

1.實作ID3 decision tree 二元分類方法

第一題

```
In [1]: 1 import math
2 def entropy(p1,n1):
3     if(p1==0 and n1==0):
4         return 1
5     elif(p1==0 or n1==0):
6         return 0
7     else:
8         pp=p1/(p1+n1)
9         np=n1/(p1+n1)
10        return -pp*math.log2(pp)-np*math.log2(np)
11 # print(entropy(29,35))
```

```
In [2]: 1 def IG(p1,n1,p2,n2):
2     num1 = p1+n1
3     num2 = p2+n2
4     num=num1+num2
5     return entropy(p1+p2,n1+n2)-num1/num*entropy(p1,n1)-num2/num*entropy(p2,n2)
6 #print(IG(21,5,8,30))
7 #print(IG(18,33,11,2))
```

```
In [3]: 1 def ID3DTtrain(feature,target):
2     node = dict()
3     node['data']=range(len(target)) #data為編號
4     tree = []
5     tree.append(node) #root=node
6     t=0
7     while(t<len(tree)): #編號t要生節點嗎
8         idx = tree[t]['data']
9         if(sum(target[idx])==0): #全部的target皆為0
10            tree[t]['leaf']=1 #走到底了
11            tree[t]['decision']=0 #以後走到這就是0
12        elif(sum(target[idx])==len(idx)): #全部target都是1
13            tree[t]['leaf']=1
14            tree[t]['decision']=1 #以後走到這都是1
15        else: #沒切乾淨要再多切幾刀
16            bestIG=0
17            for i in range(feature.shape[1]): #對每個feature都看一下
18                pool = list(set(feature[idx,i])) #看看編號idx的node裡的数据 它們的feature i都長啥樣
19                pool.sort()
20                for j in range (len(pool)-1):
21                    thres = (pool[j]+pool[j+1])/2 #區分j和j+1的中間值 e.g. 1 7 9 > 4 6(共n-1)
22                    G1=[]
23                    G2=[]
24                    for k in idx:
25                        if(feature[k][i]<thres): #如果小於thres分到group1
26                            G1.append(k)
27                        else: #不是的話就分到 group2
28                            G2.append(k)
29                    p1=sum(target[G1]==1)
30                    n1=sum(target[G1]==0)
31                    p2=sum(target[G2]==1)
32                    n2=sum(target[G2]==0)
33                    thisIG=IG(p1,n1,p2,n2)
34                    if(thisIG>bestIG): #把這個切法記起來
35                        bestIG=thisIG
36                        bestG1=G1
37                        bestG2=G2
38                        bestthres=thres
39                    bestf = i #記起來最好的切分的feature是誰
40            if(bestIG>0): #如果切分更細有幫助的話
41                tree[t]['leaf']=0
42                tree[t]['selectf']=bestf
43                tree[t]['threshold']=bestthres
44                tree[t]['child']=[len(tree),len(tree)+1]
45                node=dict()
46                node['data']=bestG1 #小於切分值得放右邊
47                tree.append(node)
48            else: #多切無益 姑且挑個比較好的值
49                tree[t]['leaf']=1 #走到底了
50                if(sum(target[idx]==1)>sum(target[idx]==0)): #因為1比0多
51                    tree[t]['decision']=1 #就讓他當1吧
52                else:
53                    tree[t]['decision']=0
54            t+=1
55        return tree
```

```
In [4]: 1 from sklearn import datasets
2 data = datasets.load_iris()
3 feature = data['data']
4 target = data['target']
5 T=ID3DTtrain(feature[0:100,:],target[0:100])
```

2.完成單一樣本測試函式

第二題

```
In [6]: 1 def ID3DTest(Tree,feature1):
2         now=0
3         while(Tree[now]['leaf']==0):
4             bestf = Tree[now]['selectf']
5             thres = Tree[now]['threshold']
6             if(feature1[bestf]<thres):
7                 now = Tree[now]['child'][0]
8             else:
9                 now = Tree[now]['child'][1]
10        return Tree[now]['decision']

In [7]: 1 # print(ID3DTest(T,feature[130,:]))
```

3. IRIS 資料庫共有 150 筆資料，共有三類花各 50 筆，今僅取後 100 筆（即 **target** 為 1 及 2 兩類）來實驗，請用這 100 筆資料建樹，再用相同資料測試，試算其預測正確率。

Ans: 100%

4. 承上，若每類改用前 30 筆建樹，後 20 筆測試，試算其測試資料的預測正確率。

Ans: 95%