# Information

$$I(A) = -\log \mathbb{P}[A]$$

- lower prob. events $\longrightarrow$ more info
- if $A, B$ independent, then

$$I(A \cap B) = -\log \mathbb{P}[A \cap B]$$
$$= -\log \mathbb{P}[A]\mathbb{P}[B]$$
$$= -\log \mathbb{P}[A] - \log \mathbb{P}[B]$$
$$= I(A) + I(B)$$

---

## Info for a random variable?

$X$ is a discrete r.v. w/ p.m.f. $f(x)$, supp. $S$

for any $x \in S$, we can write $I(X=x) = -\log f(x)$

Def. The entropy of disc. r.v. $X$ is given by

$$H(X) = \sum_{x \in S} -f(x) \log f(x) = \mathbb{E}[-\log f(x)]$$

$$\left(= \sum_i -P_i \log P_i\right)$$

$I(f(x))$

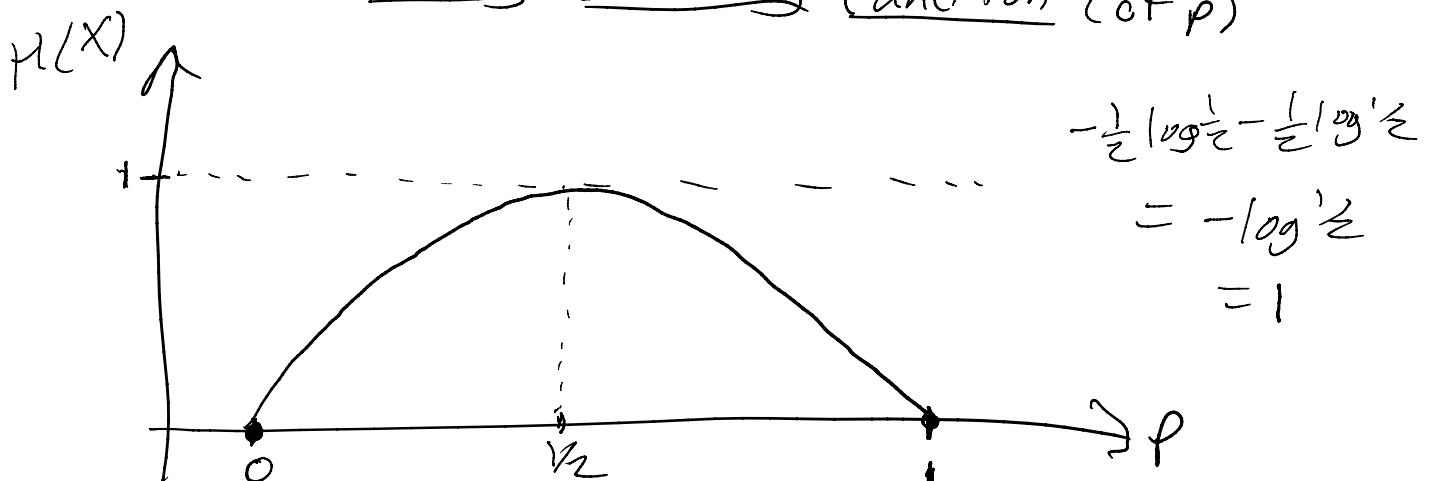-the idea: $I :$ event $\rightarrow$ real number (information)

$H : r.v. \longrightarrow$ real number (entropy)
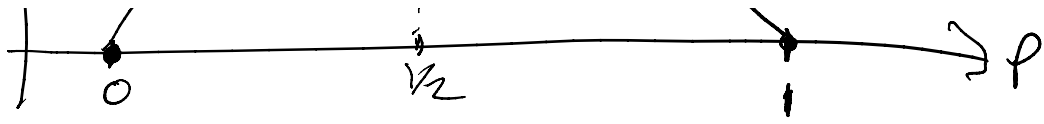
$X \longrightarrow \mathbb{E}[I(f(x))]$

note: If $f(x) = 0$, $0 \log 0 \overset{\text{in this class}}{=} 0$

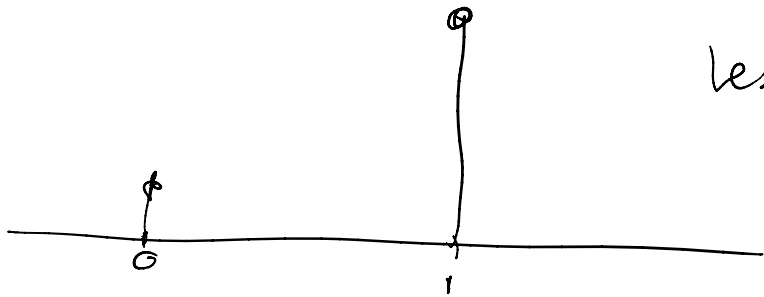Ex. What is the entropy of a Bernoulli

r.v. $X \sim Bern(p)$

$$\boxed{H(X) = -p \log p - (1-p) \log (1-p)}$$

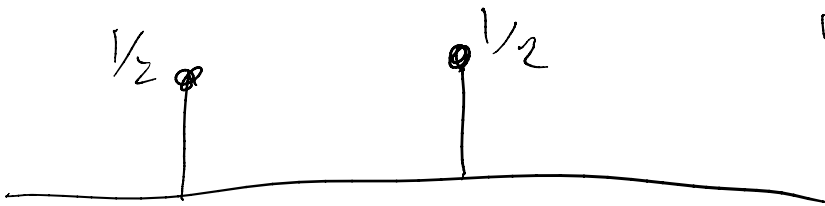Binary Entropy Function (of $p$)



$-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2}$

$= -\log \frac{1}{2}$

$= 1$

at 0, 1 $\underline{no}$ disorder (100% chance of same output)



less disorder, because you have a "better guess"

max. disorder

more entropy than Bern(½)

↑ each event has higher Info

Ex. Uniform on $1 \leq x \leq N$ "alphabet of N equiprobable events"

$$H(X) = \sum_{x=1}^{N} -\frac{1}{N} \log \frac{1}{N} = -\log \frac{1}{N} = \boxed{\log N}$$

$N \nearrow \Rightarrow H(X) \nearrow$

THM for any discrete r.v. w/ support containing $N$ pts,
$$0 \leq H(X) \leq \log N$$

Uniform is entropy-maximizing

---

Def. The joint entropy of $X$ and $Y$ is

$$H(X,Y) = -\sum_y \sum_x f(x,y) \log(f(x,y))$$

Def. The conditional entropy of $X$ given $Y$

$$H(X|Y) = -\sum_x \sum_y f(x,y) \log(f(x|y))$$

explanation at $Y=y$, $H(X|Y=y) = -\sum_x f(x|y) \log f(x|y)$

entropy of one r.v.

## more generally

$$H(X_n \mid X_1, \ldots, X_{n-1}) = \sum_{x_1, \ldots, x_n} -f(x_1, \ldots, x_n) \log f(x_n \mid x_1 \ldots x_{n-1})$$

"TPT for entropy" (directly from tpt)

$$H(X \mid Y) = \sum_i P[Y = y_i] \, H(X \mid Y = y_i)$$

## Theorem $H(X, Y) = H(Y) + H(X \mid Y)$

Proof.

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

$$= -\sum_{x, y} f(x, y) \log \left( f(y) f(x \mid y) \right)$$

$$= -\sum_{x, y} f(x, y) \log f(y) - \underbrace{\sum_{x, y} f(x, y) \log f(x \mid y)}_{H(X \mid Y)}$$

$$= -\sum_y \log f(y) \sum_x f(x, y) + H(X \mid Y)$$

$$\sim \; - \sum_y \log f(y) \sum_x f(x,y) + H(X|Y)$$

$$= -\sum_y f(y) \log f(y) + H(X|Y)$$

$$= H(Y) + H(X|Y)$$

manipulating is nice b/c of <u>logarithms</u>

<u>Generally</u>: If $X_1, \ldots, X_n$ indep

$$H(X_1, \ldots, X_n) = \sum_i H(X_i)$$

<u>Corollary</u>: $H(Y) = H(X,Y) - H(X|Y)$

<u>Def.</u> The <u>mutual information</u> between $X$ and $Y$ is

$$I(X:Y) = H(X) - H(X|Y)$$

$$I(X;Y) = H(X) - H(X|Y)$$

Cowshaw $\quad I(X;Y) = I(Y;X) = H(X) + H(Y) - H(X,Y)$

---

For cts. r.v. $X$ def. <u>differential entropy</u>

as
$$h(X) = -\int_{-\infty}^{\infty} f(x) \log f(x) \, dx$$

Ex. Entropy of $U(a,b)$

$$h(X) = -\int_a^b \frac{1}{b-a} \log \frac{1}{b-a} \, dx$$

$$= \log(b-a) \quad \leftarrow \text{ entropy-maximizing on } (a,b)$$

Ex. $X = N(0, \sigma^2)$

$$\int_{-\infty}^{\infty} \quad -x^2/2\sigma^2 \quad -x^2/2\sigma^2$$

$$h(x) = -\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-x^2/2\sigma^2} \log\left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-x^2/2\sigma^2}\right) dx$$

$$= -\underbrace{\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-x^2/2\sigma^2}}_{\text{a pdf}} \underbrace{\log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)}_{\text{x-indep}} dx + \int_{-\infty}^{\infty} \frac{x^2/2\sigma^2}{\sqrt{2\pi\sigma^2}} e^{-x^2/2\sigma^2} dx$$

$$= \frac{1}{2}\log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \overset{\sigma^2}{\overbrace{\int_{-\infty}^{\infty} (x-0)^2 f(x)\, dx}}$$

$$= \frac{1}{2}\log(2\pi\sigma^2) + \frac{1}{2}\log(e)$$

$$= \boxed{\frac{1}{2}\log(2\pi e \sigma^2)} \qquad (\text{logs here are base } e)$$

# DMS

# Discrete Memoryless Source    memoryless

A transmitter which sends $n$ independent signals from a discrete dictionary of $N$ symbols.

We call the symbols $a_1, \ldots, a_N$ and say they have probs. $P_1, \ldots, P_N$.

Over $n$ transmissions, if $n$ large

we see symbol $a_k$ approximately $n P_k$ times

So for large $n$, a sequence has Probability

$$P = \prod_{i=1}^{N} P_i^{\,n P_i} \leftarrow \text{all independent,}$$

all independent,
symbol $a_i$ occurs $n P_i$ times

$$= \prod_{i=1}^{N} 2^{\,n P_i \log P_i}$$

$$= 2^{\,n \sum P_i \log P_i}$$

$$= 2^{-nH(X)}$$

prob. of any sequence of this type is $2^{-nH(X)}$

Model: DMS output is equiprobable sequences w/ $np_i$ instances of $a_i$ each w/prob. $2^{-nH(X)}$.

$\rightarrow 2^{nH(X)}$ probable sequences (sequences of non-negligible probability)

$\rightarrow N^n$ possible sequences (ex. $\underbrace{a_1, a_1, a_1, ..., a_1}_{n \text{ times}}$)

of which only a small subset are <u>probable</u>

Real sequence space: big, not uniform in prob.
Approximation: much smaller, uniform in prob.
↖ valid w/ prob. $1 - \varepsilon$
where $\varepsilon$ can be made arbitrarily small <u>by</u> increasing $n$

To represent the output of a DMS transmitting $n$ symbols from dictionary of size $N$, need to rep $2^{nH(X)}$ seq, need

$$ \boxed{n H(X) \text{ bits}} $$

* * * * * *

$P_a = 0.75, P_b = 0.25$   $P_a = .5 = P_b$

$n=4 \begin{cases} a \ a \ a \ b \\ a \ a \ b \ a \\ a \ b \ a \ a \\ b \ a \ a \ a \end{cases}$   $\left. \begin{array}{l} a \ a \ b \ b \\ a \ b \ a \ b \\ b \ a \ a \ b \\ b \ a \ b \ a \\ b \ b \ a \ a \end{array} \right\} n = 4$

2 bits is ok   $\longrightarrow$ need 3 bits

more uniform $\longrightarrow$ more distinct probable sequences
$\longrightarrow$ higher H