

机器学习总结: Hands-on machine learning
+ 设计学习方法 CART 算法 P171

决策树:

CART training algorithm:

Classification and regression tree:

cost function the algorithm tries to minimize:

$$J(k, t_k) = \left(\frac{m_{\text{left}}}{m} \right) G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}}$$

$G_{\text{left/right}}$: measures the impurity of the left/right subset.

$\frac{\text{左边 data 个数}}{\text{整个 data 个数}}$

决策树的基石:

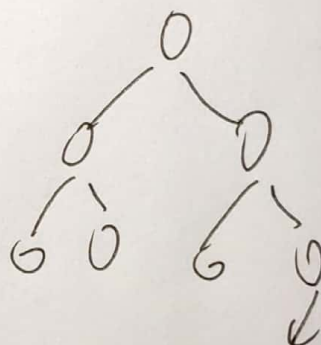
Gini impurity

$$G_i = 1 - \sum_{k=1}^n P_{i,k}^2$$

the ration

$$G = 1 - \left[\left(\frac{0}{54} \right)^2 + \left(\frac{49}{54} \right)^2 + \left(\frac{5}{54} \right)^2 \right]$$

$$= 0.168$$



里面有54个元素

属于类1的: 0

类2: 49

3: 5

↓
不纯的度量 Gini impurity is 0.168.

CART 算法: 先把 training set 分成两个 subset.

怎么分呢, 优化 $J(k, t_k)$ 见上. 使其 minimize.

→ 再分别将左右两颗树分成 2 个 subsets.

直到达到最大深度和或不能减少 impurity 了.



另一种衡量 impurity 的方法: entropy

$$H_i = - \sum_{k=1}^n P_{i,k} \log(p_{i,k})$$

$P_{i,k} \neq 0$

比较 entropy 和 Gini impurity:

区别不大.

Gini 倾向于将最频繁出现的 class 分出到一个单独的分支中.

entropy 倾向于产生均衡一点的 tree.

两边 number 差不多,

