

Entropy

In information theory, **entropy** is the average amount of information contained in each message received. Here, *message* stands for an event, sample or character drawn from a distribution or data stream. Entropy thus characterizes our uncertainty about our source of information. (Entropy is best understood as a measure of uncertainty rather than certainty as entropy is larger for more random sources). The source is also characterized by the probability distribution of the samples drawn from it. The idea here is that the less likely an event is, the more information it provides when it occurs. For some other reasons (explained below) it makes sense to define information as the negative of the logarithm of the probability distribution. The probability distribution of the events, coupled with the information amount of every event, forms a random variable whose average (a.k.a. expected) value is the average amount of information, a.k.a. entropy, generated by this distribution. The units of entropy are commonly referred to as bits, but entropy is also measured in shannons, nats, or hartleys, depending on the base of the logarithm used to define it.

The logarithm of the probability distribution is useful as a measure of information because it is additive

Entropy is a measure of *unpredictability* of *information content*

If a compression scheme is lossless—that is, you can always recover the entire original message by decompressing—then a compressed message has the same quantity of information as the original, but communicated in fewer characters. That is, it has more information, or a higher entropy, per character. This means a compressed message has less redundancy.

$$\begin{aligned} H(x) &= H(p) = - \sum p(x) \log_2 p(x) \\ &= \sum p(x) \log_2 \frac{1}{p(x)} \\ &= E \left[\log_2 \frac{1}{p(x)} \right] \end{aligned}$$

To understand the meaning of $\sum p_i \log \frac{1}{p_i}$, at first, try to define an information function, I , in terms of an event i with probability p_i . How much information is acquired due to the observation of event i ? Shannon's solution follows from the fundamental properties of information:^[8]

1. $I(p) \geq 0$ – information is a non-negative quantity
2. $I(1) = 0$ – events that always occur do not communicate information
3. $I(p_1 p_2) = I(p_1) + I(p_2)$ – information due to independent events is additive

The last is a crucial property. It states that joint probability communicates as much information as two individual events separately. Particularly, if the first event can yield one of n equiprobable outcomes and another has one of m equiprobable outcomes then there are mn possible outcomes of the joint event. This means that if $\log_2(n)$ bits are needed to encode the first value and $\log_2(m)$ to encode the second, one needs $\log_2(mn) = \log_2(m) + \log_2(n)$ to encode both. Shannon discovered that the proper choice of function to quantify information, preserving this additivity, is logarithmic, i.e.,

$$I(p) = \log(1/p)$$

The base of the logarithm can be any fixed real number greater than 1. The different units of information (bits for \log_2 , trits for \log_3 , nats for the natural logarithm \ln and so on) are just constant multiples of each other. (In contrast, the entropy would be negative if the base of the logarithm were less than 1.) For instance, in case of a fair coin toss, heads provides $\log_2(2) = 1$ bit of information, which is approximately 0.693 nats or 0.631 trits. Because of additivity, n tosses provide n bits of information, which is approximately $0.693n$ nats or $0.631n$ trits.

Now, suppose we have a distribution where event i can happen with probability p_i . Suppose we have sampled it N times and outcome i was, accordingly, seen $n_i = Np_i$ times. The total amount of information we have received is

$$\sum_i n_i I(p_i) = \sum N p_i \log(1/p_i)$$

The average amount of information that we receive with every event is therefore

$$\sum_i p_i \log \frac{1}{p_i}$$

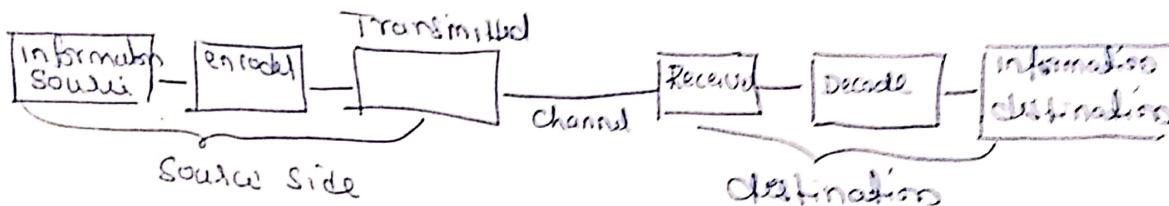
Information Theory and Coding

Information: It is the intelligence / ideas / messages in information theory.

Messages:

- 1) Electrical Signals.
- 2) Speech / Voice
- 3) Picture / Image
- 4) Video
- 5) Text

In communication system information is transmitted from source to destination



Uncertainty:

We have message $X = \{x_1, x_2, \dots, x_n\}$

corresponding probabilities $P(X) = \{P_1, P_2, \dots, P_n\}$

Total probability $P = \sum_{i=1}^n P_i$

Now Symbol x_0

$P_0 = 0$	event never happens
$P_0 = 1$	event happens 100%.

$x_0 = P_0 = 0.5 \}$ uncertainty
 $P_0 = 0.1 \}$ (NO uncertainty
in both cases)

As probability of message decreased then uncertainty increases.

Message & Information

- 1) It is information content of message
- 2) Considers an information source emitting independent messages.

$X = \{x_1, x_2, \dots, x_n\}$ with probability of occurrence

$$P = \{P_1, P_2, \dots, P_n\}$$

$$\therefore P_1 + P_2 + \dots + P_n = 1$$

The amount of information is given by

$$I_K = \log_2 \left(\frac{1}{P_K} \right) = \frac{\log \left(\frac{1}{P_K} \right)}{\log 2} \text{ bits/message}$$

Property of Information

Thm More uncertainty of message then information is

Let $P_1 = \frac{1}{4}$ having information U_1 ,

$$P_2 = \frac{1}{2} \quad - \quad - \quad - \quad U_2$$

\therefore Prove that $U_1 > U_2$

The amount of information if b_1 and b_2 are

$$I_1 = \log_2 \left(\frac{1}{P_1} \right)$$

$$= \log_2 \left(\frac{1}{\frac{1}{4}} \right)$$

$$= \log_2 2^2 = 2 \log_2 2$$

$$= 2 \text{ bits}$$

$$I_2 = \log_2 \left(\frac{1}{P_2} \right)$$

$$= \log_2 \left(\frac{1}{\frac{1}{2}} \right)$$

$$= \log_2 2$$

$$= 1 \text{ bit.}$$

$$\therefore I_1 > I_2 \Rightarrow U_1 > U_2$$

If Received Knowns message is transmitted, then there is no uncertainty in information is zero.

$$\therefore P=1 \quad \therefore I = \log_2 \left(\frac{1}{P} \right) = \log_2 (1) = 0 \text{ bit.}$$

Let I_1 is the information carried by message M_1 ,

$$I_2 - - - - - I_2 - - - - M_2$$

\therefore The combined information $I_{\text{total}} = I_1 + I_2$ (T.P.T)

$$\text{Let } I_1 = \log_2 \left(\frac{1}{P_1} \right) \quad I_2 = \log_2 \left(\frac{1}{P_2} \right)$$

\therefore Individual amount of message M_1 and M_2

new area

Ex Since messages m_1 and m_2 are independent

∴ combined probability $P = P_1 P_2$

$$\begin{aligned} \therefore I &= \log_2 \left(\frac{1}{P} \right) = \log_2 \left(\frac{1}{P_1 P_2} \right) = \log_2 \left(\frac{1}{P_1} \cdot \frac{1}{P_2} \right) \\ &= \log_2 \left(\frac{1}{P_1} \right) + \log_2 \left(\frac{1}{P_2} \right) \\ &= I_1 + I_2 \end{aligned}$$

If there are $m = 2^N$ equally likely messages, the amount of information carried by each message will be equal to N bits.

The probability of each message = $\frac{1}{m}$

$$\therefore I = \log_2 \left(\frac{1}{P} \right) = \log_2 \left(\frac{1}{1/m} \right) = \log m = \log_2 2^N = N \text{ bits}$$

Ex Calculate the amount of information if binary digits occur with equal likelihood in binary p.c.m.

Soln In a binary p.c.m. there are only two binary levels i.e. 1 or 0

Since, they occur with equal likelihood, their probabilities of occurrence will be

$$P_1 ('0' \text{ level}) = P_2 ('1' \text{ level}) = \frac{1}{2}$$

Hence amount of information carried

$$I_1 = \log_2 \left(\frac{1}{P_1} \right) \& \log_2 \left(\frac{1}{P_2} \right) = I_2$$

$$I_1 = \log_2 \left(\frac{1}{1/2} \right) \quad I_2 = \log_2 \left(\frac{1}{1/2} \right)$$

$$= \log_2 2 \quad I_2 = \log_2 2$$

$$I_1 = 1 \text{ bit} \quad I_2 = 1 \text{ bit}$$

$$\therefore I_1 = I_2 = 1 \text{ bit of information}$$

Ex: A fair coin is tossed until the first head occurs. Let X denotes the number of tosses required. Find the entropy $H(X)$ in bits.

$$\text{Let } \sum_{n=0}^{\infty} p^n = \frac{1}{1-p} \text{ and } \sum_{n=0}^{\infty} n p^n = \frac{p}{(1-p)^2}$$

Soln: Let X be a head occurs

$\therefore X$ is a geometric random variable

$$\therefore f(x) = f(n) = (1-p)^{n-1} p = q^{n-1} p.$$

$$\therefore H(X) = - \sum p(x) \log p(x) = - \sum_{n=1}^{\infty} (q^{n-1} p) \log_2 (q^{n-1} p)$$

$$= - \sum q^{n-1} p \left[(n-1) \log_2 q + \log_2 p \right]$$

$$= - \sum_{n=1}^{\infty} q^{n-1} p (n-1) \log_2 q + \sum_{n=1}^{\infty} q^{n-1} p \log_2 p$$

Put $n-1=t$ if $n=1, t=0$

$$\therefore H(X) = - \sum_{t=0}^{\infty} p q^t t \log_2 q - \sum_{t=0}^{\infty} q^t p \log_2 p$$

$$= - \cancel{p \log_2 q} \cancel{\left(\sum_{t=0}^{\infty} q^t t \right)} - p \sum_{t=0}^{\infty} q^t$$

$$= - p \log_2 q \sum_{t=0}^{\infty} q^t t - p \log_2 p \sum_{t=0}^{\infty} q^t$$

$$= - p \log_2 q \cdot \frac{q}{(1-q)^2} - p \log_2 p \cdot \frac{1}{1-q}$$

$$= - \frac{p q \log_2 q}{p} - \frac{p \log_2 p}{1-p} \quad \because 1-q=p$$

$$= - \frac{1}{p} \left[p \log_2 p + \sum q \log_2 q \right]$$

$$= - \frac{1}{p} H(p) \text{ bits}$$

<u>Rem.</u>	If $p = \frac{1}{2}$	then	$H(X) = 2$
	$p = \frac{1}{4}$	then	$H(X) = 3.2451$
	$p = \frac{1}{8}$		$H(X) = 4.38$

Ex: In a binary P.C.M if '0' occur probability $\frac{1}{4}$ and '1' occur with probability $\frac{3}{4}$, then calculate amount of information conveyed by each digits (bit)

Soln:

digit '0' has prob - $P_1 = \frac{1}{4}$

'1' has prob. $P_2 = \frac{3}{4}$

$$\therefore I_K = \log_2 \left(\frac{1}{P_K} \right)$$

$$\therefore I_1 = \log_2 \left(\frac{1}{\frac{1}{4}} \right) = \log_2 4 = \log_2 2^2 = 2 \text{ bits}$$

$$I_2 = \log_2 \left(\frac{4}{3} \right) = 0.415 \text{ bits}$$

This shows that if probability of occurrence is less, information carried is more and vice versa

Ex:

Let

$$X = \begin{cases} a & \text{with prob. } \frac{1}{2} \\ b & .. \quad \frac{1}{4} \\ c & .. \quad \frac{1}{8} \\ d & .. \quad \frac{1}{8} \end{cases}$$

Find the entropy $H(X)$.

$$\begin{aligned} \text{Soln: By definition } H(X) &= \sum P(x_i) \log_2 \frac{1}{P(x_i)} \\ &= \frac{1}{2} \log_2 \left(\frac{1}{\frac{1}{2}} \right) + \frac{1}{4} \log_2 \left(\frac{1}{\frac{1}{4}} \right) + \frac{1}{8} \log_2 \left(\frac{1}{\frac{1}{8}} \right) + \frac{1}{8} \log_2 \left(\frac{1}{\frac{1}{8}} \right) \\ &= \frac{1}{2} \log_2 2 + \frac{1}{4} \log_2 2 + \frac{1}{8} \log_2 2 + \frac{1}{8} \log_2 2 \\ &= \frac{1}{2} + \frac{2}{4} + \frac{3}{8} + \frac{3}{8} = \frac{4+4+3+3}{8} = \frac{14}{8} = \frac{7}{4} \text{ bits} \end{aligned}$$

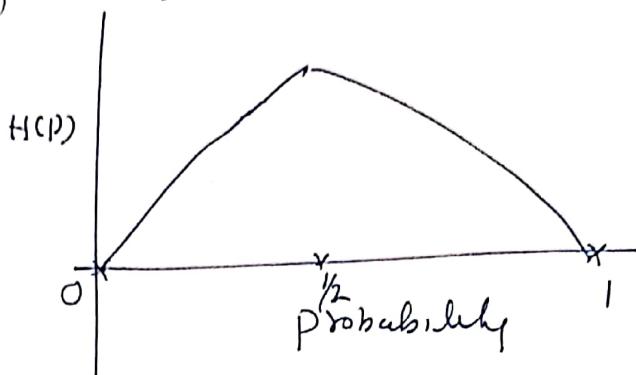
Ex: Find the entropy of the random variable

$$X = \begin{cases} 1 & \text{with prob. } \frac{1}{2} \\ 0 & \text{with prob. } \frac{1}{2} \end{cases}$$

Sol: By definition

$$H(X) = \sum p(x) \log_2 \frac{1}{p(x)} = \frac{1}{2} \log_2 2 + \frac{1}{2} \log_2 2 = 1 \text{ bit.}$$

The graph of the function $H(p)$ is as shown in figure



Basic properties of entropy

It is a concave function of the distribution and equals 0 when $p=0$ or $p=1$, i.e. the variable is not random and there is no uncertainty.

Uncertainty is maximum when $p=\frac{1}{2}$ which corresponds to the maximum value of the entropy.

Joint entropy: Let x and y are two discrete random variables with the joint probability $P(x,y)$ then the joint entropy of x, y is denoted by $H(x,y)$ and is defined as

$$H(x,y) = - \sum_{x \in X} \sum_{y \in Y} p(x,y) \log_2 p(x,y)$$

which can also be expressed as

$$H(x,y) = - E[\log_2 p(x,y)] = E\left[\log_2 \frac{1}{p(x,y)}\right]$$

Conditional entropy Let x and y are two discrete random variable with probability $P(y/x)$ or $P(x/y)$. Then the condition entropy $H(Y/X)$ is defined as

$$\begin{aligned} H(Y/X) &= - \sum_x \sum_y p(x,y) \log_2 p(y/x) \\ &= - E[\log_2 P(Y/x)] \end{aligned}$$

Thm: The entropy of a pair of a random variables is the entropy of one plus the conditional entropy of the other.

$$\therefore H(X, Y) = H(X) + H(Y|X) \text{ or } H(Y) + H(X|Y)$$

Proof: $\left[\stackrel{\text{L.L.K.T}}{P(A \cap B) = P(A)P(B|A) = P(B)P(A|B)} \right]$

Take from the definition

$$\begin{aligned} H(X, Y) &= - \sum_x \sum_y p(x, y) \log p(x, y) \\ &= - \sum_x \sum_y p(x, y) \log [p(x)p(y|x)] \\ &= - \sum_x \sum_y p(x, y) [\log p(x) + \log p(y|x)] \\ &= - \sum_x \sum_y p(x, y) \log p(x) - \sum_x \sum_y p(x, y) \log p(y|x) \\ &= - \sum_x p(x) \log p(x) - \sum_x \sum_y p(x, y) \log p(y|x) \end{aligned}$$

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

III
Hence

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

Ex? The input source to a noisy communication channel is a random variable X over the four symbols a, b, c, d. The output from this channel is a random variable Y over the same four symbols. The joint distribution of these two random variables is

$y \setminus x$	a	b	c	d
a	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{4}$
b	$\frac{1}{16}$	$\frac{1}{8}$	$\frac{1}{16}$	0
c	$\frac{1}{32}$	$\frac{1}{32}$	$\frac{1}{16}$	0
d	$\frac{1}{32}$	$\frac{1}{32}$	$\frac{1}{16}$	0

- 7
 Q) Find the Marginal distributions for X and Y
 and compute marginal entropy $H(X)$ and $H(Y)$
 in bits

The joint distribution of X and Y are

x	$a=1$	$b=2$	$c=3$	$d=4$
$P(x)$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$

y	a	b	c	d
$P(y)$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$

$$\text{Marginal entropy of } X \text{ is } H(X) = \frac{1}{4} \left[\frac{1}{4} \log_2 \frac{1}{4} \right] \\ = 4 \left[\frac{1}{2} \right] = 2 \text{ bits}$$

$$Y \in H(Y) = \frac{1}{2} \log_2 2 + \frac{1}{4} \log_2^2 \frac{1}{8} + \frac{1}{8} \log_2^3 \frac{1}{8} \\ = \frac{1}{2} + \frac{3}{4} + \frac{3}{4} = \frac{7}{4} \text{ bits}$$

- b) Find the joint entropy $H(X, Y)$

$$H(X, Y) = 1 \times \frac{1}{4} \log_2 \frac{1}{4} + 2 \left[\frac{1}{8} \log_2 \frac{1}{8} \right] + 6 \times \left[\frac{1}{16} \log_2 \frac{1}{16} \right] + 4 \left[\frac{1}{32} \log_2 \frac{1}{32} \right] \\ = \frac{1}{4} \cdot 2 + 2 \times \frac{1}{8} \times 3 + 6 \times \frac{1}{16} \times 4 + 4 \times \frac{1}{32} \times 5 \\ = \frac{1}{2} + \frac{3}{4} + \frac{3}{2} + \frac{5}{8} = \frac{4+6+12+5}{8} = \frac{27}{8} \text{ bits}$$

- c) Find $H(Y/X)$ and $H(X/Y)$

We have

$$H(X, Y) = H(X) + H(Y/X) \quad & \quad H(X, Y) = H(Y) + H(X/Y) \\ \therefore H(Y/X) = H(X, Y) - H(X) \quad & \quad H(X/Y) = H(X, Y) - H(Y) \\ = \frac{27}{8} - 2 \quad & \quad = \frac{27}{8} - \frac{7}{4} \\ = \frac{13}{8} \text{ bits} \quad & \quad = \frac{13}{8} \text{ bits}$$

Ex: Given the joint probability $P(x,y)$

$x \backslash y$	0	1	$f(x,y)$
0	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{2}{3}$
1	0	$\frac{1}{3}$	$\frac{1}{3}$
$g(u)$	$\frac{1}{3}$	$\frac{2}{3}$	1

Find $H(x)$, $H(y)$, $H(x,y)$, $H(y/x)$, $H(x/y)$

Soln: The marginal probability distributions of random variable x and y are

x	0	1
$P(x)$	$\frac{2}{3}$	$\frac{1}{3}$

y	0	1
$P(y)$	$\frac{1}{3}$	$\frac{2}{3}$

\therefore The Marginal entropy of x and y are

$$H(x) = -\sum P(x) \log_2 \frac{1}{P(x)} = -\left[\frac{2}{3} \log_2 \frac{3}{2} + \frac{1}{3} \log_2 \frac{1}{2}\right] = 0.9183$$

$$H(y) = -\sum P(y) \log_2 \frac{1}{P(y)} = -\left[\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3}\right] = 0.9183$$

Joint entropy

$$H(x,y) = -3 \times \left[\frac{1}{3} \log_2 \frac{1}{3} \right] = 1.5849$$

The conditional entropy

$$H(y/x) = H(y) + H(x,y) - H(x) = 1.5849 - 0.9183 = 0.6666$$

$$H(x/y) = H(x,y) - H(y) = 0.6666$$

Ex: Calculate the joint entropy probability $P(X \cap Y)$

X \ Y	a	b	c	$f_{(x,y)}$
x	$\frac{1}{6}$	$\frac{1}{12}$	$\frac{1}{12}$	$\frac{1}{3}$
1	$\frac{1}{12}$	$\frac{1}{6}$	$\frac{1}{12}$	$\frac{1}{3}$
2	$\frac{1}{12}$	$\frac{1}{6}$	$\frac{1}{12}$	$\frac{1}{3}$
3	$\frac{1}{12}$	$\frac{1}{6}$	$\frac{1}{12}$	$\frac{1}{3}$
$P(Y)$	$\frac{1}{3}$	$\frac{1}{4}$	$\frac{1}{4}$	1

Find $H(X)$, $H(Y)$, $H(X,Y)$, $H(X|Y)$ & $H(Y|X)$

Soln: The Marginal distributions of X and Y are

X	1	2	3
$P(x)$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

Y	a	b	c
$P(y)$	$\frac{1}{3}$	$\frac{1}{4}$	$\frac{1}{4}$

The Marginal entropy

$$H(X) = -3 \left[\frac{1}{3} \log_2 \frac{1}{3} + \frac{1}{12} \log_2 \frac{1}{12} + \frac{1}{12} \log_2 \frac{1}{12} \right] = 1.5849$$

$$H(Y) = - \frac{1}{3} \log_2 \frac{1}{3} - \frac{5}{12} \log_2 \frac{5}{12} - \frac{1}{4} \log_2 \frac{1}{4}$$

$$\text{Joint Entropy} = \frac{1}{3} \log_2 3 + \frac{5}{12} \log_2 \frac{12}{5} + \frac{1}{4} \log_2 4 = 1.5546$$

$$H(X,Y) = 3 \left[\frac{1}{6} \log_2 6 \right] + 6 \times \left[\frac{1}{12} \log_2 12 \right] = 3.08496$$

$$H(Y|X) = H(X,Y) - H(X) =$$

$$H(X|Y) = H(X,Y) - H(Y) =$$

Note: $H(X,Y) = H(X|Y) + H(Y) = H(Y|X) + H(X)$

$H(X), H(Y) \rightarrow$ uncertainty without knowing anything about Reenu

$H(Y|X), H(X|Y) \rightarrow$ uncertainty after knowing about Reenu

$H(X,Y) \rightarrow$ uncertainty towards complete Reenu

Relative Entropy and Mutual Information

The relative entropy is a measure of the distance between distributions. It arises as an expected logarithm of the likelihood ratio.

The relative entropy denoted by $D(p||q)$

Defn: The relative entropy or Kullback-Leibler distance between two probability mass functions $p(x)$ and $q(x)$ is defined as

$$\begin{aligned} D(p||q) &= \sum_x p(x) \log_2 \left[\frac{p(x)}{q(x)} \right] \\ &= E_p \log_2 \frac{p(x)}{q(x)} = - E_q \left[\log_2 \frac{q(x)}{p(x)} \right] \end{aligned}$$

Defn: Two random variables X and Y with a joint probability mass function $p(x,y)$ and marginal probability mass functions $p(x)$ and $p(y)$. The mutual information $I(X:Y)$ is the relative entropy between the joint distribution and the product distribution $p(x)p(y)$.

$$\begin{aligned} I(X:Y) &= \sum_x \sum_y p(x,y) \log_2 \frac{p(x,y)}{p(x)p(y)} \\ &= D[p(x,y)||p(x)p(y)] \\ &= E_{p(x,y)} \log_2 \frac{p(x,y)}{p(x)p(y)} \end{aligned}$$

Properties of mutual information

i) Mutual is symmetric

$$I(X:Y) = I(Y:X)$$

ii) Mutual information is always non-negative

$$I(X:Y) \geq 0$$

iii) Mutual information may be written as entropy

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

iv) Mutual information is related to joint entropy $H(X,Y)$

$$I(X:Y) = H(X) + H(Y) - H(X,Y)$$

Relationships:

The mutual information is

$$\begin{aligned}
 I(X:Y) &= \sum_{x,y} p(x,y) \log_2 \frac{p(x,y)}{p(x)p(y)} \\
 &= \sum_{x,y} p(x,y) \log_2 \frac{p(x|y)}{p(x)} \quad \left| \begin{array}{l} \therefore p(x_y) = \frac{p(x,y)}{p(y)} \\ = \frac{p(x,y)}{p(x)} \end{array} \right. \\
 &= \sum_{x,y} p(x,y) \left[\log_2 p(x|y) - \log_2 \frac{1}{p(x)} \right] \\
 &= - \sum_{x,y} p(x,y) \log_2 p(x) + \sum_{x,y} p(x,y) \log_2 p(x|y) \\
 &= - \sum_x p(x) \log_2 p(x) - \left[- \sum_{x,y} p(x,y) \log_2 p(x|y) \right]
 \end{aligned}$$

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

Thus, the mutual information $I(X:Y)$ is the reduction in the uncertainty of X due to the knowledge of Y .

By symmetry, it is also prove that

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

Thus, X says as much about Y as Y says about X .

Since $H(x,y) = H(x) + H(y/x)$, then

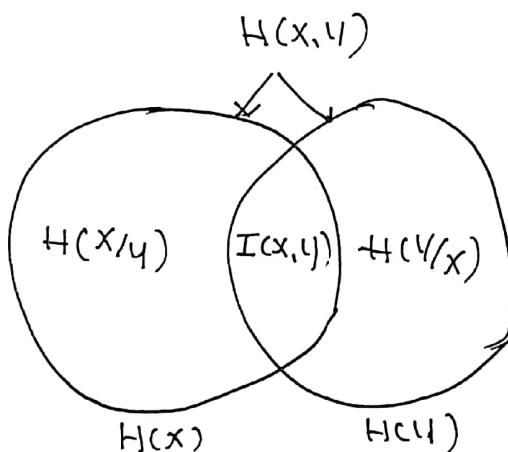
$$\begin{aligned} I(x:y) &= H(y) - H(y/x) \\ &= H(y) - [H(x,y) - H(x)] \end{aligned}$$

$$I(x:y) = H(x) + H(y) - H(x,y)$$

and

$$I(x,x) = H(x) \quad \therefore H(x/x) = 0$$

Thus, the mutual information of a random variable with itself is the entropy of the random variable. Therefore, such entropy is called self-information.



Chain rule for entropy:

The entropy of a collection of random variables is the sum of the conditional entropies

Thm: The messages x_1, x_2, \dots, x_n be drawn according

to $P(x_1, x_2, \dots, x_n)$. Then

$$H(x_1, x_2, \dots, x_n) = \sum_{i=1}^n H\left(\frac{x_i}{x_{1,2}, x_{2,3}, \dots, x_{i-1}}\right)$$

Proof: By repeated applications of the two-variable expansion rule of entropies, we have

$$H(x_1, x_2) = H(x_1) + H(x_2/x_1)$$

$$H(x_1, x_2, x_3) = H(x_1) + H(x_2, x_3/x_1) = H(x_1) + H(x_2/x_1) + H\left(\frac{x_3}{x_2, x_1}\right)$$

$$\text{so } H(x_1, x_2, \dots, x_n) = H(x_1) + H(x_2/x_1) + H\left(\frac{x_3}{x_2, x_1}\right) + \dots + H\left(\frac{x_n}{x_{n-1}, x_{n-2}, \dots, x_1}\right)$$

$$= \sum_{i=1}^n H\left(\frac{x_i}{x_{i-1}, x_{i-2}, \dots, x_1}\right)$$

E3: Let $X = \{0, 1\}$ and two distributions p and q on X . Let $p = \omega = p(0) = 1 - r$, $p(1) = r$, and $q(0) = 1 - s$, $q(1) = s$. Then

$$D(p||q) = (1-r)\log \frac{(1-r)}{(1-s)} + r \log \frac{r}{s} \quad \dots \quad (1)$$

and

$$D(q||p) = (1-s)\log \frac{(1-s)}{(1-r)} + s \log \frac{s}{r} \quad \dots \quad (2)$$

Verify that $D(p||q) = D(q||p)$.

Soln: If $r = s$, then $p = q$

$$\therefore D(p||q) = D(q||p) = 0,$$

If $r \neq s$, let $r = \frac{1}{2}$, $s = \frac{1}{4}$, then

$$\begin{aligned} D(p||q) &= \frac{1}{2} \log_2 \left(\frac{\frac{1}{2}}{\frac{3}{4}} \right) + \frac{1}{2} \log_2 \left(\frac{\frac{1}{2}}{\frac{1}{4}} \right) \\ &= \frac{1}{2} \log_2 \left(\frac{4}{6} \right) + \frac{1}{2} \log_2 \left(\frac{4}{2} \right) \\ &= \frac{1}{2} \log_2 \left(\frac{2}{3} \right) + \frac{1}{2} \log_2 2 \\ &= \frac{1}{2} + \frac{1}{2} \left[\log_2 2 - \log_2 \left| \frac{1}{2} + \frac{1}{2} \right| [1 - \log_2 3] \right] \\ &= 1 - \frac{1}{2} \log_2 3 = 0.2073 \text{ bits} \end{aligned}$$

$$D(q||p) = \frac{3}{4} \log_2 \left(\frac{\frac{3}{4}}{\frac{1}{2}} \right) + \frac{1}{4} \log_2 \left(\frac{\frac{1}{4}}{\frac{1}{2}} \right)$$

$$= \frac{3}{4} \log_2 \left(\frac{3}{2} \right) + \frac{1}{4} \log_2 \left(\frac{1}{2} \right)$$

$$= \frac{3}{4} \left[\log_2 3 - \log_2 2 \right] + \frac{1}{4} \left[\log_2 1 - \log_2 2 \right]$$

$$= -\frac{3}{4} + \frac{3}{4} \log_2 3 - \frac{1}{4} = \frac{3}{4} \log_2 3 - 1$$

$$= 0.1887 \text{ bits}$$

$$\therefore D(p||q) = 0.2073 \neq 0.1887 = D(q||p)$$

Defⁿ: The conditional mutual information of random variables x and y given z is defined as

$$I(x:y|z) = H(x|z) - H(x|y|z)$$

$$= E_{p(x,y,z)} \left[\log \frac{p(x,y|z)}{p(x|z)p(y|z)} \right]$$

which is also satisfying a chain rule:

Thmⁿ: Prove that

$$I(x_1, x_2, \dots, x_n; y) = \sum_{i=1}^n I(x_i; y | x_{i-1}, x_{i-2}, \dots, x_1)$$

Proof: By definition, we have

$$\begin{aligned} I(x_1, x_2, \dots, x_n; y) &= H(x_1, x_2, \dots, x_n) - H(x_1, x_2, \dots, x_n | y) \\ &= \sum_{i=1}^n H(x_i | x_{i-1}, x_{i-2}, \dots, x_1) - \sum_{i=1}^n H(x_i | y) \\ &= \sum_{i=1}^n I(x_i; y | x_1, x_2, \dots, x_{i-1}) \end{aligned}$$

Defⁿ: Let $p(n,y)$ and $q(n,y)$ are joint probability mass functions, the conditional relative entropy $D[p(y|x) || q(y|x)]$ is the average of the relative entropies between the conditional probability mass functions $p(y|x)$ and $q(y|x)$ averaged over the probability mass function $p(x)$

$$\begin{aligned} D[p(y|x) || q(y|x)] &= \sum_n p(x) \sum_y p(y|x) \log \frac{p(y|x)}{q(y|x)} \\ &= E_{p(x,y)} \log \frac{p(y|x)}{q(y|x)} \end{aligned}$$

Thm: Prove that $D[p(n,y) || q(n,y)] = D[p(n) || q(n)] + D[p(y|x) || q(y|x)]$

Proof: By definition, we have

$$D[p(n,y) || q(n,y)] = \sum_n \sum_y p(n,y) \log \frac{p(n,y)}{q(n,y)}$$

$$= \sum_n \sum_y p(n,y) \log \frac{p(n) p(y|x)}{q(n) q(y|x)} = \sum_n \sum_y p(n,y) \left[\log \frac{p(n)}{q(n)} + \log \frac{p(y|x)}{q(y|x)} \right]$$

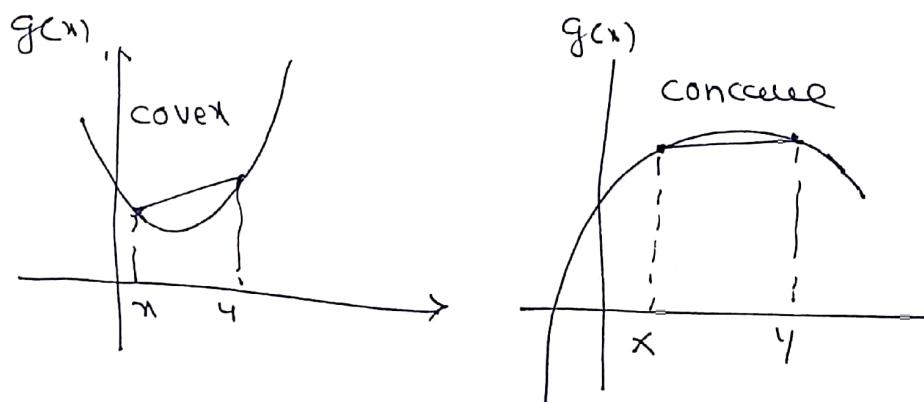
$$= \sum_n \sum_y p(n,y) \left[\log \frac{p(n)}{q(n)} + \log \frac{p(y|x)}{q(y|x)} \right]$$

$$= \sum_n p(n) \left[\log \frac{p(n)}{q(n)} + \sum_y p(y|x) \log \frac{p(y|x)}{q(y|x)} \right]$$

$$\begin{aligned} &= \sum_n \sum_y p(x,y) \log \frac{p(x)}{q(x)} + \sum_n \sum_y p(x,y) \log \frac{p(y|x)}{q(y|x)} \\ &= D[p(x) || q(x)] + D[p(y|x) || q(y|x)] \end{aligned}$$

Jensen's inequality and its consequences

A function is convex, if when you pick any two points on the graph of the function and draw a line segment between two points, the entire segment lies above the graph. If ~~does~~ the line segment always lies below the graph, the function is said to be concave. If $g(x)$ is convex if and only if $-g(x)$ is ~~convex~~ concave.



OR A function $f(x)$ is said to be convex over an interval (a,b) if for every $x_1, x_2 \in (a,b)$ and $0 \leq \lambda \leq 1$

$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$$

A function f is said to be strictly convex if equality holds if $\lambda=0$ or $\lambda=1$.

Thm: If the function f has a second derivative that is non-negative (+ve) over an interval, the function is convex over that interval.

proof: The Taylor series expansion of the function $f(x)$ around x_0 .

$$f(x) = f(x_0) + f'(x_0)(x-x_0) + \frac{f''(x_0)}{2!} (x-x_0)^2 + \dots$$

17

where x^* lies between x_0 and x_1 by hypothesis
 $f''(x^*) > 0$ and thus the last term is non-negative
 for all x .

$$\text{Let } x_0 = \lambda x_1 + (1-\lambda)x_2$$

and take $x = x_1$, then

$$f(x_1) \geq f(x_0) + f'(x_0)(1-\lambda)(x_1 - x_2) \quad \dots (1)$$

$$\text{Hence } x = x_2$$

$$f(x_2) \geq f(x_0) + f'(x_0)\lambda(x_2 - x_1) \quad \dots (2)$$

Multiplying (1) by λ and (2) by $(1-\lambda)$ and adding, we get

$$\begin{aligned} \lambda f(x_1) + (1-\lambda)f(x_2) &\geq \lambda f(x_0) + \lambda f'(x_0)(1-\lambda)(x_1 - x_2) \\ &\quad + (1-\lambda)f(x_0) + (1-\lambda)f'(x_0)\lambda(x_2 - x_1) \\ &= \cancel{\lambda f(x_0)} + f(x_0) - \cancel{\lambda f'(x_0)} + \lambda(1-\lambda)f'(x_0)(x_1 - x_2) \\ &\quad - \lambda(1-\lambda)f'(x_0)(x_1 - x_2) \\ &\geq f(x_0) \end{aligned}$$

$$\therefore \lambda f(x_1) + (1-\lambda)f(x_2) \geq f(\lambda x_1 + (1-\lambda)x_2)$$

Thm: If f is a convex function and X is a random variable

$$E[f(x)] \geq f(E(x))$$

If f is strictly convex [$f''(x) > 0$], then $X = E(X)$ with probability 1 (ie X is a constant)

proof: It can be prove by mathematical induction on the number of mass point for discrete distribution

For a two-mass-point distribution, the inequality is

$$P_1[f(x_1)] + P_2[f(x_2)] \geq f(b_1x_1 + b_2x_2) \quad \dots (1)$$

by definition of convex function by J.E

Suppose that the theorem is true for distribution with $(k-1)$ mass points. Then

$$\therefore \sum_{i=1}^{k-1} P_i f(x_i) \geq f\left[\sum_{i=1}^{k-1} P_i x_i\right] \quad \dots (2)$$

Then To prove that it is true for k mass points distribution

$$\text{Let } P_i^1 = \frac{P_i}{(1-P_k)} \quad \text{for } i = 1, 2, \dots, k-1$$

$$\therefore \sum_{i=1}^k P_i f(x_i) \geq P_k f(x_k) + (1-P_k) \sum_{i=1}^{k-1} P_i^1 f(x_i)$$

$$\geq P_k f(x_k) \\ f\left[P_k x_k + (1-P_k) \sum_{i=1}^{k-1} P_i^1 x_i\right]$$

$$\geq f\left[\sum_{i=1}^k P_i x_i\right]$$

Expected value of the function is greater than or equal to function of expected value.

Thm: Let $f(x), g(x)$, $x \in X$ be two probability mass functions, then

$$D(f||g) \geq 0$$

with equality if and only if $f(x) = g(x)$, $\forall x$.

Proof: Let $A = \{x : f(x) > 0\}$ be the support set of $f(x)$, then

$$\begin{aligned} D(f||g) &= - \sum f(x) \log \frac{f(x)}{g(x)} \\ &= \sum f(x) \log \frac{g(x)}{f(x)} \end{aligned}$$

$$\begin{aligned}
 - D(p||q) &\leq \log \sum_{x \in A} p(x) \cdot \frac{q(x)}{p(x)} \\
 &\leq \log \sum q(x) \\
 &\leq \log 1 \\
 &= 0
 \end{aligned}$$

$$\therefore D(p||q) \geq 0$$

OR From Jensen's inequality, we have

$$D(p||q) \leq \log \sum p(x) \frac{q(x)}{p(x)} \quad \dots (1)$$

Since $\log t$ is a strictly concave function of t ,
Equality holds good iff $\frac{q(x)}{p(x)}$ is constant

i.e. $q(x) = c p(x)$ everywhere

Thus $\sum_{x \in A} q(x) = c \sum_{x \in A} p(x) = c$

~~$\log \sum q(x) = \sum q(x) = 1 \Rightarrow c = 1$~~

~~Note~~ $\sum_{x \in A} q(x) = c =$

Here, we have

$$D(p||q) = 0 \text{ iff } p(x) = q(x), \forall x.$$

cor: [Non-negativity of mutual information]

For any two random variables X, Y

$$I(X, Y) \geq 0$$

if and only if X and Y are independent

Proof: By Definition

$$I(X, Y) = D(p_{(X,Y)} || p_X p_Y) \geq 0$$

Since X and Y are independent

$$p_{(X,Y)} = p(X)p(Y)$$

Cor 2: $D[p(y/x) | I(y/x)] \geq 0$ iff $p(y/x) = I(y/x)$

Cor 3: $I(x; y/z) \geq 0$ iff x and y are conditionally independent given z .

Ex Let x, y have the following joint distribution

$x \backslash y$	1	2	$f(x)$
1	0	$\frac{1}{8}$	$\frac{1}{8}$
2	$\frac{3}{4}$	$\frac{1}{8}$	$\frac{7}{8}$
$f(y)$	$\frac{3}{4}$	$\frac{1}{4}$	1

$$\begin{aligned} H(x) &= H(\frac{3}{4}, \frac{1}{4}) \quad H(\frac{1}{8}, \frac{7}{8}) = +\frac{1}{8} \log_2 8 + \frac{7}{8} \log_2 \frac{8}{7} \\ &= \frac{1}{8} \log_2 2^3 + \frac{7}{8} [\log_2 8 - \log_2 7] \\ &= \frac{3}{8} + \frac{7}{8} [3 - \log_2 7] \end{aligned}$$

$$\begin{aligned} H(y) &= \frac{3}{4} \log_2 \frac{4}{3} + \frac{1}{4} \log_2 4 = \frac{3}{4} [\log_2 4 - \log_2 3] + \frac{1}{2} \times 1 \\ &= \frac{1}{2} + \frac{3}{4} [2 - \log_2 3] = 3. \end{aligned}$$

$$\begin{aligned} H(x,y) &= 0 + \frac{1}{8} \log_2 8 + \frac{3}{4} \log_2 \frac{4}{3} + \frac{1}{8} \log_2 8 \\ &= \frac{3}{8} + \frac{3}{8} + \frac{3}{4} [\log_2 4 - \log_2 3] \\ &= \frac{3}{4} + \frac{3}{4} [2 - \log_2 3] = \frac{3}{4} [3 - \log_2 3] \end{aligned}$$

$$\begin{aligned} H(x/y) &= H(x,y) - H(y) = \frac{3}{4} [3 - \log_2 3] - \frac{3}{4} [2 - \log_2 3] - \frac{1}{2} \\ &= \frac{9}{4} - \frac{6}{4} - \frac{1}{2} = \frac{9-6-2}{4} = +\frac{1}{4} = 0.25 \end{aligned}$$

$$\begin{aligned} H(y/x) &= H(x,y) - H(x) \\ &= \frac{3}{4} [3 - \log_2 3] - \frac{3}{8} - \frac{7}{8} [3 - \log_2 7] \end{aligned}$$

Theorem: [Independent bound on entropy]

Let x_1, x_2, \dots, x_n be drawn according to $P(x_1, x_2, \dots, x_n)$, then

$$H(x_1, x_2, \dots, x_n) \leq \sum_{i=1}^n H(x_i)$$

If and only if, x_i are independent.

Proof: By Chain rule of entropy

$$\begin{aligned} H(x_1, x_2, \dots, x_n) &= \sum_{i=1}^n H\left(\frac{x_i}{x_{i-1}, x_{i-2}, \dots, x_1}\right) \\ &\leq \sum_{i=1}^n H(x_i) \end{aligned}$$

Since x_i are independent, then

$$H(x_1, x_2) = H(x_1)$$

Log Sum Inequality and its applications

Theorem: For non-negative numbers a_1, a_2, \dots, a_n and b_1, b_2, \dots, b_n

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \left(\sum_{i=1}^n a_i \right) \log \frac{\sum_{i=1}^n a_i}{\sum_{j=1}^n b_j}$$

Iff $\frac{a_i}{b_i} = \text{constant}$

Proof: Without loss of generality $a_i \geq 0$ and $b_i \geq 0$.

The function $f(t) = t \log t$ is strictly convex,

$$\text{Since } f''(t) = \frac{1}{t} \log e > 0, \forall t > 0.$$

Hence by Jensen's inequality, we have

$$\sum a_i f(t_i) \geq f(\sum a_i t_i)$$

for $a_i \geq 0$, $\sum a_i = 1$, let $x_i = \frac{b_i}{\sum_{j=1}^n b_j}$ and $t_i = \frac{a_i}{b_i}$

We obtain

$$\sum \frac{a_i}{\sum b_j} \log \frac{a_i}{b_i} \geq \sum \frac{a_i}{\sum b_j} \log \sum \frac{a_i}{\sum b_j}$$

which is the LogSum inequality.

Corollary: Using Logsum inequality, prove that $D(p||q) \geq 0$

Proof: By defn, we have

$$D(p||q) = \sum p(x) \log \frac{p(x)}{q(x)}$$

$$\geq \sum p(x) \log \frac{\sum p(x)}{\sum q(x)} = (\log 1) = 0 \quad \because \frac{h}{2} = c$$

Since both p and q are probability mass functions, $c=1$, hence, we have

$$D(p||q) = 0 \quad \text{iff} \quad p(x) = q(x), \forall x.$$

Data Processing Inequality

The data processing inequality can be used to show that no clever manipulation of the data can improve the inference that can be made from the data.

Defn: Random variables X, Y, Z are said to form a Markov chain in that order $X \rightarrow Y \rightarrow Z$, if the conditional distribution of Z depends on Y and is conditionally independent of X .

If X, Y, Z form a Markov chain $X \rightarrow Y \rightarrow Z$ if the joint probability mass function is

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

Consequences

i) $X \rightarrow Y \rightarrow Z$ iff X and Z are conditionally independent given Y . Markovity implies conditional independence because

$$p(x, z|y) = \frac{p(x, y, z)}{p(y)} = \frac{p(x|y)p(y)p(z|y)}{p(y)} = p(x|y)p(z|y)$$

This is the characterization of Markov chains that can be extended to define Markov fields, which are n -dimensional random processes in which the interior and exterior are independent given the values on the boundary.

ii) $X \rightarrow Y \rightarrow Z$ implies that $Z \rightarrow Y \rightarrow X$. Thus, the condition is sometimes written $X \leftarrow Y \leftarrow Z$.

iii) If $Z = f(Y)$, then $X \rightarrow Y \rightarrow Z$.

Sufficient Statistics

A family of probability mass functions $\{f_{\theta}(x)\}$ indexed by Θ , and let X be a sample from a distribution family. Let $T(X)$ be any statistic [μ or σ^2]. Then $\Theta \rightarrow X \rightarrow T(X)$ and by data-processing inequality, we have

$$I(\Theta; T(X)) \leq I(\Theta; X)$$

for any distributions on Θ . However, if equality holds, no information is lost.

A statistic $T(X)$ is called sufficient for Θ , if it contains all the information in X about Θ .

Defn: A function $T(x)$ is said to be a sufficient statistic relative to the family $\{f_{\theta}(x)\}$, if X is independent of Θ given $T(X)$ for any distribution on Θ [$\Theta \rightarrow T(X) \rightarrow X$ form a Markov chain]

Ex: Let $x_1, x_2, \dots, x_n, x_i \in \{0, 1\}$ be an independent and identically distributed (i.i.d.) sequence of coin tosses of a coin with unknown parameter $\Theta = P(X_i = 1)$. Given n , the number of 1's is a sufficient statistic for Θ .

$$T(x_1, x_2, \dots, x_n) = \sum_{i=1}^n x_i \quad (\text{TS.T})$$

In all sequences having that many 1's are equally likely and independent of the parameter Θ , then

$$\begin{aligned} P_T(x_1, x_2, \dots, x_n) &= P(x_1, x_2, \dots, x_n) / \sum_{\sum x_i = K} x_i \\ &= \begin{cases} \frac{1}{\binom{n}{K}} & \text{if } \sum x_i = K \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

Thus $\Theta \rightarrow \sum x_i \rightarrow (x_1, x_2, \dots, x_n)$ form a Markov chain and T is a sufficient statistic for Θ .