

# 决策树

刘杰  
人工智能学院



1

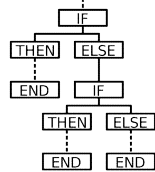
## 目录

- 决策树模型
- 特征选择
- 决策树的生成
- 决策树的修剪

2

## 决策树模型

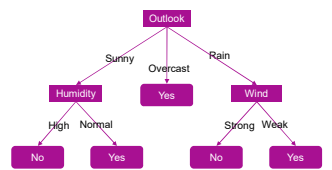
我们如何做“决策”？  
If Then规则



3

## Play Tennis?

- 分类决策树模型是一种描述对实例进行分类的树形结构。
- 决策树由结点(node) 和有向边 (directed edge) 组成。
- 结点有两种类型: 内部结点( internal node )和叶结点(leaf node)。
- 内部结点表示一个特征或属性, 叶结点表示一个类。



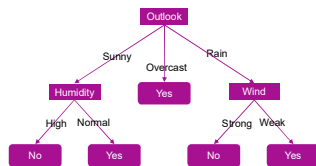
4

## Play Tennis?

如图的决策可以看做一个IF-THEN规则的集合

根结点到叶结点的每一条路径构成一条规则:

- 路径上的内部结点对应着规则的条件
- 叶结点的类对应着规则的结论



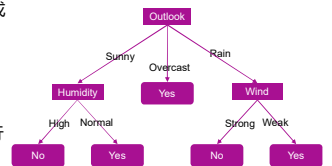
5

## 决策树

决策树是一种典型的分类方法

- 首先对数据进行处理, 利用归纳算法生成可读的规则和决策树,
- 然后使用决策对新数据进行分析。

本质上决策树是通过一系列规则对数据进行分类的过程。



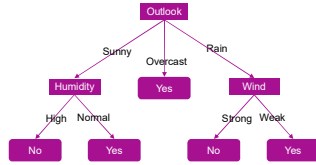
6



## 决策树

决策树的优点

- 1、推理过程容易理解，决策推理过程可以表示成为If-Then形式；
- 2、推理过程完全依赖于属性变量的取值特点；
- 3、可自动忽略目标变量没有贡献的属性变量，也为判断属性变量的重要性，减少变量的数目提供参考。

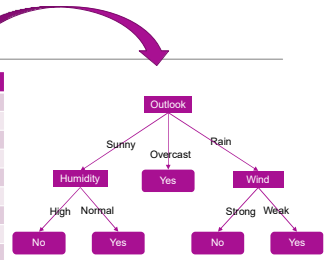


7



## 决策树

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No



8

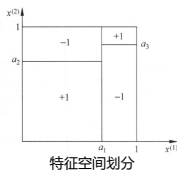


## 决策树与条件概率分布

决策树还表示给定特征条件下类的条件概率分布： $P(Y|X)$

这一条件概率分布定义在特征空间的一个划分 (partition) 上。

将特征空间划分为互不相交的单元(cell)或区域(region)，并在每个单元定义一个类的概率分布就构成了一个条件概率分布。



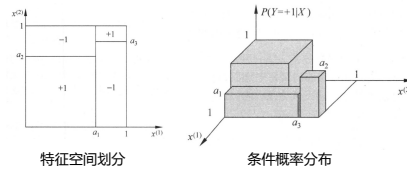
特征空间划分

9



## 决策树与条件概率分布

假设  $X$  为表示特征的随机变量， $Y$  为表示类的随机变量，那么这个条件概率分布可以表示为  $P(Y|X)$ 。 $X$  取值于给定划分下单元的集合， $Y$  取值于类的集合。各叶结点(单元)上的条件概率往往偏向某一类，即属于某一类的概率较大。决策树分类时将该结点的实例分到条件概率大的那一类去。



特征空间划分

条件概率分布

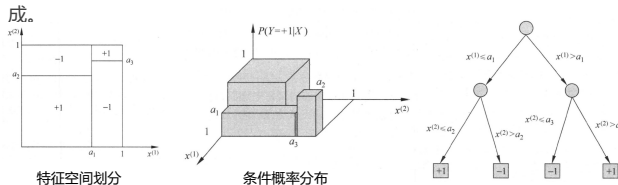
10



## 决策树与条件概率分布

决策树的一条路径对应于划分中的一个单元。

决策树所表示的条件概率分布由各个单元给定条件下类的条件概率分布组成。



特征空间划分

条件概率分布

11



## 决策树学习算法

假设给定训练数据集

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

其中， $x_i = (x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(n)})^T$  为输入实例 (特征向量)， $n$  为特征个数， $y_i \in \{1, 2, \dots, K\}$  为类标记， $i = 1, 2, \dots, N$ ， $N$  为样本容量。决策树学习的目标是根据给定的训练数据集构建一个决策树模型，使它能够对实例进行正确的分类。

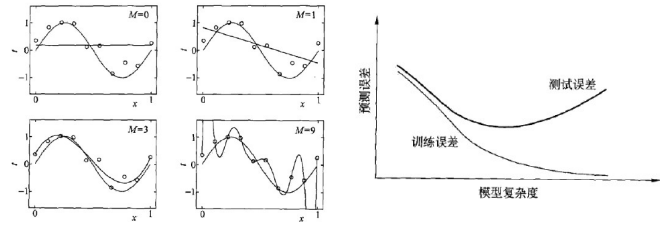
12



## 决策树学习算法

- 决策树学习本质上是从训练数据集中归纳出一组分类规则，与训练数据集不相矛盾的决策树。
- 能对训练数据进行正确分类的决策树可能多个，也可能一个也没有。我们需要的是一个与训练数据矛盾较小的决策树，同时具有很好的泛化能力。
- 决策树学习是由训练数据集估计条件概率模型。基于特征空间划分的类的条件概率模型有无穷多个。
- 我们选择的条件概率模型应该不仅对训练数据有很好的拟合，而且对未知数据有很好的预测。

13



14



## 决策树学习算法

**决策树学习策略：**以损失函数为目标函数的最小化。

**特征选择：**递归地选择最优特征，并根据该特征对训练数据进行分割，使得对各个子数据集有一个最好的分类

**决策树生成：**与特征选择过程相对应，选择的特征依次形成决策树的结点，直至所有训练数据子集都被基本正确分类，或无合适特征可选。

**决策树的剪枝：**为避免过拟合，对已生成的树自下而上进行剪枝，将树变得更简单，从而使它具有更好的泛化能力。

15



## 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

特征选择是决定用哪个特征来划分特征空间。

特征选择的准则是信息增益或信息增益比。

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

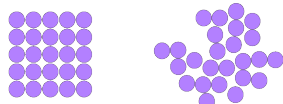
16



## 信息增益

- Shannon 1948年提出的信息论理论：
- 熵(entropy)：信息量大小的度量，即表示随机变量不确定性的度量。

$$I(a_i) = p(a_i) \log_2 \frac{1}{p(a_i)}$$



Low Entropy      High Entropy

$$I(a_1, a_2, \dots, a_n) = \sum_{i=1}^n I(a_i) = \sum_{i=1}^n p(a_i) \log_2 \frac{1}{p(a_i)}$$



17



## 信息增益

- Shannon 1948年提出的信息论理论：
- 熵(entropy)：信息量大小的度量，即表示随机变量不确定性的度量。



High Knowledge  
Low Entropy



Medium Knowledge  
Medium Entropy



Low Knowledge  
High Entropy



18



## 信息增益

- 设  $X$  是一个取有限个值的离散随机变量，其概率分布为：

$$P(X = x_i) = p_i, i = 1, 2, \dots, n$$

- 信息量：随机变量  $x_i$  的信息量  $I(x_i)$  可如下度量：

$$I(x_i) = \log \frac{1}{p(x_i)}$$

- 其中  $p(x_i)$  表示事件  $x_i$  发生的概率。

19



## 信息增益

- 熵为各事件信息量的“数学期望”：

$$H(X) = E[I(x_i)] = - \sum_{i=1}^n p(x_i) \log p(x_i)$$

对数以2为底或以  $e$  为底(自然对数)，这时熵的单位分别称作比特(bit)或纳特(nat)，熵只依赖于  $X$  的分布，与  $X$  的取值无关。

$$H(p) = - \sum_{i=1}^n p_i \log p_i$$

20



## 信息增益

- 熵的**理论解释**：

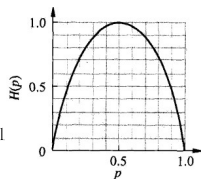
熵越大，随机变量的不确定性越大：

$$0 \leq H(p) \leq \log n$$

当  $X$  为0,1分布时：

$$P(X = 1) = p, P(X = 0) = 1 - p, 0 \leq p \leq 1$$

熵： $H(p) = -p \log_2 p - (1-p) \log_2 (1-p)$



21



## 信息增益

设有随机变量  $(X, Y)$ ，其联合概率分布为：

$$P(X = x_i, Y = y_j) = p_{ij}, i = 1, 2, \dots, n; j = 1, 2, \dots, m$$

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，定义为  $X$  给定条件下  $Y$  的条件概率分布的熵对  $X$  的数学期望：

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

当熵和条件熵中的概率由数据估计(特别是极大似然估计)得到时，所对应的熵与条件熵分别称为经验熵(empirical entropy)和经验条件熵(empirical conditional entropy)。

22



## 信息增益

- 定义(信息增益):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

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

- (Information gain)表示得知特征  $X$  的信息而使类  $Y$  的信息的不确定性减少的程度。
- 一般地, 熵  $H(Y)$  与条件熵  $H(Y|X)$  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。

23



## 信息增益的算法

- 设训练数据集为  $D$
- $|D|$  表示其样本容量, 即样本个数
- 设有  $K$  个类  $C_k, k = 1, 2, \dots, K$ ,
- $|C_k|$  为属于类  $C_k$  的样本个数
- 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  $D_1, \dots, D_n$
- $|D_i|$  为  $D_i$  的样本个数
- 记子集  $D_i$  中属于类  $C_k$  的样本集合为  $D_{ik}$
- $|D_{ik}|$  为  $D_{ik}$  的样本个数

24



## 信息增益的算法

- 输入：训练数据集D和特征A；
- 输出：特征A对训练数据集D的信息增益g(D,A)

### 1、计算数据集D的经验熵H(D)

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

### 2、计算特征A对数据集D的经验条件熵H(D|A)

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

### 3、计算信息增益

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

25



## 信息增益比

- 以信息增益作为划分训练数据集的特征，存在偏向于选择取值较多的特征的问题
- 使用信息增益比可以对这一问题进行校正
- 定义（信息增益比）特征A对训练数据集D的信息增益比定义为信息增益与训练数据集D关于特征A的值的熵之比

$$g_R(D, A) = \frac{g(D, A)}{H_A(D)}$$

$$H_A(D) = - \sum_{i=1}^n \frac{|D_i|}{|D|} \log_2 \frac{|D_i|}{|D|}, n \text{ 是特征A取值的个数}$$

26



## 决策树ID3算法

输入：训练数据集D，特征集A，阈值ε；

输出：决策树T。

- 若D中所有实例属于同一类 $C_k$ ，则T为单结点树，并将类 $C_k$ 作为该结点的类标记，返回T；
- 若 $A = \emptyset$ ，则T为单结点树，并将D中实例数最大的类 $C_k$ 作为该结点的类标记，返回T；
- 否则，根据公式计算A中各特征对D的信息增益，选择信息增益最大的特征 $A_g$ ；
- 如果 $A_g$ 的信息增益小于阈值ε，则置T为单结点树，并将D中实例数最大的类 $C_k$ 作为该结点的类标记，返回T；
- 否则，对 $A_g$ 的每一可能值 $a_i$ ，依 $A_g = a_i$ 将D分割为若干非空子集 $D_i$ ，将 $D_i$ 中实例数最大的类作为标记，构建子结点，由结点及其子结点构成树T，返回T；
- 对第i个子结点，以 $D_i$ 为训练集，以 $A - \{A_g\}$ 为特征集，递归地调用步(1)~步(5)，得到子树 $T_i$ ，返回 $T_i$ 。

27



### 第1步计算决策属性的熵

决策属性“PlayTennis”。  
该属性分两类：Yes/No

$$\begin{aligned} |C_1|(Yes) &= 9 \\ |C_2|(No) &= 5 \\ |D| &= |C_1| + |C_2| = 14 \end{aligned}$$

$$\begin{aligned} P_1 &= 9/14 = 0.643 \\ P_2 &= 5/14 = 0.357 \end{aligned}$$

$$\begin{aligned} H(D) &= -P_1 \log_2 P_1 - P_2 \log_2 P_2 \\ &= -(P_1 \log_2 P_1 + P_2 \log_2 P_2) \\ &= 0.9401 \end{aligned}$$

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

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

28



### 第2步计算条件属性的熵

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

条件属性共有4个：  
Outlook、Temperature、Humidity、Wind。  
分别计算不同属性的信息增益。

29



### 第2步 计算Outlook中各个属性的条件熵

Outlook共分三个组：  
Sunny、Overcast、Rain

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

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

30



Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

### 第2-1步计算Outlook中Sunny的熵

Outlook共分三个组：  
Sunny、Overcast、Rain  
其中当Outlook为Sunny时，Yes与No比例为2/3

$$|D_{11}|(Yes) = 2$$

$$|D_{12}|(No) = 3$$

$$|D_1| = 5$$

$$P_1 = 2/5 = 0.4$$

$$P_2 = 3/5 = 0.6$$

$$H(D_1) = -P_1 \log_2 P_1 - P_2 \log_2 P_2$$

$$= -(P_1 \log_2 P_1 + P_2 \log_2 P_2)$$

$$= 0.971$$

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

31



Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

### 第2-2步计算Outlook中Overcast的熵

Outlook共分三个组：  
Sunny、Overcast、Rain  
其中当Outlook为Overcast时，Yes与No比例为4/0

$$|D_{21}|(Yes) = 4$$

$$|D_{22}|(No) = 0$$

$$|D_2| = 4$$

$$P_1 = 1$$

$$P_2 = 0$$

$$H(D_2) = -P_1 \log_2 P_1 - P_2 \log_2 P_2$$

$$= -(P_1 \log_2 P_1 + P_2 \log_2 P_2)$$

$$= 0$$

32



Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

### 第2-3步计算计算Outlook中Rain的熵

Outlook共分三个组：  
Sunny、Overcast、Rain  
其中当Outlook为Rain时，Yes与No比例为3/2

$$|D_{31}|(Yes) = 3$$

$$|D_{32}|(No) = 2$$

$$|D_3| = 5$$

$$P_1 = 0.6$$

$$P_2 = 0.4$$

$$H(D_3) = -P_1 \log_2 P_1 - P_2 \log_2 P_2$$

$$= -(P_1 \log_2 P_1 + P_2 \log_2 P_2)$$

$$= 0.971$$

33



Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

### 第2-4步计算Outlook的熵

Outlook共分三个组：  
Sunny、Overcast、Rain  
所占比例

$$\text{Sunny: } 5/14 = 0.357$$

$$\text{Overcast: } 4/14 = 0.286$$

$$\text{Rain: } 5/14 = 0.357$$

$$\text{计算Outlook的平均信息期望}$$

$$E(\text{Outlook}) = 0.357 * 0.971 + 0.286 * 0 + 0.357 * 0.971 = 0.6932$$

$$\text{Outlook信息增益} = 0.9401 - 0.6932 = 0.2469 \quad (1)$$

34



Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

### 第3步 计算Temperature的熵

收入共分三个组：  
Hot、Mild、Cool  
 $E(\text{Temperature}) = 0.9108$   
Temperature信息增益 =  $0.9401 - 0.9108 = 0.0293 \quad (2)$

35



Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

### 第4步计算Humidity的熵

Humidity共分二组：  
High、Normal  
 $E(\text{Humidity}) = 0.7885$   
Humidity信息增益 =  $0.9401 - 0.7885 = 0.1516 \quad (3)$

36



### 第5步计算Wind的熵

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Wind共分二组:  
Weak, Strong  
 $E(Wind) = 0.892$   
Wind信息增益 =  $0.9401 - 0.892$   
= 0.0481 (4)

37



### 第6步计算选择节点

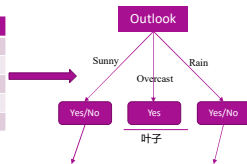
Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Outlook信息增益 =  $0.9401 - 0.6932$   
= 0.2469 (1)  
Temperature信息增益 =  $0.9401 - 0.9108$   
= 0.0293 (2)  
Humidity信息增益 =  $0.9401 - 0.7885$   
= 0.1516 (3)  
Wind信息增益 =  $0.9401 - 0.892$   
= 0.0481 (4)

38



Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes



39



Sunny中Yes与No比例为2/3

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes

$|C_1|(Yes) = 2$   
 $|C_2|(No) = 3$   
 $|D| = 5$   
 $P_1 = 2/5 = 0.4$   
 $P_2 = 3/5 = 0.6$   
 $H(DS_{Sunny}) = -P_1 \log_2 P_1 - P_2 \log_2 P_2$   
=  $-(P_1 \log_2 P_1 + P_2 \log_2 P_2)$   
= 0.9710

40



如果选择Temperature作为节点  
分为Hot, Mild, Cool

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes

$H(D_1) = 0$   
比例: 2/5  
 $H(D_2) = 1$   
比例: 2/5  
 $H(D_3) = 0$   
比例: 1/5

平均信息期望(加权总和):  
 $E(H(D_{Temperature|Sunny})) = 0.4 * 0 + 0.4 * 1 + 0.2 * 0 = 0.4$   
 $Gain(Temperature) = H(D_{Sunny}) - E(H(D_{Temperature|Sunny})) = 0.9710 - 0.4 = 0.5710$

注意

41



如果选择Humidity作为节点  
分为High, Normal

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes

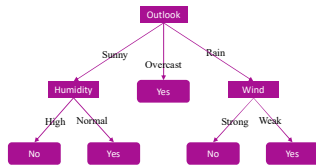
$H(D_1) = 0$   
比例: 3/5  
 $H(D_2) = 0$   
比例: 2/5

平均信息期望(加权总和):  
 $E(H(D_{Humidity|Sunny})) = 0$   
 $Gain(Humidity) = H(D_{Sunny}) - E(H(D_{Humidity|Sunny})) = 0.9710 - 0 = 0.9710$   
由于信息增益已然达到最极值, 所以选择Humidity。  
其他结点的增益也按照上述的方法计算。

42



Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No



43



## 决策树ID3算法-小结

ID3算法的基本思想是，以信息熵为度量，用于决策树节点的属性选择

每次优先选取信息量最多的属性，亦即使熵值变为最小的属性，以构造一棵熵值下降最快的决策树，到叶子节点处的熵值为0。

此时，每个叶子节点对应的实例集中的实例属于同一类。

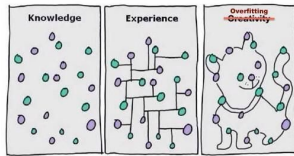
44



## 决策树面临的问题

### 过度拟合

- 决策树算法增长树的每一个分支的深度，直到恰好能对训练样例比较完美地分类。实际应用中，当数据中有噪声或训练样例的数量太少以至于不能产生目标函数的有代表性的采样时，该策略可能会遇到困难。
- 在以上情况发生时，这个简单的算法产生的树会过拟合训练样例(过拟合: Over Fitting)。



45



## 决策树面临的问题

### 理想的决策树:

- (1)叶子结点数最少;
- (2)叶子结点深度最小;
- (3)叶子结点数最少且叶子结点深度最小。

46



## 决策树的剪枝

通过极小化决策树整体的损失函数或代价函数来实现。

设树 $T$ 的叶结点个数为 $|T|$ ， $t$ 是树 $T$ 的叶结点，该叶结点有 $N_t$ 个样本点，其中 $k$ 类的样本点有 $N_{tk}$ 个， $k = 1, 2, \dots, K$ ，

$H_t(T)$ 为叶结点 $t$ 上的经验熵， $\alpha \geq 0$ 为参数，损失函数:

$$C_\alpha(T) = \sum_{t=1}^{|T|} N_t H_t(T) + \alpha |T|$$

47



## 决策树的剪枝

通过极小化决策树整体的损失函数或代价函数来实现。

设树 $T$ 的叶结点个数为 $|T|$ ， $t$ 是树 $T$ 的叶结点，该叶结点有 $N_t$ 个样本点，其中 $k$ 类的样本点有 $N_{tk}$ 个， $k = 1, 2, \dots, K$ ，

$H_t(T)$ 为叶结点 $t$ 上的经验熵， $\alpha \geq 0$ 为参数，损失函数:

$$C_\alpha(T) = \sum_{t=1}^{|T|} N_t H_t(T) + \alpha |T|$$

预测误差损失

模型复杂度损失

48





## 决策树的剪枝

通过极小化决策树整体的损失函数或代价函数来实现。

设树 $T$ 的叶结点个数为 $|T|$ ,  $t$ 是树 $T$ 的叶结点, 该叶结点有 $N_t$ 个样本点, 其中 $k$ 类的样本点有 $N_{tk}$ 个,  $k = 1, 2, \dots, K$ ,

$H_t(T)$ 为叶结点 $t$ 上的经验熵,  $\alpha \geq 0$ 为参数, 损失函数:

$$C_\alpha(T) = \sum_{t=1}^{|T|} N_t H_t(T) + \alpha |T|$$

经验熵:  $H_t(T) = - \sum_k \frac{N_{tk}}{N_t} \log \frac{N_{tk}}{N_t}$

原式第一项:  $C(T) = \sum_{t=1}^{|T|} N_t H_t(T) = - \sum_{t=1}^{|T|} \sum_{k=1}^K N_{tk} \log \frac{N_{tk}}{N_t}$

则:  $C_\alpha(T) = C(T) + \alpha |T|$

49



## 决策树的剪枝

• 树的剪枝算法:

输入: 生成算法产生的整个树  $T$ , 参数  $\alpha$ ;

输出: 修剪后的子树  $T_\alpha$ 。

(1) 计算每个结点的经验熵。

(2) 递归地从树的叶结点向上回缩。

• 设一组叶结点回缩到其父结点之后与之前的损失函数分别为:

$$C_\alpha(T_A) \text{ 与 } C_\alpha(T_B)$$

如果:  $C_\alpha(T_A) \leq C_\alpha(T_B)$  则进行剪枝

(3) 返回(2), 直至不能继续为止, 得到损失函数最小的子树  $T_\alpha$ 。

50



## CART树

CART: 分类与回归树, Classification And Regression Tree

算法由两部分组成:

- 决策树生成
- 决策树剪枝

回归树: 平方误差最小化

分类树: Gini Index

51



## CART算法

- CART算法采用一种二分递归分割的技术, 算法总是将当前样本集分割为两个子样本集
- 生成的决策树的每个非叶结点都只有两个分枝, 因此CART算法生成的决策树是结构简洁的二叉树。
- CART算法适用于样本特征的取值为是或非的场景。
- 剪枝过程特别重要, 所以在最优决策树生成过程中占有重要地位。有研究表明, 剪枝过程的重要性要比树生成过程更为重要

52



## CART的生成

### 回归树的生成

设 $Y$ 是连续变量, 给定训练数据集:

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

假设已将输入空间划分为 $M$ 个单元 $R_1, R_2, \dots, R_M$ , 并且每个单元 $R_m$ 上有一个固定的输出 $C_m$ , 回归树表示为:

$$f(x) = \sum_{m=1}^M c_m I(x \in R_m)$$

平方误差来表示预测误差, 用平方误差最小准则求解每个单元上的最优输出值

$$\sum_{x_i \in R_m} (y_i - f(x_i))^2$$

$R_m$ 上的 $C_m$ 的最优值:  $\hat{c}_m = ave(y_i | x_i \in R_m)$

53



## CART的生成

### 输入空间的划分

启发式: 选择第 $j$ 维变量 $x^{(j)}$ 和它取的值 $s$ , 作为切分变量和切分点, 定义两个区域:

$$R_1(j, s) = \{x | x^{(j)} \leq s\} \text{ 和 } R_2(j, s) = \{x | x^{(j)} > s\}$$

然后寻找最优切分变量和切分点:

$$\min_{j,s} [\min_{c_1} \sum_{x_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j,s)} (y_i - c_2)^2]$$

且:  $\hat{c}_1 = ave(y_i | x_i \in R_1(j, s))$  和  $\hat{c}_2 = ave(y_i | x_i \in R_2(j, s))$

再对两个区域重复上述划分, 直到满足停止条件。

54



## CART的生成

最小二乘回归树生成算法

输入：训练数据集 $D$ ;

输出：回归树 $f(x)$ 。

在训练数据集所在的输入空间中，递归地将每个区域划分成两个子区域并决定每个子区域上地输出值，构建二叉决策树；

(1)选择最优切分变量 $j$ 与切分点 $s$  求解

$$\min_{j,s} [\min_{c_1} \sum_{x_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j,s)} (y_i - c_2)^2]$$

遍历变量 $j$ ，对固定地切分变量 $j$ 扫描切分点 $s$ ，选择使上式达到最小值地对 $(j,s)$ 。

55



## CART的生成

•最小二乘回归树生成算法

(2)用选定地对 $(j,s)$ 划分区域并决定相应地输出值：

$$R_1(j,s) = \{x|x^{(j)} \leq s\}, R_2(j,s) = \{x|x^{(j)} > s\}$$

$$\hat{c}_m = \frac{1}{N_m} \sum_{x_i \in R_m(j,s)} y_i, m = 1, 2$$

(3)继续对两个子区域调用(1),(2),直至满足停止条件。

(4)将输出空间划分为 $M$ 个区域 $R_1, R_2, \dots, R_M$  生成决策树；

$$f(x) = \sum_{m=1}^M \hat{c}_m I(x \in R_m)$$

56



## CART的生成

•分类树的生成：

•基尼指数

•分类问题中，假设有 $k$ 个类，样本点属于 $k$ 的概率 $p_k$ ，则概率分布的基尼指数：

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

•二分类问题：

$$Gini(p) = 2p(1-p)$$

•对给定的样本集合 $D$ ，基尼指数

$$Gini(D) = 1 - \sum_{k=1}^K \left(\frac{|C_k|}{|D|}\right)^2$$

57



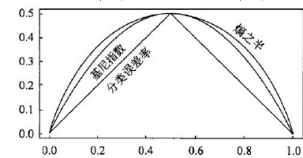
## CART的生成

如果样本集合 $D$ 根据特征 $A$ 是否为 $a$ 被分割成 $D_1$ 和 $D_2$ ，即

$$D_1 = \{(x,y) \in D | A(x) = a\}, D_2 = D - D_1$$

则在特征 $A$ 的条件下，集合 $D$ 的基尼指数：

$$Gini(D, A) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2)$$



58



## CART的生成

•CART分类树生成算法

输入：训练数据集 $D$

输出：CART分类树

从根节点开始，递归对每个结点操作

1、设结点数据集为 $D$ ，对每个特征 $A$ ，对其每个值 $a$ ，根据样本点对 $A=a$ 的测试为是或否，将 $D$ 分为 $D_1, D_2$ ，计算 $A=a$ 的基尼指数

2、在所有的特征 $A$ 以及所有可能的切分点 $a$ 中，选择基尼指数最小的特征和切分点，将数据集分配到两个子结点中。

3、对两个子结点递归调用1，2步骤

4、生成CART分类树

59



## CART树剪枝

•CART剪枝

两步

1、从生成算法产生的决策树 $T_0$ 底端开始不断剪枝，直到 $T_0$ 的根结点，形成子树序列 $\{T_0, T_1, \dots, T_n\}$ ，

2、通过交叉验证法在独立的验证数据集上对子树序列进行测试，从中选择最优子树

60



## CART树剪枝

### •CART剪枝

1、剪枝，形成子树序列

剪枝过程中，计算子树的损失函数：

$$C_{\alpha}(T) = C(T) + \alpha|T|$$

对固定的 $\alpha$ 一定存在损失函数最小的子树，表示为 $T_{\alpha}$ 。

当 $\alpha$ 变大时，最优子树 $T_{\alpha}$ 偏小，

$\alpha = 0$ 时，整体树最优， $\alpha$ 趋近无穷大，单结点最优

将 $\alpha$ 从小增大， $0 = \alpha_0 < \alpha_1 < \dots < \alpha_n < +\infty$

最优子树序列  $\{T_0, T_1, \dots, T_n\}$

61



## CART算法剪枝

### •CART剪枝

1、剪枝，形成子树序列

具体：从整树 $T_0$ 开始剪枝，以其中任意内部结点 $t$ 为单结点树的损失函数：

$$C_{\alpha}(t) = C(t) + \alpha$$

以 $t$ 为根结点的子树 $T_t$ 的损失函数：

$$C_{\alpha}(T_t) = C(T_t) + \alpha|T_t|$$

当 $\alpha = 0$ 及 $\alpha$ 很小时， $C_{\alpha}(T_t) < C_{\alpha}(t)$

不断增大 $\alpha$ ，当  $C_{\alpha}(T_t) = C_{\alpha}(t)$   $\alpha = \frac{C(t) - C(T_t)}{|T_t| - 1}$

$T_t$ 与 $t$ 有相同损失函数值，但 $t$ 结点更少，所以剪枝 $T_t$ 。

62



## CART的生成

### •CART剪枝

1、剪枝，形成子树序列

对整树 $T_0$ 中每个内部结点 $t$ ，计算：

$$g(t) = \frac{C(t) - C(T_t)}{|T_t| - 1}$$

在 $T_0$ 中剪去 $g(t)$ 最小的 $T_t$ ，将得到的子树作为 $T_1$ ，同时将最小的 $g(t)$

设为 $\alpha_1$ ， $T_1$ 为区间 $[\alpha_1, \alpha_2)$ 的最优子树

如此剪枝下去，直到根节点，不断增加 $\alpha$ 的值，产生新的区间。

63



## CART的生成

### •CART剪枝

2、在剪枝得到的子树序列 $\{T_0, T_1, \dots, T_n\}$ 中通过交叉验证选取最优子树 $T_{\alpha}$

利用独立的验证数据集，测试子树序列中各子树的平方误差或基尼指数，最小的决策树就是最优决策树。

64



## Homework

### 习题5.2

PlayTennis用CART的流程做一遍

65



END

66