**Def**: Let $P$ be a prob. distribution on a discrete $\Omega$. The (Shannon) entropy

$$H(P) := \sum_{\omega \in \Omega} P(\omega) \log \frac{1}{P(\omega)}$$

For a discrete R.V. $X: \Omega \to \mathcal{X}$

$$H(X) := H(P_X) = \sum_{x \in \mathcal{X}} P_X(x) \log \frac{1}{P_X(x)} = \mathbb{E}\left(\log \frac{1}{P_X(X)}\right)$$

↑ real R.V.

$\log \frac{1}{P_X(x)}$ = the suprisal of $X=x$ happens, $H(X)$: the uncertainty/Randomness of $P_X$.

**Rem** 1. Basis of log

$$\boxed{\log_2 \longleftrightarrow \text{bits}}$$

$\log_{256} \longleftrightarrow$ bytes

8 bits = 1 byte

2. We agree $0 \log \frac{1}{0} = 0$ by $\lim_{x \to 0} x \log \frac{1}{x} = 0$.

**Example** (Bernoulli): $X \in \{0,1\}$. $P(X=1) = p$ $P(X=0) = 1-p$

$$H(X) = p \log \frac{1}{p} + (1-p) \log \frac{1}{(1-p)} := h(p)$$

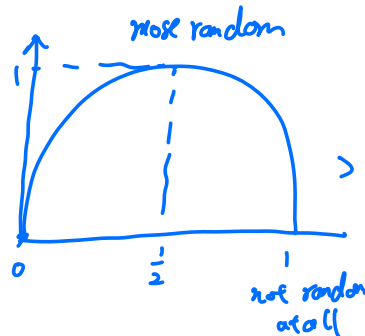where $h(\cdot)$ is called binary entropy function

$h(\frac{1}{2}) = 1$ $\qquad h(0) = h(1) = 0$

In $\log_2$ basis, $\qquad h(p) \leq 1$ and

$\qquad\qquad h(p) = 1$ iff $p = \frac{1}{2}$



most random / not random at all

**Example** (Geometric): $X \in \{0,1,2,\dots\}$ $\qquad P(X=i) = p(1-p)^i$ $i = 0,1,2 \cdots$

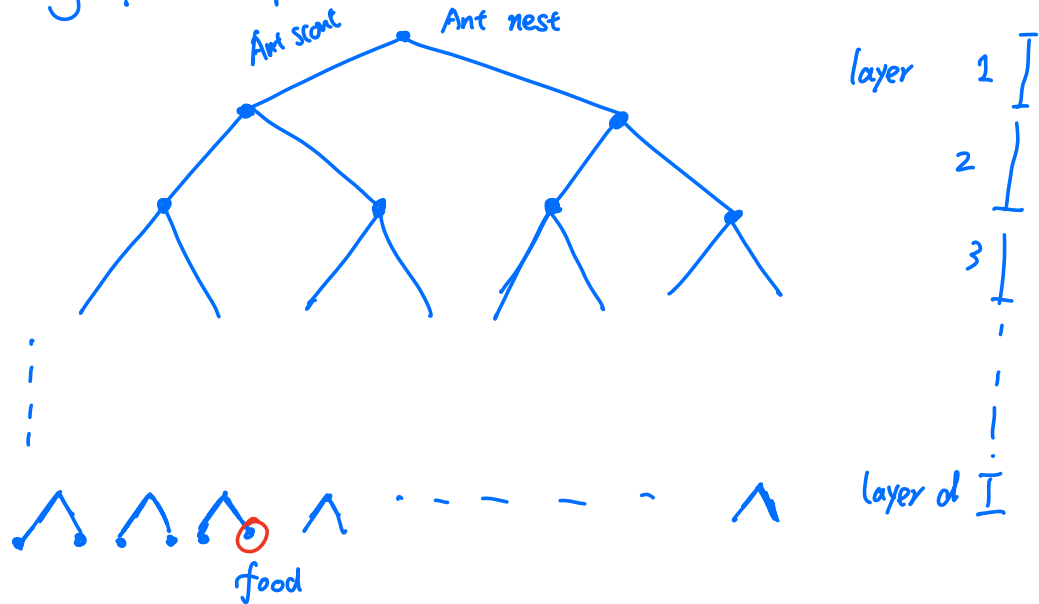$$H(X) = \sum_{i=0}^{\infty} p(1-p)^i \log \frac{1}{p(1-p)^i}$$

$$= \sum_{i=0}^{\infty} p(1-p)^i \left(\log \frac{1}{p} + i \log \frac{1}{(1-p)}\right)$$

$$= \log \frac{1}{p} \sum_{i=0}^{\infty} p(1-p)^i + p \log \frac{1}{1-p} \sum_{i=0}^{\infty} i \, p(1-p)^i$$

$$= \log \frac{1}{p} + p \log \frac{1}{(1-p)} \cdot \frac{1-p}{p^2} = \frac{h(p)}{p} \to +\infty$$

as $p \to 0$

**Example** ($\infty$ entropy): Can $H(X) = +\infty$? Yes, $P(X=k) = \frac{c}{k \ln^2 k}$, $k = 2,3 \cdots$
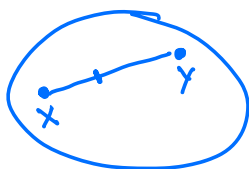
Why "log": A experiment

Ant scout · Ant nest

layer 1

2

3

layer of $\mathbb{I}$

food

Time for ant scout to describe the location of food $\sim \log_2 2^d = d$

left, right left $\cdots$ $\quad$ $d$ binary digit

$\quad$ ant communication $\quad$ $0.7 - 1$ bit/min

Convexity

$V$ a vector space $(V \cong \mathbb{R}^n)$,

A subset $S \subseteq V$ is convex if

$\quad \forall \; x, y \in S$, $\quad$ $\underline{\lambda x + (1-\lambda) y \in S}$ for $\lambda \in [0,1]$

$\qquad\qquad\qquad\qquad \downarrow$
$\qquad\qquad\qquad$ Convex combination

Convex $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ not convex

$x \quad y$

Example : ① $\mathbb{R}^n$ is convex

$$[0,1] \subseteq \mathbb{R} \qquad , \qquad (a.b) \subseteq \mathbb{R}$$

② $\mathcal{P}(X) = \{ \text{prob. distribution on } X \}$

③ $\mathcal{P}_0(\mathbb{R}) = \{ P_X \mid \mathbb{E}(X) = 0 \} \subseteq \mathcal{P}(\mathbb{R})$

$$\mathbb{E}(\lambda X + (1-\lambda) Y) = \lambda \mathbb{E}X + (1-\lambda) \mathbb{E}Y = 0$$

A function $f : S \to \mathbb{R}$ is

(i) convex if $f(\lambda x + (1-\lambda) y) \leq \lambda f(x) + (1-\lambda) f(y), \quad \forall x, y \in S, \lambda \in [0,1]$
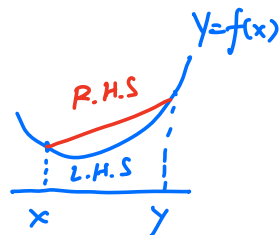
(ii) strictly convex if $f(\lambda x + 1-\lambda) < \lambda f(x) + (1-\lambda) f(y), \quad \forall x \neq y \in S$

$\qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \lambda \in (0,1)$

(iii) (strictly) concave if $-f$ is (strictly) convex

Example: ① $x \longmapsto x \log x \quad$ convex $\quad$ strictly

$\qquad \qquad x \longmapsto \log x \quad$ concave $\quad$ strictly

② $X \longmapsto \mathbb{E}X \quad$ convex but not strictly $\quad$ ( proof ? )



---

Jensen inequality : $\forall X : \Omega \to S \subseteq \mathbb{R}^n$ vector valued R.V.

$$f \text{ convex} \implies f(\mathbb{E}X) \leq \mathbb{E}f(X)$$

If $\quad$ strictly convex , $\quad$ then $\quad f(\mathbb{E}X) = \mathbb{E}f(X)$ iff

$$X = \mathbb{E}X \quad \text{a.s.}$$

$$\text{constant R.V.}$$

Pf: Convexity $\Rightarrow$ $f(\lambda_1 x_1 + \lambda_2 x_2 + \cdots + \lambda_n x_n)$
$$\leq \lambda_1 f(x_1) + \lambda_2 f(x_2) + \cdots \lambda_n f(x_n)$$
$$\lambda_i \geq 0 \qquad \sum_{i=1}^{n} \lambda_i = 1$$

$$f(\mathbb{E} X) = f\left(\sum_w P(w) X(w)\right) \leq \sum_w P(w) \, f(X(w))$$
$$\boxed{P(w) \geq 0 \quad \sum P(w) = 1} \qquad = \mathbb{E} \, f(x)$$

Properties of H

① $H(X) \geq 0$. $H(X) = 0$ iff $X$ is constant

② If $X$ is finite, $H(X) \leq \log |X|$ with equality iff $P_X$ is uniform on $X$

③ For any bijective $f$, $H(X) = H(f(X))$

④ $P \mapsto H(P)$ is strictly concave

Pf: ① $H(X) = \mathbb{E}\left[\log \frac{1}{P_X}\right] \geq 0$ $\qquad P_X(x) \leq 1, \log \frac{1}{P_X(x)} \geq 0$

② $H(X) = \mathbb{E}\left[\log \frac{1}{P_X}\right] \leq \log \mathbb{E}\left(\frac{1}{P_X}\right)$
$$= \log \sum_X P(x) \frac{1}{P(x)} = \log |X|$$

equality iff $\log \frac{1}{P_X}$ is constant

$\Leftrightarrow P_X$ constant

$$\sum_{x \in X} P_X(x) = 1 \Rightarrow P_X(x) = \frac{1}{|X|}$$

③ $P_X(x) = P(\{w \mid X(w) = x\}) = P(\{w \mid f \circ X(w) = f(x)\}) = P_{f(X)}(f(x))$

$H(X) = \sum_X P_X(x) \log \frac{1}{P_X(x)} = \sum_X P_{f(X)}(f(x)) \log \frac{1}{P_{f(X)}(f(x))} = H(X)$

④: $H(\lambda P_1 + (1-\lambda)P_2) = \sum_\omega f(\lambda P_1(\omega) + (1-\lambda)P_2(\omega))$    $f(t) = t \log \frac{1}{t}$

$$\geq \sum_\omega \lambda f(P_1(\omega)) + (1-\lambda) f(P_2(\omega)) \qquad = -t \log t$$
$$\text{concave}$$
$$= \lambda \sum_\omega f(P_1(\omega)) + 1-\lambda \sum_\omega f(P_2(\omega))$$
$$= \lambda H(P_1) + \lambda H(P_2)$$