

8.3

$$\begin{aligned} P(H(\mathbf{x}) \neq f(\mathbf{x})) &= \sum_{k=0}^{\lfloor T/2 \rfloor} \binom{T}{k} (1-\epsilon)^k \epsilon^{T-k} \\ &\leq \exp\left(-\frac{1}{2}T(1-2\epsilon)^2\right) \end{aligned}$$

[推导]: 由基分类器相互独立, 设 X 为 T 个基分类器分类正确的次数, 因此 $X \sim B(T, 1-\epsilon)$

$$\begin{aligned} P(H(\mathbf{x}) \neq f(\mathbf{x})) &= P(X \leq \lfloor T/2 \rfloor) \\ &\leq P(X \leq T/2) \\ &= P\left[X - (1-\epsilon)T \leq \frac{T}{2} - (1-\epsilon)T\right] \\ &= P\left[X - (1-\epsilon)T \leq -\frac{T}{2}(1-2\epsilon)\right] \end{aligned}$$

根据Hoeffding不等式 $P(X - (1-\epsilon)T \leq -kT) \leq \exp(-2k^2T)$ 令 $k = \frac{(1-2\epsilon)}{2}$ 得

$$\begin{aligned} P(H(\mathbf{x}) \neq f(\mathbf{x})) &= \sum_{k=0}^{\lfloor T/2 \rfloor} \binom{T}{k} (1-\epsilon)^k \epsilon^{T-k} \\ &\leq \exp\left(-\frac{1}{2}T(1-2\epsilon)^2\right) \end{aligned}$$

8.5-8.8

[推导]: 由式(8.4)可知 $H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x})$

又由式(8.11)可知

$$\alpha_t = \frac{1}{2} \ln\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$$

该分类器的权重只与分类器的错误率负相关(即错误率越大, 权重越低)

(1)先考虑指数损失函数 $e^{-f(\mathbf{x})H(\mathbf{x})}$ 的含义: f 为真实函数, 对于样本 \mathbf{x} 来说, $f(\mathbf{x}) \in \{-1, +1\}$ 只能取两个值, 而 $H(\mathbf{x})$ 是一个实数; 当 $H(\mathbf{x})$ 的符号与 $f(\mathbf{x})$ 一致时, $f(\mathbf{x})H(\mathbf{x}) > 0$, 因此 $e^{-f(\mathbf{x})H(\mathbf{x})} = e^{-|H(\mathbf{x})|} < 1$, 且 $|H(\mathbf{x})|$ 越大指数损失函数 $e^{-f(\mathbf{x})H(\mathbf{x})}$ 越小(这很合理: 此时 $|H(\mathbf{x})|$ 越大意味着分类器本身对预测结果的信心越大, 损失应该越小; 若 $|H(\mathbf{x})|$ 在零附近, 虽然预测正确, 但表示分类器本身对预测结果信心很小, 损失应该较大); 当 $H(\mathbf{x})$ 的符号与 $f(\mathbf{x})$ 不一致时, $f(\mathbf{x})H(\mathbf{x}) < 0$, 因此 $e^{-f(\mathbf{x})H(\mathbf{x})} = e^{|H(\mathbf{x})|} > 1$, 且 $|H(\mathbf{x})|$ 越大指数损失函数越大(这很合理: 此时 $|H(\mathbf{x})|$ 越大意味着分类器本身对预测结果的信心越大, 但预测结果是错的, 因此损失应该越大; 若 $|H(\mathbf{x})|$ 在零附近, 虽然预测错误, 但表示分类器本身对预测结果信心很小, 虽然错了, 损失应该较小); (2)符号 $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\cdot]$ 的含义: \mathcal{D} 为概率分布, 可简单理解为在数据集 \mathcal{D} 中进行一次随机抽样, 每个样本被取到的概率; $\mathbb{E}[\cdot]$ 为经典的期望, 则综合起来 $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\cdot]$ 表示在概率分布 \mathcal{D} 上的期望, 可简单理解为对数据集 \mathcal{D} 以概率 \mathcal{D} 进行加权后的期望。

$$\begin{aligned} \ell_{\text{exp}}(H|\mathcal{D}) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[e^{-f(\mathbf{x})H(\mathbf{x})} \right] \\ &= P(f(\mathbf{x}) = 1|\mathbf{x}) * e^{-H(\mathbf{x})} + P(f(\mathbf{x}) = -1|\mathbf{x}) * e^{H(\mathbf{x})} \end{aligned}$$

由于 $P(f(\mathbf{x}) = 1|\mathbf{x})$ 和 $P(f(\mathbf{x}) = -1|\mathbf{x})$ 为常数

故式(8.6)可轻易推知

$$\frac{\partial \ell_{\text{exp}}(H|\mathcal{D})}{\partial H(\mathbf{x})} = -e^{-H(\mathbf{x})} P(f(\mathbf{x}) = 1|\mathbf{x}) + e^{H(\mathbf{x})} P(f(\mathbf{x}) = -1|\mathbf{x})$$

令式(8.6)等于0可得

式(8.7)

$$H(\mathbf{x}) = \frac{1}{2} \ln \frac{P(f(\mathbf{x}) = 1|\mathbf{x})}{P(f(\mathbf{x}) = -1|\mathbf{x})}$$

式(8.8)显然成立

$$\begin{aligned} \text{sign}(H(\mathbf{x})) &= \text{sign}\left(\frac{1}{2} \ln \frac{P(f(\mathbf{x}) = 1|\mathbf{x})}{P(f(\mathbf{x}) = -1|\mathbf{x})}\right) \\ &= \begin{cases} 1, & P(f(\mathbf{x}) = 1|\mathbf{x}) > P(f(\mathbf{x}) = -1|\mathbf{x}) \\ -1, & P(f(\mathbf{x}) = 1|\mathbf{x}) < P(f(\mathbf{x}) = -1|\mathbf{x}) \end{cases} \\ &= \arg \max_{y \in \{-1, 1\}} P(f(\mathbf{x}) = y|\mathbf{x}) \end{aligned}$$

8.16

$$\begin{aligned} h_t(\mathbf{x}) &= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\frac{e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]} f(\mathbf{x})h(\mathbf{x}) \right] \\ &= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [f(\mathbf{x})h(\mathbf{x})] \end{aligned}$$

[推导]：假设x的概率分布是f(x) (注:本书中概率分布全都是 $\mathcal{D}(x)$)

$$\mathbb{E}(g(\mathbf{x})) = \sum_{i=1}^{|D|} f(\mathbf{x}_i)g(\mathbf{x}_i)$$

故可得

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H(\mathbf{x})}] = \sum_{i=1}^{|D|} \mathcal{D}(\mathbf{x}_i) e^{-f(\mathbf{x}_i)H(\mathbf{x}_i)}$$

由式(8.15)可知

$$\mathcal{D}_t(\mathbf{x}_i) = \mathcal{D}(\mathbf{x}_i) \frac{e^{-f(\mathbf{x}_i)H_{t-1}(\mathbf{x}_i)}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]}$$

所以式(8.16)可以表示为

$$\begin{aligned} &\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\frac{e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]} f(\mathbf{x})h(\mathbf{x}) \right] \\ &= \sum_{i=1}^{|D|} \mathcal{D}(\mathbf{x}_i) \frac{e^{-f(\mathbf{x}_i)H_{t-1}(\mathbf{x}_i)}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]} f(\mathbf{x}_i)h(\mathbf{x}_i) \\ &= \sum_{i=1}^{|D|} \mathcal{D}_t(\mathbf{x}_i) f(\mathbf{x}_i) h(\mathbf{x}_i) \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [f(\mathbf{x})h(\mathbf{x})] \end{aligned}$$

【注】：由下式(*)也可推至式(8.16)

$$P(f(x) = 1|x)e^{-H(x)} + P(f(x) = -1|x)e^{H(x)}(*)$$

首先式(*)可以拆成n个式子,n的个数为x的取值个数

$$P(f(x_i) = 1|x_i)e^{-H(x_i)} + P(f(x_i) = -1|x_i)e^{H(x_i)}(i = 1, 2, \dots, n)(**)$$

当 x_i 确定的时候 $P(f(x_i) = 1|x_i)$ 与 $P(f(x_i) = -1|x_i)$ 其中有一个为0, 另一个为1

则式(**)可以化简成

$$e^{-f(x_i)H(x_i)}(i = 1, 2, \dots, n)(***)$$

拆成n个式子是根据不同的x来拆分的, 可以把 $x = x_i$ 看成一个事件, 设为事件 A_i 。

当事件 A_i 发生时, 事件 A_j 一定不发生, 即各事件互斥, 而且各个事件发生的概率是 $P(A_i) = \mathcal{D}(x_i)$

此时可以考虑成原来的x被分成了n叉树, 每个路径的概率是 $\mathcal{D}(x_i)$, 叶子结点的值是 $e^{-f(x_i)H(x_i)}$ 相乘再相加即为期望, 同式(8.16)