

Chapter 4

4.1 看图 4.2 P74 页的流程图,

程序终止的条件, 有 5 个,

① D 为空, 则训练误差在这个子集树上为 0

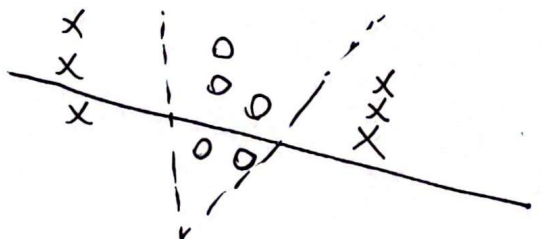
② D 样本全为一个类别, 训练误差也为 0, ~~因为类别一样~~

③ 属性集 $A = \emptyset$ or ~~D 的类取值一样~~, ^表
 A 的取值相同, 即属性相同.

~~假设~~ 因为不存在冲突数据, 属性值相同 \Rightarrow 类相同,
 \Rightarrow 训练误差为 0.

所以一定要用剪枝防止过拟合, 有极高概率过拟合.
_{否则}

4.2. 会出现以下可能, 用最小训练误差做为贪心策略,



会导致本来应该用两条线分类 ~~数据~~ 的数据,

用一条线分割后, 再分割, 不符合样本规律, 泛化能力下降.

答、4.5. 为了避免出现 4.2 问题的情况,

需要让损失函数不为“最小误差”,

而是“~~最小~~ minimal entropy”, 不知道怎么用 logistic

loss function of:

① logistic regression :
$$l(\beta) = \sum_{i=1}^m y_i \cdot p(y=1 | \hat{x}_i, \beta) + (1-y_i) \cdot p(y=0 | \hat{x}_i, \beta)$$
 实现

② logistic decision :
$$l(\beta | D, A) = \sum_{x_i \in D} \min_{y \in \{0,1\}} p(y=1 | \hat{x}_i, \beta) \cdot w_i$$

$$e(\beta | D, A) = h \left(\frac{\sum_{x_i \in \tilde{D}} w_i \cdot p(y=1 | \hat{x}_i, \beta)}{\sum_{x_i \in D} w_i} \right) + h \left(\frac{\sum_{x_i \in D} w_i - \sum_{x_i \in \tilde{D}} w_i \cdot p_i}{\sum_{x_i \in D} w_i} \right)$$

$h(x) = x \cdot \log x.$

用熵值做为 loss functions, 去避免过拟合.

(这是我的方法, 但是导数复杂度很高, 效果可能不好)

4.9. ① $Gini(D) = 1 - \sum_{k=1}^{|Y|} p_k^2$

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

②
$$p = \frac{\sum_{x \in \tilde{D}} w_x}{\sum_{x \in D} w_x}$$

$$\tilde{p}_k = \frac{\sum_{x \in \tilde{D}_k} w_x}{\sum_{x \in \tilde{D}} w_x}$$

($1 \leq k \leq |Y|$) class

$$\tilde{r}_v = \frac{\sum_{x \in \tilde{D}^v} w_x}{\sum_{x \in \tilde{D}} w_x}$$

($1 \leq v \leq V$) attribute.

$$Gini(\tilde{D}) = 1 - \sum_{k=1}^{|Y|} \tilde{p}_k^2,$$

$$Gini_index(D, \alpha) = \sum_{v=1}^V \tilde{r}_v \cdot Gini(\tilde{D}^v)$$

for $x \in D \setminus \tilde{D}$, $w_x \rightarrow \tilde{r}_v \cdot w_x$, $\forall v \leq V$. \square