

1. 随机变量的自信息

信息是对不确定性的消除。不确定性在成倍，消除后则包含的信息量就越大。概率越大则信息量就越小。相反，概率越小则信息量就越大。同样，如果事件的概率为1，那么则信无信息量。如果概率趋近于0，那么信息量就无穷大。

那么，假如有两个独立事件，设为 a_1, a_2 。

(1) 若 $P(a_1) > P(a_2)$ ，则 $f(a_1) < f(a_2)$

(2) 若 $P(a_1) = 1$ ，则 $f(a_1) = 0$

(3) 若 $P(a_1) = 0$ ，则 $f(a_1) = \infty$

(4) 若 a_1, a_2 相互独立，则 $f(a_1, a_2) = f(a_1) + f(a_2)$

通过上述数学函数值的分析，可以发现，对数函数可以满足上述要求，可以记作 $I(a_i) = -\log P(a_i) = \log \frac{1}{P(a_i)}$

我们可以从数学上分析一下，设指数为 k 。

$$\text{那么 } k \frac{1}{P(a_i)} = I(a_i)$$

可以认为是，用 k 进制的符号表来，只需要 $\frac{1}{P(a_i)}$ 位。那么进制与进制相对应。例如 $k=2$ 时，则可以用理解与 bit。

以 e 为底，单位则为 nat (nature unit)，以 10 为底，单位为 Hart (Hartley)

$$1 \text{ nat} = 1.44 \text{ bit}, \quad 1 \text{ Hart} = 3.32 \text{ bit}$$

$$= \log_2 e \text{ bit}$$

$$= \log_2 10 \text{ bit}$$

$$\log_a b = \frac{\log_c b}{\log_c a}, \quad \log_e e = \frac{\log_2 e}{\log_2 e}$$

那么，我们现在还需要证明他的唯一性。即仅有这种形式可以满足需求。

2. 唯一性的证明

对于单一变量，并不存在信息量的准确数学形式，但对于一个随机变量 X ，它的信息量 $I(X)$ 在 X 取不同值时，应该满足以下条件。

(1) 连续性条件： $f(P_1, P_2, \dots, P_n)$ 是 P_n 的连续函数。当 P_n 有微小变化时， $f(P_1, \dots, P_n)$ 也应有微小变化。

(2) 等概率时为单调增函数： $f(\frac{1}{N}, \dots, \frac{1}{N}) = g(N)$ 应是 N 的增函数。

(3) 可加性条件：当随机变量的取值不通过一次试验而通过若干次试验才取得时，随机变量在各次实验中的不确定性应该可加，且不影响最后的结果。

$$f(P_1, P_2, \dots, P_N) = f((P_1+P_2+\dots+P_k), P_{k+1}, \dots, P_N) + (P_1+P_2+\dots+P_k) f(P_1', P_2', \dots, P_k')$$

其中

$$P_i' = \frac{P_i}{P_1+P_2+\dots+P_k}, \text{ 即归一化后的 } P_i', \text{ (这个等式也是符合我们之前觉得举一个便于理解的例子)}$$

举个例子，可以将 1, 2, 3 是做一个整体，那么整个概率的信息量为 $g(\frac{1}{3})$ ，事实上等于 $f(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ 外加 "1, 2, 3" 的不确定度以及他们做为整体所发生的不确定性。这个概率可以理解为他们归一化的代价。即 P_i 中有 $\frac{1}{3} f(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ 。

下面我们开始证明唯一性。首先，假设一个 $M \times N$ 的均匀分布。

$$g(MN) = f(\frac{1}{MN}, \frac{1}{MN}, \dots, \frac{1}{MN})$$

我们利用可加性

$$g(MN) = f(\underbrace{\frac{1}{M}, \frac{1}{M}, \dots, \frac{1}{M}}_{M \uparrow}) + \sum_{i=1}^M \frac{1}{M} f(\underbrace{\frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N}}_{N \uparrow}) = g(M) + g(N)$$

$$g(s^m) = g(s \cdot s^{m-1}) = g(s) + g(s^{m-1}) \Rightarrow m g(s)$$

那么我们在继续假设取 $s, m, t \in \mathbb{N}$ ，使 $s^m \leq t^k < s^{m+1}$ 。

由上述条件 (2)，我们可以得：

$$g(s^m) \leq g(t^k) < g(s^{m+1})$$

再根据上述的推导。

$$m g(s) \leq k g(t) < (m+1) g(s)$$

同除以 $m \cdot g(s)$ 得

$$\frac{m}{n} \leq \frac{g(t)}{g(s)} < \frac{m+1}{n}$$

上述公式变形 (可根据数平均理解)

$$\left| \frac{m}{n} - \frac{g(t)}{g(s)} \right| < \frac{1}{n} \implies 0 \leq \frac{g(t)}{g(s)} - \frac{m}{n} < \frac{1}{n}$$

$$\left| \frac{m}{n} - \frac{g(t)}{g(s)} \right| < \frac{1}{n} \implies 0 \leq \frac{g(t)}{g(s)} - \frac{m}{n} < \frac{1}{n}$$

对于底数大于1的对数函数, 我们有:

$$m \log s \leq n \log t < (m+1) \log s,$$

那么同理, 我们也可得到不等式

$$\left| \frac{m}{n} - \frac{\log t}{\log s} \right| < \frac{1}{n} \implies 0 \leq \frac{\log t}{\log s} - \frac{m}{n} < \frac{1}{n}$$

将上述两式相加:

$$\left| \frac{m}{n} - \frac{g(t)}{g(s)} \right| + \left| \frac{m}{n} - \frac{\log t}{\log s} \right| < \frac{2}{n} \implies \left| \frac{\log t}{\log s} - \frac{g(t)}{g(s)} \right| < \frac{2}{n}.$$

(这里利用了 $|a+b| \leq |a| + |b|$. 由几何意义, 我们还可以得到 $< \frac{1}{n}$ 的结论)

考虑到 n 可任意大, 那么 n 为无穷大时, 上述不等式也需满足.

$$\lim_{n \rightarrow \infty} \left| \frac{g(t)}{g(s)} - \frac{\log t}{\log s} \right| = 0 \implies g(t) = C \log t.$$

上面我们考虑了均匀分布, 下面考虑 X 是非均匀分布, 但是有理数取值的情况. 即 $X = P_X(x)$, (这里可以简称为 $P(x)$)

令 $P(x) = \frac{m_x}{\sum_{i=1}^N m_i} = \frac{m_x}{M}$. (有理数可用分数表示), 那么 $g(M) = f\left(\frac{1}{M}, \dots, \frac{1}{M}\right)$, 从而

因为是有理数, 可以将 $g(M)$ 进一步分组 $g(M) = f\left(\underbrace{\frac{1}{M}, \dots, \frac{1}{M}}_{m_1 \uparrow}, \underbrace{\frac{1}{M}, \dots, \frac{1}{M}}_{m_2 \uparrow}, \dots, \underbrace{\frac{1}{M}, \dots, \frac{1}{M}}_{m_N \uparrow}\right)$

根据可加性的条件:

$$g(M) = f\left(\frac{m_1}{M}, \frac{m_2}{M}, \dots, \frac{m_N}{M}\right) + \sum_{i=1}^N \frac{m_i}{M} g(m_i)$$

$$\implies f\left(\frac{m_1}{M}, \dots, \frac{m_N}{M}\right) = g(M) - \sum_{i=1}^N \frac{m_i}{M} g(m_i)$$

$$= C \log M - \sum_{i=1}^N \frac{m_i}{M} C \log m_i$$

$$= C \left[\sum_{i=1}^N \frac{m_i}{M} (\log \frac{M}{m_i}) \right] = -C \sum_{i=1}^N P_i \log P_i.$$

我们可以发现, 在均匀分布中也符合这个结论.

最后, 我们考虑无理数的情况. 无理数的证明, 可以利用有理数逼近, 然后利用连续性即可得证.

于是, 我们已经讨论了所有情况下的数学形式, 从而证明了不确定度的唯一性. 即

$$f(P_X(x)) = - \sum_{i=1}^N P_i \log P_i$$

3. 信息熵.

不确定度的数学形式和热力学熵的数学形式是一致的, 因为不确定度也被定义为信息熵. 在推导下-D 分布中 (固体物理初步-散后-部分), 我们简单提到了熵的不确定性. 世界是从有序到无序的, 也就是熵的增加. 而我们所处理的信息也是自然界的一部分, 从无噪声信号变为有噪声, 这也符合熵增加理论. 所以用熵来定义是很合理的. 数学上我们也证明了这一点.

在热力学中有一个著名的理论即麦克斯韦佯谬. Landau 提出: "麦克斯韦佯谬在减少绝热系统熵的同时所需的信息熵不少于热力学熵的减少量". 从而解决了这个佯谬. 这个例子也说明了信息熵和热力学熵的一致性.

3.1 定义

我们这里用较为严格的数学语言定义他:

离散随机变量 X 的信息熵 $H(X)$ 定义为

$$H(X) = - \sum_{x \in X} P(x) \log P(x)$$

对于 $0 \log 0$, 认为 $x \log x = 0$. 显然并非所有相同信息熵, 而均匀分布时, 信息熵是最大.

3.2 联合熵

对于一对随机变量 X, Y , 我们定义 $H(X, Y)$. 很自然的我们可以推广到联合的随机事件. 那么就得到了联合熵

例如对于一对离散随机变量 (X, Y) 的联合熵定义为

$$H(X, Y) = - \sum_{x \in X, y \in Y} P(x, y) \log P(x, y)$$

1. 熵的公式: $H(X, Y) = -\sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x, y)$

$$H(X, Y) = -\sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x, y)$$

联合熵的大小与两个随机变量的熵之和不一定相等。如两个事件相互独立, 熵则相等。即

$$H(X, Y) = H(X) + H(Y) \quad (\text{这个在数学上的证明, 只要将 } P(x, y) = P(x) \cdot P(y) \text{ 代入即可})$$

如两个事件不独立的话, 则存在一定的相关性 (需要注意这是与相关系数没有关系, 相关系数表示的是线性相关), 那么

$$H(X, Y) < H(X) + H(Y), \text{ 因为当获取到 } X \text{ 时, 就已经获取到一部分 } Y, \text{ 相当于减小了 } Y \text{ 的信息量, 所以联合熵会偏小。}$$

如如果我们想说明等号的话 $P(x, y) = P(x) \cdot P(y|x)$ 或 $P(x, y) = P(y) \cdot P(x|y)$ 则一直成立。这时我们可以引出条件熵的概念。

3.3 条件熵

我们希望得到 $H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y)$ (这样的结果, 那么我们可以根据条件熵来描述条件熵的数学形式。

$$H(X, Y) = -\sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x, y), \quad \text{代入 } p(x, y) = p(x) \cdot p(y|x)$$

$$= -\sum_{x \in X} \sum_{y \in Y} p(x, y) [\log p(x) + \log p(y|x)]$$

$$= -\sum_{x \in X} \sum_{y \in Y} p(x) \cdot p(y|x) \log p(x) - \sum_{x \in X} \sum_{y \in Y} p(x) \cdot p(y|x) \log p(y|x)$$

$$= -\sum_{x \in X} p(x) \log p(x) \left[\sum_{y \in Y} p(y|x) \right] - \sum_{x \in X} \sum_{y \in Y} p(x) \cdot p(y|x) \log p(y|x), \quad \sum_{y \in Y} p(y|x) = 1$$

那么第一项即为 $H(X)$, 第二项即为 $H(Y|X)$

我们给出其他的精确定义:

$$\text{若 } (X, Y) \sim p(x, y), \text{ 则条件熵定义为 } H(Y|X) = -\sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(y|x) = -\sum_{x \in X} p(x) \sum_{y \in Y} p(y|x) \log p(y|x)$$

我们可以把 $-\sum_{y \in Y} p(y|x) \log p(y|x)$ 记作 $H(Y|X=x)$,

而 $H(Y|X)$ 并不是 $H(Y|X=x)$, 需要注意。

有最基本的条件熵的极限, 我们可以用结论进行推导。

$$(1) \quad H(X, Y|Z) = H(X|Z) + H(Y|X, Z)$$

简单证明一下: $H(X) = -E(\log p(x))$

$$p(x, y|z) = p(x|z) \cdot p(y|x, z)$$

$$\log p(x, y|z) = \log p(x|z) + \log p(y|x, z)$$

$$\Rightarrow H(X, Y|z) = H(X|z) + H(Y|X, z)$$

[说明] 通过对数函数, 即使概率分布的数学结论从乘法变为加法, 而信息熵则可以认为是对其加权平均。

这个加权平均则可以认为是联合熵的加权平均, 所以类似的结论都是成立的, 这样的满足可加性。联合观察必须等于逐次观察。

(2) 设两个变量 X_1, X_2, \dots, X_n 满足分布 $p(x_1, x_2, \dots, x_n)$, 则

$$H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i | X_{i-1}, X_{i-2}, \dots, X_1)$$

这也满足可加性。

在通信中, X_1, \dots, X_n 通常也是依次到达的, 我们也可以根据可加性做次序性的检测, 每次只观察最新到达的信息即可。

(3) 若 X, Y 相互独立, 那么 $p(x, y) = p(x) \cdot p(y)$, 则有 $H(X, Y) = H(X) + H(Y)$, 这已经证明过。

3.4 信息熵的性质

1) 对称性

$$H(p_1, p_2, \dots, p_n) = H(p_{k(1)}, p_{k(2)}, \dots, p_{k(n)})$$

信息熵与排列顺序无关, 并不在于信息本身, 这也是经典信息论的核心思想

2) 非负性

熵的最小值为 0, 即为概率为 1 的时候

3) 可加性

$$H(X, Y) = H(X) + H(Y|X), \text{ 即为上述的 } (1)$$

4) 条件减小熵

$$H(X|Y) \leq H(X)$$

因为条件熵总是小于等于

上述的结论对任意事件 X, Y 都成立, 所以减小熵总是成立的, 这也可以理解为 X, Y 总是有关系的

4). 条件减小熵:

$H(X|Y) \leq H(X)$,
 这符合我们直觉. 知道了统计相关性变量, 则可以减小不确定性. 绝对平均意义上减少不确定性
 但是针对某一个Y取值不一. 从哲学上来说, 事物是广泛联系的, 但是不意味着 $H(X|Y=y) \leq H(X)$
 因为对于某一个Y取值来说, 有可能使 $P(X|y)$ 比无条件熵减小从而使信度增大. 例如不可能事件.

5) 离散熵和变量 X 在等概率分布时, 熵取极大值, 也就是熵增加. 之所以叫熵增加, 用熵也一定满足这个结论的.
 我们还是要给数学上证明. 即:

$$H(p_1, p_2, \dots, p_N) \leq H\left(\frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N}\right) = \log N = \log |X| = G(N)$$

我们可以借助数学归纳法. 一凸集, 这部分内容含在数学部分. 这是新直接应用数学结论了.

证明过程如下:

我们已经推导出熵分布是凸函数, 所以只需证明 $H(x)$ 是凸函数即可.

$$\begin{aligned} & H(\lambda p_1 + (1-\lambda)p_2) - \lambda H(p_1) - (1-\lambda)H(p_2) \\ &= -\sum_{n=1}^N \left[\lambda p_{1n} + (1-\lambda)p_{2n} \right] \log \left[\lambda p_{1n} + (1-\lambda)p_{2n} \right] + \lambda \sum_{n=1}^N p_{1n} \log p_{1n} + (1-\lambda) \sum_{n=1}^N p_{2n} \log p_{2n} \\ &= \sum_{n=1}^N \lambda p_{1n} \log \frac{p_{1n}}{\lambda p_{1n} + (1-\lambda)p_{2n}} + \sum_{n=1}^N (1-\lambda)p_{2n} \log \frac{p_{2n}}{\lambda p_{1n} + (1-\lambda)p_{2n}} \end{aligned}$$

由 $\lim_{x \rightarrow 0} \frac{\log a^{1+x}}{x} = \lim_{x \rightarrow 0} \log a \frac{1+x}{x} = \log a^e = e / \ln a \Rightarrow \log x \sim \frac{1}{\ln^2} (x-1)$ 或 $\log x \leq \frac{1}{\ln^2} (x-1)$
 代入上式

$$\begin{aligned} & \geq \sum_{n=1}^N \lambda p_{1n} \cdot \frac{1}{\ln^2} \left[1 - \frac{\lambda p_{1n} + (1-\lambda)p_{2n}}{p_{1n}} \right] + \sum_{n=1}^N (1-\lambda)p_{2n} \cdot \frac{1}{\ln^2} \left[1 - \frac{\lambda p_{1n} + (1-\lambda)p_{2n}}{p_{2n}} \right] \\ & \geq \frac{1}{\ln^2} \sum_{n=1}^N \lambda p_{1n} + (1-\lambda)p_{2n} - (\lambda + 1 - \lambda) [\lambda p_{1n} + (1-\lambda)p_{2n}] \end{aligned}$$

≥ 0

于是我们证明了熵函数具有凸性. 即 $H(p)$ 是凸函数, 那么我们就可以直接利用凸函数的 Jensen 不等式

$$H(E(X)) \geq E[H(X)]$$

等于在等根时取到极大值

3.1 从离散到连续:

这件事我们在本原论与数理统计已经帮我们完成了. 现在采用黎曼积分的方法来求解. 即所谓的定积分

$$\begin{aligned} H(X) &= -\sum_k P(x_k) \Delta x \log P(x_k) \Delta x = -\sum_k P(x_k) \log P(x_k) \cdot \Delta x - \sum_k [P(x_k) \log \Delta x] \cdot \Delta x \\ &= -\int P(x) \log P(x) dx - \int P(x) \lim_{\Delta x \rightarrow 0} \log \Delta x dx = h(X) + \infty \end{aligned}$$

(这也与直觉一致. 例如又通信量无穷大. 无论用多少个比特我们都无法表示清楚它.
 香农针对这一问题提出了微分熵的概念

微分熵
$$h(x) = \int_{-\infty}^{+\infty} -P(x) \log P(x) dx$$

从这个定义还可以推广得到:

联合熵
$$h(X, Y) = -\iint P(x, y) \log P(x, y) dx dy$$

条件熵
$$h(X|Y) = -\iint P(x, y) \log P(x|y) dx dy = -\int P(y) \int P(x|y) \log P(x|y) dx dy$$

以及一系列关系:

$$h(X, Y) = h(X) + h(Y|X) = h(Y) + h(X|Y)$$

$$h(X|Y) \leq h(X), \quad h(Y|X) \leq h(Y), \quad h(X, Y) \leq h(X) + h(Y)$$

但微分熵可能为负值. 因为他不是信息熵. 例如.

$$h(X|Y) \leq h(X), \quad h(Y|X) \leq h(Y), \quad h(X, Y) \leq h(X) + h(Y)$$

但是微分熵可能为负值。因为他并不是信息度量。例如。

$$X: p(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b, \quad b-a < 1. \\ 0, & x > b, \quad x < a. \end{cases}$$

$$\text{则 } h(X) = - \int_a^b \frac{1}{b-a} \log \frac{1}{b-a} dx = \log(b-a) < 0.$$

我们可以发现连续型随机变量 $p(x)$ 并不一定大于等于 1。

还有一个重要的分布是高斯分布。

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left[-\frac{(x-m)^2}{2\sigma^2}\right].$$

那么我们就代入公式计算。

$$\begin{aligned} h(X) &= - \int_{-\infty}^{+\infty} p(x) \log \left\{ \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left[-\frac{(x-m)^2}{2\sigma^2}\right] \right\} dx \\ &= - \int_{-\infty}^{+\infty} p(x) \log \frac{1}{\sqrt{2\pi}\sigma^2} dx - \int_{-\infty}^{+\infty} p(x) \log \exp\left[-\frac{(x-m)^2}{2\sigma^2}\right] dx \\ &= \log \sqrt{2\pi}\sigma^2 + \log e \int_{-\infty}^{+\infty} p(x) \frac{(x-m)^2}{2\sigma^2} dx \\ &= \log \sqrt{2\pi}\sigma^2 + \frac{1}{2} \log e = \frac{1}{2} \log 2\pi e \sigma^2 \end{aligned}$$

这个公式的重要性会在信道容量中体现。连续型高斯分布，那么它所带来的信息量与均值无关，只与方差有关。在通信系统中，方差和功率谱密度是紧密相关的。那么一个连续型高斯分布的随机变量与他的功率有关。

由此我们还有一个定理

若给定连续型随机变量 X 的均值 m 与方差 σ^2 ，则当 X 服从高斯分布时，微分熵最大。

证明如下：

设 $p(x)$ 为均值 m ，方差 σ^2 的高斯分布的概率密度函数 (PDF)，设为随机变量 Y 。

$q(x)$ 为均值 m ，方差 σ^2 的任意分布的 PDF。设为随机变量 Z 。

我们直接做差。这实际上也体现了信息论之间的差距的思想。尽量让写的能够体现其思路。

$$H(Z) = - \int_{-\infty}^{+\infty} q(x) \log q(x) dx = H(Y) - H(Y) + H(Z)$$

$$= H(Y) - [H(Y) - H(Z)]$$

$$H(Y) - H(Z) = - \int_{-\infty}^{+\infty} q(x) \log q(x) dx - \int_{-\infty}^{+\infty} p(x) \log p(x) dx.$$

对于高斯分布，我们可以证明。

$$- \int_{-\infty}^{+\infty} q(x) \log p(x) dx = - \int_{-\infty}^{+\infty} q(x) \log \frac{1}{\sqrt{2\pi}\sigma^2} dx + \log e \int_{-\infty}^{+\infty} q(x) \frac{(x-m)^2}{2\sigma^2} dx = \frac{1}{2} \log 2\pi e \sigma^2 = H(Y)$$

于是。

$$H(Y) - H(Z) = \int_{-\infty}^{+\infty} q(x) \log \frac{p(x)}{q(x)}, \leq \int_{-\infty}^{+\infty} q(x) \left(\frac{p(x)}{q(x)} - 1 \right) dx = 0.$$

证毕。

我们令 $\int_{-\infty}^{+\infty} q(x) \log \frac{p(x)}{q(x)}$ 记为 $-D(q(x)||p(x))$ 即为相对熵。

证明.

我们作 $\int_{-\infty}^{+\infty} q(x) \log \frac{p(x)}{q(x)} dx$ 记为 $-D(q||p)$, 即为相对熵.

这也说明了 $H(x) = \log |x| - D(p||\omega)$, 对于连续变量的推广. 这个道理也与热力学第二定律有相同的地方. 值得深思.