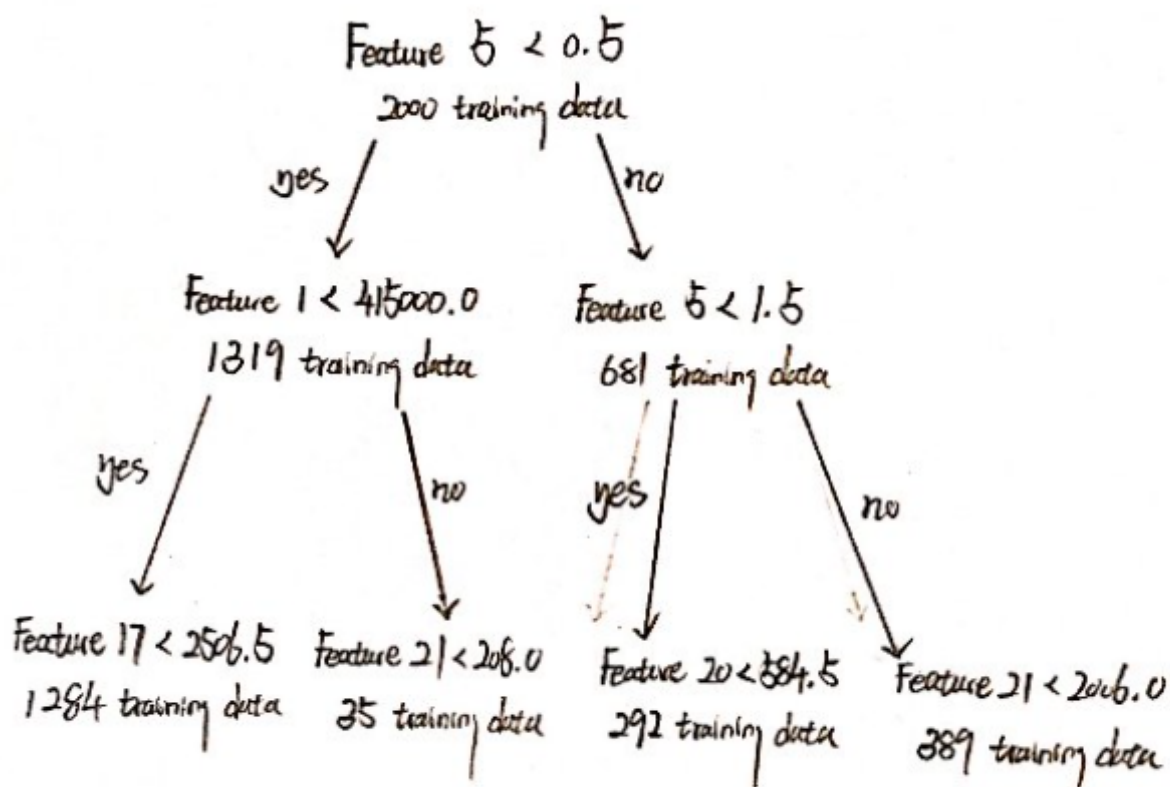**1. First, build an ID3 Decision Tree classifier based on the data in pa2train.txt. Do not use pruning. Draw the first three levels decision tree that you obtain. For each node that you draw, if it is a leaf node, write down the label that will be predicted for this node, as well as how many of the training data points lie in this node. If it is an internal node, write down the splitting rule for the node, as well as how many of the training data points lie in this node. (Hint: If your code is correct, the root node will involve the rule Feature 5 < 0.5.)**

Feature 5 < 0.5
2000 training data

yes /     \ no

Feature 1 < 415000.0      Feature 5 < 1.5
1319 training data       681 training data

yes /    \ no      yes /    \ no

Feature 17 < 2506.5   Feature 21 < 208.0    Feature 20 < 884.5   Feature 21 < 2006.0
1284 training data    35 training data      292 training data     389 training data

```python
1    import math
2    import numpy
3    import operator
4
5    def loadData():
6        # open the training data file
7        data_file = open("pa2train.txt")
8        # open the feature file
9        features_file = open("pa2features.txt")
10       # list that stores the training data
11       data = []
12       # list that stores the features
13       features = []
14       # read the training data file
15       for line in data_file:
16           d = line.split()
17           for i in range(0, 23):
18               d[i] = float([i])
19           if d[22] == 1:
20               d[22] = "yes"
21           else:
22               d[22] = "no"
23           data.append(data)
24       # read the features file
25       for line in features_file:
26           features.append(line)
27       return data, features
28
29   def calculateEntropy(data):
30       # dictionary that maps label to its count
31       labels = {}
32       # initialize the dictionary
33       for i in data:
34           labels[i[-1]] = 0
```

```python
35          # update the dictionary
36          for i in data:
37              labels[i[-1]] = labels[i[-1]] + 1
38          h = 0.0
39          for i in labels:
40              # calculate the probability
41              p = float(labels[i]) / float(len(data))
42              if p != 0:
43                  # calculate the entropy
44                  h = h - p * math.log(p, 2)
45          return h

47      def split(data, feature, value):
48          splitedData = []
49          for i in data:
50              # split the feature vector
51              if i[feature] == value:
52                  splitedFeatureVector = i[:feature]
53                  splitedFeatureVector = i[feature + 1:]
54                  splitedData.append(splitedFeatureVector)
55          return splitedData

57      def getBestSplitFeature(data):
58          h = calculateEntropy(data)
59          bestFeature = -1
60          bestInformationGain = 0.0
61          # iterate through the features
62          for feature in range(0, 22):
63              # get all the values for each feature
64              values = []
65              for j in data:
66                  values.append(j[feature])
67              valuesSet = set(values)
```

```
68              conditionalH = 0.0
69              for value in valuesSet:
70                  # make the split
71                  splitedData = split(data, feature, value)
72                  p = float(len(splitedData)) / float(len(data))
73                  # calculate the conditional entropy
74                  conditionalH = conditionalH + p * calculateEntropy(splitedData)
75                  # get the information gain
76                  informationGain = h - conditionalH
77                  # get the best feature for splitting
78                  if (informationGain > bestInformationGain):
79                      bestFeature = feature
80                      bestInformationGain = informationGain
81          return bestFeature
82
83  def getMajority(labels):
84      # dictionary that maps label to its count
85      labelCount = {}
86      # initialize the dictionary
87      for i in labels:
88          labelCount[i] = 0
89      # update the dictionary
90      for i in labels:
91          labelCount[i] = labelCount[i] + 1
92      # sort the dictionary
93      labelCount = sorted(labelCount.items(), key = operator.itemgetter(1))
94      # get the label
95      label = labelCount[len(labelCount) - 1][0]
96      return label
97
98  def buildTree(data,features):
99      # list that stores all the labels
100     labels = []
```

```
101         for i in data:
102             labels.append(i[-1])
103         # if there is only one label in the list
104         length = len(labels)
105         if labels.count(labels[0]) == length:
106             return labels[0]
107         # if there is no features left
108         if len(data[0]) == 1:
109             return getMajority(labels)
110         # get the best feature
111         bestFeature = getBestSplitFeature(data)
112         # get the label of the best feature
113         bestLabel = features[bestFeature]
114         # build the tree
115         tree = {bestLabel:{}}
116         # get the values for the best feature
117         featureValues = []
118         for i in data:
119             featureValues.append(i[bestFeature])
120         featureValues = set(featureValues)
121         new_features = []
122         # remove the best feature from the features list
123         for i in features:
124             if i != features[bestFeature]:
125                 new_features.append(i)
126         # recursion
127         for values in featureValues:
128             # make the split
129             splitedData = split(data,bestFeature,values)
130             tree[bestLabel][values] = buildTree(splitedData,new_features)
131         return tree
132
133     def prune(tree, validation_data):
134         queue = []
```

```python
135             queue.append(tree[0])
136         # go through the tree with BFS order
137         while len(queue) != 0:
138             key = queue.pop(0)
139             # calculate the current validation error
140             old_error = getError(tree[key][0], validation_data)
141             getMajority(tree[key][0])
142             new_error = getError(tree[key][0], validation_data)
143             if new_error <= old_error:
144                 return new_error
145             for i in tree[key][1]:
146                 queue.append(i)
147         # the tree can not be pruned
148         return -1
149
150
151     if __name__ == '__main__':
152         # load the data
153         data, features = loadData()
154         # build the decision tree
155         tree = buildTree(data, features)
156         # print the tree to see the top three levels
157         print(tree)
158         # read the validation data
159         validation_file = open("pa2validation.txt")
160         validation_data = []
161         for line in validation_file:
162             d = line.split()
163             for i in range(0, 23):
164                 d[i] = float([i])
165             if d[22] == 1:
166                 d[22] = "yes"
167             else:
```

```
168              d[22] = "no"
169          validation_data.append(d)
170      # first round of pruning
171      prune(tree, validation_data)
172      # second round of pruning
173      prune(tree, validation_data)
174      # print the pruned tree
175      print(tree)
176
```