

HW3.

3.2 证明对于参数 w .

(1) $y = \frac{1}{1 + e^{-(w^T x + b)}}$ 是非凸函数.

(2) $l(\beta) = \sum_{i=1}^m (-y_i \beta^T \hat{x}_i + \ln(1 + e^{\beta^T \hat{x}_i}))$ 是凸函数.

首先关于凸函数的定理:

若 $f(x)$ 二阶连续可微, 那么 $f(x)$ 是 D 上的凸函数的充要条件是, $f(x)$ 的 Hesse 矩阵在 D 上是半正定的.

(1) 对 y , 求 Hessian 矩阵. 要证 y 非凸, 证明该矩阵非半正定即可.

$$\begin{aligned} \frac{\partial y}{\partial w} &= -[1 + e^{-(w^T x + b)}]^{-2} \cdot e^{-(w^T x + b)} \cdot (-x) \\ &= -\frac{(1-y)x}{\frac{1}{y^2}} = y(1-y)x = x(y-y^2) \end{aligned}$$

$$\begin{aligned} \frac{\partial^2 y}{\partial w^2} &= \frac{\partial}{\partial w^T} \left(\frac{\partial y}{\partial w} \right) = \frac{\partial}{\partial w^T} (x(y-y^2)) = \left(\frac{\partial y}{\partial w} \right)^T \left[\frac{\partial (xy - xy^2)}{\partial y} \right] \\ &= x^T (y - y^2) (x - 2xy) \\ &= x x^T (1-2y) (y-y^2) \end{aligned}$$

$\therefore y$ 的值域为 $(0, 1)$, $x x^T$ 相当于 k 倍单位阵

且当 $y \in (0.5, 1)$ 时, $y(1-y)(1-2y) < 0$, $\frac{d}{dw^T} \left(\frac{\partial y}{\partial w} \right)$ 半负定

综上 y 为非凸函数.

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

$$\frac{\partial^2 l}{\partial \beta^2} = \frac{\partial (\frac{\partial l}{\partial \beta})}{\partial \beta^T} = \sum_{i=1}^m \hat{x}_i \hat{x}_i^T p_i(\hat{x}_i; \beta) (1 - p_i(\hat{x}_i; \beta))$$

$$= X P X^T$$

其中 X 为 (n, m) 矩阵, 每一列对应一个样本, P 为对角矩阵

$$p_{ii} = p_i(\hat{x}_i; \beta) (1 - p_i(\hat{x}_i; \beta))$$

$X P X^T$ 对任意向量 z 有:

$$z^T X P X^T z = (x^T z)^T P (x^T z) = V^T P V$$

$$= \sum_i p_{ii} v_i^2 \geq 0$$

因此此 Hessian 矩阵半正定

$l(\beta)$ 对于参数 w 是凸函数

3.7 码长为 9, 类别数为 4. 求最佳 EOC = 元码.

好的纠错码应该达到行分离, 列分离, 且两个分类器的编码不互为反码.

	f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9
$C_1 \rightarrow$	+1	+1	+1	+1	+1	+1	+1	-1	-1
$C_2 \rightarrow$	-1	-1	-1	-1	+1	+1	+1	+1	+1
$C_3 \rightarrow$	-1	-1	+1	+1	-1	-1	+1	-1	+1
$C_4 \rightarrow$	-1	+1	-1	+1	-1	+1	-1	+1	+1

如上的 ECCC 二元码, 不存在编码互为反码的分类器

任意两组分类器之间的海明距离均大于等于 5

4 类中不存在 f_i 和 f_j , 使得 f_i 和 f_j 总是相等 (这种情况, 即为过剩的码长, 对分类效果没有贡献)

下面证明这种情况为最优, 分为 2 步:

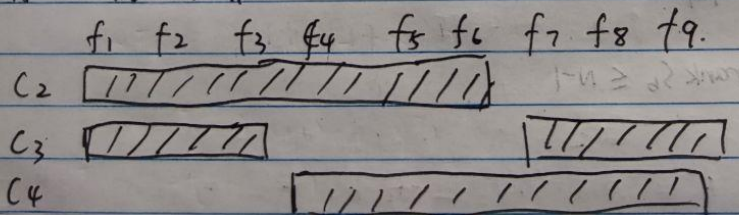
- ① 证明不存在任意两组分类器间海明距离均大于等于 7 的纠错码。
- ② 证明存在任意两组分类器间海明距离均为 6 的纠错码, 但其不满足列分离。

①: 以 C_1 为标准, 若 C_2 与 C_1 海明距离为 7, 则有 7 位编码相异, 不失一般性, 不妨设为 $f_1 \sim f_7$ 位编码相异。

由 C_3 和 C_1 海明距离为 7, 码长为 9 可知, C_3 和 C_2 必有至少 $(7+7-9)=5$ 位编码相同, 但此时 C_3 和 C_2 的海明距离至多为 4, 无法达到 7。

故命题得证。

- ②: 同 ① 的分析, 以 C_1 为标准, 若海明距离为 6, 则 C_1 与 C_2, C_3, C_4 分别有 6 位编码相异。



上图为 1 例, 阴影部分表示与 C_1 相异的 6 位编码位置。

可以看出,其满足任意两组分类器编码的海明距离为6。
 但其中 $\{f_1, f_2, f_3\}$, $\{f_4, f_5, f_6\}$, $\{f_7, f_8, f_9\}$ 三组中,3个互
 不互相分离,不是一个好的 ELOC 纠错码。
 综上,如最开始的任意两组分类器编码的海明距离大于等于5
 的纠错码为最优。

补充题1: 在 LDA 多分类情形下,计算类间散度矩阵 S_b 的秩并证明。

$$S_b = \sum_{i=1}^N m_i (u_i - u)(u_i - u)^T$$

$$= [u_1 - u, u_2 - u, \dots, u_N - u] \begin{pmatrix} m_1 & & \\ & \ddots & \\ & & m_N \end{pmatrix} \begin{pmatrix} (u_1 - u)^T \\ (u_2 - u)^T \\ \vdots \\ (u_N - u)^T \end{pmatrix}$$

$$\text{记 } M = \text{diag}(m_1, m_2, \dots, m_N)$$

$$A = (u_1 - u, u_2 - u, \dots, u_N - u)^T$$

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

$$= \text{rank } (A^T M^{\frac{1}{2}}) = \text{rank } A^T = \text{rank } (u_1 - u, u_2 - u, \dots, u_N - u)$$

$$\text{其中,因为 } \sum_{i=1}^N m_i u_i = \left(\sum_{i=1}^N m_i \right) u.$$

$$\text{即 } \sum_{i=1}^N m_i (u_i - u) = 0, \text{ 故 } \text{rank } (u_1 - u, u_2 - u, \dots, u_N - u) \leq N - 1$$

$$\text{rank } S_b \leq N - 1$$

补充题2. 给出公式 3.45 的推导过程.
已知式 3.44, 优化目标为 $\max_W \frac{\text{tr}(W^T S_b W)}{\text{tr}(W^T S_w W)}$

固定此目标的分母, 即 $\text{tr}(W^T S_w W)$ 为 1, 则等价于

优化问题是:
$$\begin{aligned} \min_W & -\text{tr}(W^T S_b W) \\ \text{s.t.} & \text{tr}(W^T S_w W) = 1 \end{aligned}$$

根据拉格朗日乘子法, 可定义上述优化问题的拉格朗日函数

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

根据矩阵微分式

$$\frac{\partial \text{tr}(X^T B X)}{\partial X_{ij}} = \sum_{q=1}^m b_{iq} x_{qj} + \sum_{p=1}^m b_{pi} x_{pj} = [BX + B^T X]_{ij}$$

$$\frac{\partial \text{tr}(X^T B X)}{\partial X} = [BX + B^T X]$$

$L(W, \lambda)$ 对 W 求偏导则有:

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

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

$$= -2S_b W + 2\lambda S_w W$$

令上式等于 0, 即有:

$$-2S_b W + 2\lambda S_w W = 0$$

$$S_b W = \lambda S_w W$$

补充题3. 证明 $X(X^T X)^{-1} X^T$ 是投影矩阵, 对线性回归模型从

投影角度解释.

① 设向量 b 在矩阵 X 空间上的投影为 p , 则

$$p = Xk = [x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} k_1 \\ \vdots \\ k_n \end{bmatrix} = x_1 k_1 + \dots + x_n k_n$$

其中 $x_1 \dots x_n$ 为 X 的列向量.

令误差向量 $e = b - p$.

则由于 e 垂直于列空间的平面, 则有:

$$\begin{cases} x_1^T (b - p) = 0 \\ x_2^T (b - p) = 0 \\ \vdots \\ x_n^T (b - p) = 0 \end{cases}$$

$$\text{即有 } X^T (b - p) = 0$$

$$\because p = Xk$$

$$\therefore X^T (b - Xk) = 0$$

$$X^T b = X^T X k$$

$$k = (X^T X)^{-1} X^T b$$

$$p = Xk = X(X^T X)^{-1} X^T b$$

$$\text{令 } W = X(X^T X)^{-1} X^T$$

$$p = Wb$$

得出 W 为投影矩阵.

② 解释线性回归模型

可以把特征矩阵 X 看作一个向量组, 每一列(特征)都是一个 n 维向量, 我们有

d 个这样的向量. 我们假设 $d < n$ 且所有特征的线性无关. 则 X 张成的空间是个 d 维空间. 真实值 y 是一个 $n \times 1$ 的向量, 处于 n 维空间中. 多元线性回归就是在 X 张成的 d 维空间中, 寻找 n 维空间中 y 的投影.

HW 4

4.1 决策树的递归停止条件为：(1) 当前结点包含的样本属于同一类别，无需划分。(2) 当前的属性集为空，或是当前样本在所有属性上的取值相同，无法划分。(3) 当前结点的样本集合为空，不能划分。

假设对于训练得到的决策树存在结点，该结点中有无法划分的数据，即存在冲突数据，训练误差不为0，与原假设矛盾。

因此，对于不含冲突数据的训练集，必存在与训练集一致的决策树。

4.9 给定训练集 D 和属性 a ，令 \tilde{D} 表示 D 中在属性 a 上没有缺失值的样本子集。假设 a 有 V 个可取值 $\{a^1, a^2, \dots, a^V\}$ 。令 \tilde{D}^v 表示 \tilde{D} 上 a 属性取值为 a^v 的样本子集， \tilde{D}_k 表示 \tilde{D} 中属于第 k 类 ($k=1, \dots, |Y|$) 的样本子集。

$$\text{显然有 } \tilde{D} = \bigcup_{k=1}^{|Y|} \tilde{D}_k \quad \tilde{D} = \bigcup_{v=1}^V \tilde{D}^v$$

为每个样本赋予一个权重 w_x

$$\text{令 } p = \frac{\sum_{x \in \tilde{D}} w_x}{\sum_{x \in D} w_x} \quad \text{表示在 } a \text{ 上无缺失值的样本比例}$$

$$\tilde{p}_k = \frac{\sum_{x \in \tilde{D}_k} w_x}{\sum_{x \in \tilde{D}} w_x} \quad \text{表示无缺失值样本中第 } k \text{ 类权重比例}$$

$$\tilde{r}_v = \frac{\sum_{x \in \tilde{D}^v} w_x}{\sum_{x \in \tilde{D}} w_x} \quad \text{表示无缺失值样本中 } a \text{ 属性取 } a^v \text{ 的样本权重比例}$$

则由 Gini 指数定义：

$$Gini(D) = 1 - \sum_{k=1}^{|Y|} \tilde{p}_k^2 \quad \text{即排除了所有随机抽取两样本后，其类别标记一致的情况的概率。}$$

采用类似的符号表示，属性 a 的基尼指数定义为

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

(相当于原定义中的 $\frac{|D^v|}{|D|}$ 换为了 \tilde{r}_v)

补充题 1. $X = \{1, 2, \dots, k\}$ $P(X=k) = P_k$

熵 $H(P) = -\sum_k P_k \log_2 P_k$

写为规划问题, 即有:

$$\max -\sum_k P_k \log_2 P_k$$

$$\text{s.t. } \sum_{i=1}^k P_i = 1$$

拉格朗日函数: $L(P, \lambda) = -\sum_{i=1}^k P_i \log_2 P_i + \lambda \left(\sum_{i=1}^k P_i - 1 \right)$

求偏导: $\frac{\partial L(P, \lambda)}{\partial P_i} = -\log_2 P_i - \frac{1}{\ln 2} + \lambda$

$$\frac{\partial L(P, \lambda)}{\partial \lambda} = \sum_{i=1}^k P_i - 1 = 0$$

令 $\frac{\partial L(P, \lambda)}{\partial P_i} = 0$, 则 $P_i = 2^{\lambda - \frac{1}{\ln 2}}$

又 $\sum_{i=1}^k P_i = 1$

$$\therefore P_i = \frac{1}{k} \quad \lambda = \log_2 k + \frac{1}{\ln 2} = -\log_2 k + \frac{1}{\ln 2}$$

由此求出的 P_i 就是 $L(P, \lambda)$ 在条件 $\sum_{i=1}^k P_i = 1$ 下的可能的极值点。又因为这样的点只有一个, 可直接确定这是所求的点。因此使熵最大的分布是均匀分布。

补充题 2.

$$(a). P(\text{类别} = "+") = \frac{1}{2}$$

$$P(\text{类别} = "-") = \frac{1}{2}$$

$$\text{Ent}(D) = -(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2}) = 1$$

(b) 以属性 A 来分析, 有 "T", "F" 两个取值.

$$D^1 = \{1, 2, 3, 8\} \quad D^1 \text{ 中有 3 个 "+" , 1 个 "-"}$$

$$D^2 = \{4, 5, 6, 7, 9, 10\} \quad D^2 \text{ 中有 2 个 "+" , 4 个 "-"}$$

$$\text{Ent}(D^1) = -(\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4})$$

$$= 2 - \frac{3}{4} \log_2 3$$

$$\text{Ent}(D^2) = -(\frac{2}{6} \log_2 \frac{2}{6} + \frac{4}{6} \log_2 \frac{4}{6})$$

$$= \log_2 3 - \frac{2}{3}$$

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

$$= 1 - \frac{2}{5} (2 - \frac{3}{4} \log_2 3) - \frac{3}{5} (\log_2 3 - \frac{2}{3})$$

$$= \frac{3}{5} - \frac{3}{10} \log_2 3$$

以属性 B 来分析, 有 "T", "F" 两个取值.

$$D^1 = \{1, 2, 5, 6, 9\} \quad D^1 \text{ 中有 2 个 "+" , 3 个 "-"}$$

$$D^2 = \{3, 4, 7, 8, 10\} \quad D^2 \text{ 中有 3 个 "+" , 2 个 "-"}$$

$$\text{Ent}(D^1) = -(\frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5})$$

$$= \log_2 5 - \frac{2}{5} - \frac{3}{5} \log_2 3$$

$$\text{Ent}(D^2) = -(\frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5})$$

$$= \log_2 5 - \frac{2}{5} - \frac{3}{5} \log_2 3$$

$$\text{Gain}(D, B) = \text{Ent}(D) - \sum_{v=1}^2 \frac{|D^v|}{|D|} \text{Ent}(D^v)$$

$$= 1 - (\log_2 5 - \frac{2}{5} \log_2 3)$$

$$= \frac{7}{5} + \frac{2}{5} \log_2 3 - \log_2 5$$

(c) 统计 C 类, 共 8 种取值, 1.0, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 8.0

不妨取划分点分别为 1.5, 2.5, 3.5, 4.5, 5.5, 6.5, 7.5, 1

① 划分点为 1.5

$D^1 = \{1\}$ 为 "+"

$D^2 = \{2, 3, 4, 5, 6, 7, 8, 9\}$ 4 个为 "+", 5 个为 "-"

$$\text{Ent}(D^1) = -\log_2 1 = 0$$

$$\text{Ent}(D^2) = -(\frac{4}{9} \log_2 \frac{4}{9} + \frac{5}{9} \log_2 \frac{5}{9}) = -\frac{8}{9} - \frac{5}{9} \log_2 5 + 2 \log_2 3$$

$$\text{Gain}(D, C)_1 = \text{Ent}(D) - \sum_{v=1}^2 \frac{|D^v|}{|D|} \text{Ent}(D^v)$$

$$= 1 - \frac{1}{10} \times 0 + \frac{9}{10} [\frac{8}{9} + \frac{5}{9} \log_2 5 - 2 \log_2 3]$$

$$= \frac{9}{5} + \frac{1}{5} \log_2 5 - \frac{9}{5} \log_2 3$$

② 划分点为 2.5

$D^1 = \{1, 10\}$ 2 个均为 "+"

$D^2 = \{2, 3, 4, 5, 6, 7, 8, 9\}$ 3 个为 "+", 5 个为 "-"

$$\text{Ent}(D^1) = 0$$

$$\text{Ent}(D^2) = -(\frac{3}{8} \log_2 \frac{3}{8} + \frac{5}{8} \log_2 \frac{5}{8}) = 3 - \frac{3}{8} \log_2 3 - \frac{5}{8} \log_2 5$$

$$\text{Gain}(D, C)_2 = \text{Ent}(D) - \sum_{v=1}^2 \frac{|D^v|}{|D|} \text{Ent}(D^v)$$

$$= 1 - \frac{2}{10} (3 - \frac{3}{8} \log_2 3 - \frac{5}{8} \log_2 5) = \frac{3}{10} \log_2 3 + \frac{1}{5} \log_2 5 - \frac{7}{5}$$

③ 划分点为 3.5

$$D^1 = \{1, 6, 10\}$$

2个为"+", 1个为"-"

$$D^2 = \{2, 3, 4, 5, 7, 8, 9\}$$

3个为"+", 4个为"-"

$$\text{Ent}(D^1) = -\left[\frac{2}{3}\log_2\frac{2}{3} + \frac{1}{3}\log_2\frac{1}{3}\right] = \log_2 3 - \frac{2}{3}$$

$$\text{Ent}(D^2) = -\left[\frac{3}{7}\log_2\frac{3}{7} + \frac{4}{7}\log_2\frac{4}{7}\right] = \log_2 7 - \frac{3}{7}\log_2 3 - \frac{8}{7}$$

$$\text{Gain}(D, C)_3 = \text{Ent}(D) - \sum_{v=1}^2 \frac{|D^v|}{|D|} \text{Ent}(D^v)$$

$$0.8 - 0.5 \cdot 0.8 - 0.5 \cdot 0.2 = 0.4 - \frac{1}{10}(\log_2 3 - \frac{2}{3}) - \frac{7}{10}(\log_2 7 - \frac{3}{7}\log_2 3 - \frac{8}{7})$$

$$= 0.4 - \frac{1}{10}\log_2 3 + \frac{2}{30} - \frac{7}{10}\log_2 7 + \frac{21}{70}\log_2 3 + \frac{56}{70} = 2 - \frac{7}{10}\log_2 7$$

④ 划分点为 4.5

$$D^1 = \{1, 4, 6, 10\}$$

3个为"+", 1个为"-"

$$D^2 = \{2, 3, 5, 7, 8, 9\}$$

2个为"+", 4个为"-"

$$\text{Ent}(D^1) = -\left[\frac{3}{4}\log_2\frac{3}{4} + \frac{1}{4}\log_2\frac{1}{4}\right] = 2 - \frac{3}{4}\log_2 3$$

$$\text{Ent}(D^2) = -\left[\frac{2}{6}\log_2\frac{2}{6} + \frac{4}{6}\log_2\frac{4}{6}\right] = \log_2 3 - \frac{2}{3}$$

$$\text{Gain}(D, C)_4 = \text{Ent}(D) - \sum_{v=1}^2 \frac{|D^v|}{|D|} \text{Ent}(D^v)$$

$$= \frac{3}{5} - \frac{3}{10}\log_2 3$$

⑤ 划分点为 5.5

$$D^1 = \{1, 3, 4, 6, 9, 10\}$$

3个为"+", 3个为"-"

$$D^2 = \{2, 5, 7, 8\}$$

2个为"+", 2个为"-"

$$\text{Ent}(D^1) = -\log_2\frac{1}{2} = 1$$

$$\text{Ent}(D^2) = -\log_2\frac{1}{2} = 1$$

$$\text{Gain}(D, C)_5 = \text{Ent}(D) - \sum_{v=1}^2 \frac{|D^v|}{|D|} \text{Ent}(D^v) = 0$$

⑥ 划分点为 6.5

$$D' = \{1, 2, 3, 4, 6, 9, 10\}$$

4个为"+", 3个为="-"

$$D'' = \{5, 7, 8\}$$

1个为"+", 2个为="-"

$$\text{Ent}(D') = -\left[\frac{4}{7}\log_2\frac{4}{7} + \frac{3}{7}\log_2\frac{3}{7}\right] = \log_2 7 - \frac{3}{7}\log_2 3 - \frac{8}{7}$$

$$\text{Ent}(D'') = -\left[\frac{1}{3}\log_2\frac{1}{3} + \frac{2}{3}\log_2\frac{2}{3}\right] = \log_2 3 - \frac{2}{3}$$

$$\begin{aligned}\text{Gain}(D, C)_6 &= \text{Ent}(D) - \sum_{v=1}^2 \frac{|D^v|}{|D|} \text{Ent}(D^v) \\ &= 2 - \frac{7}{10} \log_2 7\end{aligned}$$

⑦ 划分点为 7.5

$$D' = \{1, 2, 3, 4, 5, 6, 8, 9, 10\}$$

5个为"+", 4个为="-"

$$D'' = \{7\}$$

为 "-"

$$\text{Ent}(D') = -\left(\frac{5}{9}\log_2\frac{5}{9} + \frac{4}{9}\log_2\frac{4}{9}\right) = -\frac{8}{9} - \frac{5}{9}\log_2 5 + 2\log_2 3$$

$$\text{Ent}(D'') = 0$$

$$\begin{aligned}\text{Gain}(D, C)_7 &= \text{Ent}(D) - \sum_{v=1}^2 \frac{|D^v|}{|D|} \text{Ent}(D^v) \\ &= \frac{9}{5} + \frac{1}{2}\log_2 5 - \frac{9}{5}\log_2 3\end{aligned}$$

(d) 对于属性 A:

$$\text{Gini}(D') = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = \frac{3}{8}$$

$$\text{Gini}(D'') = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = \frac{4}{9}$$

$$\text{Gini-index}(D, A) = \frac{3}{8} \times \frac{2}{5} + \frac{4}{9} \times \frac{3}{5} = \frac{5}{12}$$

对于属性B:

$$Gini(D') = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = \frac{12}{25}$$

$$Gini(D'') = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = \frac{12}{25}$$

$$Gini_index(D, B) = \frac{12}{25}$$

$$\therefore Gini_index(D, A) < Gini_index(D, B)$$

\therefore 属性B为最优划分

(e) 计算可知, *C的划分中信息增益最高的为

以 2.5 为划分点

$$\text{且 } Gain(D, C)_2 \approx 0.3645$$

$$\text{信息增益比 } Gain_R(D, C)_2 = \frac{\frac{3}{10} \log_2 3 + \frac{1}{2} \log_2 5 - \frac{7}{5}}{-\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5}} \approx 0.3275$$

$$Gain_R(D, A) = \frac{\frac{3}{5} - \frac{3}{10} \log_2 3}{-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}} \approx 0.1282$$

$$Gain_R(D, B) = \frac{\frac{7}{5} + \frac{3}{5} \log_2 3 - \log_2 5}{-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}} \approx 0.0291$$

故 $Gain_R(D, C)_2 > Gain_R(D, A) > Gain_R(D, B)$

选择以 C 属性的情况为根结点

决策树

