
320     # Base case: if the node is a leaf, it will not have a child
321     if node.left is None or node.right is None:
322         print("-" * depth + "Leaf, estimate: " + LABEL_STR[node.point_estimate])
323         return

324     # Recursive case: print the current node's split information
325     if feature_names and node.feature_to_split in feature_names:
326         feature_name = feature_names[node.feature_to_split]
327     else:
328         feature_name = str(node.feature_to_split)

329     print("-" * depth + f"Node: Split Feature={feature_name}, Info Gain={node.info_gain:.10f}")

330     # Recursively print the left subtree
331     print("-" * depth + "L(w/ feature):")
332     print_tree(node.left, depth + 1, feature_names)

333     # Recursively print the right subtree
334     print("-" * depth + "R(wo/ feature):")
335     print_tree(node.right, depth + 1, feature_names)

342 ##### TESTING #####

344 # Use decision tree to predict the label for a single document
345 def predict_label(node, document_word_array):
346     # If we have reached a leaf node, return its point estimate
347     if node.left is None and node.right is None:
348         return node.point_estimate
349     # If the document contains the word_id at the current node, go left
350     elif node.feature_to_split in document_word_array:
351         return predict_label(node.left, document_word_array)
352     # If the document does not contain the word_id, go right
353     else:
354         return predict_label(node.right, document_word_array)

356 # Function to calculate the accuracy of the decision tree
357 def calculate_accuracy(tree, data, labels):
358     correct_predictions = 0
359     # Iterate over all documents in the test data
360     for doc_id, document_word_array in data.items():
361         # Use the tree to predict the label for the current document
362         predicted_label = predict_label(tree, document_word_array)
363         # If the predicted label matches the actual label, increment the correct predictions count
364         if predicted_label == labels[doc_id]:
365             correct_predictions += 1
366     # Calculate the percentage of correctly classified samples

```

```

367     accuracy = (correct_predictions / len(data)) * 100
368     return accuracy

371 ##### CONCURRENCY #####

373 # Function to build a tree and calculate accuracies for a given number of nodes
374 def compute_accuracies_for_nodes(max_nodes):
375     words = read_word_data('./words.txt')
376     train_data = read_document_data('./trainData.txt')
377     train_labels = read_label_data('./trainLabel.txt')
378     test_data = read_document_data('./testData.txt')
379     test_labels = read_label_data('./testLabel.txt')

381     # Build trees using both methods
382     tree1 = build_decision_tree(train_data, train_labels, method=1, subreddit_dict=↵
        train_labels, words=words, max_nodes=max_nodes)
383     tree2 = build_decision_tree(train_data, train_labels, method=2, subreddit_dict=↵
        train_labels, words=words, max_nodes=max_nodes)

385     # Calculate accuracies for both trees
386     accuracy_train_tree1 = calculate_accuracy(tree1, train_data, train_labels)
387     accuracy_test_tree1 = calculate_accuracy(tree1, test_data, test_labels)
388     accuracy_train_tree2 = calculate_accuracy(tree2, train_data, train_labels)
389     accuracy_test_tree2 = calculate_accuracy(tree2, test_data, test_labels)

391     # Return a tuple of accuracies
392     print(f"Done_max_node={str(max_nodes)}".)
393     return (max_nodes, accuracy_train_tree1, accuracy_test_tree1, accuracy_train_tree2, ↵
        accuracy_test_tree2)

396 ##### MAIN #####

399 if __name__ == '__main__':

401     # Reading data from files
402     words = read_word_data('./words.txt')
403     train_data = read_document_data('./trainData.txt')
404     train_labels = read_label_data('./trainLabel.txt')
405     test_data = read_document_data('./testData.txt')
406     test_labels = read_label_data('./testLabel.txt')

408     ##### b) #####

410     print("---_building_tree_1_---\n")
411     tree1 = build_decision_tree(train_data, train_labels, method=1, subreddit_dict=↵
        train_labels, words=words, max_nodes=10)
412     print("\n---_method_1_tree_---\n")
413     print_tree(tree1, feature_names=words)

415     print("\n---_building_tree_2_---\n")
416     tree2 = build_decision_tree(train_data, train_labels, method=2, subreddit_dict=↵
        train_labels, words=words, max_nodes=10)
417     print("\n---_method_2_tree_---\n")
418     print_tree(tree2, feature_names=words)

420     # Validate the decision trees tree1 and tree2

```



```

421 accuracy_tree1 = calculate_accuracy(tree1, test_data, test_labels)
422 accuracy_tree2 = calculate_accuracy(tree2, test_data, test_labels)

424 # Print the accuracies
425 print(f"\nAccuracy of tree1 (Method 1): {accuracy_tree1:.5f}%\n")
426 print(f"\nAccuracy of tree2 (Method 2): {accuracy_tree2:.5f}%\n")

428 ##### C) #####

431 # Number of processes to use
432 num_processes = 100 # For school server, laptop will explode with this

434 # Create a pool of processes
435 with Pool(processes=num_processes) as pool:
436     # Map the compute_accuracies_for_nodes function to each number of max_nodes
437     results = pool.map(compute_accuracies_for_nodes, range(1, 101))

439 # Now we will process the results to plot them
440 # Initialize lists to hold accuracies for plotting
441 nodes = []
442 accuracy_train_method1 = []
443 accuracy_test_method1 = []
444 accuracy_train_method2 = []
445 accuracy_test_method2 = []

447 # Populate the lists with data
448 i = 1
449 for result in results:
450     print(f"{str(i)} node, train_m1={result[1]:.5f}, test_m1={result[2]:.5f}, train_m2={result[3]:.5f}, test_m2={result[4]:.5f}")
451     nodes.append(result[0])
452     accuracy_train_method1.append(result[1])
453     accuracy_test_method1.append(result[2])
454     accuracy_train_method2.append(result[3])
455     accuracy_test_method2.append(result[4])
456     i += 1

458 # Plotting the results
459 plt.figure(figsize=(14, 7))
460 # Plot for method 1
461 plt.subplot(1, 2, 1)
462 plt.scatter(nodes, accuracy_train_method1, label='Training Accuracy (Method 1)')
463 plt.scatter(nodes, accuracy_test_method1, label='Testing Accuracy (Method 1)')
464 plt.xlabel('Number of Nodes')
465 plt.ylabel('Accuracy (%)')
466 plt.title('Training and Testing Accuracies (Method 1)')
467 plt.legend()
468 # Plot for method 2
469 plt.subplot(1, 2, 2)
470 plt.scatter(nodes, accuracy_train_method2, label='Training Accuracy (Method 2)')
471 plt.scatter(nodes, accuracy_test_method2, label='Testing Accuracy (Method 2)')
472 plt.xlabel('Number of Nodes')
473 plt.ylabel('Accuracy (%)')
474 plt.title('Training and Testing Accuracies (Method 2)')
475 plt.legend()
476 plt.savefig('method_accuracies.png')

```

b) --- method 1 tree ---

```
3 Node: Split Feature = christian, Info Gain = 0.5002058812
4 L (w/ feature):
5 -Leaf, estimate: Atheism
6 R (wo/ feature):
7 -Node: Split Feature = atheism, Info Gain = 0.5001930379
8 -L (w/ feature):
9 --Leaf, estimate: Atheism
10 -R (wo/ feature):
11 --Node: Split Feature = christians, Info Gain = 0.4998761311
12 --L (w/ feature):
13 ---Leaf, estimate: Atheism
14 --R (wo/ feature):
15 ---Node: Split Feature = beliefs, Info Gain = 0.4995539340
16 ---L (w/ feature):
17 ----Leaf, estimate: Atheism
18 ---R (wo/ feature):
19 ----Node: Split Feature = atheists, Info Gain = 0.4987661687
20 ----L (w/ feature):
21 -----Leaf, estimate: Atheism
22 ----R (wo/ feature):
23 -----Node: Split Feature = brain, Info Gain = 0.4982474614
24 -----L (w/ feature):
25 -----Leaf, estimate: Atheism
26 -----R (wo/ feature):
27 -----Node: Split Feature = aa, Info Gain = 0.4976745555
28 -----L (w/ feature):
29 -----Leaf, estimate: Atheism
30 -----R (wo/ feature):
31 -----Node: Split Feature = murder, Info Gain = 0.4973448755
32 -----L (w/ feature):
33 -----Leaf, estimate: Atheism
34 -----R (wo/ feature):
35 -----Node: Split Feature = proof, Info Gain = 0.4969295610
36 -----L (w/ feature):
37 -----Leaf, estimate: Atheism
38 -----R (wo/ feature):
39 -----Node: Split Feature = logic, Info Gain = 0.4966008344
40 -----L (w/ feature):
41 -----Leaf, estimate: Atheism
42 -----R (wo/ feature):
43 -----Leaf, estimate: Books
```

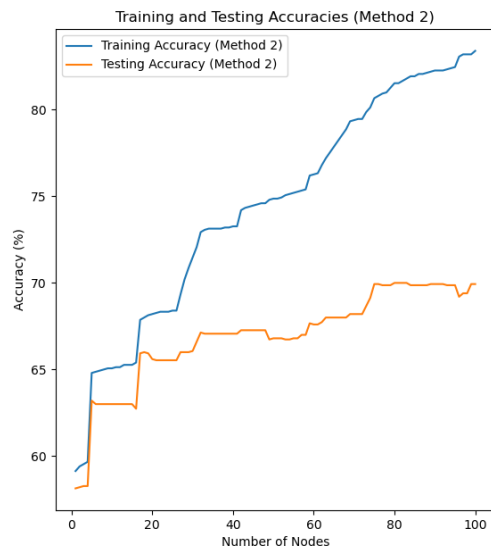
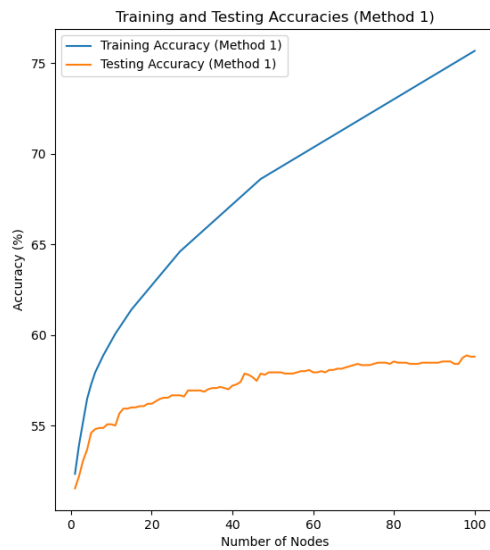
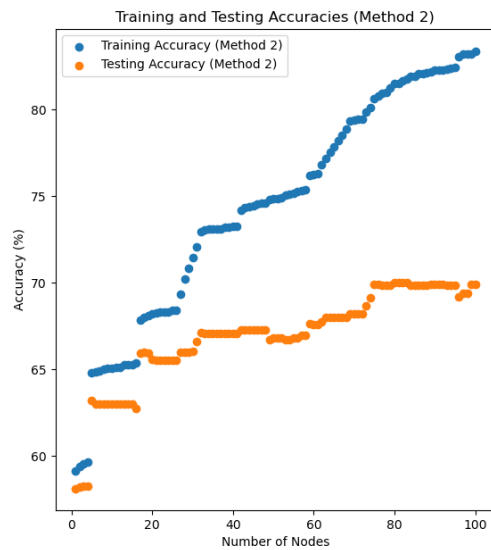
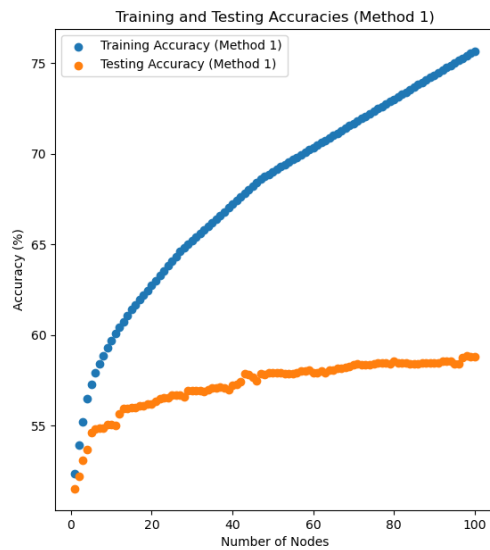
```

1      --- method 2 tree ---

3  Node: Split Feature = book, Info Gain = 0.0770187045
4  L (w/ feature):
5  -Node: Split Feature = bible, Info Gain = 0.1153579525
6  -L (w/ feature):
7  --Leaf, estimate: Atheism
8  -R (wo/ feature):
9  --Node: Split Feature = call, Info Gain = 0.0694586299
10 --L (w/ feature):
11 ---Leaf, estimate: Atheism
12 --R (wo/ feature):
13 ---Node: Split Feature = sent, Info Gain = 0.0842240762
14 ---L (w/ feature):
15 ----Leaf, estimate: Atheism
16 ----R (wo/ feature):
17 ----Node: Split Feature = controlling, Info Gain = 0.0587581571
18 ----L (w/ feature):
19 -----Leaf, estimate: Atheism
20 -----R (wo/ feature):
21 -----Leaf, estimate: Books
22 R (wo/ feature):
23 -Node: Split Feature = books, Info Gain = 0.0599261848
24 -L (w/ feature):
25 --Node: Split Feature = sure, Info Gain = 0.0976949605
26 --L (w/ feature):
27 ---Node: Split Feature = soon, Info Gain = 0.9182958341
28 ---L (w/ feature):
29 ----Leaf, estimate: Books
30 ---R (wo/ feature):
31 ----Leaf, estimate: Atheism
32 --R (wo/ feature):
33 ---Node: Split Feature = spirit, Info Gain = 0.0969446061
34 ---L (w/ feature):
35 ----Leaf, estimate: Atheism
36 ---R (wo/ feature):
37 ----Leaf, estimate: Books
38 -R (wo/ feature):
39 --Node: Split Feature = religion, Info Gain = 0.0356734070
40 --L (w/ feature):
41 ---Leaf, estimate: Atheism
42 --R (wo/ feature):
43 ---Leaf, estimate: Atheism

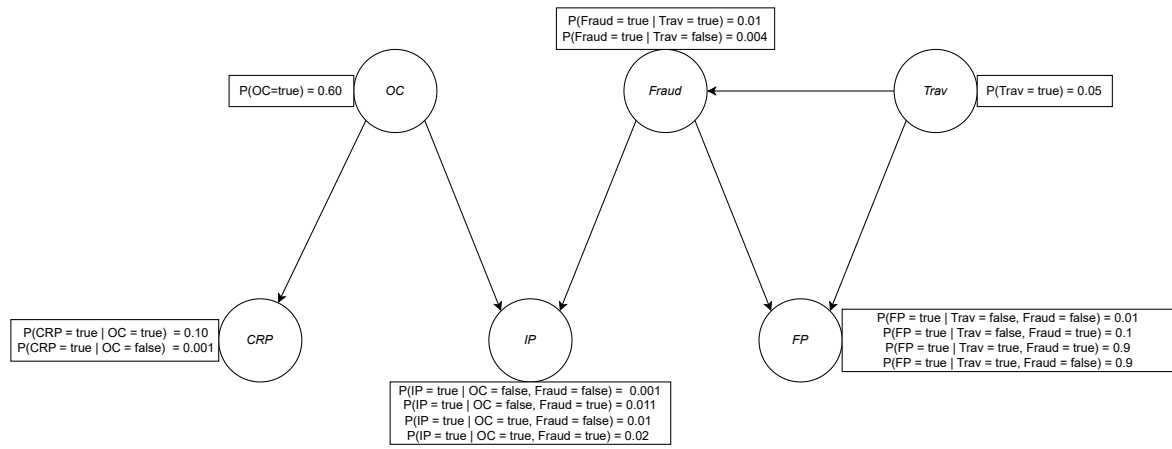
```

c) Traing & Testing Accuracies:



Problem 2

a) Bayes Network:



b) The prior probability that the current transaction is a fraud is:

$$P(Fraud = true) = \sum_{Trav} P(Fraud = true|Trav)P(Trav) \\ = \sum_{Trav} f_1(Trav)f_0(Trav)$$

where $f_0(Trav) = P(Trav)$

| $Trav$ | $f_0(Trav)$ |
|--------|-------------|
| t | 0.05 |
| f | 0.95 |

where $f_1(Trav) = P(Fraud = true|Trav)$

| $Trav$ | $P(Fraud = true)$ |
|--------|-------------------|
| t | 0.01 |
| f | 0.004 |

Thus we get $P(Fraud = true) \propto \sum_{Trav} f_1(Trav)f_0(Trav)$
 $= 0.05 \times 0.01 + 0.95 \times 0.004 = 0.0043$

$f_0(Trav) = P(Trav) =$

| $Trav$ | $f_0(Trav)$ |
|--------|-------------|
| t | 0.05 |
| f | 0.95 |

$f_1(Fraud, Trav) = P(Fraud|Trav) =$

| $Fraud$ | $Trav$ | $f_1(Fraud Trav)$ |
|---------|--------|---------------------|
| t | t | 0.01 |
| t | f | 0.004 |
| f | t | $1 - 0.01 = 0.99$ |
| f | f | $1 - 0.004 = 0.996$ |

$f_2(OC) = P(OC) =$

| OC | $f_2(OC)$ |
|------|-----------|
| t | 0.6 |
| f | 0.4 |

$f_3(Fraud, Trav) = P(FP = true|Fraud, Trav) =$

| $Fraud$ | $Trav$ | $f_3(Fraud, Trav)$ |
|---------|--------|--------------------|
| t | t | 0.9 |
| t | f | 0.1 |
| f | t | 0.9 |
| f | f | 0.01 |

$f_4(OC) = P(CRP = true|OC) =$

| OC | $f_4(OC)$ |
|------|-----------|
| t | 0.1 |
| f | 0.001 |

$f_5(Fraud, OC) = P(IP = false|Fraud, OC) =$

| $Fraud$ | OC | $f_5(Fraud, OC)$ |
|---------|------|---------------------|
| t | t | $1 - 0.02 = 0.98$ |
| t | f | $1 - 0.011 = 0.989$ |
| f | t | $1 - 0.01 = 0.99$ |
| f | f | $1 - 0.001 = 0.999$ |

So that $P(Fraud|FP = true, IP = false, CRP = true)$

$$\propto \sum_{OC} \sum_{Trav} f_0(Trav)f_1(Fraud, Trav)f_2(OC)f_3(Fraud, Trav)f_4(OC)f_5(Fraud, OC) \\ = \sum_{OC} f_2(OC)f_4(OC)f_5(Fraud, OC) \sum_{Trav} f_0(Trav)f_1(Fraud, Trav)f_3(Fraud, Trav) \\ = \sum_{OC} f_2(OC)f_4(OC)f_5(Fraud, OC)f_6(Fraud)$$

Where $f_6(Fraud) = \sum_{Trav} f_0(Trav)f_1(Fraud, Trav)f_3(Fraud, Trav)$

$$=$$

| <i>Fraud</i> | $f_6(Fraud)$ |
|--------------|--|
| <i>t</i> | $0.05 \times 0.01 \times 0.9 + 0.95 \times 0.004 \times 0.1 = 0.00083$ |
| <i>f</i> | $0.05 \times 0.99 \times 0.9 + 0.95 \times 0.996 \times 0.01 = 0.054012$ |

Then, we get $\sum_{OC} f_2(OC)f_4(OC)f_5(Fraud, OC)f_6(Fraud) = f_7(Fraud)f_6(fraud)$

Where $f_7(Fraud) = \sum_{OC} f_2(OC)f_4(OC)f_5(Fraud, OC)$

$$=$$

| <i>Fraud</i> | $f_7(Fraud)$ |
|--------------|---|
| <i>t</i> | $0.6 \times 0.1 \times 0.98 + 0.4 \times 0.001 \times 0.989 = 0.0591956$ |
| <i>f</i> | $0.6 \times 0.1 \times 0.99 + 0.4 \times 0.001 \times 0.999 = 0.059799599999999994$ |

Thus, we get $f_8(Fraud) = f_7(Fraud)f_6(fraud)$

$$=$$

| <i>Fraud</i> | $f_8(Fraud)$ |
|--------------|--|
| <i>t</i> | $0.00083 \times 0.0591956 = 0.000049132348$ |
| <i>f</i> | $0.054012 \times 0.059799599999999994 = 0.0032298959951999997$ |

Thus, we get $P(Fraud = true|FP = true, IP = false, CRP = true)$ as

$$\frac{0.000049132348}{0.000049132348 + 0.0032298959951999997} = 0.014983813147541082 \approx 0.01498$$

c) $f_0(Fraud) = P(Fraud|Trav = true) =$

| <i>Fraud</i> | $f_0(Fraud)$ |
|--------------|---------------------|
| <i>t</i> | 0.01 |
| <i>f</i> | $(1 - 0.01) = 0.99$ |

$f_1(OC) = P(OC) =$

| <i>OC</i> | $f_1(OC)$ |
|-----------|-----------|
| <i>t</i> | 0.6 |
| <i>f</i> | 0.4 |

$f_2(Fraud) = P(FP = true|Fraud, Trav = true) =$

| <i>Fraud</i> | $f_2(Fraud)$ |
|--------------|--------------|
| <i>t</i> | 0.9 |
| <i>f</i> | 0.9 |

$f_3(OC) = P(CRP = true|OC) =$

| <i>OC</i> | $f_3(OC)$ |
|-----------|-----------|
| <i>t</i> | 0.1 |
| <i>f</i> | 0.001 |

$f_5(Fraud, OC) = P(IP = false|Fraud, OC) =$

| <i>Fraud</i> | <i>OC</i> | $f_5(Fraud, OC)$ |
|--------------|-----------|---------------------|
| <i>t</i> | <i>t</i> | $1 - 0.02 = 0.98$ |
| <i>t</i> | <i>f</i> | $1 - 0.011 = 0.989$ |
| <i>f</i> | <i>t</i> | $1 - 0.01 = 0.99$ |
| <i>f</i> | <i>f</i> | $1 - 0.001 = 0.999$ |

So, $P(Fraud|FP = true, IP = false, CRP = true, Trav = true)$
 $= \sum_{OC} f_0(Fraud)f_1(OC)f_2(Fraud)f_3(OC)f_5(Fraud, OC)$
 $= f_0(Fraud)f_2(Fraud)\sum_{OC} f_1(OC)f_3(OC)f_5(Fraud, OC)$
 $= f_0(Fraud)f_2(Fraud)f_6(Fraud)$

Where $f_6(Fraud) = \sum_{OC} f_1(OC)f_3(OC)f_5(Fraud, OC)$

$=$

| <i>Fraud</i> | $f_6(Fraud)$ |
|--------------|--|
| <i>t</i> | $0.6 \times 0.1 \times 0.98 + 0.4 \times 0.001 \times 0.989 = 0.0591956$ |
| <i>f</i> | $0.6 \times 0.1 \times 0.99 + 0.4 \times 0.001 \times 0.999 = 0.0597996$ |

Thus, we get $f_7(Fraud) = f_0(Fraud)f_2(Fraud)f_6(Fraud)$

$=$

| <i>Fraud</i> | $f_7(Fraud)$ |
|--------------|---|
| <i>t</i> | $0.01 \times 0.9 \times 0.0591956 = 0.0005327604$ |
| <i>f</i> | $0.99 \times 0.9 \times 0.0597996 = 0.0532814436$ |

Thus, we get $P(Fraud = true|FP = true, IP = false, CRP = true, Trav = true)$

$= \frac{0.0005327604}{0.0005327604 + 0.0532814436} = 0.00989999591929298 \approx 0.0099$

d) We will need to calculate $P(Fraud|IP = true)$.

$$f_0(Trav) = P(Trav) =$$

| <i>Trav</i> | $f_0(Trav)$ |
|-------------|-------------|
| <i>t</i> | 0.05 |
| <i>f</i> | 0.95 |

$$f_1(Fraud, Trav) = P(Fraud|Trav) =$$

| <i>Fraud</i> | <i>Trav</i> | $f_1(Fraud Trav)$ |
|--------------|-------------|---------------------|
| <i>t</i> | <i>t</i> | 0.01 |
| <i>t</i> | <i>f</i> | 0.004 |
| <i>f</i> | <i>t</i> | $1 - 0.01 = 0.99$ |
| <i>f</i> | <i>f</i> | $1 - 0.004 = 0.996$ |

$$f_2(OC) = P(OC) =$$

| <i>OC</i> | $f_2(OC)$ |
|-----------|-----------|
| <i>t</i> | 0.6 |
| <i>f</i> | 0.4 |

$$f_3(FP, Fraud, Trav) = P(FP|Fraud, Trav) =$$

| <i>FP</i> | <i>Fraud</i> | <i>Trav</i> | $f_3(Fraud, Trav)$ |
|-----------|--------------|-------------|--------------------|
| <i>t</i> | <i>t</i> | <i>t</i> | 0.9 |
| <i>t</i> | <i>t</i> | <i>f</i> | 0.1 |
| <i>t</i> | <i>f</i> | <i>t</i> | 0.9 |
| <i>t</i> | <i>f</i> | <i>f</i> | 0.01 |
| <i>f</i> | <i>t</i> | <i>t</i> | $1 - 0.9 = 0.1$ |
| <i>f</i> | <i>t</i> | <i>f</i> | $1 - 0.1 = 0.9$ |
| <i>f</i> | <i>f</i> | <i>t</i> | $1 - 0.9 = 0.1$ |
| <i>f</i> | <i>f</i> | <i>f</i> | $1 - 0.01 = 0.99$ |

$$f_4(CRP, OC) = P(CRP|OC) =$$

| <i>CRP</i> | <i>OC</i> | $f_4(CRP, OC)$ |
|------------|-----------|---------------------|
| <i>t</i> | <i>t</i> | 0.10 |
| <i>t</i> | <i>f</i> | 0.001 |
| <i>f</i> | <i>t</i> | $1 - 0.10 = 0.9$ |
| <i>f</i> | <i>f</i> | $1 - 0.001 = 0.999$ |

$$f_5(Fraud, OC) = P(IP = true|Fraud, OC) =$$

| <i>Fraud</i> | <i>OC</i> | $f_5(Fraud, OC)$ |
|--------------|-----------|------------------|
| <i>t</i> | <i>t</i> | 0.02 |
| <i>t</i> | <i>f</i> | 0.011 |
| <i>f</i> | <i>t</i> | 0.01 |
| <i>f</i> | <i>f</i> | 0.001 |

Thus we get $P(Fraud|IP = true)$

$$\begin{aligned} &\propto \sum_{CRP} \sum_{OC} f_2(OC) f_4(CRP, OC) f_5(Fraud, OC) \sum_{FP} \sum_{Trav} f_0(Trav) f_1(Fraud, Trav) f_3(FP, Fraud, Trav) \\ &= \sum_{CRP} \sum_{OC} f_2(OC) f_4(CRP, OC) f_5(Fraud, OC) \sum_{FP} f_6(Fraud, FP) \end{aligned}$$

Where $f_6(Fraud, FP) = \sum_{Trav} f_0(Trav) f_1(Fraud, Trav) f_3(FP, Fraud, Trav)$

$$=$$

| <i>Fraud</i> | <i>FP</i> | $f_6(Fraud, FP)$ |
|--------------|-----------|--|
| <i>t</i> | <i>t</i> | $0.05 \times 0.01 \times 0.9 + 0.95 \times 0.004 \times 0.1 = 0.00083$ |
| <i>t</i> | <i>f</i> | $0.05 \times 0.01 \times 0.1 + 0.95 \times 0.004 \times 0.9 = 0.00347$ |
| <i>f</i> | <i>t</i> | $0.05 \times 0.99 \times 0.9 + 0.95 \times 0.996 \times 0.01 = 0.054012$ |
| <i>f</i> | <i>f</i> | $0.05 \times 0.99 \times 0.1 + 0.95 \times 0.996 \times 0.99 = 0.941688$ |

$$\begin{aligned} &\text{So, } \sum_{CRP} \sum_{OC} f_2(OC) f_4(CRP, OC) f_5(Fraud, OC) \sum_{FP} f_6(Fraud, FP) \\ &= \sum_{CRP} \sum_{OC} f_2(OC) f_4(CRP, OC) f_5(Fraud, OC) f_7(Fraud) \end{aligned}$$

Where $f_7(Fraud) = \sum_{FP} f_6(Fraud, FP)$

| <i>Fraud</i> | $f_7(Fraud)$ |
|--------------|--------------------------------|
| <i>t</i> | $0.00083 + 0.00347 = 0.0043$ |
| <i>f</i> | $0.054012 + 0.941688 = 0.9957$ |

So, $\sum_{CRP} \sum_{OC} f_2(OC) f_4(CRP, OC) f_5(Fraud, OC) f_7(Fraud) = \sum_{CRP} f_8(Fraud, CRP)$

Where $f_8(Fraud, CRP) = \sum_{OC} f_2(OC) f_4(CRP, OC) f_5(Fraud, OC) f_7(Fraud)$

| <i>Fraud</i> | <i>CRP</i> | $f_8(Fraud, CRP)$ |
|--------------|------------|--|
| <i>t</i> | <i>t</i> | $0.6 \times 0.1 \times 0.02 \times 0.0043 + 0.4 \times 0.001 \times 0.011 \times 0.0043 = 0.00000517892$ |
| <i>t</i> | <i>f</i> | $0.6 \times 0.9 \times 0.02 \times 0.0043 + 0.4 \times 0.999 \times 0.011 \times 0.0043 = 0.00006534108$ |
| <i>f</i> | <i>t</i> | $0.6 \times 0.1 \times 0.01 \times 0.9957 + 0.4 \times 0.001 \times 0.001 \times 0.9957 = 0.00059781828$ |
| <i>f</i> | <i>f</i> | $0.6 \times 0.9 \times 0.01 \times 0.9957 + 0.4 \times 0.999 \times 0.001 \times 0.9957 = 0.00577466172$ |

So, $f_9(Fraud) = \sum_{CRP} f_8(Fraud, CRP)$

| <i>Fraud</i> | $f_9(Fraud)$ |
|--------------|--|
| <i>t</i> | $0.00000517892 + 0.00006534108 = 0.00007052$ |
| <i>f</i> | $0.00059781828 + 0.00577466172 = 0.00637248$ |

So, we get $P(Fraud = true | IP = true) = \frac{0.00007052}{0.00007052 + 0.00637248} = 0.010945211857830203 \approx 0.01095$.

Now we proceed with $P(Fraud = true | IP = true, CRP = true)$:

| <i>Trav</i> | $f_0(Trav)$ |
|-------------|-------------|
| <i>t</i> | 0.05 |
| <i>f</i> | 0.95 |

| <i>Fraud</i> | <i>Trav</i> | $f_1(Fraud Trav)$ |
|--------------|-------------|---------------------|
| <i>t</i> | <i>t</i> | 0.01 |
| <i>t</i> | <i>f</i> | 0.004 |
| <i>f</i> | <i>t</i> | $1 - 0.01 = 0.99$ |
| <i>f</i> | <i>f</i> | $1 - 0.004 = 0.996$ |

| <i>OC</i> | $f_2(OC)$ |
|-----------|-----------|
| <i>t</i> | 0.6 |
| <i>f</i> | 0.4 |

| <i>FP</i> | <i>Fraud</i> | <i>Trav</i> | $f_3(Fraud, Trav)$ |
|-----------|--------------|-------------|--------------------|
| <i>t</i> | <i>t</i> | <i>t</i> | 0.9 |
| <i>t</i> | <i>t</i> | <i>f</i> | 0.1 |
| <i>t</i> | <i>f</i> | <i>t</i> | 0.9 |
| <i>t</i> | <i>f</i> | <i>f</i> | 0.01 |
| <i>f</i> | <i>t</i> | <i>t</i> | $1 - 0.9 = 0.1$ |
| <i>f</i> | <i>t</i> | <i>f</i> | $1 - 0.1 = 0.9$ |
| <i>f</i> | <i>f</i> | <i>t</i> | $1 - 0.9 = 0.1$ |
| <i>f</i> | <i>f</i> | <i>f</i> | $1 - 0.01 = 0.99$ |

| <i>OC</i> | $f_4(OC)$ |
|-----------|-----------|
| <i>t</i> | 0.10 |
| <i>f</i> | 0.001 |

$$f_5(Fraud, OC) = P(IP = true | Fraud, OC) =$$

| <i>Fraud</i> | <i>OC</i> | $f_5(Fraud, OC)$ |
|--------------|-----------|------------------|
| <i>t</i> | <i>t</i> | 0.02 |
| <i>t</i> | <i>f</i> | 0.011 |
| <i>f</i> | <i>t</i> | 0.01 |
| <i>f</i> | <i>f</i> | 0.001 |

Thus we get $P(Fraud | IP = true, CRP = true)$

$$\propto \sum_{OC} f_2(OC) f_4(OC) f_5(Fraud, OC) \sum_{FP} \sum_{Trav} f_0(Trav) f_1(Fraud, Trav) f_3(FP, Fraud, Trav)$$

$$= \sum_{OC} f_2(OC) f_4(OC) f_5(Fraud, OC) \sum_{FP} f_6(Fraud, FP)$$

Where $f_6(Fraud, FP) = \sum_{Trav} f_0(Trav) f_1(Fraud, Trav) f_3(FP, Fraud, Trav)$

| <i>Fraud</i> | <i>FP</i> | $f_6(Fraud, FP)$ |
|--------------|-----------|--|
| <i>t</i> | <i>t</i> | $0.05 \times 0.01 \times 0.9 + 0.95 \times 0.004 \times 0.1 = 0.00083$ |
| <i>t</i> | <i>f</i> | $0.05 \times 0.01 \times 0.1 + 0.95 \times 0.004 \times 0.9 = 0.00347$ |
| <i>f</i> | <i>t</i> | $0.05 \times 0.99 \times 0.9 + 0.95 \times 0.996 \times 0.01 = 0.054012$ |
| <i>f</i> | <i>f</i> | $0.05 \times 0.99 \times 0.1 + 0.95 \times 0.996 \times 0.99 = 0.941688$ |

So, $\sum_{OC} f_2(OC) f_4(OC) f_5(Fraud, OC) \sum_{FP} f_6(Fraud, FP)$

$$= \sum_{OC} f_2(OC) f_4(OC) f_5(Fraud, OC) f_7(Fraud)$$

Where $f_7(Fraud) = \sum_{FP} f_6(Fraud, FP)$

| <i>Fraud</i> | $f_7(Fraud)$ |
|--------------|--------------------------------|
| <i>t</i> | $0.00083 + 0.00347 = 0.0043$ |
| <i>f</i> | $0.054012 + 0.941688 = 0.9957$ |

So, $\sum_{OC} f_2(OC) f_4(OC) f_5(Fraud, OC) f_7(Fraud) = f_8(Fraud)$

| <i>Fraud</i> | $f_8(Fraud)$ |
|--------------|--|
| <i>t</i> | $0.6 \times 0.1 \times 0.02 \times 0.0043 + 0.4 \times 0.001 \times 0.011 \times 0.0043 = 0.00000517892$ |
| <i>f</i> | $0.6 \times 0.1 \times 0.01 \times 0.9957 + 0.4 \times 0.001 \times 0.001 \times 0.9957 = 0.00059781828$ |

So, we get $P(Fraud = true | IP = true, CRP = true) = \frac{0.00000517892}{0.00000517892 + 0.00059781828} = 0.008588630262296408 \approx 0.00859$.

So by having $CRP = true$ prior to my internet purchase, the chance of the transaction being rejected as a possible fraud is reduced by $0.010945211857830203 - 0.008588630262296408 = 0.00235658 \approx 0.00236$.