

3.2

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

求导

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

二阶导

$$\frac{\partial}{\partial w^T} \left(\frac{\partial y}{\partial w^T} \right) = x(1 - 2y) \left(\frac{\partial y}{\partial w} \right)^T = x(1 - 2y)x^T(y - y^2)$$

对于 xPx^T , 对任意向量 Z 有 $Z^T xPx^T Z = \sum_i P_{ii}v_i^2 \geq 0$, 因此其海森矩阵半正定
在二阶导中, xx^T 的秩为1, 非零特征向量为1, y 的值域为 $(0, 1)$, 当 $y \in (0.5, 1)$ 时

$$y(y - 1)(1 - 2y) < 0$$

二阶导半负定, 则公式(3.18)非凸

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

$$\frac{\partial}{\partial \beta^T} \left(\frac{\partial l}{\partial \beta} \right) = \sum_{i=1}^m \hat{x}_i \hat{x}_i^T p_1(\hat{x}_i; \beta)(1 - p_1(\hat{x}_i; \beta)) \geq 0$$

所以公式(3.27)为凸函数

3.7

在类别为4时, 其可行的编码有7种

	f0	f1	f2	f3	f4	f5	f6
c1	1	1	1	1	1	1	1
c2	0	0	0	0	1	1	1
c3	0	0	1	1	0	0	1
c4	0	1	0	1	0	1	0

当码长为9时, 那么 f_6 之后加任意两个编码, 即为最优编码, 因为此时再加任意的编码都是先有编码的反码, 此时, 类别之间最小的海明距离都为4, 不会再增加。

▲在LDA多分类情形下, 试计算类间散度矩阵 S_b 的秩并证明

$$\text{rank}(S_b) \leq k - 1$$

其中 k 为属性的数量

因为.

$$S_b = \sum_{j=1}^k N_j (\mu_j - \mu)(\mu_j - \mu)^T$$

对于单独的 $\mu_j - \mu$, 它的秩为1, 因此协方差矩阵相加后的秩 $\leq k$, 但是由于均值的性质, 最后一个 μ_k 能够用前 $k - 1$ 个 μ_j 线性表示, 因此

$$\text{rank}(S_b) \leq k - 1$$

▲给出公式3.45的推导过程

$$\max_W \frac{\text{tr}(W^T S_b W)}{\text{tr}(W^T S_w W)} = \max_W \frac{\sum w_i^T S_b w_i}{\sum w_i^T S_w w_i}$$

上式为广义瑞利熵，其中 w_i 表示 W 只有第 i 行不为0的矩阵， $W = \sum w_i$

$$w_i^T = \begin{bmatrix} 0 & 0 & 0 & \cdots & 0 \\ & & & \cdots & \\ w_{i1} & w_{i2} & w_{i3} & \cdots & w_{in} \\ & & & \cdots & \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}$$

此时最优化问题同式(3.35)，利用拉格朗日乘子法，即等价于求

$$S_b W = \lambda S_w W$$

▲证明 $X(X^T X)^{-1} X^T$ 是投影矩阵，并对线性回归模型从投影角度解释

假设向量 b 在空间 X 上的投影为 p ，则

$$p = [X_1 \quad X_2 \quad \cdots \quad X_n] \begin{bmatrix} x_1 \\ x_2 \\ \cdots \\ x_n \end{bmatrix} = Xx$$

根据投影的性质， $b - p$ 应该和空间 X 垂直

$$\begin{aligned} X^T(b - p) &= 0 \\ X^T(b - Xx) &= 0 \\ X^T b &= X^T Xx \\ x &= (X^T X)^{-1} X^T b \end{aligned}$$

代回

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

即

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

为投影矩阵

我们可以将特征矩阵 X 看作是一个向量组，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 $d < n$ 且所有特征都线性无关，那 X 张成的空间是个 d 维度空间。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。

4.1

假设不存在与训练集一致的决策树，那么训练集训练得到的决策树至少有一个节点上存在无法划分的多个数据（若节点上没有冲突数据，那么总是能够将数据分开的）这与前提-不含冲突数据 矛盾，因此必存在与训练集一致的决策树

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. 对问题(1), 显然我们仅可根据 \tilde{D} 来判断属性 a 的优劣. 假定属性 a 有 V 个可取值 $\{a^1, a^2, \dots, a^V\}$, 令 \tilde{D}^v 表示 \tilde{D} 中在属性 a 上取值为 a^v 的样本子集, \tilde{D}_k 表示 \tilde{D} 中属于第 k 类 ($k = 1, 2, \dots, |\mathcal{Y}|$) 的样本子集, 则显然有 $\tilde{D} = \bigcup_{k=1}^{|\mathcal{Y}|} \tilde{D}_k$, $\tilde{D} = \bigcup_{v=1}^V \tilde{D}^v$. 假定我们为每个样本 \mathbf{x} 赋予一个权重 $w_{\mathbf{x}}$, 并定义

$$\rho = \frac{\sum_{\mathbf{x} \in \tilde{D}} w_{\mathbf{x}}}{\sum_{\mathbf{x} \in D} w_{\mathbf{x}}}, \quad (4.9)$$

$$\tilde{p}_k = \frac{\sum_{\mathbf{x} \in \tilde{D}_k} w_{\mathbf{x}}}{\sum_{\mathbf{x} \in \tilde{D}} w_{\mathbf{x}}} \quad (1 \leq k \leq |\mathcal{Y}|), \quad (4.10)$$

$$\tilde{r}_v = \frac{\sum_{\mathbf{x} \in \tilde{D}^v} w_{\mathbf{x}}}{\sum_{\mathbf{x} \in \tilde{D}} w_{\mathbf{x}}} \quad (1 \leq v \leq V). \quad (4.11)$$

将基尼指数推广为:

$$Gini_{index}(D, a) = \rho \times \left(\sum_{v=1}^V \tilde{r}_v Gini(\tilde{D}^v) \right)$$

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

▲ 假设离散随机变量 $X \in \{1, \dots, K\}$, 其取值为 k 的概率 $P(X = k) = p_k$, 其熵为 $H(p) = -\sum_k p_k \log_2 p_k$, 试用拉格朗日乘子法证明熵最大的分布为均匀分布

写出拉格朗日函数:

$$L(p_1, \dots, p_K) = -\sum_i p_i \log_2 p_i - \lambda \left(\sum_i p_i - 1 \right)$$

对 p_1 到 p_K 求偏导可得方程

$$\begin{cases} -\sum_i p_i \log_2 p_i - \lambda (\sum_i p_i - 1) = 0 \\ -\log_2 p_i - \ln 2 - \lambda = 0 \\ \sum_i p_i - 1 = 0 \end{cases}$$

可得到 p_i 平均分布时, 可以得到熵的最大值

▲ 习题

- 下表表示的二分类数据集，具有三个属性A,B,C，样本标记为两类“+”，“-”。请运用你学过的知识完成如下问题：

实例	A	B	C	类别
1	T	T	1.0	+
2	T	T	6.0	+
3	T	F	5.0	-
4	F	F	4.0	+
5	F	T	7.0	-
6	F	T	3.0	-
7	F	F	8.0	-
8	T	F	7.0	+
9	F	T	5.0	-
10	F	F	2.0	+

- 整个训练样本关于类属性的熵是多少
- 数据集中A，B两个属性的信息增益各是多少
- 对于属性C，计算所有可能划分的信息增益
- 根据Gini指数，A和B两个属性哪个是最优划分
- 采用算法C4.5，构造决策树

(a)

$$Ent(D) = -\frac{5}{10}\log\frac{5}{10} - \frac{5}{10}\log\frac{5}{10} = 1$$

(b)

$$\begin{aligned} Gain(D, A) &= Ent(D) - [\frac{4}{10}Ent(D_1) + \frac{6}{10}Ent(D_2)] \\ &= 1 - [\frac{4}{10}(-\frac{3}{4}\log\frac{3}{4} - \frac{1}{4}\log\frac{1}{4}) + \frac{6}{10}(-\frac{2}{6}\log\frac{2}{6} - \frac{4}{6}\log\frac{4}{6})] \\ &= 0.1245 \end{aligned}$$

$$\begin{aligned} Gain(D, B) &= Ent(D) - [\frac{5}{10}Ent(D_1) + \frac{5}{10}Ent(D_2)] \\ &= 1 - [\frac{5}{10}(-\frac{2}{5}\log\frac{2}{5} - \frac{3}{5}\log\frac{3}{5}) + \frac{5}{10}(-\frac{3}{5}\log\frac{3}{5} - \frac{2}{5}\log\frac{2}{5})] \\ &= 0.0290 \end{aligned}$$

(c)

$$t = 1.5 \quad Gain(D, C, t) = 1 + \frac{9}{10}\frac{4}{9}\log\frac{4}{9} + \frac{5}{10}\frac{5}{9}\log\frac{5}{9} = 0.1080$$

$$t = 2.5 \quad Gain(D, C, t) = 1 + \frac{8}{10}\frac{3}{8}\log\frac{3}{8} + \frac{8}{10}\frac{5}{8}\log\frac{5}{8} = 0.2365$$

$$t = 3.5 \quad Gain(D, C, t) = 1 + \frac{3}{10}\frac{2}{3}\log\frac{2}{3} + \frac{3}{10}\frac{1}{3}\log\frac{1}{3} + \frac{7}{10}\frac{3}{7}\log\frac{3}{7} + \frac{7}{10}\frac{4}{7}\log\frac{4}{7} = 0.0349$$

$$t = 4.5 \quad Gain(D, C, t) = 1 + \frac{4}{10}\frac{3}{4}\log\frac{3}{4} + \frac{4}{10}\frac{1}{4}\log\frac{1}{4} + \frac{6}{10}\frac{2}{6}\log\frac{2}{6} + \frac{6}{10}\frac{4}{6}\log\frac{4}{6} = 0.1245$$

$$t = 5.5 \quad Gain(D, C, t) = 0$$

$$t = 7.5 \quad Gain(D, C, t) = 1 + \frac{9}{10}\frac{5}{9}\log\frac{5}{9} + \frac{9}{10}\frac{4}{9}\log\frac{4}{9} = 0.1080$$

(d)

$$Gini_index(D, A) = \frac{4}{10}(1 - \frac{3^2}{4} - \frac{1^2}{4}) + \frac{6}{10}(1 - \frac{2^2}{6} - \frac{4^2}{6}) = 0.4167$$

$$Gini_index(D, B) = \frac{5}{10} \left(1 - \frac{2^2}{5} - \frac{3^2}{5} \right) + \frac{5}{10} \left(1 - \frac{2^2}{5} - \frac{3^2}{5} \right) = 0.48$$

A是最优化分

(e)

