

Hw 3

1 (3.2)

$$y = \frac{1}{1 + e^{-(w^T x + b)}}$$

$$\frac{\partial y}{\partial w} = \frac{x e^{-(w^T x + b)}}{(1 + e^{-(w^T x + b)})^2} = \vec{x} y (1 - y)$$

$$\frac{\partial^2 y}{\partial w^T \partial w} = \frac{\partial}{\partial w^T} \frac{\partial y}{\partial w} = \frac{\partial}{\partial w^T} (\vec{x} y (1 - y)) = y'^T (1 - y) \vec{x} - y y'^T \vec{x} = x^T x y (1 - y) (1 - 2y)$$

由 $x^T x > 0$, 故 $y \in (0.5, 1)$ 时 $\frac{\partial^2 y}{\partial w^T \partial w} < 0$, 此时函数非凸

$$\text{而 } l(\beta) = \sum_{i=1}^m (-y_i \beta^T x_i + \ln(1 + e^{\beta^T x_i}))$$

$$\frac{\partial l(\beta)}{\partial \beta} = \sum_{i=1}^m (-y_i \vec{x}_i + \frac{1}{1 + e^{\beta^T x_i}} \cdot \vec{x}_i \cdot e^{\beta^T x_i})$$

$$\frac{\partial^2 l(\beta)}{\partial \beta^T \partial \beta} = \sum_{i=1}^m \frac{\beta^T x_i}{(1 + e^{\beta^T x_i})^2} x_i^T x_i \geq 0$$

所以 $l(\beta)$ 为凸函数

2(3.7) 最优 EOC 要求 任意两 编码之间 最小距离最大,

任意两 类别之间 海明距离大, 且反码之间 距离最大

可采取如下编码

	f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9
C_1	1	1	1	1	1	1	1	X	X
C_2	0	0	0	0	1	1	1	X	X
C_3	0	0	1	1	0	0	1	X	X
C_4	0	1	0	1	0	1	0	X	X

类别间 海明距离为 4, f_8, f_9 可为任意编码

$$3. \text{ 令 } S_b = \sum_i m_i (\mu_i - \mu)(\mu_i - \mu)^T = A M A^T = (\mu_1 - \mu, \dots, \mu_N - \mu) \begin{pmatrix} m_1 & & \\ & m_2 & \\ & & \ddots \\ & & & m_N \end{pmatrix} \begin{pmatrix} \mu_1 - \mu \\ \vdots \\ \mu_N - \mu \end{pmatrix}$$

$$\Rightarrow \text{rank}(S_b) = \text{rank}(A M A^T) = \text{rank}(A M^{\frac{1}{2}} (A M^{\frac{1}{2}})^T) = \text{rank}(A M^{\frac{1}{2}}) = \text{rank}(A)$$

而 $\sum m_i (\mu_i - \mu) = 0$ 线性相关,

故 $\text{rank}(S_b) = \text{rank } A \leq N - 1$

4. 证明 $S_b W = \lambda S_w W$

解: 原问题: $\max_w \frac{\text{Tr}(W^T S_b W)}{\text{Tr}(W^T S_w W)}$

$$\Leftrightarrow \max_w \text{Tr}(W^T S_b W) \quad \text{s.t.} \quad \text{tr}(W^T S_w W) = 1.$$

$$\Leftrightarrow L(W, \lambda) = -\text{tr}(W^T S_b W) + \lambda (\text{tr}(W^T S_w W) - 1)$$

$$\frac{\partial L}{\partial W} = 0$$

$$\Rightarrow \lambda \left(\frac{\partial \text{tr}(W^T S_w W)}{\partial W} \right) = \frac{\partial \text{tr}(W^T S_b W)}{\partial W}$$

$$\Rightarrow \lambda (S_w + S_w^T) W = (S_b + S_b^T) W$$

$$\Rightarrow S_b W = \lambda S_w W$$

5. 令 $P = X(X^T X)^{-1} X^T$

$$P^T = X (X (X^T X)^{-1})^T = X ((X^T X)^{-1})^T X^T$$

$$= X (X^T X)^{-1} X^T = P.$$

$$P^2 = X (X^T X)^{-1} X^T X (X^T X)^{-1} X^T = X (X^T X)^{-1} X^T = P.$$

所以 P 是投影矩阵

解释: 对任意向量 y , $\hat{y} = X (X^T X)^{-1} X^T y$ 是在 y 空间的线性投影

Hw 4.

1. (4.1). 反证: 假设不存在, 那么决策树必然有与训练集冲突的数据, 这与假设矛盾,

$$\begin{aligned} 2. (4.9). \quad Gini(D, a) &= p \times Gini-index(\tilde{D}, a) \\ &= p \sum_{v=1}^{|\tilde{V}|} \tilde{r}_v Gini(\tilde{D}^v) \\ &= p \sum_{v=1}^{|\tilde{V}|} \tilde{r}_v (1 - \sum_{i=1}^K \tilde{r}_i^2) \end{aligned}$$

3. 设随机变量 $X \in \{1, \dots, N\}$, 取值为 k 的概率: $P(X=k) = p_k$

$$\text{熵为: } H = -\sum_{k=1}^N p_k \log_2 p_k \quad \text{s.t. } \sum_k p_k = 1.$$

$$\text{Lagrange: } L = -\sum_{k=1}^N p_k \log_2 p_k + \lambda (\sum_{k=1}^N p_k - 1)$$

$$\Rightarrow \frac{\partial L}{\partial p_i} = -\log_2 p_i - \frac{1}{\ln 2} + \lambda = 0$$

$$\Rightarrow \left. \begin{aligned} p_i &= 2^{\lambda - \frac{1}{\ln 2}} \quad (i=1, \dots, N) \\ \frac{\partial L}{\partial \lambda} &= \sum_{i=1}^N p_i - 1 = 0 \end{aligned} \right\} \Rightarrow p_i = \frac{1}{N}$$

$$4. (1). \quad p^+ = \frac{1}{2}, \quad p^- = \frac{1}{2}$$

$$\Rightarrow Ent(D) = -p^+ \log p^+ - p^- \log p^- = 1$$

$$(2). \quad Gain(D, A) = Ent(D) - \sum_{v=1}^{|D|} \frac{|D^v|}{|D|} Ent(D^v)$$

$$= 1 - \left[\frac{4}{10} \left(-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right) + \frac{6}{10} \left(-\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \right) \right]$$

$$= 0.125$$

$$Gain(D, B) = 1 - \left[\frac{5}{10} \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) + \frac{5}{10} \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) \right]$$

$$= 0.029$$

(3) 选取: $C \in \{1.5 + k \mid k=0\}$ 作为二分节点

$$\text{则 } G(D, C, C_k) \text{ 分别为: } \begin{matrix} 0.108, & 0.236, & 0.035, & 0.125, & 0, & 0.035, & 0.108 \\ k=0 & 1 & 2 & 3 & 4 & 5 & 6 \end{matrix}$$

$$(4) \quad G_{ini}(D, A) = \frac{4}{10} \left(1 - \left(\frac{3}{4} \right)^2 - \left(\frac{1}{4} \right)^2 \right) + \frac{6}{10} \left(1 - \left(\frac{2}{6} \right)^2 - \left(\frac{4}{6} \right)^2 \right) \\ = 0.417$$

$$G_{ini}(D, B) = 0.48.$$

因此 A 更优。

$$(5) \quad \begin{cases} Gain_r(D, A) = \frac{Gain(D, A)}{IV(D, A)} \approx 0.129. \\ Gain_r(D, B) = \frac{Gain(D, B)}{IV(D, B)} \approx 0.029 \\ Gain_r(D, C, 2.5) = \frac{Gain(D, C, 2.5)}{IV(D, C)} \approx 0.326. \end{cases}$$

First. 以 $(C, 2.5)$ 为划分。

对 $D'_{C \leq 2.5}$ 为正例, 不再细分

对 $D'_{C > 2.5}$ 有 继续上述操作, 得 C4.5 决策树。

