

决策树自主实践与SKlearn调用

TA: 吕松霖

E-mail: lvsl@lamda.nju.edu.cn

实验内容

1. 从0构建决策树对简单数据分类
2. 可视化构建出的决策树
3. 调用SKlearn包中的决策树直接处理数据

实验目标

1. 从代码层面了解决策树的构建与计算
2. 学会直接调用现有决策树函数解决问题

实验内容

1 导入实验中的相关包



```
1 """
2 pickle包可以将决策树保存下来，方便下次直接调用
3 """
4 from matplotlib.font_manager import FontProperties
5 import matplotlib.pyplot as plt
6 from math import log
7 import operator
8 import pickle
```

2 构建简单的分类数据集



```
1
```

```

2  """
3  函数说明：创建测试数据集
4
5  Parameters:
6      None
7
8  Returns:
9      dataSet - 数据集
10     labels - 分类属性
11
12     用字典存储决策树结构：
13     {'有自己的房子':{0:{'有工作':{0:'no', 1:'yes'}}, 1:'yes'}}
14     年龄：0代表青年，1代表中年，2代表老年
15     有工作：0代表否，1代表是
16     有自己的房子：0代表否，1代表是
17     信贷情况：0代表一般，1代表好，2代表非常好
18     类别（是否给贷款）：no代表否，yes代表是
19  """
20  def createDataSet():
21      # 数据集
22      dataSet = [[0, 0, 0, 0, 'no'],
23                  [0, 0, 0, 1, 'no'],
24                  [0, 1, 0, 1, 'yes'],
25                  [0, 1, 1, 0, 'yes'],
26                  [0, 0, 0, 0, 'no'],
27                  [1, 0, 0, 0, 'no'],
28                  [1, 0, 0, 1, 'no'],
29                  [1, 1, 1, 1, 'yes'],
30                  [1, 0, 1, 2, 'yes'],
31                  [1, 0, 1, 2, 'yes'],
32                  [2, 0, 1, 2, 'yes'],
33                  [2, 0, 1, 1, 'yes'],
34                  [2, 1, 0, 1, 'yes'],
35                  [2, 1, 0, 2, 'yes'],
36                  [2, 0, 0, 0, 'no']]
37      # 分类属性

```

```
38     labels = ['年龄', '有工作', '有自己的房子', '信贷情况']
39     # 返回数据集和分类属性
40     return dataSet, labels
```

3 计算香农熵函数

```
1  """
2  函数说明：计算给定数据集的经验熵（香农熵）
3      Ent (D)  = -SUM (kp*Log2 (kp) )
4
5  Parameters:
6      dataSet - 数据集
7
8  Returns:
9      shannonEnt - 经验熵（香农熵）
10 """
11 def calcShannonEnt(dataSet):
12     # 返回数据集的行数
13     numEntires = len(dataSet)
14     # 保存每个标签（Label）出现次数的“字典”
15     labelCounts = {}
16     # 对每组特征向量进行统计
17     for featVec in dataSet:
18         # 提取标签（Label）信息
19         currentLabel = featVec[-1]
20         # 如果标签（Label）没有放入统计次数的字典，添加进去
21         if currentLabel not in labelCounts.keys():
22             # 创建一个新的键值对，键为currentLabel值为0
23             labelCounts[currentLabel] = 0
24         # Label计数
25         labelCounts[currentLabel] += 1
26     # 经验熵（香农熵）
27     shannonEnt = 0.0
28     # 计算香农熵
29     for key in labelCounts:
```

```

30         # 选择该标签 (Label) 的概率
31         prob = float(labelCounts[key]) / numEntires
32         # 利用公式计算
33         shannonEnt -= prob*log(prob, 2)
34     # 返回经验熵 (香农熵)
35     return shannonEnt
36
37 """
38 函数说明: 选择最优特征
39     Gain(D,g) = Ent(D) - SUM(|Dv|/|D|)*Ent(Dv)
40
41 Parameters:
42     dataSet - 数据集
43
44 Returns:
45     bestFeature - 信息增益最大的 (最优) 特征的索引值
46 """
47 def chooseBestFeatureToSplit(dataSet):
48     # 特征数量
49     numFeatures = len(dataSet[0]) - 1
50     # 计算数据集的香农熵
51     baseEntropy = calcShannonEnt(dataSet)
52     # 信息增益
53     bestInfoGain = 0.0
54     # 最优特征的索引值
55     bestFeature = -1
56     # 遍历所有特征
57     for i in range(numFeatures):
58         # 获取dataSet的第i个所有特征存在featList这个列表中 (列表生成式)
59         featList = [example[i] for example in dataSet]
60         # 创建set集合{}, 元素不可重复, 重复的元素均被删掉
61         # 从列表中创建集合是python语言得到列表中唯一元素值得最快方法
62         uniqueVals = set(featList)
63         # 经验条件熵
64         newEntropy = 0.0
65         # 计算信息增益

```

```

66         for value in uniqueVals:
67             # subDataSet划分后的子集
68             subDataSet = splitDataSet(dataSet, i, value)
69             # 计算子集的概率
70             prob = len(subDataSet) / float(len(dataSet))
71             # 根据公式计算经验条件熵
72             newEntropy += prob * calcShannonEnt(subDataSet)
73             # 信息增益
74             infoGain = baseEntropy - newEntropy
75             # 打印每个特征的信息增益
76             print("第%d个特征的信息增益为%.3f" % (i, infoGain))
77             # 计算信息增益
78             if(infoGain > bestInfoGain):
79                 # 更新信息增益，找到最大的信息增益
80                 bestInfoGain = infoGain
81                 # 记录信息增益最大的特征的索引值
82                 bestFeature = i
83             # 返回信息增益最大的特征的索引值
84             return bestFeature
85

```

4 必要的工具函数

```

1  """
2  函数说明：按照给定特征划分数据集
3
4  Parameters:
5      dataSet - 待划分的数据集
6      axis - 划分数据集的特征
7      values - 需要返回的特征的值
8
9  Returns:
10     None
11 """
12 def splitDataSet(dataSet, axis, value):

```

```

13     # 创建返回的数据集列表
14     retDataSet = []
15     # 遍历数据集的每一行
16     for featVec in dataSet:
17         if featVec[axis] == value:
18             # 去掉axis特征
19             reducedFeatVec = featVec[:axis]
20             # 将符合条件的添加到返回的数据集
21             # extend() 函数用于在列表末尾一次性追加另一个序列中的多个值（用新
列表扩展原来的列表）。
22             reducedFeatVec.extend(featVec[axis+1:])
23             # 列表中嵌套列表
24             retDataSet.append(reducedFeatVec)
25     # 返回划分后的数据集
26     return retDataSet
27
28 """
29 函数说明：统计classList中出现次数最多的元素（类标签）
30     服务于递归第二个终止条件
31
32 Parameters:
33     classList - 类标签列表
34
35 Returns:
36     sortedClassCount[0][0] - 出现次数最多的元素（类标签）
37 """
38 def majorityCnt(classList):
39     classCount = {}
40     # 统计classList中每个元素出现的次数
41     for vote in classList:
42         if vote not in classCount.keys():
43             classCount[vote] = 0
44             classCount[vote] += 1
45     # 根据字典的值降序排序
46     # operator.itemgetter(1)获取对象的第1列的值

```

```
47     sortedClassCount = sorted(classCount.items(), key =  
    operator.itemgetter(1), reverse = True)  
48     # 返回classList中出现次数最多的元素  
49     return sortedClassCount[0][0]
```

5 决策树相关函数

```
1  """  
2  函数说明：创建决策树（ID3算法）  
3      递归有两个终止条件：1、所有的类标签完全相同，直接返回类标签  
4      2、用完所有标签但是得不到唯一类别的分组，即特征不够  
    用，挑选出现数量最多的类别作为返回  
5  
6  Parameters:  
7      dataSet - 训练数据集  
8      labels - 分类属性标签  
9      featLabels - 存储选择的最优特征标签  
10  
11 Returns:  
12     myTree - 决策树  
13 """  
14 def createTree(dataSet, labels, featLabels):  
15     # 取分类标签（是否放贷：yes or no）  
16     classList = [example[-1] for example in dataSet]  
17     # 如果类别完全相同则停止继续划分  
18     if classList.count(classList[0]) == len(classList):  
19         return classList[0]  
20     # 遍历完所有特征时返回出现次数最多的类标签  
21     if len(dataSet[0]) == 1:  
22         return majorityCnt(classList)  
23     # 选择最优特征  
24     bestFeat = chooseBestFeatureToSplit(dataSet)  
25     # 最优特征的标签  
26     bestFeatLabel = labels[bestFeat]  
27     featLabels.append(bestFeatLabel)
```

```

28     # 根据最优特征的标签生成树
29     myTree = {bestFeatLabel: {}}
30     # 删除已经使用的特征标签
31     del(labels[bestFeat])
32     # 得到训练集中所有最优解特征的属性值
33     featValues = [example[bestFeat] for example in dataSet]
34     # 去掉重复的属性值
35     uniqueVals = set(featValues)
36     # 遍历特征, 创建决策树
37     for value in uniqueVals:
38         myTree[bestFeatLabel][value] =
createTree(splitDataSet(dataSet, bestFeat, value), labels,
featLabels)
39     return myTree
40
41 """
42 函数说明: 获取决策树叶子结点的数目
43
44 Parameters:
45     myTree - 决策树
46
47 Returns:
48     numLeafs - 决策树的叶子结点的数目
49 """
50 def getNumLeafs(myTree):
51     # 初始化叶子
52     numLeafs = 0
53     # python3中myTree.keys()返回的是dict_keys, 不是list, 所以不能用
54     # myTree.keys()[0]的方法获取结点属性, 可以使用list(myTree.keys())
[0]
55     # next() 返回迭代器的下一个项目 next(iterator[, default])
56     firstStr = next(iter(myTree))
57     # 获取下一组字典
58     secondDict = myTree[firstStr]
59     for key in secondDict.keys():
60         # 测试该结点是否为字典, 如果不是字典, 代表此节点为叶子结点

```



```

61         if type(secondDict[key]).__name__ == 'dict':
62             numLeafs += getNumLeafs(secondDict[key])
63         else:
64             numLeafs += 1
65     return numLeafs
66 """
67 函数说明：获取决策树的层数
68
69 Parameters:
70     myTree – 决策树
71
72 Returns:
73     maxDepth – 决策树的层数
74 """
75 def getTreeDepth(myTree):
76     # 初始化决策树深度
77     maxDepth = 0
78     # python3中myTree.keys()返回的是dict_keys,不是list,所以不能用
79     # myTree.keys()[0]的方法获取结点属性,可以使用list(myTree.keys())
80     [0]
81     # next() 返回迭代器的下一个项目 next(iterator[, default])
82     firstStr = next(iter(myTree))
83     # 获取下一个字典
84     secondDict = myTree[firstStr]
85     for key in secondDict.keys():
86         # 测试该结点是否为字典, 如果不是字典, 代表此节点为叶子结点
87         if type(secondDict[key]).__name__ == 'dict':
88             thisDepth = 1 + getTreeDepth(secondDict[key])
89         else:
90             thisDepth = 1
91     # 更新最深层数
92     if thisDepth > maxDepth:
93         maxDepth = thisDepth
94     # 返回决策树的层数
95     return maxDepth


```

```

96 """
97 函数说明：使用决策树分类
98
99 Parameters:
100     inputTree - 已经生成的决策树
101     featLabels - 存储选择的最优特征标签
102     testVec - 测试数据列表，顺序对应最优特征标签
103
104 Returns:
105     classLabel - 分类结果
106 """
107 def classify(inputTree, featLabels, testVec):
108     # 获取决策树结点
109     firstStr = next(iter(inputTree))
110     # 下一个字典
111     secondDict = inputTree[firstStr]
112     featIndex = featLabels.index(firstStr)
113     for key in secondDict.keys():
114         if testVec[featIndex] == key:
115             if type(secondDict[key]).__name__ == 'dict':
116                 classLabel = classify(secondDict[key], featLabels,
testVec)
117             else:
118                 classLabel = secondDict[key]
119     return classLabel

```

6 利用pickle包储存及读取决策树



```

1 """
2 函数说明：存储决策树
3
4 Parameters:
5     inputTree - 已经生成的决策树
6     filename - 决策树的存储文件名
7

```

```

8 Returns:
9     None
10 """
11 def storeTree(inputTree, filename):
12     with open(filename, 'wb') as fw:
13         pickle.dump(inputTree, fw)
14
15
16 """
17 函数说明：读取决策树
18
19 Parameters:
20     filename - 决策树的存储文件名
21
22 Returns:
23     pickle.load(fr) - 决策树字典
24 """
25 def grabTree(filename):
26     fr = open(filename, 'rb')
27     return pickle.load(fr)
28

```

7 主函数


```

1
2 dataSet, features = createDataSet()
3 featLabels = []
4 myTree = createTree(dataSet, features, featLabels)
5 # storeTree(myTree, 'classifierStorage.txt')
6 # myTree = grabTree('classifierStorage.txt')
7 # print(myTree)
8 # 测试数据
9 testVec = [0, 1, 1, 1]
10 result = classify(myTree, featLabels, testVec)
11 if result == 'yes':

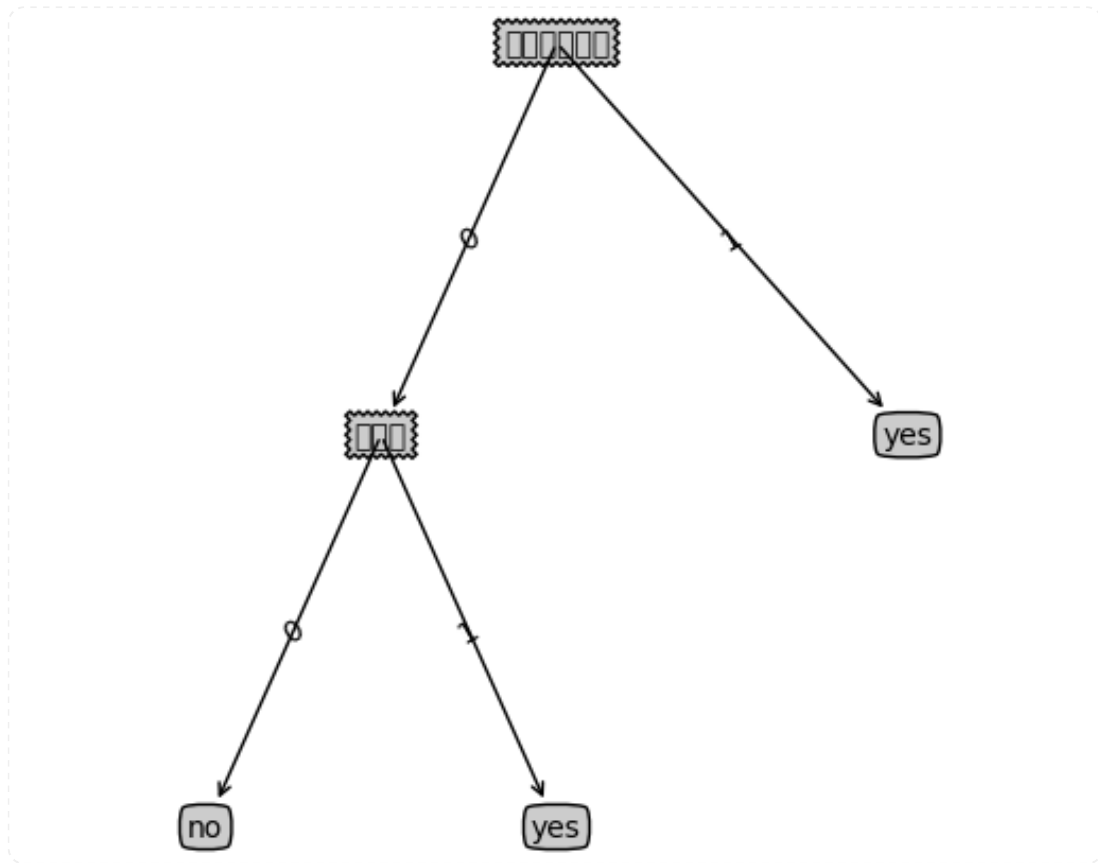
```

```
12     print('放贷')
13 if result == 'no':
14     print('不放贷')
15 print(myTree)
16 createPlot(myTree)
17 # print(dataSet)
18 # print(calcShannonEnt(dataSet))
19 print("最优特征索引值:" + str(chooseBestFeatureToSplit(dataSet)))
```

输出如下:



```
1 第0个特征的增益为0.083
2 第1个特征的增益为0.324
3 第2个特征的增益为0.420
4 第3个特征的增益为0.363
5 放贷
6 {'有自己的房子': {0: {'有工作': {0: 'no', 1: 'yes'}}}, 1: 'yes'}}
7
8 第0个特征的增益为0.083
9 第1个特征的增益为0.324
10 第2个特征的增益为0.420
11 第3个特征的增益为0.363
12 最优特征索引值:2
```



8 额外内容——SKlearn直接调用

安装SKlearn包指令：



```
1 pip install scikit-learn
```

数据及处理代码：



```
1 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
2 from sklearn import tree
3
4 dataSet = [[0, 0, 0, 0, 'no'],
5            [0, 0, 0, 1, 'no'],
6            [0, 1, 0, 1, 'yes'],
7            [0, 1, 1, 0, 'yes'],
8            [0, 0, 0, 0, 'no'],
9            [1, 0, 0, 0, 'no'],
10           [1, 0, 0, 1, 'no'],
11           [1, 1, 1, 1, 'yes'],
12           [1, 0, 1, 2, 'yes'],
```

```

13         [1, 0, 1, 2, 'yes'],
14         [2, 0, 1, 2, 'yes'],
15         [2, 0, 1, 1, 'yes'],
16         [2, 1, 0, 1, 'yes'],
17         [2, 1, 0, 2, 'yes'],
18         [2, 0, 0, 0, 'no']]
19
20 # 提取特征和标签
21 X = [row[:-1] for row in dataSet] # 特征（去掉最后一列）
22 y = [row[-1] for row in dataSet] # 目标标签（最后一列）
23
24 # 用 LabelEncoder 把 'no' 和 'yes' 转换为数值
25 le = LabelEncoder()
26 y = le.fit_transform(y) # 'no' -> 0, 'yes' -> 1
27
28 clf = tree.DecisionTreeClassifier(max_depth = 4) #创建
           DecisionTreeClassifier()类
29 clf = clf.fit(X, y) #使用数据，构建决策树
30
31 print(clf.predict([[1,1,1,0]])) #预测
32 # 结果: 1
33

```

DecisionTreeClassifier 是 scikit-learn 中的决策树分类器，其构造函数参数如下：

1. 基本参数

参数	说明	默认值
criterion	选择用于划分的标准 ，可选 "gini"（基尼系数）、"entropy"（信息增益）、"log_loss"（对数损失，用于概率估计）	"gini"
splitter	选择节点划分策略 ，"best" 表示选择最优划分，"random" 表示随机选择一个划分	"best"
max_depth	树的最大深度 ，None 表示不限制深度（树会生长到所有叶子纯度为止）	None

2. 控制树的形状

参数	说明	默认值
min_samples_split	内部节点再划分所需的最小样本数 （整数：至少多少个样本，浮点数：按比例）	2
min_samples_leaf	叶子节点最少包含的样本数 （防止过拟合）	1
min_weight_fraction_leaf	叶子节点最小权重比例 （适用于带权重数据，避免小样本极端分割）	0
max_features	用于寻找最佳划分的特征数量 ，可选整数（指定数量）、浮点数（比例）、"auto"（全部）、"sqrt"（平方根个）、"log2"（对数个）	None
max_leaf_nodes	最大叶子节点数 （None 表示不限制）	None
min_impurity_decrease	分裂所需的最小不纯度减少量 （防止过拟合）	0

3. 其他参数

参数	说明	默认值
random_state	随机数种子 ，用于保证结果可复现	None
class_weight	类别权重 ，可选 None（均等权重）、"balanced"（自动调整）、字典 {class_label: weight}	None
ccp_alpha	剪枝参数（复杂度惩罚） ，值越大剪枝越多	0

问题

-
1. 能否可视化库函数中的决策树？
 2. 能否从网上搜索表格数据并用决策树完成分类或回归任务？
 3. 自行尝试其他信息增益算法