

• 在概率论中，该不等式给出随机变量的和与其期望值偏差的概率上限。

伯努利随机变量的特例

- 霍夫丁不等式经常被用于一些有独立分布的伯努利随机变量的重要特例中。

- 硬币 A 面 p , B 面 $1-p$, n 次后 A 面朝上期望 np .

A 面次数不超过 k 次的概率为:

$$P(H(n) \leq k) = \sum_{i=0}^k \binom{n}{i} p^i (1-p)^{n-i}.$$

\wedge 表示 A 面朝上的次数.

对某 $\varepsilon > 0$, 当 $k = (p - \varepsilon)n$ 时, 上面不等式确定的霍夫丁上界将会按照指数级变化:

$$P(H(n) \leq (p - \varepsilon)n) \leq e \exp(-2\varepsilon^2 n).$$

相似的, 对某 $\varepsilon > 0$, 当 $k = (p + \varepsilon)n$ 时, 霍夫丁不等式的概率边界同样可以确定为:

$$P(H(n) \geq (p + \varepsilon)n) \leq e \exp(-2\varepsilon^2 n).$$

根据上两式可以得到:

$$P((p - \varepsilon)n \leq H(n) \leq (p + \varepsilon)n) \geq 1 - 2e \exp\{-2\varepsilon^2 n\}.$$

现在令 $\varepsilon = \sqrt{\ln n / n}$

那么可以得到:

$$P(|H(n) - pn| \leq \sqrt{n \ln n}) \geq 1 - 2e \exp\{-2 \ln n\} = 1 - \frac{2}{n^2},$$

普遍情况

- 现在令 X_1, X_2, \dots, X_n 为 $[0, 1]$ 的独立随机变量, 即 $0 \leq X_i \leq 1$, 我们定义变量的经验均值为:

$$\bar{x} = \frac{1}{n} (X_1 + \dots + X_n)$$

- 霍夫丁定理-中的不等式为:

$$P(\bar{x} - E[\bar{x}] \geq t) \leq e^{-2nt^2}$$

当知道 X_i 严格有界范围 a_i, b_i 时, 霍夫丁定理-中更加严格:

$$P(\bar{x} - E[\bar{x}] \geq t) \leq e^{-2nt^2 / \sum (b_i - a_i)^2}$$

有界区间,

$$P(|\bar{x} - E[\bar{x}]| \geq t) \leq 2e^{-2nt^2 / \sum (b_i - a_i)^2}$$

证明

霍夫丁引理

- 假设 X 为均值为 0 的实数随机变量并且满足,

$$P(X \in [a, b]) = 1$$

那么有如下不等式.

$$E[e^{sX}] \leq e^{tP\left\{\frac{1}{8}s^2(b-a)^2\right\}}$$

- 假如 X_1, \dots, X_n 为几个独立分布随机变量并且.

$$P(X_i \in [a_i, b_i]) = 1 \quad 1 \leq i \leq n$$

令

$$S_n = X_1 + \dots + X_n$$

对于 $s, t \geq 0$, 马尔可夫不等式以及独立性质说明:

$$\begin{aligned} P(S_n - E[S_n] \geq t) &= P(e^{s(S_n - E[S_n])} \geq e^{st}) \\ &\leq e^{-st} \cdot E[e^{s(S_n - E[S_n])}] \\ &= e^{-st} \prod E[e^{s(X_i - E[X_i])}] \\ &\leq e^{-st} \prod e^{\frac{s^2(b_i - a_i)^2}{8}} \\ &= e^{tP\left\{-st + \frac{1}{8}s^2 \sum (b_i - a_i)^2\right\}} \end{aligned}$$

为了得到最好界限上限, 将不等式右边等于为一个关于 s 的函数,

$$\begin{cases} g: \mathbb{R}_+ \rightarrow \mathbb{R} \\ g(s) = -st + \frac{s^2}{8} \sum (b_i - a_i)^2 \end{cases}$$

注意到是二次函数, 要取最小值需满足

$$s = \frac{4t}{\sum (b_i - a_i)^2}$$

这样我们有.

$$P(S_n - E[S_n] \geq t) \leq e^{tP\left\{-\frac{2t^2}{\sum (b_i - a_i)^2}\right\}}$$

证明

置信区间 • 霍夫丁不等式被用来分析样本的置信区间,

$$P(\bar{X} - E[\bar{X}] \geq t) \leq e^{-2nt^2}$$

这个不等式说明估计值大 t 的概率被指数边界控制,

$$P(-\bar{X} + E[\bar{X}] \geq t) \leq e^{-2nt^2}$$

将两式合并

$$P(|\bar{X} - E[\bar{X}]| \geq t) \leq 2e^{-2nt^2}$$

上述不等式可以理解为:

$$\alpha = P(\bar{X} \notin [E[\bar{X}] - t, E[\bar{X}] + t]) \leq 2e^{-2nt^2}$$

即真值估计范围。其中:

$$n \geq -\frac{\log(\alpha/2)}{2t^2}$$

所以我们要减少上述不等式右边式子的样本数量从而使估计区间更加靠近真值。





