

§4 决策树

1. 基本流程: divide and conquer 分裂属性, 使各分裂点尽可能纯
递归返回: (1) 当前结点包含的样本全属同一类别, 无需划分 (2) 当前结点样本为空, 即所有样本在所有属性上取值相同, 无法划分 (3) 当前结点包含的样本为空

2. 划分选择

信息熵 (information entropy): 度量样本纯度

$$Ent(D) = - \sum_{k=1}^{|D|} p_k \log_2 p_k$$

信息增益 (information gain): $Gain(D, a) = Ent(a) - \sum_{v=1}^V \frac{|D_v|}{|D|} Ent(D_v)$ — ID3

信息增益率 (gain ratio): $Gain_ratio(D, a) = \frac{Gain(D, a)}{IV(a)}$ — C4.5

$$IV(a) = - \sum_{v=1}^V \frac{|D_v|}{|D|} \log_2 \frac{|D_v|}{|D|}$$

基尼系数 (Gini index): $Gini(D) = \sum_{k=1}^{|D|} \sum_{k' \neq k} p_k p_{k'} = 1 - \sum_{k=1}^{|D|} p_k^2$ — CART

$$Gini_index(D, a) = \sum_{v=1}^V \frac{|D_v|}{|D|} Gini(D_v)$$

$Gini(D)$ 越小, D 纯度越高

3. 剪枝处理 pruning

预剪枝 (pre-pruning): 决策树情况

后剪枝 (post pruning)

预剪枝: ID3

优点	缺点	特点	衡量指标
建树时间短	单变量决策树	X 与目标变量	信息增益
建树简单	只运用小样本	构造决策树	
	依赖于特征取值		
	数目较多的特征		
	抗噪性差		

C4.5

优点	缺点	特点	衡量指标
规则易于理解	需进行多次扫描	多分支决策树	信息增益率
对噪声敏感	和树深, 效率较低	选信息增益率最大的	
实现简单	只运用小样本	建树	

CART		
优点	特点	衡量指标
可用于分类和回归	构造二叉树	Gini

算法总结

算法	支持模型	树结构	划分特征选择	连续值处理	缺失值处理	不平衡	特征属性多义性
ID3	分类	多叉树	信息增益	×	×	×	×
C4.5	分类	多叉树	信息增益率	✓	✓	✓	×
CART	分类, 回归	二叉树	Gini, 均方差	✓	✓	✓	✓

4. 多变量决策树

算法: OC1 - 贪心地寻找每个属性的最佳取值 在局部最优的基础上再对分类边界进行随机扰动以试图找到更优的边界

第 i 个神经元的连接权重

w_1, w_2, \dots, w_n



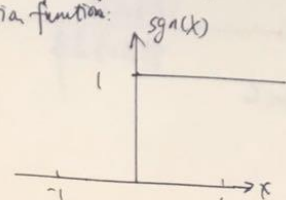
神经网络

输出 $y = f(\sum_{i=1}^n w_i x_i - \theta)$

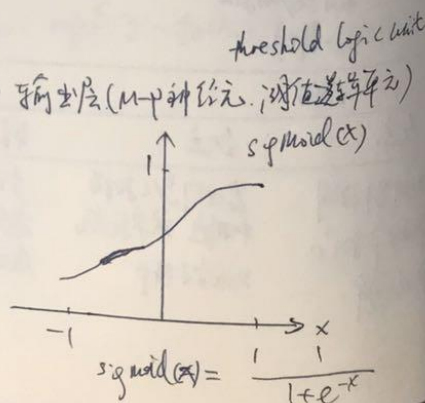
神经网络

1. 感知机 (perceptron): 输入层接收信号 \rightarrow 输出层 ($M-P$ 神经元, 阈值逻辑单元)

activation function:



$$\text{sgn}(x) = \begin{cases} 1 & x \geq 0 \\ 0 & x < 0 \end{cases}$$



$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$