

证明: 定义: $\text{Gain}(D, a) = \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v)$.

由Ent定义: $\text{Ent}(D) = - \sum_{k=1}^n p_k \log_2 p_k = - \sum_{k=1}^n \frac{|D^k|}{|D|} \log_2 \frac{|D^k|}{|D|}$.

其中: n 为类别个数.

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

$$= - \sum_{k=1}^n \frac{|D^k|}{|D|} \log_2 \frac{|D^k|}{|D|} + \sum_{v=1}^V \frac{|D^v|}{|D|} \left[\sum_{k=1}^n \frac{|D^{vk}|}{|D^v|} \log_2 \frac{|D^{vk}|}{|D^v|} \right]$$

这两项均表示权重, 因此可交换, 且形式为乘积, 因此可交换, 即表示两权重之积.

换个字母

$$= - \sum_{k=1}^n \frac{|D^k|}{|D|} \log_2 \frac{|D^k|}{|D|} + \sum_{k=1}^n \frac{|D^k|}{|D|} \sum_{v=1}^V \frac{|D^{vk}|}{|D^k|} \log_2 \frac{|D^{vk}|}{|D^v|}$$

提出因子

$$= \sum_{k=1}^n \frac{|D^k|}{|D|} \left(\sum_{v=1}^V \frac{|D^{vk}|}{|D^k|} \log_2 \frac{|D^{vk}|}{|D^k|} - \log_2 \frac{|D^k|}{|D|} \right).$$

故欲证 $\text{Gain}(D, a) \geq 0$,

只需证: $\sum_{v=1}^V \frac{|D^{vk}|}{|D^k|} \log_2 \frac{|D^{vk}|}{|D^k|} \geq \log_2 \frac{|D^k|}{|D|}$ for $\forall k \in \{1, 2, \dots, n\}$ 成立.

注意到: $|D|$ 与 $|D^k|$ 均表示样本量, 即 $|D| = |D^k|$

↓

$$\therefore \text{即证 } \sum_{v=1}^V \frac{|D^{vk}|}{|D^k|} \log_2 \frac{|D^{vk}|}{|D^k|} \geq \log_2 \frac{|D^k|}{|D|} \text{ for } \forall k \in \{1, 2, \dots, n\} \text{ 成立.}$$

$f(x) = \log_2(x)$ 为凹函数 ($f''(x) < 0$), $\therefore \forall \theta_1, \theta_2, \dots, \theta_V \geq 0$, subject to $\theta_1 + \theta_2 + \dots + \theta_V = 1$, 有 $f(\theta_1 x_1 + \theta_2 x_2 + \dots + \theta_V x_V) \geq \theta_1 f(x_1) + \theta_2 f(x_2) + \dots + \theta_V f(x_V)$.

取 $\theta_i = \frac{|D^{ik}|}{|D^k|}$ ($i=1, 2, \dots, V$), 取

即证 $\sum_{v=1}^V \frac{|D^{vk}|}{|D^k|} \left(-\log_2 \frac{|D^{vk}|}{|D^k|} \right) \geq \log_2 \frac{|D^k|}{|D|}$ for $\forall k \in \{1, 2, 3, \dots, n\}$ 成立.

由 $f(x) = -\log_2(x)$ 为凸函数, $\therefore \forall \theta_1, \theta_2, \dots, \theta_V \geq 0$, subject to $\theta_1 + \theta_2 + \dots + \theta_V = 1$, 有 $f(\theta_1 x_1 + \theta_2 x_2 + \dots + \theta_V x_V) \leq \theta_1 f(x_1) + \theta_2 f(x_2) + \dots + \theta_V f(x_V)$.

取 $\theta_i = \frac{|D^{ik}|}{|D^k|}$ ($i=1, 2, 3, \dots, V$), 取 $x_i = \frac{|D^k|}{|D^{ik}|}$ ($i=1, 2, \dots, V$). (有 $\theta_1 + \theta_2 + \dots + \theta_V = 1$)

\therefore 左式 $\geq -\log_2 (1 + 1 + \dots + 1) = \log_2 \frac{1}{V}$.

\therefore 只需证 $V \cdot \frac{|D^k|}{|D^k|} \leq |D|$ 即 $V \leq \frac{|D|}{|D^k|}$.

只需证 $V \leq \frac{|D|}{|D^k|}$ min. 考虑本身意义, 当且仅当类别 (输出种类) 个数为 1 时, $\frac{|D|}{|D^k|} = 1$.



上述证明不能进行下去, 不妨利用条件概率记:

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

$$= - \sum_{k=1}^n p_k \log_2 p_k + \sum_{x=1}^n p_x \sum_{k=1}^n p(k|x) \log_2 (k|x)$$

简记作

$$= - \sum_k p_k \log_2 p_k + \sum_x p_x \sum_k p(k|x) \log_2 (k|x)$$

由条件概率性质

$$= - \sum_k \sum_x p(x, k) \log_2 p_k + \sum_k \sum_x p(k, x) \log_2 (k|x)$$

$$= \sum_x \sum_k p(x, k) \log_2 \frac{p(k|x)}{p(k)}$$

~~$-f(x) = \log_2 x$~~

$$= \sum_x \sum_k p(x, k) \log_2 \frac{p(x)p(k)}{p(x, k)}$$

由 $f(x) = \log_2 x$

$$\geq - \log_2 p(x, k) \cdot \frac{p(x)p(k)}{p(x, k)} = - \log_2 p(x)p(k)$$

$\because p(x), p(k) \leq 1, \therefore -\log_2 p(x)p(k) \leq 0, \therefore$ 原式 ≥ 0 . 得证.



①方法也可证明:

原因:之前写错一个字母.

求证: $\sum_{v=1}^V \frac{|D^v K|}{|DK|} \log_2 \frac{|D^v K|}{|D^v|} \geq \log_2 \frac{|DK|}{|D|}$

由 $-\log x$ 为凸函数得: (凸函数性质)

$$\sum_{v=1}^V \left(\frac{|D^v K|}{|DK|} \right) \left(-\log_2 \frac{|D^v|}{|D^v K|} \right)$$

→ 系数正, 且和为1.

$$\geq -\log_2 \left(\frac{|D^1|}{|DK|} + \frac{|D^2|}{|DK|} + \dots + \frac{|D^V|}{|DK|} \right) = \log_2 \frac{|DK|}{|D|}$$

∴得证,

