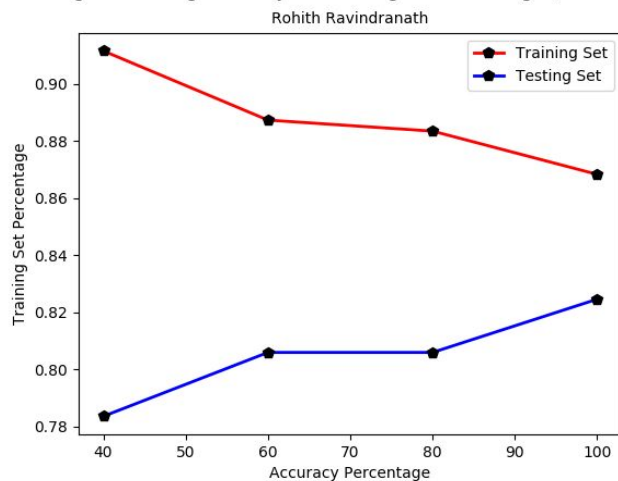Rohith Ravindranath
PUID: 0028822977
Dan Goldwasser
CS 37300
3rd March 2018

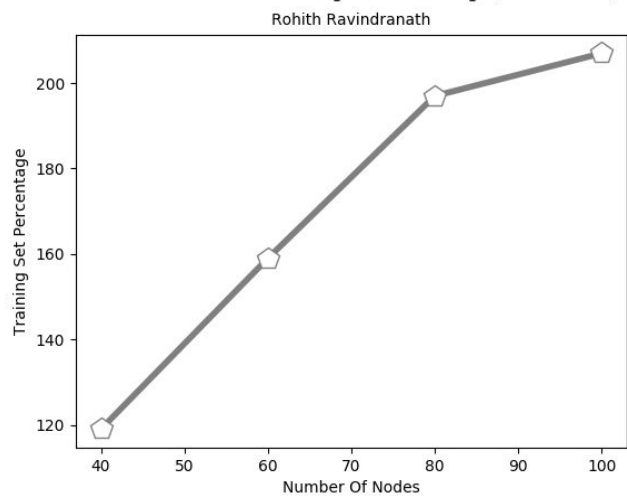<p style="text-align:center">Homework #2</p>

I collaborated with Vishal Vasan, Meera Haridasa, Jullian Haresco, and Rohan Saxena. I affirm that I wrote the solutions in my own words and that I understand the solutions I am submitting.

1. Vanilla Tree Plots



Training and Testing Accuracy vs. Training Set Percentage (Vanilla Tree)
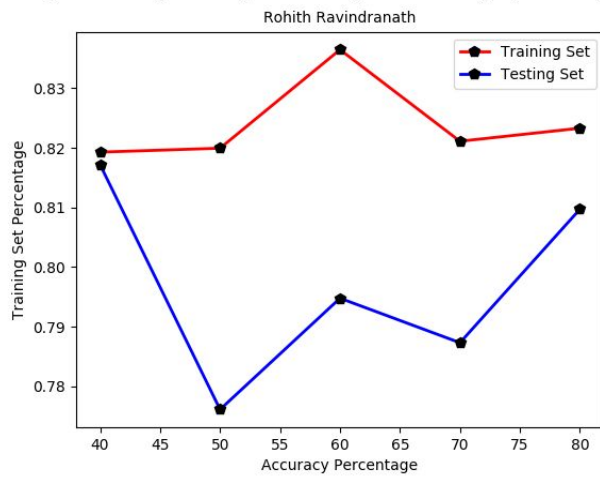Rohith Ravindranath



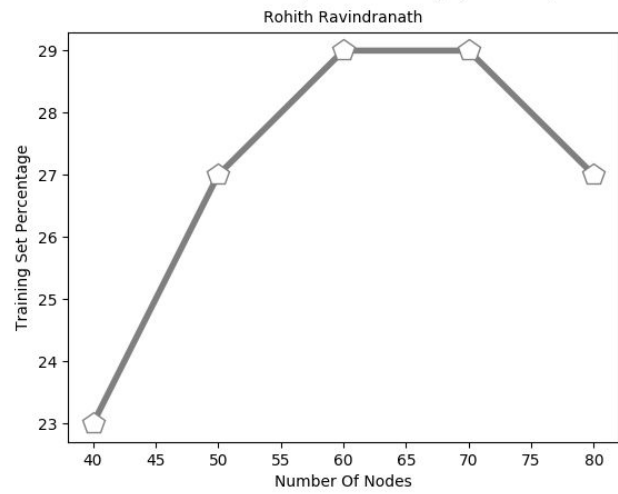Number Of Nodes vs. Training Set Percentage (Vanilla Tree)
Rohith Ravindranath

Question 1 Code Snippet

## 2. Depth Tree Plots

**Training and Testing Accuracy vs. Training Set Percentage (Static Depth Tree**
Rohith Ravindranath



**Number Of Nodes vs. Training Set Percentage (Static Depth Tree)**
Rohith Ravindranath



**Depth vs. Training Set Percentage (Static Depth Tree)**
Rohith Ravindranath

# Question 2 Code Snippet

```python
76  ################## QUESTION 2 - STATIC DEPTH TREE ANALYSIS #####################
77  print('################## QUESTION 2 - STATIC DEPTH TREE ANALYSIS #####################')
78  training_set_percent_lis = [40,50,60,70,80]
79  max_depth_lis = [5,10,15,20]
80  number_of_nodes = []
81  training_set_accuracy = []
82  testing_set_accuracy = []
83  optional_choice_of_depth = []
84  validation_percent = 20
85
86  for training_percent in training_set_percent_lis:
87
88      max_depth = 0
89      training_acc = 0
90      testing_acc = 0
91      validation_acc  =  0
92      max_nodes = 0
93
94      for depth in max_depth_lis:
95
96          depth_tree = ID3.Tree(TRAINING_FILE_DATA,TRAINING_FILE_LABEL,TESTING_FILE_DATA,TESTING_FILE_LABEL,'depth')
97          data = depth_tree.load_data(TRAINING_FILE_DATA,TRAINING_FILE_LABEL)
98          test_data = depth_tree.load_data(TESTING_FILE_DATA,TESTING_FILE_LABEL)
99
100         """
101         split = float((training_percent) / 100)
102         v_split = float((validation_percent) / 100)
103         sub_data = data.sample(frac=split)
104         v_sub_data = data.sample(frac=v_split)
105         """
106
107         sub_data = data[: (int(len(data) * int(training_percent) / 100))]
108         v_sub_data = data[ (int(len(data) * int(100 - validation_percent) / 100)): ]
109
110         tree= depth_tree.construct_tree(sub_data,max_depth)
111         num_nodes = depth_tree.count_nodes(tree)
112
113         training_predicted_vals = depth_tree.predict(sub_data,tree)
114         training_accuracy = depth_tree.analyze(training_predicted_vals)
115
116         validation_predicted_vals = depth_tree.predict(v_sub_data,tree)
117         validation_accuracy = depth_tree.analyze(validation_predicted_vals)
118
119         testing_predicted_vals = depth_tree.predict(test_data ,tree)
120         testing_accuracy = depth_tree.analyze(testing_predicted_vals)
121         #print(str(validation_accuracy) + '   '  + str(depth))
122         if validation_accuracy > validation_acc:
123             max_depth = depth
124             training_acc  = training_accuracy
125             testing_acc = testing_accuracy
126             validation_acc = validation_accuracy
127             max_nodes = num_nodes
128
129     print('Best depth for ' + str(training_percent) + '% is ' + str(max_depth))
130     print('Number of Nodes: ' + str(max_nodes))
131     print("Training set accuracy: %.4f" % training_acc)
132     print("Validation set accuracy: %.4f" % validation_acc)
133     print("Testing set accuracy: %.4f" % testing_acc)
134     print()
135
136     number_of_nodes.append(max_nodes)
137     training_set_accuracy.append(training_acc)
138     testing_set_accuracy.append(testing_acc)
139     optional_choice_of_depth.append(max_depth)
140
141  #Beging plotting
142  plt.figure()
143  plt.plot(training_set_percent_lis,training_set_accuracy,'-p',  color='red', markersize=7, linewidth=2, markerfacecolor='black', markeredgecolor='black', label='Traini
144  plt.plot(training_set_percent_lis,testing_set_accuracy, '-p', color='blue', markersize=7, linewidth=2, markerfacecolor='black', markeredgecolor='black', label='Testin
145  plt.xlabel('Accuracy Percentage')
146  plt.legend()
147  plt.ylabel('Training Set Percentage')
148  plt.title('Rohith Ravindranath' , fontsize = 10)
149  plt.suptitle('Training and Testing Accuracy vs. Training Set Percentage (Static Depth Tree)',  fontsize = 12)
150  #plt.show()
151  plt.savefig('Training_v_Testing_Plot_SD.png')
152
153  plot_graph(number_of_nodes, training_set_percent_lis,'Number Of Nodes', 'Training Set Percentage', 'Number Of Nodes vs. Training Set Percentage (Static Depth Tree)','
154  plot_graph(optional_choice_of_depth, training_set_percent_lis,'Depth', 'Training Set Percentage', 'Depth vs. Training Set Percentage (Static Depth Tree)','Depth_SD' )
155
```

## 3. Prune Tree Plots

### Training and Testing Accuracy vs. Training Set Percentage (Prune Tree)
Rohith Ravindranath



### Number Of Nodes vs. Training Set Percentage (Prune Tree)
Rohith Ravindranath



## Question 3 Code Snippet

```python
157  print('################ QUESTION 3 - PRUNE TREE ANALYSIS ######################')
158  training_set_percent_lis = [40,50,60,70,80]
159  number_of_nodes = []
160  training_set_accuracy = []
161  testing_set_accuracy = []
162  validation_percent = 20
163
164  for training_percent in training_set_percent_lis:
165      prune_tree = ID3.Tree(TRAINING_FILE_DATA,TRAINING_FILE_LABEL,TESTING_FILE_DATA,TESTING_FILE_LABEL,'prune')
166      data = prune_tree.load_data(TRAINING_FILE_DATA,TRAINING_FILE_LABEL)
167      test_data = prune_tree.load_data(TESTING_FILE_DATA,TESTING_FILE_LABEL)
168
169      max_depth = float('inf')
170      """
171      split = float((training_set_percent) / 100)
172      v_split = float((validation_set_percent) / 100)
173      sub_data = data.sample(frac=split)
174      v_sub_data = data.sample(frac=v_split)
175      """
176
177      sub_data = data[: (int(len(data) * int(training_percent) / 100))]
178      v_sub_data = data[ (int(len(data) * int(100 - validation_percent) / 100)): ]
179
180      tree = prune_tree.construct_tree(sub_data,max_depth, v_data = v_sub_data)
181      num_nodes = prune_tree.count_nodes(tree)
182      number_of_nodes.append(num_nodes)
183      print('Number of Nodes: ' + str(num_nodes))
184
185      training_predicted_vals = prune_tree.predict(sub_data,tree)
186      training_accuracy = prune_tree.analyze(training_predicted_vals)
187      training_set_accuracy.append(training_accuracy)
188      print("Training set accuracy: %.4f" % training_accuracy)
189
190      testing_predicted_vals = prune_tree.predict(test_data ,tree)
191      testing_accuracy = prune_tree.analyze(testing_predicted_vals)
192      testing_set_accuracy.append(testing_accuracy)
193      print("Testing set accuracy: %.4f" % testing_accuracy)
194      print()
195
196  #Beging plotting
197  plt.figure()
198  plt.plot(training_set_percent_lis,training_set_accuracy,'-p',  color='red', markersize=7, linewidth=2, markerfacecolor='black', markeredgecolor='black', label='Traini
199  plt.plot(training_set_percent_lis,testing_set_accuracy, '-p', color='blue', markersize=7, linewidth=2, markerfacecolor='black', markeredgecolor='black', label='Testin
200  plt.xlabel('Accuracy Percentage')
201  plt.legend()
202  plt.ylabel('Training Set Percentage')
203  plt.title('Rohith Ravindranath' , fontsize = 10)
204  plt.suptitle('Training and Testing Accuracy vs. Training Set Percentage (Prune Tree)',  fontsize = 12)
205  #plt.show()
206  plt.savefig('Training_v_Testing_Plot_P.png')
207
208  plot_graph(number_of_nodes, training_set_percent_lis,'Number Of Nodes', 'Training Set Percentage', 'Number Of Nodes vs. Training Set Percentage (Prune Tree)','Num_Nod
```

Question 4

We don't want to prune directly on the test set as this will make the decision optimized for that specific data set. In other words, pruning using the test data set will create a biased towards that specific data set and give a higher accuracy percentage than what it should be. Another reason that, we the testing set is too big and will not be efficient to use it as pruning mechanism.

Question 5

For a normal decision tree (vanilla tree) instead of classifying the terminal node with only one class label, you would calculate the probability distribution of all the class labels from the sub-set of data that corresponds to this terminal node. When we are using the tree to predict for a given test data set, we will get a ranking model based on a probability distribution of all class labels rather than only one class labe when we reach a terminal node. With respect to a depth decision tree, when we reach the max depth we will simply convert that node into a terminal node and calculate the probability distribution of all the class labels from the sub-set of data that corresponds to this node. When we are pruning the tree, if we decide to prune the child nodes of a certain node, we would simple calculate the probability distribution for the node based on the sub-set of data of both the child nodes.