

《机器学习实战》学习笔记（三）：决策树

原创 我是管小亮 2019-08-22 09:28:48 3168 收藏 24

版权

分类专栏: Machine Learning 文章标签: 机器学习 机器学习实战 读书笔记 决策树

欢迎关注WX公众号：【程序员管小亮】

【机器学习】《机器学习实战》读书笔记及代码 总目录

- <https://blog.csdn.net/TeFuirnever/article/details/99701256>

GitHub代码地址:

- <https://github.com/TeFuirnever/Machine-Learning-in-Action>

目录

欢迎关注WX公众号：【程序员管小亮】

本章内容

- 1、决策树
 - 2、决策树的构造
 - 1) 特征选择
 1. 香农熵
 2. 信息增益
 - 2) 决策树的生成和修剪
 - 3、决策树的可视化
 - 4、使用决策树执行分类
 - 5、决策树的存储
 - 6、Sklearn构建决策树分类器预测隐形眼镜类型
 - 7、sklearn.tree.DecisionTreeClassifier
 - 8、总结
- 参考文章

本章内容

- 决策树简介
- 在数据集中度量一致性
- 使用递归构造决策树
- 使用Matplotlib绘制树形图

关于决策树的西瓜书笔记和课后习题的博客如下:

- 读书笔记 <https://blog.csdn.net/TeFuirnever/article/details/98944801>
- 习题参考答案 <https://blog.csdn.net/TeFuirnever/article/details/99056057>

1、决策树

你是否玩过二十个问题的游戏，游戏的规则很简单：参与游戏的一方在脑海里想某个事物，其他参与者向他提问题，只允许提20个问题，问题的答案也只能用对或错回答。问问题的人通过推断分解，逐步缩小待猜测事物的范围。或者是酒桌上的猜数游戏，游戏规则有点类似，参与游戏的一方在脑海里想一个固定的数值，需要在固定的范围内，其他参与者进行猜测，他会先给出猜测的正确与否，如果正确直接喝酒，如果错误，就在下一个参与者开始之前给出新的猜数范围，逐步缩小待猜测事物的范围。

决策树的工作原理与上面两个游戏类似，用户输入一系列数据，然后给出游戏的答案。我们经常使用决策树处理分类问题，近来的调查表明决策树也是最经常使用的数据挖掘算法。它之所以如此流行，一个很重要的原因就是不需要了解机器学习知识，就能搞明白决策树是如何工作的。

如果以前没有接触过决策树，也完全不用担心，它的概念非常简单。通过简单的图形就可以了解其工作原理，图3-1所示的流程图就是一个决策树，长方形代表 **判断模块**（decision block），椭圆形代表 **终止模块**（terminating block），表示已经得出结论，可以终止运行。从 **判断模块** 引出的左右箭头称作 **分支**（branch），它可以到达另一个判断模块或者终止模块。

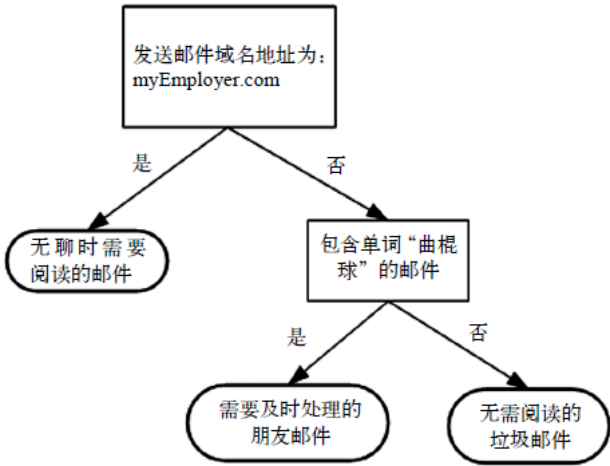


图3-1 流程图形式的决策树

上面构造的是一个假想的邮件分类系统，首先它检测发送邮件域名地址：

- 如果地址为myEmployer.com，则将其放在分类“无聊时需要阅读的邮件”中。
- 如果邮件不是来自这个域名，则检查邮件内容里是否包含单词曲棍球。
 - 如果包含则将邮件归类到“需要及时处理的朋友邮件”。
 - 如果不包含则将邮件归类到“无需阅读的垃圾邮件”。

第二章介绍的k-近邻算法可以完成很多分类任务，详见博客：【机器学习】《机器学习实战》读书笔记及代码：第2章 - k-近邻算法，但是它最大的缺点就是无法给出数据的内在含义，决策树的主要优势就在于数据形式非常容易理解，因此决策树可以使用不熟悉的数据集合，并从中提取出一系列规则，在这些机器根据数据集创建规则时，就是 **机器学习** 的过程。

决策树给出结果往往可以匹敌在当前领域具有几十年工作经验的人类专家。（有点厉害0-0）

那么如何从一堆原始数据中构造决策树呢？别着急，过程如下：

- 首先讨论构造决策树的方法，以及如何编写构造树的Python代码；
- 接着提出一些度量算法成功率的方法；
- 最后使用递归建立分类器，并且使用Matplotlib绘制决策树图。

决策树
优点：计算复杂度不高，输出结果易于理解，对中间值的缺失不敏感，可以处理不相关特征数据。
缺点：可能会产生过度匹配问题。
适用数据类型：数值型和标称型。

2、决策树的构造

决策树的一般流程

- (1) 收集数据：可以使用任何方法。
- (2) 准备数据：树构造算法只适用于标称型数据，因此数值型数据必须离散化。
- (3) 分析数据：可以使用任何方法，构造树完成之后，我们应该检查图形是否符合预期。
- (4) 训练算法：构造树的数据结构。
- (5) 测试算法：使用经验树计算错误率。
- (6) 使用算法：此步骤可以适用于任何监督学习算法，而使用决策树可以更好地理解数据的内在含义。

决策树要如何构建呢？通常，这一过程可以概括为3个步骤：特征选择、决策树的生成和决策树的修剪。

1) 特征选择

特征选择在于选取对训练数据具有分类能力的特征，这样可以提高决策树学习的效率；如果利用一个特征进行分类的结果与随机分类的结果没有很大差别，则称这个特征是没有分类能力的。经验上扔掉这样的特征对决策树学习的精度影响不大。

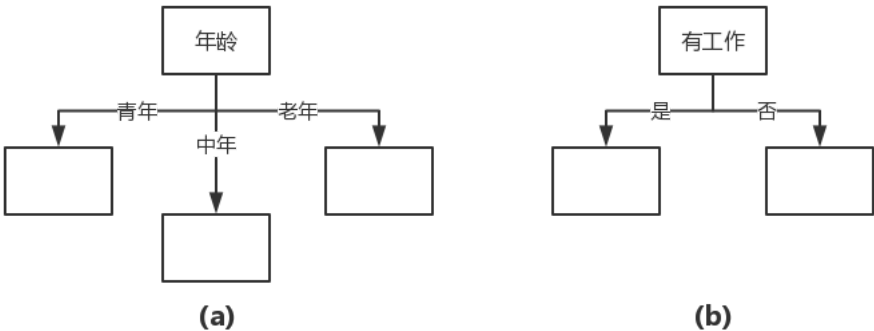
通常特征选择的标准是信息增益(information gain)或信息增益比，为了简单，本文章使用信息增益作为选择特征的标准。那么，什么是信息增益？

在继续讲解之前，先看一组实例，贷款申请样本数据表。

ID	年龄
1	青年
2	青年
3	青年
4	青年
5	青年
6	中年
7	中年
8	中年
9	中年
10	中年
11	老年
12	老年
13	老年
14	老年
15	老年

希望通过所给的训练数据学习一个贷款申请的决策树，用以对未来的贷款申请进行分类，即当新的客户提出贷款申请时，根据申请人的特征利用决策树决定是否批准贷款申请。

特征选择就是决定用哪个特征来划分特征空间。比如，我们通过上述数据表得到两个可能的决策树，分别由两个不同特征的根结点构成。



<http://blog.csdn.net/c406495762>

现在我们要决定选择第一个特征还是第二个特征，但是问题是：究竟选择哪个特征更好些？这就要求确定选择特征的准则。直观上，如果一个特征具有更好的分类能力，或者说，按照这一特征将训练数据集分割成子集，使得各个子集在当前条件下有最好的分类，那么就更应该选择这个特征。信息增益就能够很好地表示这一直观的准则。那么什么是信息增益呢？

在划分数据集之前之后信息发生的变化成为信息增益，知道如何计算信息增益，我们就可以计算每个特征值划分数据集获得的信息增益，获得信息增益最高的特征就是最好的选择。

1. 香农熵

在划分数据集之前之后信息发生的变化称为 **信息增益**，知道如何计算信息增益，就可以计算每个特征值划分数据集获得的 **信息增益**，获得 **信息增益** 最高的特征就是最好的选择。在可以评测哪种数据划分方式是最好的数据划分之前，我们必须学习如何计算 **信息增益**。集合信息的度量方式称为 **香农熵** 或者简称为 **熵**，这个名字来源于信息论之父克劳德·香农。

如果看不明白什么是信息增益 (information gain) 和熵 (entropy), 请不要着急——它们自诞生的那一天起, 就注定会令世人十分费解。克劳德·香农写完信息论之后, 约翰·冯·诺依曼建议使用“熵”这个术语, 因为大家都不知道它是什么意思。

熵定义为信息的期望值, 在明晰这个概念之前, 我们必须知道信息的定义。如果待分类的事务可能划分在多个分类之中, 则符号 x_i 的信息定义为

$$I(x_i) = -\log_2 p(x_i)$$

其中 $p(x_i)$ 是选择该分类的概率, 就是个固定公式, 记住就ok了。

为了计算熵, 我们需要计算所有类别, 所有可能值包含的信息期望值, 通过下面的公式得到:

$$H = -\sum_{i=1}^n p(x_i) \log_2 p(x_i)$$

其中 n 是分类的数目。

当熵中的概率由数据估计(特别是最大似然估计)得到时, 所对应的熵称为 **经验熵**(empirical entropy)。定义贷款申请样本数据表中的数据为训练数据集D, 则训练数据集D的经验熵为H(D), |D|表示其样本容量, 及样本个数。设有K个类 C_k , $k = 1, 2, 3, \dots, K$, $|C_k|$ 为属于类 C_k 的样本个数, 这经验熵公式可以写为:

$$H(D) = -\sum_{k=1}^K \frac{|C_k|}{|D|} \log_2 \frac{|C_k|}{|D|}$$

根据此公式计算经验熵H(D), 分析贷款申请样本数据表中的数据。最终分类结果只有两类, 即放贷和不放贷。根据表中的数据统计可知, 在15个数据中, 9个数据的结果为放贷, 6个数据的结果为不放贷。所以数据集D的经验熵H(D)为:

$$H(D) = -\frac{9}{15} \log_2 \frac{9}{15} - \frac{6}{15} \log_2 \frac{6}{15} = 0.971$$

经过计算可知, 数据集D的经验熵H(D)的值为0.971。

下面就可以准备编写代码来进行验证, 在这之前, 先对数据集进行属性标注。

- 年龄: 0代表青年, 1代表中年, 2代表老年;
- 有工作: 0代表否, 1代表是;
- 有自己的房子: 0代表否, 1代表是;
- 信贷情况: 0代表一般, 1代表好, 2代表非常好;
- 类别(是否给贷款): no代表否, yes代表是。

```
1  from math import log
2
3  """
4  Parameters:
5      无
6  Returns:
7      dataSet - 数据集
8      labels - 分类属性
9  """
10 # 函数说明: 创建测试数据集
11 def createDataSet():
12     dataSet = [[0, 0, 0, 0, 'no'], #数据集
13                [0, 0, 0, 1, 'no'],
14                [0, 1, 0, 1, 'yes'],
15                [0, 1, 1, 0, 'yes'],
16                [0, 0, 0, 0, 'no'],
17                [1, 0, 0, 0, 'no'],
18                [1, 0, 0, 1, 'no'],
19                [1, 1, 1, 1, 'yes'],
20                [1, 0, 1, 2, 'yes'],
21                [1, 0, 1, 2, 'yes'],
22                [2, 0, 1, 2, 'yes'],
23                [2, 0, 1, 1, 'yes'],
24                [2, 1, 0, 1, 'yes'],
```

```

24         [2, 1, 0, 2, 'yes'],
25         [2, 0, 0, 0, 'no']]
26     labels = ['年龄', '有工作', '有自己的房子', '信贷情况'] # 分类属性
27     return dataSet, labels # 返回数据集和分类属性
28
29 """
30 Parameters:
31     dataSet - 数据集
32 Returns:
33     shannonEnt - 经验熵(香农熵)
34 """
35 # 函数说明: 计算给定数据集的经验熵(香农熵)
36 def calcShannonEnt(dataSet):
37     numEntires = len(dataSet) # 返回数据集的行数
38     labelCounts = {} # 保存每个标签(Label)出现次数的字典
39     for featVec in dataSet: # 对每组特征向量进行统计
40         currentLabel = featVec[-1] # 提取标签(Label)信息
41         if currentLabel not in labelCounts.keys(): # 如果标签(Label)没有放入统计次数的字典, 添加进去
42             labelCounts[currentLabel] = 0
43         labelCounts[currentLabel] += 1 # Label计数
44     shannonEnt = 0.0 # 经验熵(香农熵)
45     for key in labelCounts: # 计算香农熵
46         prob = float(labelCounts[key]) / numEntires # 选择该标签(Label)的概率
47         shannonEnt -= prob * log(prob, 2) # 利用公式计算
48     return shannonEnt # 返回经验熵(香农熵)
49
50 if __name__ == '__main__':
51     dataSet, features = createDataSet()
52     print(dataSet)
53     print(calcShannonEnt(dataSet))
54
55
56 >>>
57 [[0, 0, 0, 0, 'no'], [0, 0, 0, 1, 'no'], [0, 1, 0, 1, 'yes'], [0, 1, 1, 0, 'yes'], [0, 0, 0, 0, 'no'], [1, 0, 0, 0, 'no'], [1, 0, 0, 1, 'yes'], [1, 0, 0, 0, 'no'], [1, 0, 0, 1, 'yes']]
58 0.9709505944546686

```

2. 信息增益

如何选择特征, 需要看 **信息增益**。也就是说, **信息增益** 是相对于特征而言的, **信息增益** 越大, 特征对最终的分类结果影响也就越大, 所以我们应该选择对最终分类结果影响最大的那个特征作为分类特征。在讲解信息增益定义之前, 还需要明确一个概念, **条件熵**。熵我们知道是什么, **条件熵** 又是个什么鬼? **条件熵** $H(Y|X)$ 表示在已知随机变量X的条件下随机变量Y的不确定性, 随机变量X给定的条件下随机变量Y的 **条件熵**(conditional entropy) $H(Y|X)$, 定义X在给定Y条件下的条件概率分布的熵对X的数学期望:

$$H(Y|X) = \sum_{i=1}^n p_i H(Y|X = x_i)$$

<http://blog.csdn.net/c406495762>

这里, $p_i = P(X = x_i)$, $i = 1, 2, \dots, n$ 。同理, 当 **条件熵** 中的概率由数据估计(特别是极大似然估计)得到时, 所对应的 **条件熵** 成为 **条件经验熵**(empirical conditional entropy)。

接下来说 **信息增益**, 前面也提到了, **信息增益** 是相对于特征而言的。所以, 特征A对训练数据集D的信息增益 $g(D, A)$, 定义为集合D的经验熵 $H(D)$ 与特征A给定条件下D的经验条件熵 $H(D|A)$ 之差, 即

$$g(D, A) = H(D) - H(D|A)$$

一般地, 熵 $H(D)$ 与条件熵 $H(D|A)$ 之差称为 **互信息**(mutual information)。决策树学习中的 **信息增益** 等价于训练数据集中类与特征的 **互信息**。

设特征A有n个不同的取值 a_1, a_2, \dots, a_n , 根据特征A的取值将D划分为n个子集 D_1, D_2, \dots, D_n , $|D_i|$ 为 D_i 的样本个数。记子集 D_i 中属于 C_k 的样本的集合为 D_{ik} , 即 $D_{ik} = D_i \cap C_k$, $|D_{ik}|$ 为 D_{ik} 的样本个数。于是经验条件熵的公式可以写为:

$$H(D|A) = \sum_{i=1}^n \frac{|D_i|}{|D|} H(D_i) = - \sum_{i=1}^n \frac{|D_i|}{|D|} \sum_{k=1}^K \frac{|D_{ik}|}{|D_i|} \log_2 \frac{|D_{ik}|}{|D_i|}$$

啥，听不懂，哈哈，正常，下面我们看例子说话，回头你再看概念就懂了。还是以贷款申请样本数据表为例，年龄这一行为特征 A_1 ，一共三个类别，分别是：青年、中年和老年。

从数据表格中可以看出年龄是青年的数据在训练数据集出现的概率是5 / 15，同理，年龄是中年的数据在训练数据集出现的概率也都是5 / 15。年龄是青年的数据的最终得到贷款的概率为2 / 5，因为在五个数据中，只有两个数据显示拿到了最终的贷款，同理，年龄是中年的数据最终得到贷款的概率分别为3 / 5、4 / 5。

所以计算年龄的信息增益，过程如下：

$$\begin{aligned} g(D, A_1) &= H(D) - \left[\frac{5}{15} H(D_1) + \frac{5}{15} H(D_2) + \frac{5}{15} H(D_3) \right] \\ &= 0.971 \\ &\quad - \left[\frac{5}{15} \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) + \frac{5}{15} \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) \right. \\ &\quad \left. + \frac{5}{15} \left(-\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} \right) \right] = 0.971 - 0.888 = 0.083 \end{aligned}$$

同理，计算其余特征的信息增益 $g(D, A_2)$ 、 $g(D, A_3)$ 和 $g(D, A_4)$ 。分别为：

$$\begin{aligned} g(D, A_2) &= H(D) - \left[\frac{5}{15} H(D_1) + \frac{10}{15} H(D_2) \right] \\ &= 0.971 - \left[\frac{5}{15} \times 0 + \frac{10}{15} \left(-\frac{4}{10} \log_2 \frac{4}{10} - \frac{6}{10} \log_2 \frac{6}{10} \right) \right] \\ &= 0.971 - 0.647 = 0.324 \end{aligned}$$

$$\begin{aligned} g(D, A_3) &= H(D) - \left[\frac{6}{15} H(D_1) + \frac{9}{15} H(D_2) \right] \\ &= 0.971 - \left[\frac{6}{15} \times 0 + \frac{9}{15} \left(-\frac{3}{9} \log_2 \frac{3}{9} - \frac{6}{9} \log_2 \frac{6}{9} \right) \right] \\ &= 0.971 - 0.551 = 0.420 \end{aligned}$$

$$g(D, A_4) = 0.971 - 0.608 = 0.363$$

最后，比较特征的信息增益大小，发现特征 A_3 (有自己的房子) 的信息增益值最大，所以选择 A_3 作为最优特征。

```

1  from math import log
2
3  """
4  Parameters:
5      dataSet - 数据集
6  Returns:
7      shannonEnt - 经验熵(香农熵)
8  """
9  # 函数说明: 计算给定数据集的经验熵(香农熵)
10 def calcShannonEnt(dataSet):
11     numEntires = len(dataSet)          # 返回数据集的行数
12     labelCounts = {}                  # 保存每个标签(Label)出现次数的字典
13     for featVec in dataSet:           # 对每组特征向量进行统计
14         currentLabel = featVec[-1]     # 提取标签(Label)信息
15         if currentLabel not in labelCounts.keys(): # 如果标签(Label)没有放入统计次数的字典, 添加进去
16             labelCounts[currentLabel] = 0
17             labelCounts[currentLabel] += 1 # Label计数

```

```

18     shannonEnt = 0.0 #经验熵(香农熵)
19     for key in labelCounts: #计算香农熵
20         prob = float(labelCounts[key]) / numEntires #选择该标签(Label)的概率
21         shannonEnt -= prob * log(prob, 2) #利用公式计算
22     return shannonEnt #返回经验熵(香农熵)
23
24 """
25 Parameters:
26     无
27 Returns:
28     dataSet - 数据集
29     labels - 分类属性
30 """
31 # 函数说明: 创建测试数据集
32 def createDataSet():
33     dataSet = [[0, 0, 0, 0, 'no'], #数据集
34                [0, 0, 0, 1, 'no'],
35                [0, 1, 0, 1, 'yes'],
36                [0, 1, 1, 0, 'yes'],
37                [0, 0, 0, 0, 'no'],
38                [1, 0, 0, 0, 'no'],
39                [1, 0, 0, 1, 'no'],
40                [1, 1, 1, 1, 'yes'],
41                [1, 0, 1, 2, 'yes'],
42                [1, 0, 1, 2, 'yes'],
43                [2, 0, 1, 2, 'yes'],
44                [2, 0, 1, 1, 'yes'],
45                [2, 1, 0, 1, 'yes'],
46                [2, 1, 0, 2, 'yes'],
47                [2, 0, 0, 0, 'no']]
48     labels = ['年龄', '有工作', '有自己的房子', '信贷情况'] #分类属性
49     return dataSet, labels #返回数据集和分类属性
50
51 """
52 Parameters:
53     dataSet - 待划分的数据集
54     axis - 划分数据集的特征
55     value - 需要返回的特征的值
56 Returns:
57     无
58 """
59 # 函数说明: 按照给定特征划分数据集
60 def splitDataSet(dataSet, axis, value):
61     retDataSet = [] #创建返回的数据集列表
62     for featVec in dataSet: #遍历数据集
63         if featVec[axis] == value:
64             reducedFeatVec = featVec[:axis] #去掉axis特征
65             reducedFeatVec.extend(featVec[axis+1:]) #将符合条件的添加到返回的数据集
66             retDataSet.append(reducedFeatVec)
67     return retDataSet #返回划分后的数据集
68
69 """
70 Parameters:
71     dataSet - 数据集
72 Returns:
73     bestFeature - 信息增益最大的(最优)特征的索引值
74 """
75 # 函数说明: 选择最优特征
76 def chooseBestFeatureToSplit(dataSet):
77     numFeatures = len(dataSet[0]) - 1 #特征数量
78     baseEntropy = calcShannonEnt(dataSet) #计算数据集的香农熵
79     bestInfoGain = 0.0 #信息增益
80     bestFeature = -1 #最优特征的索引值
81     for i in range(numFeatures): #遍历所有特征
82         #获取dataSet的第i个所有特征
83         featList = [example[i] for example in dataSet]
84         uniqueVals = set(featList) #创建set集合{},元素不可重复
85         newEntropy = 0.0 #经验条件熵
86         for value in uniqueVals: #计算信息增益
87             subDataSet = splitDataSet(dataSet, i, value) #subDataSet划分后的子集
88             prob = len(subDataSet) / float(len(dataSet)) #计算子集的概率
89             newEntropy += prob * calcShannonEnt(subDataSet) #根据公式计算经验条件熵
90         infoGain = baseEntropy - newEntropy #信息增益
91         print("第%d个特征的信息增益为%.3f" % (i, infoGain)) #打印每个特征的信息增益
92         if (infoGain > bestInfoGain): #计算信息增益
93             bestInfoGain = infoGain #更新信息增益 找到最大的信息增益
94             bestFeature = i
95     return bestFeature

```



```

89         bestInfoGain = infoGain           #更新信息增益，找到最大的信息增益
90         bestFeature = i                   #记录信息增益最大的特征的索引值
91     return bestFeature                     #返回信息增益最大的特征的索引值
92
93
94 if __name__ == '__main__':
95     dataSet, features = createDataSet()
96     print("最优特征索引值:" + str(chooseBestFeatureToSplit(dataSet)))

```

```

1  >>>
2  第0个特征的增益为0.083
3  第1个特征的增益为0.324
4  第2个特征的增益为0.420
5  第3个特征的增益为0.363
6  最优特征索引值:2

```

最佳索引是2，也就是特征A3(有自己的房子)。

2) 决策树的生成和修剪

ID3算法的核心是在决策树各个结点上对应 **信息增益** 准则选择特征，递归地构建决策树，ID3相当于用极大似然法进行概率模型的选择。

具体方法是：

从根结点(root node)开始，对结点计算所有可能的特征的信息增益，选择信息增益最大的特征作为结点的特征，由该特征的不同取值建立子节点；
再对子结点递归地调用以上方法，构建决策树；
直到所有特征的信息增益均很小或没有特征可以选择为止，最后得到一个决策树。

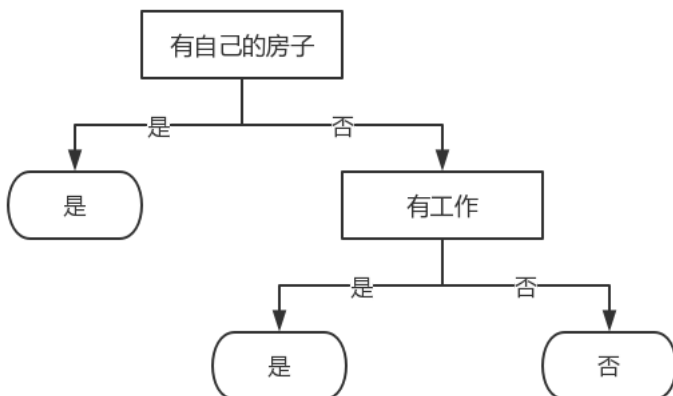
上面得到了特征A3(有自己的房子)的信息增益值最大，所以选择特征A3作为根结点的特征。它将训练集D划分为两个子集D1(A3取值为“是”)和D2(A3取值为“否”)。由于D1只有同一类的样本点，所以它成为一个叶结点，结点的类标记为“是”。对D2则需要从特征A1(年龄)，A2(有工作)和A4(信贷情况)中选择新的特征，计算各个特征的信息增益：

$$g(D2,A1) = H(D2) - H(D2 | A1) = 0.251$$

$$g(D2,A2) = H(D2) - H(D2 | A2) = 0.918$$

$$g(D2,A3) = H(D2) - H(D2 | A3) = 0.474$$

根据计算，选择信息增益最大的特征A2(有工作)作为结点的特征。由于A2有两个可能取值，从这一结点引出两个子结点：一个对应“是”(有工作)的子结点，包含3个样本，它们属于同一类，所以这是一个叶结点，类标记为“是”；另一个是对应“否”(无工作)的子结点，包含6个样本，它们也属于同一类，所以这也是一个叶结点，类标记为“否”。这样就生成了一个决策树，该决策树只用了两个特征(有两个内部结点)，生成的决策树如下图所示：



<http://blog.csdn.net/c406495762>

```

1  from math import log
2  import operator
3
4  """
5  Parameters:

```



```

6 Parameters:
7     dataSet - 数据集
8 Returns:
9     shannonEnt - 经验熵(香农熵)
10 """
11 # 函数说明: 计算给定数据集的经验熵(香农熵)
12 def calcShannonEnt(dataSet):
13     numEntires = len(dataSet) # 返回数据集的行数
14     labelCounts = {} # 保存每个标签(Label)出现次数的字典
15     for featVec in dataSet: # 对每组特征向量进行统计
16         currentLabel = featVec[-1] # 提取标签(Label)信息
17         if currentLabel not in labelCounts.keys(): # 如果标签(Label)没有放入统计次数的字典, 添加进去
18             labelCounts[currentLabel] = 0
19         labelCounts[currentLabel] += 1 # Label计数
20     shannonEnt = 0.0 # 经验熵(香农熵)
21     for key in labelCounts: # 计算香农熵
22         prob = float(labelCounts[key]) / numEntires # 选择该标签(Label)的概率
23         shannonEnt -= prob * log(prob, 2) # 利用公式计算
24     return shannonEnt # 返回经验熵(香农熵)
25
26 """
27 Parameters:
28     无
29 Returns:
30     dataSet - 数据集
31     labels - 特征标签
32 """
33 # 函数说明: 创建测试数据集
34 def createDataSet():
35     dataSet = [[0, 0, 0, 0, 'no'], # 数据集
36                [0, 0, 0, 1, 'no'],
37                [0, 1, 0, 1, 'yes'],
38                [0, 1, 1, 0, 'yes'],
39                [0, 0, 0, 0, 'no'],
40                [1, 0, 0, 0, 'no'],
41                [1, 0, 0, 1, 'no'],
42                [1, 1, 1, 1, 'yes'],
43                [1, 0, 1, 2, 'yes'],
44                [1, 0, 1, 2, 'yes'],
45                [2, 0, 1, 2, 'yes'],
46                [2, 0, 1, 1, 'yes'],
47                [2, 1, 0, 1, 'yes'],
48                [2, 1, 0, 2, 'yes'],
49                [2, 0, 0, 0, 'no']]
50     labels = ['年龄', '有工作', '有自己的房子', '信贷情况'] # 特征标签
51     return dataSet, labels # 返回数据集和分类属性
52
53 """
54 Parameters:
55     dataSet - 待划分的数据集
56     axis - 划分数据集的特征
57     value - 需要返回的特征的值
58 Returns:
59     无
60 """
61 # 函数说明: 按照给定特征划分数据集
62 def splitDataSet(dataSet, axis, value):
63     retDataSet = [] # 创建返回的数据集列表
64     for featVec in dataSet: # 遍历数据集
65         if featVec[axis] == value:
66             reducedFeatVec = featVec[:axis] # 去掉axis特征
67             reducedFeatVec.extend(featVec[axis+1:]) # 将符合条件的添加到返回的数据集
68         retDataSet.append(reducedFeatVec)
69     return retDataSet # 返回划分后的数据集
70
71 """
72 Parameters:
73     dataSet - 数据集
74 Returns:
75     bestFeature - 信息增益最大的(最优)特征的索引值
76 """
77 # 函数说明: 选择最优特征
78 def chooseBestFeatureToSplit(dataSet):
79     numFeatures = len(dataSet[0]) - 1 # 特征数量
80     baseEntropy = calcShannonEnt(dataSet) # 计算数据集的香农熵
81     bestInfoGain = 0.0 # 信息增益
82     bestFeature = -1 # 最佳特征索引值

```

```

77     bestFeature = -1 #最优特征的索引值
78     for i in range(numFeatures): #遍历所有特征
79         #获取dataSet的第i个所有特征
80         featList = [example[i] for example in dataSet]
81         uniqueVals = set(featList) #创建set集合{},元素不可重复
82         newEntropy = 0.0 #经验条件熵
83         for value in uniqueVals: #计算信息增益
84             subDataSet = splitDataSet(dataSet, i, value) #subDataSet划分后的子集
85             prob = len(subDataSet) / float(len(dataSet)) #计算子集的概率
86             newEntropy += prob * calcShannonEnt(subDataSet) #根据公式计算经验条件熵
87             infoGain = baseEntropy - newEntropy #信息增益
88             # print("第%d个特征的增益为%.3f" % (i, infoGain)) #打印每个特征的信息增益
89             if (infoGain > bestInfoGain): #计算信息增益
90                 bestInfoGain = infoGain #更新信息增益,找到最大的信息增益
91                 bestFeature = i #记录信息增益最大的特征的索引值
92     return bestFeature #返回信息增益最大的特征的索引值
93
94 """
95 Parameters:
96     classList - 类标签列表
97 Returns:
98     sortedClassCount[0][0] - 出现此处最多的元素(类标签)
99 """
100 # 函数说明:统计classList中出现此处最多的元素(类标签)
101 def majorityCnt(classList):
102     classCount = {}
103     for vote in classList: #统计classList中每个元素出现的次数
104         if vote not in classCount.keys(): classCount[vote] = 0
105         classCount[vote] += 1
106     sortedClassCount = sorted(classCount.items(), key = operator.itemgetter(1), reverse = True) #根据字典的值降序排序
107     return sortedClassCount[0][0] #返回classList中出现次数最多的元素
108
109 """
110 Parameters:
111     dataSet - 训练数据集
112     labels - 分类属性标签
113     featLabels - 存储选择的最优特征标签
114 Returns:
115     myTree - 决策树
116 """
117 # 函数说明: 创建决策树
118 def createTree(dataSet, labels, featLabels):
119     classList = [example[-1] for example in dataSet] #取分类标签(是否放贷:yes or no)
120     if classList.count(classList[0]) == len(classList): #如果类别完全相同则停止继续划分
121         return classList[0]
122     if len(dataSet[0]) == 1: #遍历完所有特征时返回出现次数最多的类标签
123         return majorityCnt(classList)
124     bestFeat = chooseBestFeatureToSplit(dataSet) #选择最优特征
125     bestFeatLabel = labels[bestFeat] #最优特征的标签
126     featLabels.append(bestFeatLabel)
127     myTree = {bestFeatLabel: {}} #根据最优特征的标签生成树
128     del(labels[bestFeat]) #删除已经使用特征标签
129     featValues = [example[bestFeat] for example in dataSet] #得到训练集中所有最优特征的属性值
130     uniqueVals = set(featValues) #去掉重复的属性值
131     for value in uniqueVals: #遍历特征, 创建决策树。
132         myTree[bestFeatLabel][value] = createTree(splitDataSet(dataSet, bestFeat, value), labels, featLabels)
133     return myTree
134
135 if __name__ == '__main__':
136     dataSet, labels = createDataSet()
137     featLabels = []
138     myTree = createTree(dataSet, labels, featLabels)
139     print(myTree)
140
141
142 1 >>>
143 2 {'有自己的房子': {0: {'有工作': {0: 'no', 1: 'yes'}}, 1: 'yes'}}

```

可见，决策树构建完成了。很明显这个决策树看着别扭，虽然能看懂，但是如果多点的结点，就不一定了。

下面我们使用强大的Matplotlib绘制决策树。

3、决策树的可视化

可视化需要用到的函数：

- getNumLeafs：获取决策树叶子结点的数目
- getTreeDepth：获取决策树的层数
- plotNode：绘制结点
- plotMidText：标注有向边属性值
- plotTree：绘制决策树
- createPlot：创建绘制面板
- 为了显示中文，需要设置FontProperties

```
1  from matplotlib.font_manager import FontProperties
2  import matplotlib.pyplot as plt
3  from math import log
4  import operator
5
6  """
7  Parameters:
8      dataSet - 数据集
9  Returns:
10     shannonEnt - 经验熵(香农熵)
11 """
12 # 函数说明: 计算给定数据集的经验熵(香农熵)
13 def calcShannonEnt(dataSet):
14     numEntires = len(dataSet) #返回数据集的行数
15     labelCounts = {} #保存每个标签(Label)出现次数的字典
16     for featVec in dataSet: #对每组特征向量进行统计
17         currentLabel = featVec[-1] #提取标签(Label)信息
18         if currentLabel not in labelCounts.keys(): #如果标签(Label)没有放入统计次数的字典,添加进去
19             labelCounts[currentLabel] = 0
20         labelCounts[currentLabel] += 1 #Label计数
21     shannonEnt = 0.0 #经验熵(香农熵)
22     for key in labelCounts: #计算香农熵
23         prob = float(labelCounts[key]) / numEntires #选择该标签(Label)的概率
24         shannonEnt -= prob * log(prob, 2) #利用公式计算
25     return shannonEnt #返回经验熵(香农熵)
26
27 """
28 Parameters:
29     无
30 Returns:
31     dataSet - 数据集
32     labels - 特征标签
33 """
34 # 函数说明: 创建测试数据集
35 def createDataSet():
36     dataSet = [[0, 0, 0, 0, 'no'], #数据集
37                [0, 0, 0, 1, 'no'],
38                [0, 1, 0, 1, 'yes'],
39                [0, 1, 1, 0, 'yes'],
40                [0, 0, 0, 0, 'no'],
41                [1, 0, 0, 0, 'no'],
42                [1, 0, 0, 1, 'no'],
43                [1, 1, 1, 1, 'yes'],
44                [1, 0, 1, 2, 'yes'],
45                [1, 0, 1, 2, 'yes'],
46                [2, 0, 1, 2, 'yes'],
47                [2, 0, 1, 1, 'yes'],
48                [2, 1, 0, 1, 'yes'],
49                [2, 1, 0, 2, 'yes'],
50                [2, 0, 0, 0, 'no']]
51     labels = ['年龄', '有工作', '有自己的房子', '信贷情况'] #特征标签
52     return dataSet, labels #返回数据集和分类属性
53
54 """
55 Parameters:
56     dataSet - 待划分的数据集
57     axis - 划分数据集的特征
58     value - 需要返回的特征的值
59 Returns:
60 """
```

```

57     无
58     """
59     # 函数说明: 按照给定特征划分数据集
60     def splitDataSet(dataSet, axis, value):
61         retDataSet = []                                # 创建返回的数据集列表
62         for featVec in dataSet:                        # 遍历数据集
63             if featVec[axis] == value:
64                 reducedFeatVec = featVec[:axis]        # 去掉axis特征
65                 reducedFeatVec.extend(featVec[axis+1:]) # 将符合条件的添加到返回的数据集
66                 retDataSet.append(reducedFeatVec)
67         return retDataSet                              # 返回划分后的数据集
68
69     """
70     Parameters:
71         dataSet - 数据集
72     Returns:
73         bestFeature - 信息增益最大的(最优)特征的索引值
74     """
75     # 函数说明: 选择最优特征
76     def chooseBestFeatureToSplit(dataSet):
77         numFeatures = len(dataSet[0]) - 1             # 特征数量
78         baseEntropy = calcShannonEnt(dataSet)         # 计算数据集的香农熵
79         bestInfoGain = 0.0                             # 信息增益
80         bestFeature = -1                               # 最优特征的索引值
81         for i in range(numFeatures):                  # 遍历所有特征
82             # 获取dataSet的第i个所有特征
83             featList = [example[i] for example in dataSet]
84             uniqueVals = set(featList)                 # 创建set集合{}, 元素不可重复
85             newEntropy = 0.0                           # 经验条件熵
86             for value in uniqueVals:                   # 计算信息增益
87                 subDataSet = splitDataSet(dataSet, i, value) # subDataSet划分后的子集
88                 prob = len(subDataSet) / float(len(dataSet)) # 计算子集的概率
89                 newEntropy += prob * calcShannonEnt(subDataSet) # 根据公式计算经验条件熵
90             infoGain = baseEntropy - newEntropy        # 信息增益
91             # print("第%d个特征的增益为%.3f" % (i, infoGain)) # 打印每个特征的信息增益
92             if (infoGain > bestInfoGain):              # 计算信息增益
93                 bestInfoGain = infoGain                # 更新信息增益, 找到最大的信息增益
94                 bestFeature = i                        # 记录信息增益最大的特征的索引值
95         return bestFeature                             # 返回信息增益最大的特征的索引值
96
97     """
98     Parameters:
99         classList - 类标签列表
100     Returns:
101         sortedClassCount[0][0] - 出现此处最多的元素(类标签)
102     """
103     # 函数说明: 统计classList中出现此处最多的元素(类标签)
104     def majorityCnt(classList):
105         classCount = {}
106         for vote in classList: # 统计classList中每个元素出现的次数
107             if vote not in classCount.keys(): classCount[vote] = 0
108             classCount[vote] += 1
109         sortedClassCount = sorted(classCount.items(), key = operator.itemgetter(1), reverse = True) # 根据字典的值降序排序
110         return sortedClassCount[0][0] # 返回classList中出现次数最多的元素
111
112     """
113     Parameters:
114         dataSet - 训练数据集
115         labels - 分类属性标签
116         featLabels - 存储选择的最优特征标签
117     Returns:
118         myTree - 决策树
119     """
120     # 函数说明: 创建决策树
121     def createTree(dataSet, labels, featLabels):
122         classList = [example[-1] for example in dataSet] # 取分类标签(是否放贷: yes or no)
123         if classList.count(classList[0]) == len(classList): # 如果类别完全相同则停止继续划分
124             return classList[0]
125         if len(dataSet[0]) == 1: # 遍历完所有特征时返回出现次数最多的类标签
126             return majorityCnt(classList)
127         bestFeat = chooseBestFeatureToSplit(dataSet) # 选择最优特征
128         bestFeatLabel = labels[bestFeat] # 最优特征的标签
129         featLabels.append(bestFeatLabel)
130         myTree = {bestFeatLabel: {}} # 根据最优特征的标签生成树
131         del(labels[bestFeat]) # 删除已经使用特征标签

```

```

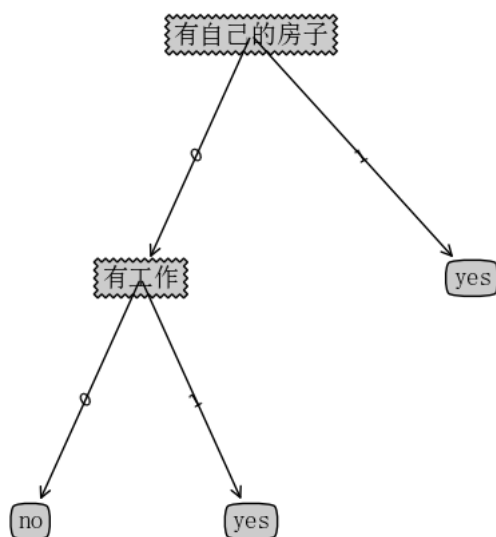
128     featValues = [example[bestFeat] for example in dataSet]#得到训练集中所有最优特征的属性值
129     uniqueVals = set(featValues)                                #去掉重复的属性值
130     for value in uniqueVals:                                    #遍历特征, 创建决策树。
131         myTree[bestFeatLabel][value] = createTree(splitDataSet(dataSet, bestFeat, value), labels, featLabels)
132     return myTree
133
134 """
135 Parameters:
136     myTree - 决策树
137 Returns:
138     numLeafs - 决策树的叶子结点的数目
139 """
140 # 函数说明: 获取决策树叶子结点的数目
141 def getNumLeafs(myTree):
142     numLeafs = 0                                #初始化叶子
143     #python3中myTree.keys()返回的是dict_keys, 不在是list,
144     #所以不能使用myTree.keys()[0]的方法获取结点属性, 可以使用list(myTree.keys())[0]
145     firstStr = next(iter(myTree))
146     secondDict = myTree[firstStr]                #获取下一组字典
147     for key in secondDict.keys():
148         if type(secondDict[key]).__name__=='dict':#测试该结点是否为字典, 如果不是字典, 代表此结点为叶子结点
149             numLeafs += getNumLeafs(secondDict[key])
150         else: numLeafs +=1
151     return numLeafs
152
153 """
154 Parameters:
155     myTree - 决策树
156 Returns:
157     maxDepth - 决策树的层数
158 """
159 # 函数说明: 获取决策树的层数
160 def getTreeDepth(myTree):
161     maxDepth = 0                                #初始化决策树深度
162     #python3中myTree.keys()返回的是dict_keys, 不在是list,
163     #所以不能使用myTree.keys()[0]的方法获取结点属性, 可以使用list(myTree.keys())[0]
164     firstStr = next(iter(myTree))
165     secondDict = myTree[firstStr]                #获取下一个字典
166     for key in secondDict.keys():
167         if type(secondDict[key]).__name__=='dict':    #测试该结点是否为字典, 如果不是字典, 代表此结点为叶子结点
168             thisDepth = 1 + getTreeDepth(secondDict[key])
169         else: thisDepth = 1
170         if thisDepth > maxDepth: maxDepth = thisDepth    #更新层数
171     return maxDepth
172
173 """
174 Parameters:
175     nodeTxt - 结点名
176     centerPt - 文本位置
177     parentPt - 标注的箭头位置
178     nodeType - 结点格式
179 Returns:
180     无
181 """
182 # 函数说明: 绘制结点
183 def plotNode(nodeTxt, centerPt, parentPt, nodeType):
184     arrow_args = dict(arrowstyle="<-")            #定义箭头格式
185     font = FontProperties(fname=r"c:\windows\fonts\simsun.ttc", size=14)    #设置中文字体
186     createPlot.ax1.annotate(nodeTxt, xy=parentPt,  xycoords='axes fraction',#绘制结点
187                             xytext=centerPt, textcoords='axes fraction',
188                             va="center", ha="center", bbox=nodeType, arrowprops=arrow_args, FontProperties=font)
189
190 """
191 Parameters:
192     cntrPt、parentPt - 用于计算标注位置
193     txtString - 标注的内容
194 Returns:
195     无
196 """
197 # 函数说明: 标注有向边属性值
198 def plotMidText(cntrPt, parentPt, txtString):
199     xMid = (parentPt[0]-cntrPt[0])/2.0 + cntrPt[0]#计算标注位置
200     yMid = (parentPt[1]-cntrPt[1])/2.0 + cntrPt[1]
201     createPlot.ax1.text(xMid, yMid, txtString, va="center", ha="center", rotation=30)
202
203 """
204 Parameters:

```

```

199     myTree - 决策树(字典)
200     parentPt - 标注的内容
201     nodeTxt - 结点名
202 Returns:
203     无
204 """
205 # 函数说明: 绘制决策树
206 def plotTree(myTree, parentPt, nodeTxt):
207     decisionNode = dict(boxstyle="sawtooth", fc="0.8") # 设置结点格式
208     leafNode = dict(boxstyle="round4", fc="0.8") # 设置叶结点格式
209     numLeafs = getNumLeafs(myTree) # 获取决策树叶结点数目,
210
211     depth = getTreeDepth(myTree) # 获取决策树层数
212     firstStr = next(iter(myTree)) # 下个字典
213     cntrPt = (plotTree.xOff + (1.0 + float(numLeafs))/2.0/plotTree.totalW, plotTree.yOff) # 中心位置
214     plotMidText(cntrPt, parentPt, nodeTxt) # 标注有向边属性值
215     plotNode(firstStr, cntrPt, parentPt, decisionNode) # 绘制结点
216     secondDict = myTree[firstStr] # 下一个字典,
217
218     plotTree.yOff = plotTree.yOff - 1.0/plotTree.totalD # y 偏移
219     for key in secondDict.keys():
220         if type(secondDict[key]).__name__ == 'dict': # 测试该结点是否为字典,
221
222             plotTree(secondDict[key], cntrPt, str(key)) # 不是叶结点,
223
224         else: # 如果是叶结点, 绘制叶结点,
225
226             plotTree.xOff = plotTree.xOff + 1.0/plotTree.totalW
227             plotNode(secondDict[key], (plotTree.xOff, plotTree.yOff), cntrPt, leafNode)
228             plotMidText((plotTree.xOff, plotTree.yOff), cntrPt, str(key))
229             plotTree.yOff = plotTree.yOff + 1.0/plotTree.totalD
230
231 """
232 Parameters:
233     inTree - 决策树(字典)
234 Returns:
235     无
236 """
237 # 函数说明: 创建绘制面板
238 def createPlot(inTree):
239     fig = plt.figure(1, facecolor='white') # 创建fig
240     fig.clf() # 清空fig
241     axprops = dict(xticks=[], yticks=[])
242     createPlot.ax1 = plt.subplot(111, frameon=False, **axprops) # 去掉x、y轴
243     plotTree.totalW = float(getNumLeafs(inTree)) # 获取决策树叶结点数目
244     plotTree.totalD = float(getTreeDepth(inTree)) # 获取决策树层数
245     plotTree.xOff = -0.5/plotTree.totalW; plotTree.yOff = 1.0; # x 偏移
246     plotTree(inTree, (0.5, 1.0), '') # 绘制决策树
247     plt.show() # 显示绘制结果
248
249 if __name__ == '__main__':
250     dataSet, labels = createDataSet()
251     featLabels = []
252     myTree = createTree(dataSet, labels, featLabels)
253     print(myTree)
254     createPlot(myTree)
255

```



<https://blog.csdn.net/TeFuirnever>

4、使用决策树执行分类

依靠训练数据构造了决策树之后，可以将它用于实际数据的分类。在执行数据分类时，需要决策树以及用于构造树的标签向量；然后，程序比较测试数据与决策树上的数值，递归执行该过程直到进入叶子结点；最后将测试数据定义为叶子结点所属的类型。

在构建决策树的代码，可以看到，有个 `featLabels` 参数。它是用来记录各个分类结点的，在用决策树做预测的时候，按顺序输入需要的分类结点的属性值即可。举个例子，比如用上述已经训练好的决策树做分类，只需要提供这个人是否有房子，是否有工作这两个信息即可，无需提供其他冗余的信息。

```

1  # -*- coding: UTF-8 -*-
2  from math import log
3  import operator
4
5  """
6  Parameters:
7      dataSet - 数据集
8  Returns:
9      shannonEnt - 经验熵(香农熵)
10 """
11 # 函数说明: 计算给定数据集的经验熵(香农熵)
12 def calcShannonEnt(dataSet):
13     numEntires = len(dataSet) #返回数据集的行数
14     labelCounts = {} #保存每个标签(Label)出现次数的字典
15     for featVec in dataSet: #对每组特征向量进行统计
16         currentLabel = featVec[-1] #提取标签(Label)信息
17         if currentLabel not in labelCounts.keys(): #如果标签(Label)没有放入统计次数的字典,添加进去
18             labelCounts[currentLabel] = 0
19         labelCounts[currentLabel] += 1 #Label计数
20     shannonEnt = 0.0 #经验熵(香农熵)
21     for key in labelCounts: #计算香农熵
22         prob = float(labelCounts[key]) / numEntires #选择该标签(Label)的概率
23         shannonEnt -= prob * log(prob, 2) #利用公式计算
24     return shannonEnt #返回经验熵(香农熵)
25
26 """
27 Parameters:
28     无
29 Returns:
30     dataSet - 数据集
31     labels - 特征标签
32 """
33 # 函数说明: 创建测试数据集
34 def createDataSet():
35     dataSet = [[0, 0, 0, 0, 'no'], #数据集
36                [0, 0, 0, 1, 'no'],
37                [0, 1, 0, 1, 'yes'],
38                [0, 1, 1, 0, 'yes'],
39                [1, 0, 0, 0, 'no'],
40                [1, 0, 0, 1, 'no'],
41                [1, 0, 1, 0, 'yes'],
42                [1, 0, 1, 1, 'yes'],
43                [1, 1, 0, 0, 'yes'],
44                [1, 1, 0, 1, 'yes'],
45                [1, 1, 1, 0, 'yes'],
46                [1, 1, 1, 1, 'yes']]
47     labels = ['s1', 's2', 's3', 's4']
48     return dataSet, labels
  
```



```

38         [0, 0, 0, 0, 'no'],
39         [1, 0, 0, 0, 'no'],
40         [1, 0, 0, 1, 'no'],
41         [1, 1, 1, 1, 'yes'],
42         [1, 0, 1, 2, 'yes'],
43         [1, 0, 1, 2, 'yes'],
44         [2, 0, 1, 2, 'yes'],
45         [2, 0, 1, 1, 'yes'],
46         [2, 1, 0, 1, 'yes'],
47         [2, 1, 0, 2, 'yes'],
48         [2, 0, 0, 0, 'no']]
49 labels = ['年龄', '有工作', '有自己的房子', '信贷情况']#特征标签
50 return dataSet, labels#返回数据集和分类属性
51
52 """
53 Parameters:
54     dataSet - 待划分的数据集
55     axis - 划分数据集的特征
56     value - 需要返回的特征的值
57 Returns:
58     无
59 """
60 # 函数说明: 按照给定特征划分数据集
61 def splitDataSet(dataSet, axis, value):
62     retDataSet = [] #创建返回的数据集列表
63     for featVec in dataSet: #遍历数据集
64         if featVec[axis] == value:
65             reducedFeatVec = featVec[:axis] #去掉axis特征
66             reducedFeatVec.extend(featVec[axis+1:])#将符合条件的添加到返回的数据集
67             retDataSet.append(reducedFeatVec)
68     return retDataSet #返回划分后的数据集
69
70 """
71 Parameters:
72     dataSet - 数据集
73 Returns:
74     bestFeature - 信息增益最大的(最优)特征的索引值
75 """
76 # 函数说明: 选择最优特征
77 def chooseBestFeatureToSplit(dataSet):
78     numFeatures = len(dataSet[0]) - 1 #特征数量
79     baseEntropy = calcShannonEnt(dataSet) #计算数据集的香农熵
80     bestInfoGain = 0.0 #信息增益
81     bestFeature = -1 #最优特征的索引值
82     for i in range(numFeatures): #遍历所有特征
83         #获取dataSet的第i个所有特征
84         featList = [example[i] for example in dataSet]
85         uniqueVals = set(featList) #创建set集合{},元素不可重复
86         newEntropy = 0.0 #经验条件熵
87         for value in uniqueVals: #计算信息增益
88             subDataSet = splitDataSet(dataSet, i, value) #subDataSet划分后的子集
89             prob = len(subDataSet) / float(len(dataSet)) #计算子集的概率
90             newEntropy += prob * calcShannonEnt(subDataSet)#根据公式计算经验条件熵
91         infoGain = baseEntropy - newEntropy #信息增益
92         # print("第%d个特征的增益为%.3f" % (i, infoGain)) #打印每个特征的信息增益
93         if (infoGain > bestInfoGain): #计算信息增益
94             bestInfoGain = infoGain #更新信息增益,找到最大的信息增益
95             bestFeature = i #记录信息增益最大的特征的索引值
96     return bestFeature #返回信息增益最大的特征的索引值
97
98 """
99 Parameters:
100     classList - 类标签列表
101 Returns:
102     sortedClassCount[0][0] - 出现此处最多的元素(类标签)
103 """
104 # 函数说明: 统计classList中出现此处最多的元素(类标签)
105 def majorityCnt(classList):
106     classCount = {}
107     for vote in classList: #统计classList中每个元素出现的次数
108         if vote not in classCount.keys():classCount[vote] = 0
109         classCount[vote] += 1
110     sortedClassCount = sorted(classCount.items(), key = operator.itemgetter(1), reverse = True)#根据字典的值降序排序
111     return sortedClassCount[0][0] #返回classList中出现次数最多的元素
112
113 """

```

```

109 Parameters:
110     dataSet - 训练数据集
111     labels - 分类属性标签
112     featLabels - 存储选择的最优特征标签
113 Returns:
114     myTree - 决策树
115 """
116 # 函数说明: 创建决策树
117 def createTree(dataSet, labels, featLabels):
118     classList = [example[-1] for example in dataSet] #取分类标签(是否放贷:yes or no)
119     if classList.count(classList[0]) == len(classList): #如果类别完全相同则停止继续划分
120         return classList[0]
121     if len(dataSet[0]) == 1: #遍历完所有特征时返回出现次数最多的类标签
122         return majorityCnt(classList)
123     bestFeat = chooseBestFeatureToSplit(dataSet) #选择最优特征
124     bestFeatLabel = labels[bestFeat] #最优特征的标签
125     featLabels.append(bestFeatLabel)
126     myTree = {bestFeatLabel: {}} #根据最优特征的标签生成树
127     del(labels[bestFeat]) #删除已经使用特征标签
128     featValues = [example[bestFeat] for example in dataSet] #得到训练集中所有最优特征的属性值
129     uniqueVals = set(featValues) #去掉重复的属性值
130     for value in uniqueVals: #遍历特征, 创建决策树。
131         myTree[bestFeatLabel][value] = createTree(splitDataSet(dataSet, bestFeat, value), labels, featLabels)
132     return myTree
133 """
134 Parameters:
135     inputTree - 已经生成的决策树
136     featLabels - 存储选择的最优特征标签
137     testVec - 测试数据列表, 顺序对应最优特征标签
138 Returns:
139     classLabel - 分类结果
140 """
141 # 函数说明: 使用决策树分类
142 def classify(inputTree, featLabels, testVec):
143     firstStr = next(iter(inputTree)) #获取决策树结点
144     secondDict = inputTree[firstStr] #下一个字典
145     featIndex = featLabels.index(firstStr)
146     for key in secondDict.keys():
147         if testVec[featIndex] == key:
148             if type(secondDict[key]).__name__ == 'dict':
149                 classLabel = classify(secondDict[key], featLabels, testVec)
150             else: classLabel = secondDict[key]
151     return classLabel
152
153 if __name__ == '__main__':
154     dataSet, labels = createDataSet()
155     featLabels = []
156     myTree = createTree(dataSet, labels, featLabels)
157     testVec = [0,1] #测试数据
158     result = classify(myTree, featLabels, testVec)
159     if result == 'yes':
160         print('放贷')
161     if result == 'no':
162         print('不放贷')

```

这里只增加了 `classify` 函数，用于决策树分类。输入测试数据[0,1]，代表没有房子，但是有工作，分类结果如下所示：

```

1 >>>
2 放贷

```

5、决策树的存储

构造决策树是很耗时的任务，即使处理很小的数据集，如前面的样本数据，也要花费几秒的时间，如果数据集很大，将会耗费很多计算时间。然而用创建好的决策树解决分类问题，则可以很快完成。因此，为了节省计算时间，最好能够在每次执行分类时调用已经构造好的决策树。

为了解决这个问题，需要使用Python模块 `pickle` 序列化对象。序列化对象可以在磁盘上保存对象，并在需要的时候读取出来。假设我们已经得到决策树{'有自己的房子': {0: {'有工作': {0: 'no', 1: 'yes'}}, 1: 'yes'}}，使用 `pickle.dump` 存储决策树。

```

1  import pickle
2
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  Returns:
8      无
9  """
10 # 函数说明: 存储决策树
11 def storeTree(inputTree, filename):
12     with open(filename, 'wb') as fw:
13         pickle.dump(inputTree, fw)
14
15 if __name__ == '__main__':
16     myTree = {'有自己的房子': {0: {'有工作': {0: 'no', 1: 'yes'}}, 1: 'yes'}}
17     storeTree(myTree, 'classifierStorage.txt')
18

```

运行代码，在该Python文件的相同目录下，会生成一个名为classifierStorage.txt的txt文件，这个是个二进制存储的文件，会生成即可，下次使用的时候：

```

1  import pickle
2
3  """
4  Parameters:
5      filename - 决策树的存储文件名
6  Returns:
7      pickle.load(fr) - 决策树字典
8  """
9  # 函数说明: 读取决策树
10 def grabTree(filename):
11     fr = open(filename, 'rb')
12     return pickle.load(fr)
13
14 if __name__ == '__main__':
15     myTree = grabTree('classifierStorage.txt')
16     print(myTree)
17

```

```

1  >>>
2  {'有自己的房子': {0: {'有工作': {0: 'no', 1: 'yes'}}, 1: 'yes'}}

```

6、Sklearn构建决策树分类器预测隐形眼镜类型

1. 准备数据

一共有24组数据，数据的Labels依次是age、prescript、astigmatic、tearRate、class，也就是第一列是年龄，第二列是症状，第三列是是否散光，第四列是眼泪数量，第五列是最终的分类标签。

```
lenses.txt - 记事本
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)
young    myope    no        reduced  no lenses
young    myope    no        normal   soft
young    myope    yes       reduced  no lenses
young    myope    yes       normal   hard
young    hyper    no        reduced  no lenses
young    hyper    no        normal   soft
young    hyper    yes       reduced  no lenses
young    hyper    yes       normal   hard
pre       myope    no        reduced  no lenses
pre       myope    no        normal   soft
pre       myope    yes       reduced  no lenses
pre       myope    yes       normal   hard
pre       hyper    no        reduced  no lenses
pre       hyper    no        normal   soft
pre       hyper    yes       reduced  no lenses
pre       hyper    yes       normal   no lenses
presbyopic myope    no        reduced  no lenses
presbyopic myope    no        normal   no lenses
presbyopic myope    yes       reduced  no lenses
presbyopic myope    yes       normal   hard
presbyopic hyper    no        reduced  no lenses
presbyopic hyper    no        normal   soft
presbyopic hyper    yes       reduced  no lenses
presbyopic hyper    yes       normal   no lenses
```

2. 测试算法

对string类型的数据序列化，需要先生成pandas数据，这样方便序列化工作。所以先原始数据->字典->pandas数据，编写代码如下：

```
1 import pandas as pd
2
3 if __name__ == '__main__':
4     with open('lenses.txt', 'r') as fr: #加载文件
5         lenses = [inst.strip().split('\t') for inst in fr.readlines()] #处理文件
6         lenses_target = [] #提取每组数据的类别，保存在列表里
7         for each in lenses:
8             lenses_target.append(each[-1])
9
10        lensesLabels = ['age', 'prescript', 'astigmatic', 'tearRate'] #特征标签
11        lenses_list = [] #保存Lenses数据的临时列表
12        lenses_dict = {} #保存Lenses数据的字典，用于生成pandas
13        for each_label in lensesLabels: #提取信息，生成字典
14            for each in lenses:
15                lenses_list.append(each[lensesLabels.index(each_label)])
16            lenses_dict[each_label] = lenses_list
17            lenses_list = []
18        print(lenses_dict) #打印字典信息
19        lenses_pd = pd.DataFrame(lenses_dict) #生成pandas.DataFrame
20        print(lenses_pd)

1 >>>
2 {'age': ['young', 'young', 'young', 'young', 'young', 'young', 'young', 'young', 'pre', 'pre', 'pre', 'pre', 'pre', 'pre', 'pre']}
3      age astigmatic prescript tearRate
4 0    young         no      myope reduced
5 1    young         no      myope normal
6 2    young         yes      myope reduced
7 3    young         yes      myope normal
8 4    young         no      hyper reduced
9 5    young         no      hyper normal
10 6    young         yes      hyper reduced
11 7    young         yes      hyper normal
12 8      pre         no      myope reduced
13 9      pre         no      myope normal
14 10     pre         yes      myope reduced
15 11     pre         yes      myope normal
```

16	12	pre	no	hyper	reduced
17	13	pre	no	hyper	normal
18	14	pre	yes	hyper	reduced
19	15	pre	yes	hyper	normal
20	16	presbyopic	no	myope	reduced
21	17	presbyopic	no	myope	normal
22	18	presbyopic	yes	myope	reduced
23	19	presbyopic	yes	myope	normal
24	20	presbyopic	no	hyper	reduced
25	21	presbyopic	no	hyper	normal
26	22	presbyopic	yes	hyper	reduced
27	23	presbyopic	yes	hyper	normal

接下来，将数据序列化，编写代码如下：

```

1 import pandas as pd
2 from sklearn.preprocessing import LabelEncoder
3
4 if __name__ == '__main__':
5     with open('lenses.txt', 'r') as fr:                #加载文件
6         lenses = [inst.strip().split('\t') for inst in fr.readlines()] #处理文件
7         lenses_target = []                             #提取每组数据的类别，保存在列表里
8         for each in lenses:
9             lenses_target.append(each[-1])
10
11     lensesLabels = ['age', 'prescript', 'astigmatic', 'tearRate'] #特征标签
12     lenses_list = []                                             #保存Lenses数据的临时列表
13     lenses_dict = {}                                             #保存Lenses数据的字典，用于生成pandas
14     for each_label in lensesLabels:                             #提取信息，生成字典
15         for each in lenses:
16             lenses_list.append(each[lensesLabels.index(each_label)])
17             lenses_dict[each_label] = lenses_list
18             lenses_list = []
19     # print(lenses_dict)                                         #打印字典信息
20     lenses_pd = pd.DataFrame(lenses_dict)                       #生成pandas.DataFrame
21     print(lenses_pd)                                            #打印pandas.DataFrame
22     le = LabelEncoder()                                         #创建LabelEncoder()对象，用于序列化
23     for col in lenses_pd.columns:                               #为每一列序列化
24         lenses_pd[col] = le.fit_transform(lenses_pd[col])
25     print(lenses_pd)

```

```

1 >>>
2      age astigmatic prescript tearRate
3 0    young        no    myope  reduced
4 1    young        no    myope  normal
5 2    young        yes    myope  reduced
6 3    young        yes    myope  normal
7 4    young        no    hyper  reduced
8 5    young        no    hyper  normal
9 6    young        yes    hyper  reduced
10 7    young        yes    hyper  normal
11 8     pre         no    myope  reduced
12 9     pre         no    myope  normal
13 10    pre         yes    myope  reduced
14 11    pre         yes    myope  normal
15 12    pre         no    hyper  reduced
16 13    pre         no    hyper  normal
17 14    pre         yes    hyper  reduced
18 15    pre         yes    hyper  normal
19 16  presbyopic    no    myope  reduced
20 17  presbyopic    no    myope  normal
21 18  presbyopic    yes    myope  reduced
22 19  presbyopic    yes    myope  normal
23 20  presbyopic    no    hyper  reduced
24 21  presbyopic    no    hyper  normal
25 22  presbyopic    yes    hyper  reduced
26 23  presbyopic    yes    hyper  normal
27

```

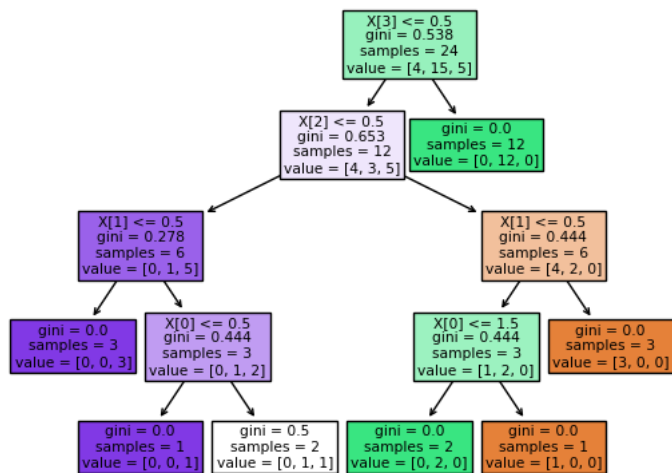
	age	astigmatic	prescript	tearRate
29	0	2	0	1
30	1	2	0	1
31	2	2	1	1
32	3	2	1	0
33	4	2	0	0
34	5	2	0	0
35	6	2	1	0
36	7	2	1	0
37	8	0	0	1
38	9	0	0	1
39	10	0	1	1
40	11	0	1	0
41	12	0	0	1
42	13	0	0	0
43	14	0	1	1
44	15	0	1	0
45	16	1	0	1
46	17	1	0	0
47	18	1	1	1
48	19	1	1	0
49	20	1	0	1
50	21	1	0	0
51	22	1	1	1
	23	1	1	0

可以看出数据已经顺利完成了序列化，可以进行fit()数据了。

```

1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4
5  from sklearn.preprocessing import LabelEncoder, OneHotEncoder
6  from sklearn.externals.six import StringIO
7  from sklearn import tree
8
9  if __name__ == '__main__':
10     with open('lenses.txt', 'r') as fr:
11         lenses = [inst.strip().split('\t') for inst in fr.readlines()] #加载文件
12         lenses_target = [] #处理文件
13         #提取每组数据的类别，保存在列表里
14         for each in lenses:
15             lenses_target.append(each[-1])
16         print(lenses_target)
17
18         lensesLabels = ['age', 'prescript', 'astigmatic', 'tearRate'] #特征标签
19         lenses_list = [] #保存Lenses数据的临时列表
20         lenses_dict = {} #保存Lenses数据的字典，用于生成pandas
21         #提取信息，生成字典
22         for each_label in lensesLabels:
23             for each in lenses:
24                 lenses_list.append(each[lensesLabels.index(each_label)])
25             lenses_dict[each_label] = lenses_list
26             lenses_list = []
27         # print(lenses_dict) #打印字典信息
28         lenses_pd = pd.DataFrame(lenses_dict) #生成pandas.DataFrame
29         # print(lenses_pd) #打印pandas.DataFrame
30         le = LabelEncoder() #创建LabelEncoder()对象，用于序列化
31         #序列化
32         for col in lenses_pd.columns:
33             lenses_pd[col] = le.fit_transform(lenses_pd[col])
34         #打印编码信息
35
36         clf = tree.DecisionTreeClassifier(max_depth = 4) #创建DecisionTreeClassifier()类
37         #使用数据，构建决策树
38         tree.plot_tree(clf.fit(lenses_pd.values.tolist(), lenses_target), filled=True)
39         plt.show()

```



<https://blog.csdn.net/TeFuimever>

官方可视化代码示例教程：https://scikit-learn.org/stable/auto_examples/tree/plot_iris_dtc.html

```
1 | print(clf.predict([[1,1,1,0]]))#预测
```

```
1 | >>>
2 | ['hard']
```

7、sklearn.tree.DecisionTreeClassifier

sklearn.tree.DecisionTreeClassifier是一个很好的模型，决策树算法就是通过它实现的，详细的看这个博客——[sklearn.tree.DecisionTreeClassifier\(\)函数解析（最清晰的解释）](#)

8、总结

决策树分类器就像带有终止块的流程图，终止块表示分类结果。开始处理数据集时，首先需要测量集合中数据的不一致性，也就是熵，然后寻找最优方案划分数数据集，直到数据集中的所有数据属于同一分类。ID3算法可以用于划分标称型数据集。构建决策树时，通常采用递归的方法将数据集转化为决策树。一般并不构造新的数据结构，而是使用Python语言内嵌的数据结构字典存储树节点信息。

使用Matplotlib的注解功能，可以将存储的树结构转化为容易理解的图形。Python语言的pickle模块可用于存储决策树的结构。隐形眼镜的例子表明决策树可能会产生过多的数据集划分，从而产生过度匹配数据集的问题。我们可以通过裁剪决策树，合并相邻的无法产生大量信息增益的叶节点，消除过度匹配问题。还有其他的决策树的构造算法，最流行的是C4.5和CART，第9章讨论回归问题时将介绍CART算法。

第2章、第3章讨论的是结果确定的分类算法，数据实例最终会被明确划分到某个分类中。

- 【机器学习】《机器学习实战》读书笔记及代码：第2章 - k-近邻算法
- 【机器学习】《机器学习实战》读书笔记及代码：第3章 - 决策树

下一章讨论的分类算法将不能完全确定数据实例应该划分到某个分类，或者只能给出数据实例属于给定分类的概率。