

CART

优点: 可用于分类和回归  
特点: 构造二叉树  
衡量指标: Gini

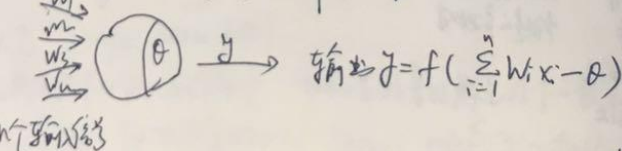
算法总结

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	剪枝	特征属性选择
ID3	分类	多叉树	信息增益	X	X	X	X
C4.5	分类	多叉树	信息增益率	✓	✓	✓	X
CART	分类、回归	二叉树	Gini 增益	✓	✓	✓	✓

4. 多变量决策树

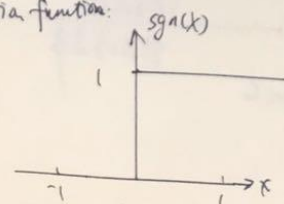
主要算法: OCI - 贪心地寻找每个属性的最佳取值, 在局部最优的基础上对分裂边界进行随机扰动以试图找到更好的边界

神经网络

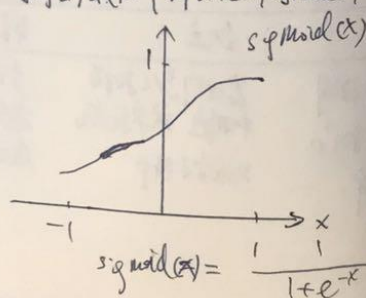


1. 感知机 (perceptron): 输入层接收信号  $\rightarrow$  输出层 ( $M$  个神经元, 阈值逻辑单元)

activation function:



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



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



1) 以多组不同参数值初始化多个神经网络, 按标注方法训练后, 取其中误差最小的作为最终参数

2) 模拟退火, 每一步以一定概率接受比当前更差的结果

3) 种群随机梯度下降

4 其他常见神经网络

1) RBF (Radial Basis Function) 径向基函数

$$\phi(x) = \sum_{i=1}^n W_i p(x, c_i) \quad p(x, c_i) = e^{-\beta_i \|x - c_i\|^2}$$

确定神经元中心  $c_i$ , 再利用BP算法确定  $W_i, \beta_i$

2) ART (adaptive resonance theory) 自适应谐振理论, 竞争型学习  
winner-take-all

3) SOM (self-organizing map) 自组织映射: 每个输出层神经元含4个连接权重, 自身携带的权向量之间的距离, 距离最近的神经元成为竞争获胜者, beat matching unit 最佳匹配单元

## §6 支持向量机

1. 感知器模型

$$y = \text{sgn}(\theta \cdot x) = \begin{cases} 1, & \theta \cdot x > 0 \\ -1, & \theta \cdot x < 0 \end{cases} \quad \begin{matrix} y \theta x > 0 & \text{正确分类} \\ y \theta x < 0 & \text{错误分类} \end{matrix}$$

Loss function, 期望使所有样本到超平面的距离之和最小

$$L = \sum_{i=1}^k \frac{y_i \theta \cdot x_i}{\|\theta\|_2} \rightarrow L = - \sum_{i=1}^k y_i \theta \cdot x_i \quad \frac{\partial L(\theta)}{\partial \theta} = - \sum_{i=1}^k y_i x_i$$

线性可分 SVM: 距离超平面最近的点, 令其距离超平面这个点 超平面 (separating hyperplane)  
非线性可分 SVM

支持向量 (support vector) 离分割超平面最近的点

间隔 (margin): 支持向量表点到分割超平面的距离