

1 HW2

1.1

Before answering the questions, let's first take a look at the structure of the code.

```

1 class DescisionTree():
2     def __init__(self,m,k):
3         self.m = m
4         self.k = k
5         self.tree = [None]*(2**self.k-1)
6         self.label = [None]*(2**self.k-1)
7
8     def GenerateDataPoint(self):
9         '''generate data according to the given pattern '''
10
11    def Partition(self,matrix,x):
12        '''partition the data matrix based on x, \
13        return two matrixs according to value of x '''
14
15    def InformationGain(self,matrix):
16        '''compute the IGs and return max IG and its corresponding x '''
17
18    def GenerateTree(self,matrix,id = 0):
19        '''recursively building tree and label'''
20
21    def GetError(self,matrix):
22        '''get error of a built tree '''
23
24    if __name__ == '__main__':

```

Then we go back to question 1. For a given value of k and m, the following part of code generates a training data set based on the given pattern

```

1 def GenerateDataPoint(self):
2     '''generate data according to the given pattern '''
3     #compute weight list
4     w = []
5     denominator = 0
6     for i in range(2,self.k+1):
7         denominator += 0.9**i
8     for i in range(1,self.k+1):

```

```

9         w_i = 0.9**i/denominator
10        w.append(w_i)
11
12    #compute data point
13    matrix = np.zeros((self.m, self.k+1), 'int ')
14    for i in range(self.m):
15        value = 0
16        matrix[i][0] = 1 if random.random() < 0.5 else 0
17        for j in range(1, self.k):
18            matrix[i][j] = matrix[i][j-1] if \
19                random.random() < 0.75 else 1-matrix[i][j-1]
20            value += matrix[i][j]*w[j]
21        matrix[i][self.k] = matrix[i][0] if value >= 1/2 else 1-matrix[i][0]
22
23    return matrix
24    print(f'generate {self.m} data points successful ')

```

1.2

The following four parts relatively achieve three functions: partitioning the data; compute information gain; building decision tree; return the error of a built tree.

```

1  def Partition(self, matrix, x):
2      '''partition the data matrix based on x, \
3      return two matrixs according to value of x '''
4      list_x_0 = []
5      list_x_1 = []
6      for i in range(matrix.shape[0]):
7          if matrix[i][x] == 0:
8              list_x_0.append(matrix[i])
9          else:
10             list_x_1.append(matrix[i])
11      mat0 = np.array(list_x_0).reshape(-1, self.k+1)
12      mat1 = np.array(list_x_1).reshape(-1, self.k+1)
13      return mat0, mat1

```

```

1  def InformationGain(self, matrix):
2      '''compute the IGs and return max IG and its corresponding x '''
3      #compute H(Y)
4      count_y = 0
5      for i in range(matrix.shape[0]):
6          count_y = count_y+1 if matrix[i][self.k] == 1 else count_y
7      p_y = count_y/matrix.shape[0]
8      H_y = -p_y*math.log(p_y, 2) - (1-p_y)*math.log((1-p_y), 2)
9
10     #compute each H(Y/Xi)
11     list_IG = []
12     for j in range(self.k):
13         count_x_1 = 0 #compute sum of Xi=1
14         count_x_1_y_1 = 0 #compute sum of yi=1 when Xi=1
15         count_x_0_y_1 = 0 #compute sum of yi=1 when Xi=0

```

```

16         for i in range(matrix.shape[0]):
17             if matrix[i][j] == 1:
18                 count_x_1 += 1
19                 if matrix[i][self.k] == 1:
20                     count_x_1_y_1 += 1
21             else:
22                 if matrix[i][self.k] == 1:
23                     count_x_0_y_1 += 1
24
25         p1 = count_x_1/matrix.shape[0]
26         p0 = 1-p1
27         p11 = count_x_1_y_1/count_x_1 if count_x_1 != 0 else 0
28         p10 = 1-p11
29         p01 = count_x_0_y_1/(matrix.shape[0]-count_x_1) if \
30         (matrix.shape[0]-count_x_1) != 0 else 0
31         p00 = 1-p01
32         #print(p1,p0,p11,p10,p01,p00)
33
34         H_yx = p1*(-p11*math.log((p11+0.0000001),2)-\
35         p10*math.log((p10+0.0000001),2)) + \
36         p0*(-p01*math.log((p01+0.0000001),2)- \
37         p00*math.log((p00+0.0000001),2))
38         IG = H_y-H_yx
39         list_IG.append([j,IG])
40     return max(list_IG, key=lambda x: x[1])

```

```

1  def GenerateTree(self,matrix,id = 0):
2      '''recursively building tree and label'''
3      #print('id',id)
4      count = 0
5      for i in range(matrix.shape[0]):
6          if matrix[i][-1]==0:
7              count += 1
8
9      #print('bizhi ',count/matrix.shape[0])
10     #input()
11     if count/matrix.shape[0]>0.95 :
12         self.label[id] = 0
13         return
14     if count/matrix.shape[0]<0.05 :
15         self.label[id] = 1
16         return
17
18     pos, IG = self.InformationGain(matrix)
19     #print('pos&IG',pos,IG)
20     #input()
21     self.tree[id] = pos+1
22     mat0, mat1 = self.Partition(matrix, pos)
23     self.GenerateTree(mat0,2*id+1)
24     #print('id', id)
25     self.GenerateTree(mat1,2*id+2)

```

```

1 def GetError(self, matrix):
2     '''get error of a built tree '''
3     if matrix.shape[1]-1 != self.k:
4         print('depth of matrix does not comply with depth of the built tree')
5         return 0
6
7     right = 0
8     for i in range(matrix.shape[0]):
9         pos = 0
10        id = self.tree[pos]
11        while id:
12            if matrix[i][id-1] == 0:
13                pos = 2*pos+1
14            else:
15                pos = 2*pos+2
16            if pos > len(self.tree):
17                id = None
18                break
19            id = self.tree[pos]
20        predict = self.label[pos]
21        if predict == matrix[i][-1]:
22            right += 1
23
24    return 1-right/matrix.shape[0]

```

1.3

For $k = 4$ and $m = 30$, here is a sample and the corresponding decision tree.

```

1 original matrix generated by k=4 m=30:
2 [[0 0 1 0 1]
3  [0 0 0 0 1]
4  [1 1 1 1 1]
5  [1 1 0 0 0]
6  [0 0 0 0 1]
7  [0 0 0 0 1]
8  [0 1 1 1 0]
9  [1 1 0 0 0]
10 [0 0 0 0 1]
11 [0 0 1 1 0]
12 [0 1 1 1 0]
13 [0 1 0 0 1]
14 [0 0 1 1 0]
15 [1 1 0 0 0]
16 [1 1 1 1 1]
17 [0 0 0 0 1]
18 [1 1 0 0 0]
19 [0 0 0 0 1]
20 [0 0 1 0 1]
21 [1 0 0 0 0]
22 [1 1 0 1 1]

```

23	[1 1 1 1 1]
24	[1 1 0 0 0]
25	[0 0 0 0 1]
26	[0 0 0 0 1]
27	[0 0 1 1 0]
28	[0 0 1 1 0]
29	[0 0 1 1 0]
30	[0 0 1 1 0]
31	[0 0 0 1 1]

Based on the data points above, we get the decision tree as follows:

```

1 decision tree of original matrix:
2 [3, 1, 1, None, 4, 4, None, None, None, None, None, None, None, None, None]
3
4 corresponding label:
5 [None, None, None, 1, None, None, 1, None, None, 0, 1, 1, 0, None, None, \
6 None, None, None, None, None, None, None, None, None, None, None, None, None, \
7 None, None, None, None]
```

Below we draw the tree:

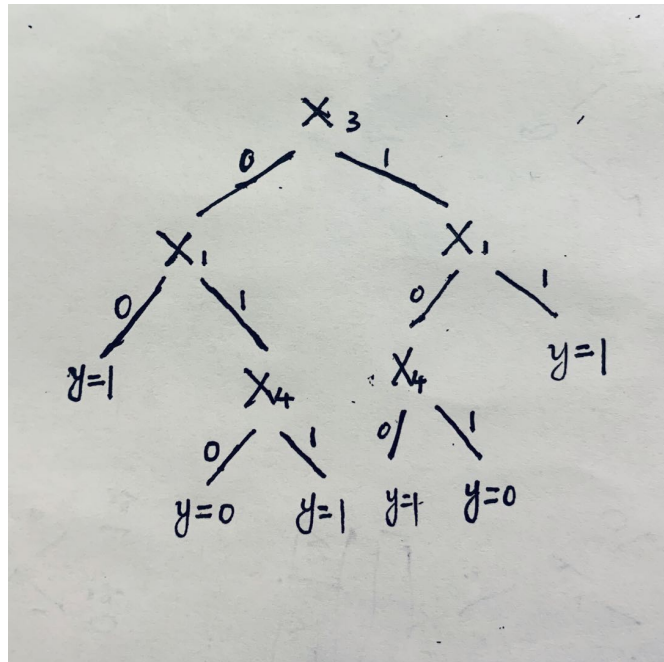


Figure 1: The Decision Tree

As we can see from above, the depth of the tree is 3, which is less than 4. So our approach helps reduce the depth of tree.

1.4

Based on $k=4$ and $m=30$, we first generate a decision tree from original 30 data points. Then we generate 2000 new data points according to the same scheme. Finally we try to fit the decision tree to new data, and get the average

error, which is around 0.06.

```
1  if __name__ == '__main__':
2      decision_tree = DescisonTree(30,4)
3      original = decision_tree.GenerateDataPoint()
4      print(f'original matrix generated by k=4 m=30:\n{original}')
5
6      decision_tree.GenerateTree(original)
7      print(f'decision tree of original matrix:\n{decision_tree.tree}\n')
8      print(f'corresponding label:\n{decision_tree.label}')
9
10     error_train = decision_tree.GetError(original)
11     print(f'training error = {error_train}')
12
13     #generate new data points to get test error
14     decision_tree.m = 2000
15     new = decision_tree.GenerateDataPoint()
16     print(f'new matrix generated by same k and m:\n{new}')
17
18     error_test = decision_tree.GetError(new)
19     print(f'test error = {error_test}')
```

1.5

When $k=10$, there are 2^{10} possible data. That is to say, if we have 2^{10} data points, we are likely to see every possible situation. So we take m from 100 to 5001(which is much bigger than 2^{10}) with step 100. For each m , we repeat for 50 times and compute the average of $|error_{train} - error_{test}|$

Below is the difference drawn according to above scheme. And from this we can see that, when number of data points increases, the difference between $error_{train}$ and $error_{test}$ decreases. That is because we are likely to see every possible data, thus the decision tree we get can generalize better.

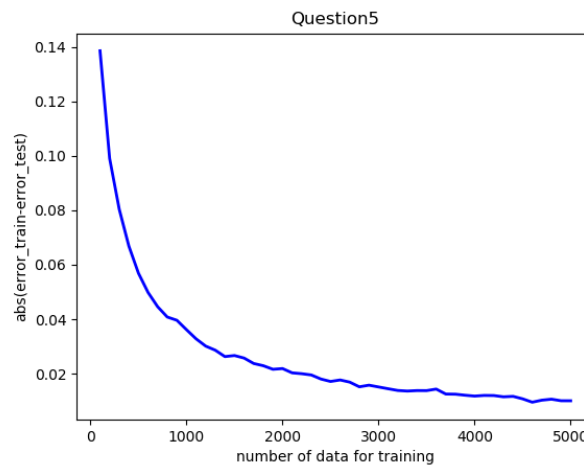


Figure 2: Difference

1.6

I implement Gini-Impurity in generating decision tree.(I heard about Gini-Impurity from others, and learned how it works from the Internet) Gini impurity is a measure of how likely a element is incorrectly labeled. Though they have different definition, they have a similar curve as mentioned on some websites.

Below is the main code for Gini-impurity.

```
1 def Gini(self ,matrix):
2     '''compute the Gini impurity and return min Gini \
3     and its corresponding x '''
4
5     list_gini = []
6     for j in range(self.k):
7         count_x_1 = 0      #compute sum of Xi=1
8         count_x_1_y_1 = 0  #compute sum of yi=1 when Xi=1
9         count_x_0_y_1 = 0  #compute sum of yi=1 when Xi=0
10        for i in range(matrix.shape[0]):
11            if matrix[i][j] == 1:
12                count_x_1 += 1
13                if matrix[i][self.k] == 1:
14                    count_x_1_y_1 += 1
15            else:
16                if matrix[i][self.k] == 1:
17                    count_x_0_y_1 += 1
18
19        p1 = count_x_1/matrix.shape[0]
20        p0 = 1-p1
21        p11 = count_x_1_y_1/count_x_1 if count_x_1 != 0 else 0
22        p10 = 1-p11
23        p01 = count_x_0_y_1/(matrix.shape[0]-count_x_1) if \
24        (matrix.shape[0]-count_x_1) != 0 else 0
25        p00 = 1-p01
26        #print(p1,p0,p11,p10,p01,p00)
27
28        if p0==0 or p1==0:    #if last time use x1 to part, \
29        this time p0=0 or p1 =0 when computing based on x1
30            continue
31        gini = p0*(1-p00**2-p01**2)+p1*(1-p10**2-p11**2)
32
33        list_gini.append([j,gini])
34    #print(list_gini)
35    #input()
36    return min(list_gini , key=lambda x: x[1])
```

Now we do what we did in question 5 (we take m from 100 to 5001 with step 100. For each m, we repeat for 50 times),we can get the difference curve in Figure 3:

When we put the difference curve of question 5 and question 6 into one graph as in Figure 4(blue line is for question 5; red line is for question 6). We can see that both approach have similar performances.

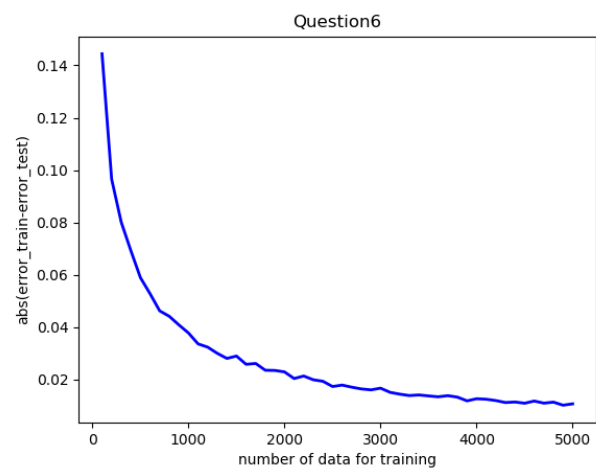


Figure 3: Difference

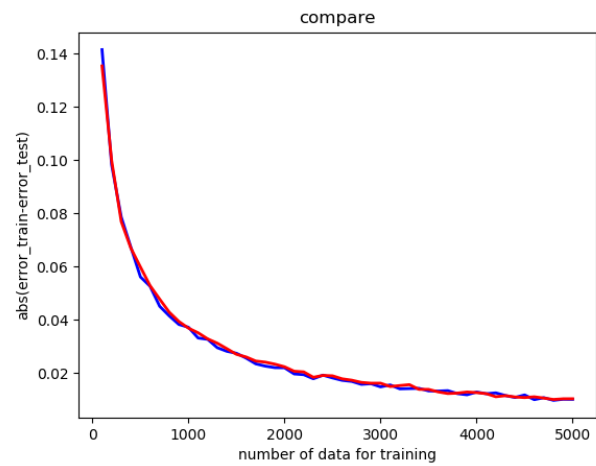


Figure 4: compare