Entropy & Information (Selected topics)

Entropy: a natural measure for the "richness"
of a distribution

Say $P$ is specified by $p_1, \cdots, p_M$ where $\sum p_j = 1$, $p_i \geq 0$;
$p_j$ is the probability that a symbol drawn from
$P$ is "$j$".

$H(P)$ = # of yes-no questions, on average, required to
determine which symbol is drawn.

Will show $\boxed{H(P) = -\sum p_j \log_2 p_j}$

Say, $P$ has $2^n$ symbols, each with $p_j = 2^{-n}$.

$n$ yes-no questions are necessary, and sufficient

necessary: only $2^n$ possible sequences of answers
sufficient: dichotomy strategy
        Say $n = 3$:                    ↗ yes: refine $\{0, 1, 2, 3\}$
                is it in $\{0, 1, 2, 3\}$?
                                        ↘ no: refine $\{4, 5, 6, 7\}$.

On, more compactly:
        Express $j$ as a binary number of $j$ digits —
                ask about each one.

$$-\sum_{j=1}^{2^n} 2^{-n} \log_2 (2^{-n}) = -\log_2 (2^{-n}) = n.$$

Say $P$ has $a$ symbols, $2^{n-1} < a < 2^n$.

$n-1$ yes-no questions will not suffice.
$n$ questions will suffice. so $(n-1) < H < n$.

Consider symbols in pairs. (e.g, $a=7$)
$\underbrace{36}\;\underbrace{40}\;\underbrace{31}\;5\;05\;221\;\cdots$

This is an alphabet of $a^2$ symbols.
More generally, considering symbols in $k$-tuples
gives an alphabet of $a^k$ symbols.

Determining a $k$-tuple $=$ determining $k$ 1-tuples.

Say $n_k$ is s.t. $2^{n_k-1} < a^k < 2^{n_k}$.

Then $\quad n_k - 1 < kH < n_k$.

$2^{n_k-1} < a^k < 2^{n_k} \iff n_k - 1 < \log_2 a^k < n_k$

$$\frac{n_k - 1}{k} < \log_2 a < \frac{n_k}{k}$$

$n_k = \lceil (\log_2 a) \cdot k \rceil$ . $\lceil u \rceil = $ legest integer $\geqslant u$.

$$\frac{1}{k}\left( \lceil (\log_2 (a) \cdot k) \rceil - 1 \right) < H < \frac{1}{k} \lceil (\log_2 a) \cdot k \rceil \text{, all } k.$$
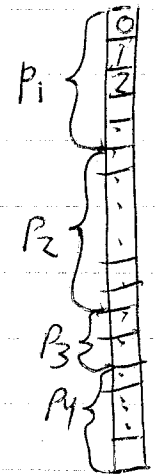
$$H = \log_2 a.$$

3.

Unequal probabilities.
[sketch]

Say each $p_j = N_j / a$.

We need, on average, $\log_2 a$ bits to determine which of $a$ "hidden" symbols is present. But, with probability $p_j$, we have $\log_2 N_j$ "excess" questions.

So
$$H = \log_2 a - \sum_j p_j \log_2 N_j$$

$$= \log_2 a - \sum_j p_j \log_2 (a p_j)$$

$$= \log_2 a - \sum_j p_j \log_2 a - \sum p_j \log_2 p_j$$

$$= - \sum_j p_j \log_2 p_j.$$

A few basic properties:

① Say $P$ & $Q$ are independent processes. Then,
$$H(P \times Q) = H(P) + H(Q).$$

P-stream $x_1 \, x_2 \, x_3 \cdots x_k \cdots$     prob $(x = x_\alpha) = p_\alpha$

Q-stream $y_1 \, y_2 \, y_3 \cdots y_k$     prob $(y = y_\beta) = q_\beta$

$R = P+Q$-stream $(x_1, y_1) \cdots \quad (x_k, y_k)$     $r_{\alpha\beta} = p_\alpha q_\beta$

prob $(x, y) = (x_\alpha, y_\beta) = r_{\alpha\beta}$

4

$$H(P+Q) = -\sum_{\alpha,\beta} r_{\alpha\beta} \log_2 r_{\alpha\beta}$$

$$= -\sum_{\alpha,\beta} p_\alpha q_\beta \log_2 p_\alpha q_\beta$$

$$= -\sum_{\alpha,\beta} p_\alpha q_\beta \left(\log_2 p_\alpha + \log_2 q_\beta\right)$$

$$= -\sum_{\alpha,\beta} p_\alpha q_\beta \log_2 p_\alpha - \sum_{\alpha,\beta} p_\alpha q_\beta \log_2 q_\beta$$

$$= -\sum_{\alpha} p_\alpha \log_2 p_\alpha - \sum_{\beta} q_\beta \log_2 q_\beta$$

$$(\sum q_\beta = 1) \qquad\qquad (\sum p_\alpha = 1)$$

$$= H(P) + H(Q).$$

② Mixing. "$R_z = (1-z)P + zQ$"  $\qquad \dfrac{d^2 H(R_z)}{dz^2} < 0$ $\forall Q \neq P$.

$P$ & $Q$ both distributions on same letters.

$$r_\alpha = (1-z)p_\alpha + z q_\alpha.$$

$$\frac{d H(R_z)}{dz} = \frac{d}{dz} \sum_\alpha \left((1-z)p_\alpha + z q_\alpha\right) \log_2 \left((1-z)p_\alpha + z q_\alpha\right)$$
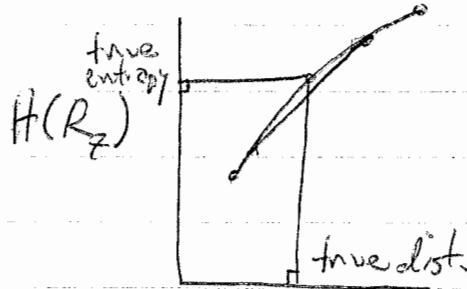
$$= -\sum_\alpha \left[ \frac{(-p_\alpha + q_\alpha)}{\ln 2} + +(-p_\alpha + q_\alpha)\log_2 \left((1-z)p_\alpha + z q_\alpha\right)\right]$$

$$= -\sum_\alpha (-p_\alpha + q_\alpha) \log_2 \left((1-z)p_\alpha + z q_\alpha\right)$$

$$\frac{d^2 H(R_z)}{dz^2} = -\sum_\alpha \frac{(-p_\alpha + q_\alpha)^2}{\ln 2} \left(\frac{1}{(1-z)p_\alpha + z q_\alpha}\right) < 0$$

5.

Consequence of mixing property for estimates of entropy



"plug-in" estimate is always downward-biased.
Amount of bias, for an estimate of $p_1, \cdots, p_K$ from $N$ samples, is

$\begin{bmatrix} \text{Miller} \\ \text{Carlton} \\ \text{Toeves} \\ \text{Panzeri} \end{bmatrix}$
$$\frac{K-1}{2 \ln 2} \cdot \frac{1}{N} + O\left(\frac{1}{N^2}\right) \quad \text{provided } K \ll N.$$

Caution — bias estimate is very wrong if $K \gtrsim N$.

③ Max Ent distribution on $(x_\alpha, y_\beta)$ subj to $\sum_\beta r_{\alpha\beta} = p_\alpha$
$\sum_\alpha r_{\alpha\beta} = q_\beta$

We know (from previous maxent analyses) that

$$r_{\alpha\beta} = C \, e^{-\lambda_\alpha - \mu_\beta}$$

so, $r_{\alpha\beta}$ must be a product, so $r_{\alpha\beta} = p_\alpha q_\beta$.

Any other dist. on $(x_\alpha, y_\beta)$ with $\sum_\beta r_{\alpha\beta} = p_\alpha$, $\sum_\alpha r_{\alpha\beta} = q_\beta$

must have lower entropy (via property ②).

So we now have a nonparametric measure of association between two variables:

Say $\quad r_{\alpha\beta} = \text{prob.}(x_\alpha, y_\beta)$ ; $p_\alpha = \sum_\beta r_{\alpha\beta}$, $q_\beta = \sum_\alpha r_{\alpha\beta}$.

Then $\quad H(P) + H(Q) - H(R) \geq 0$;

this is $0$ only for independence.

$H(R)$ can never be less than $H(P)$ or $H(Q)$

So this quantity can never be larger than $\min(H(P), H(Q))$.

This is the "mutual information" between $X$ and $Y$.

The above ideas also make sense if the sequence of symbols is __not__ independent

$$s_1 \; s_2 \; s_3 \; s_4 \cdots$$

i.e., if $p(s_1 = x_\alpha, s_2 = x_\beta) \neq p(s_1 = x_\alpha) \cdot p(s_2 = x_\beta)$.

But we can still talk about the entropy per symbol,

$$H = \lim_{k \to \infty} \frac{1}{k} \{\text{entropy of } k\text{-tuples}\}.$$

7.

Example:

Sequence of 0's & 1's, equally probable, but

$$p(0\ 0) = c$$
$$p(0, 1) = \tfrac{1}{2} - c \qquad \text{since } p(0,0) + p(0,1) = \tfrac{1}{2}$$
$$p(1, 0) = \tfrac{1}{2} - c \qquad \qquad \text{ } p(0,0) + p(1,0) = \tfrac{1}{2}$$
$$p(1, 1) = c \qquad \qquad \text{ } p(1,1) + p(0,1) = \tfrac{1}{2}$$

Recalling orig. def. of entropy ( # of yes-no questions required to specify ):

0 1 1 1 0 1 0 1 1 0 1 0 ···

can be represented by

○ △ = = △ △ △ ○ = △ △ ○ ···

but = and △ are independent, so,   $p(=) = 2c$
k - symbol entropy =   $p(△) = 1 - 2c$

$$1 + (k-1)\left[ -2c\log_2 2c - (1-2c)\log_2(1-2c) \right]$$

↓
initial
symbol

So $H = -2c\log_2 2c - (1-2c)\log_2(1-2c)$.

( maximum at $c = \tfrac{1}{4}$;   $H = 1$ ).

Extends to arbitrary-order Markov processes.

8.

Transmitted information    [Mutual Information]

P : a symbol sequence (for simplicity, independent)
    c̄ probabilities $p_1, \cdots, p_M$

Q : a symbol sequence, possibly dependent on P, but
    no other serial dependence $q_1, \cdots, q_N$.

Idea: Information that Q has about P =

$\boxed{H_1}$ # of bits, on average, required to determine a sample of P
    [without looking at Q]

$\boxed{H_2}$ — # of bits, on average, required to determine a sample of P
    [after observing  ]

Use $r_{\alpha\beta}$ to describe coupling of P & Q $\left( \begin{array}{l} p_\alpha = \sum_\beta r_{\alpha\beta} \\ q_\beta = \sum_\alpha r_{\alpha\beta} \end{array} \right)$

$$H_1 = -\sum_\alpha p_\alpha \log_2 p_\alpha.$$

$$H_2 = \sum_\beta q_\beta \left\{ -\sum_\alpha prob(\alpha|\beta) \log_2 prob(\alpha|\beta) \right\}$$

where $prob(\alpha|\beta) = r_{\alpha\beta}/q_\beta$. So

$$H_2 = -\sum_{\alpha,\beta} r_{\alpha\beta} \left( \log_2 r_{\alpha\beta} - \log_2 q_\beta \right)$$

$$= -\sum_{\alpha,\beta} r_{\alpha\beta} \log_2 r_{\alpha\beta} + \sum_\beta q_\beta \log_2 q_\beta$$

Info of Q about P = $H_1 - H_2 = -\sum_\alpha p_\alpha \log_2 p_\alpha - \sum_\beta q_\beta \log_2 q_\beta + \sum_{\alpha\beta} r_{\alpha\beta} \log_2 r_{\alpha\beta}$

$= H(P) + H(Q) - H(R)$    [mutual info of P & Q]

9.

Bias (large $N$ limit)

$$= \frac{(K_P - 1)(K_Q - 1) - (K_R - 1)}{2 \ln 2 \cdot N}$$

$$[K_R \ll N]$$

If all cells are occupiable,

$$K_R = K_P K_Q \;, \qquad \rightarrow \quad \text{bias} < 0 \quad (\text{plug-in estimate is an overestimate})$$

But not all cells need be occupiable. Bias can be $+$ or $-$.



# DATA PROCESSING THM.

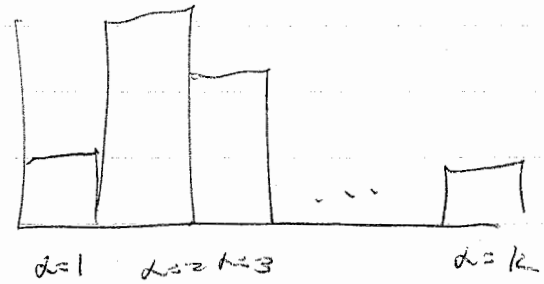$$X \longrightarrow Y \longrightarrow Z$$

$$\vdots \qquad\qquad \vdots \qquad\qquad \vdots$$

$$P_\alpha \qquad\quad q_\beta \qquad\quad s_\gamma$$

Say $r_{\alpha\beta}$ indicates relationship of $p_\alpha, q_\beta$ — indep. of $\gamma$

$t_{\beta\gamma}$ " " " $q_\beta, s_\gamma$ — indep. of $\alpha$

Information of $Z$ about $X$ $\leq$ Information of $Y$ about $X$.

Rel. of $s_\gamma, p_\alpha$ given by $u_{\alpha\gamma} = \sum_\beta r_{\alpha\beta} t_{\beta\gamma}$.

Continuum case.

Replace P.



$\alpha=1$   $\alpha=2$ $\alpha=3$      $\alpha=k$

by $p(x)$ where $\int p(x)dx = 1$.

(x may be a vector).

What is $\lim\limits_{\Delta x \to 0} -\sum\limits_{i} \left( p(x_i) \Delta x_i \right) \log_2 \left( p(x_i) \Delta x_i \right)$ ?

This is a natural notion for entropy of P, but, note that this =

$$\lim\limits_{\Delta x \to 0} -\sum\limits_{i} p(x_i) \Delta x \left( \log_2 p(x_i) + \log_2 \Delta x \right)$$

$$= -\int p(x) \log_2 p(x) dx \underbrace{\hspace{3cm}}_{\substack{\text{"Differential Entropy"} \\ \text{of } P}} - \underbrace{\log_2 \Delta x}_{\text{limit does not exist}}$$

Can still compare entropies, and, can still calculate mutual informations.

Differential entropy of a Gaussian, covariance matrix $V$

$$p(\vec{x}) = \frac{1}{(2\pi)^{d/2}} \frac{1}{\sqrt{\det V}} e^{-\vec{x}^T V^{-1} \vec{x}/2}$$

$$-\ln p(\vec{x}) = \frac{d}{2} \ln 2\pi + \frac{1}{2} \ln \det V + \frac{1}{2} x^T V^{-1} x \; .$$

Differential entropy $= \frac{1}{\ln 2} \int \left( \frac{d}{2} \ln 2\pi + \frac{1}{2} \ln \det V + \frac{1}{2} x^T V^{-1} x \right) p(\vec{x}) d\vec{x}$

$$\int x^T V^{-1} x \, p(\vec{x}) d\vec{x} = \int (y^T y) \cdot \frac{1}{(2\pi)^{d/2}} e^{-y^T y/2} d\vec{y} = d$$

$$y \text{ s.t. } y^T y = x^T V^{-1} x \; ; \quad d\vec{y} = \frac{1}{\sqrt{\det V}} d\vec{x}$$

So, diff. entropy $= \frac{1}{\ln 2} \left[ \frac{d}{2}(1 + \ln 2\pi) + \frac{1}{2} \ln \det V \right]$

## Mutual Info in the continuous setting:

No new issues. And, since $\log(\sigma_X \sigma_Y) = \log \sigma_X + \log \sigma_Y$, the annoying term drops out.

Consequence of Data Processing Inequality.

$$X \longrightarrow Y \underset{\uparrow}{\longleftrightarrow} Y' \Longrightarrow MI(X, Y) = MI(X, Y')$$
$$\text{invertible}$$

since $X \rightarrow Y \rightarrow Y' \rightarrow Y \qquad MI(XY) \leq MI(X, Y')$

and $X \rightarrow Y \rightarrow Y' \qquad \Rightarrow \qquad MI(X, Y') \leq MI(X, Y)$

Important example!

 Gaussian signal, additive Gaussian noise

$s$: possibly multivariate signal with covariance $\langle s\, s^T\rangle = V_S$

$r$: response. $r = As + x$,  $x =$ Gaussian noise;
      noise indep. of $s$; covariance $\langle x\, x^T\rangle = V_X$.

Note $\langle r\, s^T\rangle = \langle(As + x)s^T\rangle = A V_S$  $(\langle x\, s^T\rangle = 0)$

$\langle s\, r^T\rangle = \langle s(As + x)^T\rangle = V_S A^T$

$\langle r\, r^T\rangle = \langle(As+x)(As+x)^T\rangle = A V_S A^T + V_X = V_R$

So $\begin{pmatrix} s \\ r \end{pmatrix}(s^T\ r^T) = \begin{pmatrix} V_S & V_S A^T \\ A V_S & A V_S A^T + V_X \end{pmatrix} = V_{SR}$

Mutual information is

$$\frac{1}{\ln 2}\left[\frac{d_S}{2}(1 + \ln 2\pi) + \tfrac{1}{2}\ln\det V_S\right.$$

$$+\frac{d_r}{2}(1 + \ln 2\pi) + \tfrac{1}{2}\ln\det V_R$$

$$\left. -\frac{d_{SR}}{2}(1 + \ln 2\pi) + \tfrac{1}{2}\ln\det V_{SR}\right]$$

$$d_{SR} = d_S + d_R$$

$$= \frac{1}{\ln 2}\cdot\frac{1}{2}\left[\ln(\det V_S)(\det V_R)/\det V_{SR}\right]$$

13.

Elem. row op. on $V_{SR}$

$$\det \begin{pmatrix} V_S & V_S A^T \\ A V_S & A V_S A^T + V_X \end{pmatrix} = \det \begin{pmatrix} V_S & A V_S A^T \\ 0 & V_X \end{pmatrix}$$

(subtract $A(V_S \quad V_S A^T)$ from 2nd row)

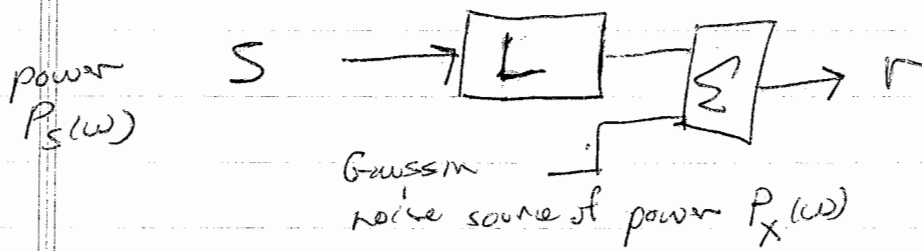So $\det V_{SR} = \det V_S \det V_X$.

Model into =

$$\frac{1}{\ln 2} \cdot \frac{1}{2} \ln \frac{\det V_R}{\det V_X} = \frac{1}{2} \log_2 \frac{\det(A V_S A^T + V_X)}{\det V_X}$$

$$= \log_2 \sqrt{\det\left(1 + \underbrace{A V_S A^T}_{\text{signal}} / \underbrace{V_X}_{\text{noise}}\right)}$$

Works in frequency domain too → each frequency
can be considered separately (if linear, stationary)

14.

Information rate for a Gaussian channel

$$\text{power } P_S(\omega) \quad S \longrightarrow \boxed{L} \longrightarrow \boxed{\Sigma} \longrightarrow r$$

Gaussian
noise source of power $P_X(\omega)$

$$\boxed{\begin{array}{c} P_R(\omega) = |L(\omega)|^2 P_S(\omega) \\ + P_X(\omega) \end{array}}$$

Over a time $T$, and sampling at $\delta t$, the relevant frequencies are $\frac{2\pi k}{T}$, $k = 1, \cdots, \frac{1}{2}\left(\frac{T}{\delta t}\right)$.

At each frequency, the response from time $0$ to $T$ has a Fourier component $\tilde{r}(\omega)$ whose real & imaginary parts are each independently Gaussian distributed, with variance $\frac{1}{2} T P_R(\omega)$

$$\text{Transmitted info} = \sum_{k=1}^{\frac{1}{2}\frac{T}{\delta t}}$$
$$2 \sum_{k=1} \log_2 \sqrt{\left( \frac{\frac{1}{2} T P_R(\omega)}{\frac{1}{2} T P_X(\omega)} \right)} \qquad \omega_k = \frac{2\pi k}{T}$$

$\uparrow$
$\cos \& \sin$

$$= \sum_{k=1}^{\frac{1}{2}\frac{T}{\delta t}} \log_2 \left( \frac{P_X(\omega_k) + |L(\omega_k)|^2 P_S(\omega_k)}{P_X(\omega_k)} \right) \qquad \Delta\omega = \frac{2\pi}{\delta t}$$

$$\longrightarrow \int_0^\infty \log_2 \left( 1 + |L(\omega)|^2 \frac{P_S(\omega)}{P_X(\omega)} \right) d\omega$$