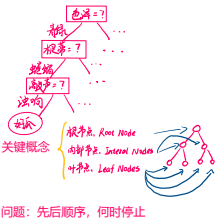


基本流程



概念

**熵** 对一种事物的不确定性称为熵 (概率和信息熵是随机事件的一体两面)

**信息** 消除不确定的事物  
调整概率  
作用: 排除干扰  
确定情况

**噪音** 不能消除某人对某件事情不确定性的事物

**数据** 噪音+信息

参照一个不确定的事件用作单位 例: 抛硬币的不确定性 概率: 50%正面, 50%反面; 相当于猜一次抛硬币的不确定性, 记为1bit

熵如何量化

**等概率** 8个等概率的不确定情况, 相当于抛8个硬币 ( $2^8$ ), 熵为3bit

**均匀分布** 假设有m种等概率的不同情况, 那么熵为  $\log_2 m$

**每种情况概率不相等** 一般分布 计算公式  $Ent(D) = -\sum_{k=1}^m p_k \log_2 p_k$

例: 假设发生A的概率为1/2, 发生B的概率为1/3, 发生C的概率为1/6, 该概率为1/6, 该如何求熵?

信息如何量化

得知信息前后, 熵的差值就是信息的量

假设我并不知道选择题是ABCD哪个选项时的熵  $\log_2 4 = 2$

有人告诉我C有一半概率正确后的熵  $\log_2 2 + \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} = 1.79$

信息 "C有一半概率正确"  $2 - 1.79 = 0.21$

划分选择

信息增益  
决策树ID3训练算法

信息增益的概念与度量  
信息熵  
信息增益

西瓜数据集

编号	A 色泽	B 根茎	C 敲声	D 纹理	E 脐部	D 触感	好瓜
1	青绿	蜷缩	浊响	清晰	凹	硬滑	是
2	青绿	蜷缩	...	...	...	...	...
...	...	...	...	...	...	...	...
17	浅白	...	...	...	...	...	否

色泽: 青绿 青绿 乌黑 浅白  
正例 3 4 1  
反例 3 2 4

目标: 分类是好瓜还是坏瓜。共17个样本, 其中8好9坏, 属于一般分布, 则总熵为  $Ent(D) = \frac{8}{17} \log_2 \frac{8}{17} + \frac{9}{17} \log_2 \frac{9}{17} \approx 0.998$

若想使得机器基于信息增益生成的决策树具有最佳的分类效果, 应当按照以下步骤执行:

1) 欲确定第一次分类的属性, 应当计算各属性下所有属性值的熵。以 "色泽" 为例, 计算出不同色泽下的熵分别为  $Ent(D_A), Ent(D_B), Ent(D_C)$ , 则色泽的熵就是它们的熵的概率和, 即  $\frac{8}{17} Ent(D_A) + \frac{9}{17} Ent(D_B) + \frac{9}{17} Ent(D_C)$ , 将其与总熵  $Ent(D)$  作差得到色泽的信息增益为0.109。以同样的方式求出其它属性的信息增益, 最终得到纹理的信息增益-0.381最大, 则第一轮应优先分类纹理。

2) 欲确定第一次分类后不同属性值的下一轮分类的属性, 注意需要对第一次分类后不同属性值的各样本集分别分析。以色泽清晰的样本为例, 数据总量为9, 则计算其它不同属性的熵时只需要针对这8个数据, 其余的数据忽略即可。若这9个样本中正例有7个, 反例有2个, 则未分类时的总熵为  $Ent(D') = \frac{7}{9} \log_2 \frac{7}{9} + \frac{2}{9} \log_2 \frac{2}{9} \approx 0.764$ 。最终算得此时信息增益最大的属性是根茎, 即对于色泽清晰的样本, 它们下一轮分类的属性是根茎。同理, 若样本的色泽是稍暗, 则只针对色泽稍暗的样本计算它们其它属性的信息增益, 最终得到信息增益最大的是触感, 即对于色泽稍暗的样本, 下一轮分类的属性是触感。

增益率  
决策树C4.5训练算法

实际上, 信息增益准则对可取值数目较多的属性有所偏好, 为减少这种偏好可能带来的不利影响, 可使用增益率来选择

最优划分属性  $Gain_{rate}(D, a) = IV(a)$

其中  $IV(a) = -\sum_{i=1}^n \frac{|D_i|}{|D|} \log_2 \frac{|D_i|}{|D|}$  称为属性a的固有值; 属性a的可能取值数目越多,  $IV(a)$  值通常越大

对可取值数目较少的属性有所偏好  
一般先从候选划分属性中找出信息增益高于平均水平的属性, 再从中选择增益最高的

基尼指数  
决策树CART训练算法  
Classification and Regression

分类树: 基尼指数最小准则

数据 (各项指标只有 "YES(1)" 或 "NO(0)" 两种状态); 根据这3个指标判断其是否为好瓜 (由于存在缺失值, 3个类别统计后总数据可能不同)

指标1	指标2	指标3	好瓜
0	0	1	1
1	1	0	1
1	0	Nan	0
...	...	...	...

纯度计算: 以指标1为例, 随机取两个样本, 它们是同一类别的概率为  $(\frac{105}{105+99})^2 + (\frac{105}{105+99})^2$

反之, 若该值越小, 则基尼系数越小, 纯度越高。数据集D的纯度可用基尼值来度量  $Gini(D) = \sum_{k=1}^m p_k(1-p_k) = 1 - \sum_{k=1}^m p_k^2$

公式  $Gini(p) = \sum_{k=1}^m p_k(1-p_k) = 1 - \sum_{k=1}^m p_k^2$

$p_k$  表示选中的样本属于k类别的概率, 则这个样本被错分的概率为  $(1-p_k)$

样本集含有k个类别, 一个随机选中的样本可以属于这k个类别中的任意一个, 对类别进行相加

当为二分类时,  $Gini(p) = 2p(1-p)$

例: 基于基尼指数构建决策树

1) 要求根节点, 首先需要计算3个指标的基尼指数。以指标1为例, 样本具有指标1时, 基尼指数为  $Gini(D_1) = 1 - (\frac{105}{105+99})^2 - (\frac{99}{105+99})^2$ , 不具有指标1时, 基尼指数为  $Gini(D_2) = 1 - (\frac{34}{34+105})^2 - (\frac{105}{34+105})^2$ , 则指标1的基尼指数就是它们的加权平均值  $\frac{105}{105+99} Gini(D_1) + \frac{99}{105+99} Gini(D_2)$ 。同理计算出另外两个指标的基尼指数, 得到指标2的基尼指数=0.36最小, 所以以指标2做第一次分叉。

2) 对于第二次及以上分叉, 以具有指标2的样本为例 ( $m=37, m=127$ )。首先计算出此时的基尼指数  $Gini(D') = 1 - (\frac{34}{34+127})^2 - (\frac{127}{34+127})^2 \approx 0.28$ , 再计算具有指标2的样本, 其指标1和指标3的基尼指数。最终得到指标3的基尼  $\frac{1}{13} \log_2 \frac{1}{13} + \frac{1}{24} \log_2 \frac{1}{24}$ , 首先计算出此时的基尼指数  $Gini(D') = 1 - (\frac{1}{13+24})^2 - (\frac{24}{13+24})^2 \approx 0.29$ , 即精度下降到4/7=57.1%, 应剪枝

3) 注意, 若分叉后得到的基尼指数大于分叉前的基尼指数, 即纯度下降, 就不应该分叉。

目的: 对付过拟合

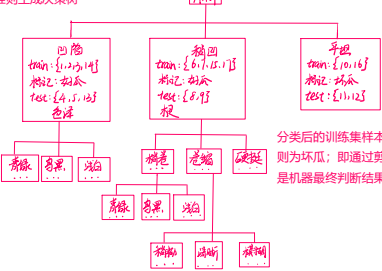
数据集  $\begin{matrix} \text{编号} & A \text{ 色泽} & B \text{ 根茎} \\ 1-8 & \dots & \text{是} \\ 9-17 & \dots & \text{否} \end{matrix}$

划分出训练集和验证集 (该样本中以测试集充当训练集)

其中训练集有5个正例, 5个反例; 验证集有3个正例, 4个反例

训练集  
信息增益准则生成决策树

训练集与验证集在决策树上的情况



剪枝处理

自上而下, 若分完后, 精度上升, 则应该分; 若精度不变或下降, 则没必要分支, 应剪枝

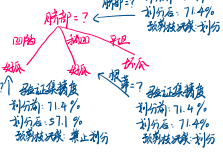
对于脐部  $\text{train: } m^+=5, m^-=5$   
 $\text{test: } m^+=3, m^-=4$   
剪枝

将部分数据划分到 "凹陷" 后, 对于色泽

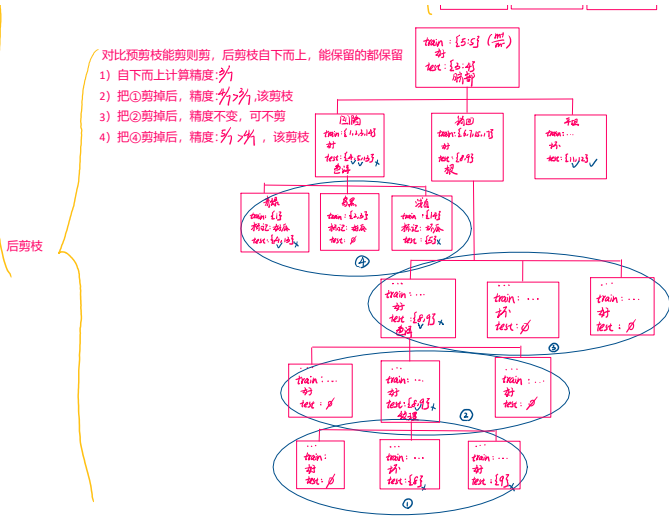
分之前: 5/7=71.4%  
分之后, 若有:

青绿  $\text{train: } \{1, 2, 3\}$   
青黑  $\text{train: } \{4, 5\}$   
平坦  $\text{train: } \{6, 7\}$

最终预剪枝决策树



由于将5号误分为坏瓜 (之前分类正确, 分支后分类错误) 即精度下降到4/7=57.1%, 应剪枝



连续与缺失值

连续值处理 二分法，使得信息增益最大的阈值均为连续值

例: 根据身高对西瓜数据集分类

- 1) 按照浓度从小到大的顺序对数据样本排序
- 2) 通过计算信息增益，查找每一次分类的最佳阈值。

继续以剪枝处理的西瓜数据集D为例，即数据集为

编号	A-色泽	决策
1-8	...	是
9-17	...	否

其中，1-8是正例，9-17是反例。

缺失值处理

- 1) 以色泽属性为例，缺失了1号，5号，13号在色泽上的数据，即无缺失值的样例子集D'中只有14个样例，计算分类前的信息熵  $Ent(D') = -\log_2 \frac{14}{28} - \log_2 \frac{14}{28} = 1.0$
- 2) 令  $D^+$ ,  $D^-$  分别表示在色泽属性上取值为青绿，乌黑和浅白的样例子集，有  $Ent(D^+) = -\log_2 \frac{14}{28} - \log_2 \frac{14}{28} = 1.0$ ,  $Ent(D^-) = -\log_2 \frac{14}{28} - \log_2 \frac{14}{28} = 1.0$ ,  $Ent(D) = -\log_2 \frac{14}{28} - \log_2 \frac{14}{28} = 1.0$
- 3) 因此，样本子集D'上色泽的信息增益为  $G_{gain}(D, 色泽) = Ent(D) - \frac{14}{28} Ent(D^+) - \frac{14}{28} Ent(D^-) = 1.0 - (\frac{14}{28} \times 1.0 + \frac{14}{28} \times 1.0) = 0.0$
- 4) 于是，样本集D上色泽属性的信息增益为  $G_{gain}(D, 色泽) = P \times G_{gain}(D', 色泽) = \frac{14}{28} \times 0.0 = 0.0$ ，类似地可计算出其它属性在D上的信息增益

多变量决策树

对于正常基于连续值生成的决策树，它们不同的属性最后往往呈线性的关系，因为每次的取值都是一次二分

例:

关于密度和含糖率生成的决策树

多变量决策树的分类边界不再是某个属性，而是对属性的线性组合进行一种测试

多变量决策树对应的分类边界

决策树对应的分类边界

决策树对复杂分类边界的分段近似