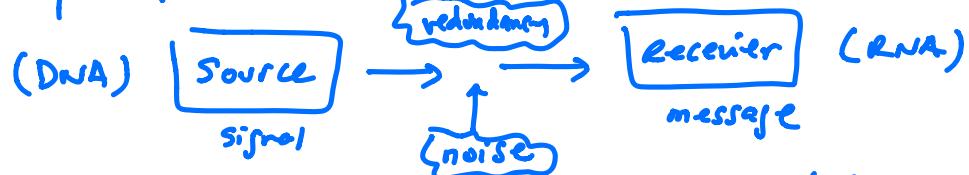


Information theory Claude Shannon (1912-2001)

"A mathematical theory of communication" (1948)

Mackay chapter 2, lectures 1, 2



how much information you need to send so  
that message is reproduced faithfully ?

need to quantify information

1. I had breakfast this morning
2. Today is not my birthday
3. Today it's my birthday
4. I run the 1<sup>st</sup> Cambridge 1/2 marathon
5. I run the 1<sup>st</sup>+2<sup>nd</sup> " "
6. I run the 1<sup>st</sup>+2<sup>nd</sup>+3<sup>rd</sup> " "
7. I have been to Antarctica

which of these messages carries more info?

→ BREAKROOMS 5'

Information content ( $I$ ) OC  
proportional  
to inverse of how frequent an event is

$p$  - probability of an event

$$I(e) \text{ OC } \frac{1}{p}$$

- the rarer an event is, the lower its prob. and the more info you would obtain by running the experiment
- the more ignorant about an outcome (many possibilities) each has lower prob., and the more info you get by listening at the end of the channel

Information content should be additive

$e_1, e_2$  independent events

$$I(e_1, e_2) = I(e_1) + I(e_2)$$

Putting both together

$$I(e) \propto -P(e)$$

$$I(e_1, e_2) = I(e_1) + I(e_2)$$

→ Shannon proposed

$$\boxed{I(e) = \log \frac{1}{P(e)}}$$

$$\text{in base 2} \quad I(e) = \log_2 \frac{1}{P} \text{ bits}$$

In shannon own words:

"useful"  
"intuitive"  
"suitable"

subjective!

$$P(e_1, e_2) = P(e_1) \cdot P(e_2) \text{ if } e_1 \perp e_2 \text{ (independent)}$$

$$\begin{aligned} I(e_1, e_2) &= \log \frac{1}{P(e_1, e_2)} = \log \frac{1}{P(e_1) P(e_2)} \\ &= \log \frac{1}{P(e_1)} + \log \frac{1}{P(e_2)} \\ &= I(e_1) + I(e_2) \end{aligned}$$

$$\log(ab) = \log(a) + \log(b)$$

$$\log\left(\frac{a}{b}\right) = \log(a) - \log(b)$$

log base

(natural)  $\log a = b \rightarrow a = e^b$  (nats)

$\log_2 a = b_2 \rightarrow a = 2^{b_2}$  (bits)

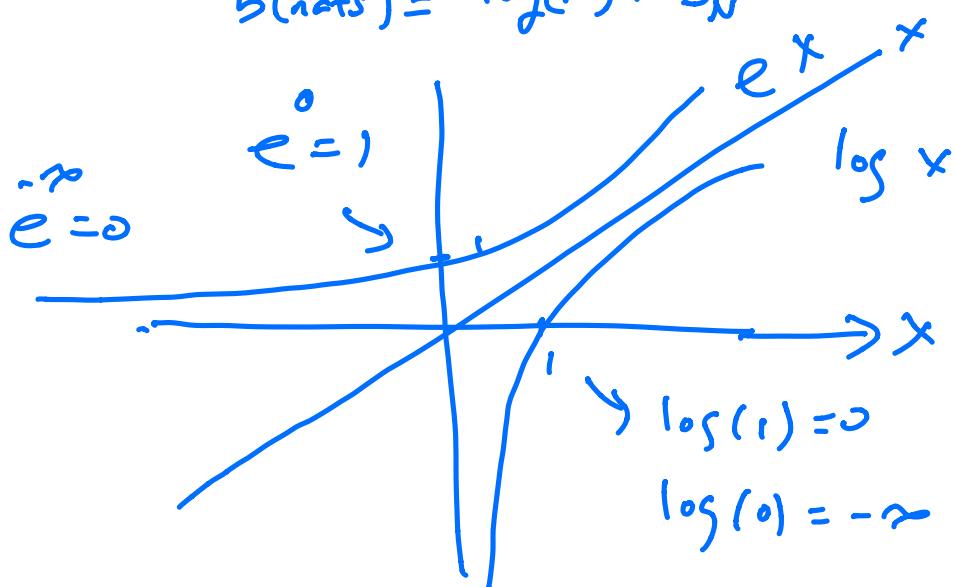
$\log_{10} a = b_{10} \rightarrow a = 10^{b_{10}}$

$\log_N a = b_N \rightarrow a = N^{b_N}$

$$e^b = 2^{b_2} \rightarrow \log(e^b) = \log(2^{b_2})$$

$$b = b_2 \log 2$$

$$b(\text{nats}) = \log(N) \cdot b_N$$



	P	I	bits
1. I had breakfast today	1	$\log 1$	0
2. Today not my b-day	$\frac{364}{365}$	$\log \frac{365}{364}$	2.5
3. " is "	$\frac{1}{365}$	$-\log 365$	8.5
4. I run 1 <sup>st</sup> Can $\frac{1}{2}$	$\frac{4500}{4.7 \cdot 10^3}$		10.0
5. I " 1 <sup>st</sup> + 2 <sup>nd</sup> Can $\frac{1}{2}$	$\frac{4500}{4.7 \cdot 10^3} \cdot \frac{6500}{4.7 \cdot 10^3}$		13.5

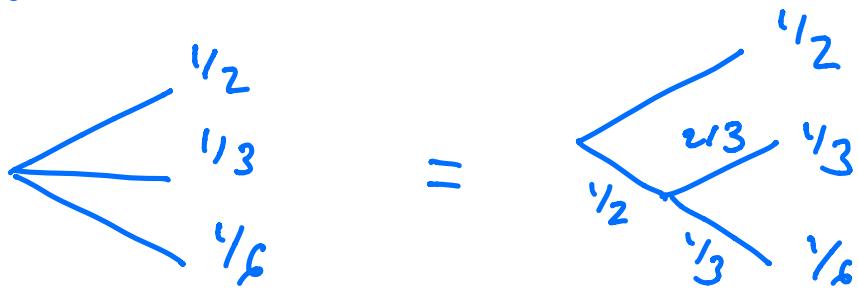
Entropy average information in a probability distribution

$X$   $P(x)$  probability density (pdf)

$$H(X) = \langle \log P(X) \rangle = \int_X \log P(X) \cdot P(X) dX$$

$$H(X) \geq 0 \quad H(X)=0 \Leftrightarrow X=0$$

A desirable property described in Shannon's paper is



$$H(p_1, p_2, p_3) = H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{1}{2} H\left(\frac{2}{3}, \frac{1}{3}\right)$$

(exercise to proof this)

## Relative entropy the Kullback-Liebler divergence

To compare 2 pdf's  $P(x)$ ,  $g(x)$

$$D_{KL}(P||g) = \int_x P(x) \log \frac{P(x)}{g(x)}$$

$$\text{i)} D_{KL}(P||g) \neq D_{KL}(g||P)$$

$$\text{ii)} D_{KL}(P||g) > 0 \quad \text{and} \quad \Rightarrow P=g$$

Mutual information  $M I$   $X, Y$   $\frac{P(X,Y)}{P(X)P(Y)} - P$

$$D_{KL}(P_{XY}, P_X P_Y) =$$

$$MI = \int_x \int_y P(x,y) \log \frac{P(x,y)}{P(x)P(y)}$$

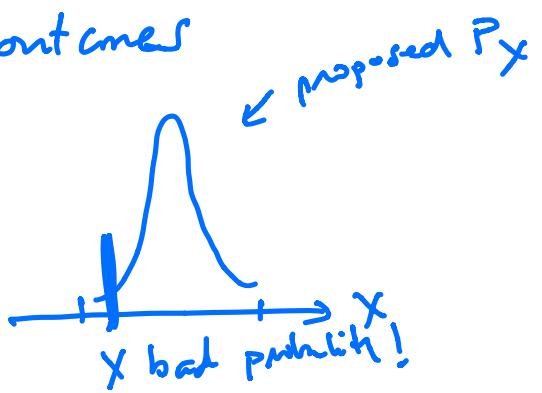
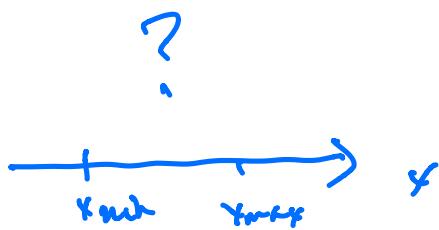
$$X \perp Y: P_{XY}(x,y) = P_X(x)P_Y(y) \Rightarrow MI=0$$

## the principle of Maximum Entropy

$X \rightarrow$  which is the  $P_x$  that has max entropy?

it is the distribution that favors the feast

any of the possible outcomes



Optimization question.

given  $X$  which is  $P_X$  such that

$H(P_X)$  is maximal

$$H(P_X) = \int_X P_X(x) \log P_X(x) dx$$

$$L = \int_X P_X(x) \log \frac{1}{P_X(x)} dx - \lambda \left[ \int_X P_X(x) dx - 1 \right]$$

$$\frac{\delta L}{\delta P_X(y)} = \log \hat{P}_X^*(y) - \frac{P_X(y)}{P_X(y)} - \lambda = 0$$

$$\log \hat{P}_X^*(y) = -\lambda - 1$$

$$\hat{P}_X^*(y) = e^{-\lambda-1} \rightarrow \text{const!}$$

$$\int_y \hat{P}_X^*(y) dy = e^{-\lambda-1} * (b-a) = 1 \quad e^{-\lambda-1} = \frac{1}{b-a}$$

$$\boxed{\hat{P}_X^*(x) = \frac{1}{b-a}}$$

for uni. fun distri.

If you know the means

$$L = - \int_y P_X(y) \log P_X(y) dy - \lambda \left[ \int_y P_X(y) dy - 1 \right] - \lambda_\mu \left[ \int_y y P_X(y) dy - \mu \right]$$

$$\frac{\delta L}{\delta P_X(x)} = -\log P_X(x) - 1 - \lambda - \lambda_\mu x = 0$$

$$\log P_X^*(x) = -(1+\lambda) - \lambda_\mu x$$

$$P_X^*(x) = e^{-(\lambda+1) - \lambda_\mu x}$$

$$\int_a^b P_X^*(x) dx = e^{-(\lambda+1)} \int_a^b e^{-\lambda_\mu x} dx = 1$$

$$P_X^*(x) = \frac{e^{-\lambda_\mu x}}{\int_y e^{-\lambda_\mu y} dy} \quad \text{an exponential distribution}$$

makes sense if  $\lambda > 0$   
 $y \in [0, +\infty)$

When you know the variance  $\sigma^2$

$$\sigma^2 = \int (x - \mu)^2 P_X(x) dx$$

$$L = - \int_y P_X(y) \log P_X(y) dy - \lambda \left[ \int_y P_X(y) dy - 1 \right] \\ - \lambda_\sigma \left[ \int_y (y - \mu)^2 P_X(y) dy - \sigma^2 \right]$$

$$\frac{\delta L}{\delta P_X(x)} = -\log P_X(x) - 1 - \lambda - \lambda_\sigma (y - \mu)^2 = 0$$

$$\log P_X^*(x) = -(1 + \lambda) - \lambda_\sigma (x - \mu)^2$$

$$P_X^*(x) = \frac{e^{-(1+\lambda)}}{e^{-\lambda_\sigma (x-\mu)^2}}$$

$$P_X^*(x) = \frac{e^{-\lambda_\sigma (x-\mu)^2}}{\int_y e^{-\lambda_\sigma (y-\mu)^2}} \rightarrow \begin{array}{l} \text{Normal} \\ \text{dist} \end{array}$$

$$\sigma^2 = \frac{1}{\lambda_\sigma}$$

makes sense if  $\lambda_\sigma > 0$   
 $x \in (-\infty, +\infty)$