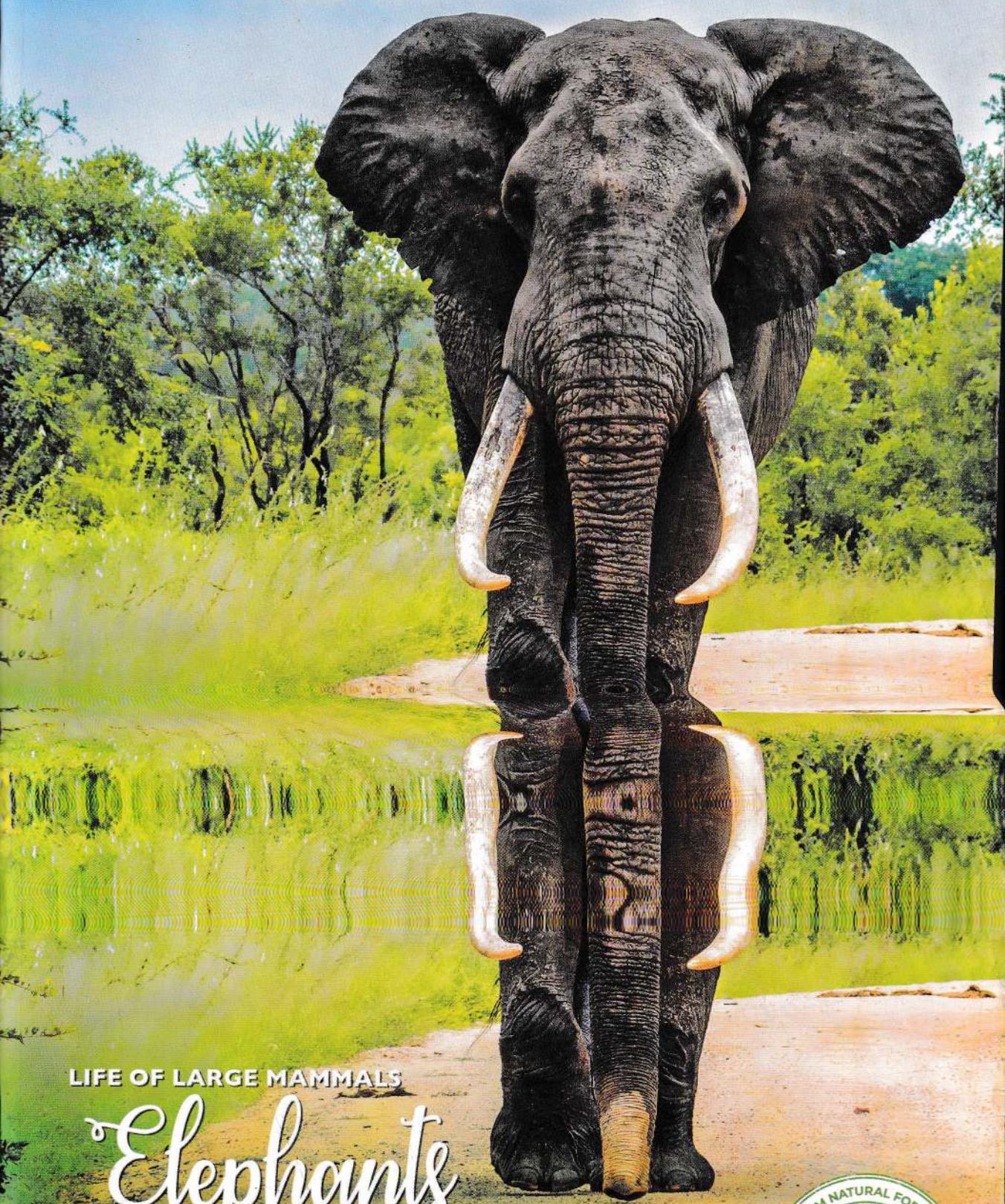




classmate



LIFE OF LARGE MAMMALS

Elephants

Customize your notebook covers at
www.classmateshop.com



The image shows five large, bold letters ('I', 'N', 'D', 'E', 'X') in a light pink color, each centered on a white square. The letters are arranged horizontally from left to right.

11-1

NAME: SOORAJ-S. STD.: _____ SEC.: _____ ROLL NO.: _____ SUB.: _____

S. No.	Date	Title	Page No.	Teacher's Sign / Remarks
		<p>QUANTUM COMPUTATION & QUANTUM INFORMATION</p> <p>— Nielsen & Chuang</p>		

11

ENTROPY & INFORMATION

Suppose we learn the value of a random variable X . The Shannon entropy of X quantifies how much information we gain, on average, when we learn the value of X .

Alternatively, the entropy of X measures the amount of uncertainty about X before we learn its value.

Conceptually,

information can be thought of as being stored in or transmitted, as variables that can take on different values.

A variable can be thought of as a unit of storage that can take on, at different times, one of several different specified values, following some process for taking on those values.

Informally, we get information from a variable by looking at its value, and the information is about the process behind the variable.

The entropy of a variable is the "amount of information" contained in the variable.

This amount is determined not just by the # of different values the variable can take on, just as the information in an email is quantified not just by the # of words in the email or the different possible words in the language of the email.

Informally, the amount of information in an email is proportional to the amount of "surprise" its reading causes.

Ex:-

If an email is simply a repeat of an earlier email, then it is not informative at all. On the other hand, if say the email reveals the outcome of a cliff-hanger election, then it is highly informative.

Similarly, the information in a variable is tied to the amount of surprise that value of the variable causes when revealed.

Shannon entropy quantifies the amount of information in a variable

Storage and transmission of information can intuitively be expected to be tied to the amount of information involved.

Ex:-

Information may be about the outcome of a coin toss. This information can be stored in a Boolean variable that can take on the values 0 or 1. We can use the variable to represent

the raw data corresponding to the coin toss
- whether the coin toss came up heads or not.

This Boolean variable can be represented in a single bit. This bit directly stores the value of the variable, i.e., the raw data corresponding to the outcome of the coin toss. It does not succinctly capture the information in the coin toss, e.g., whether the coin is biased or unbiased and, if biased, how biased.

Whereas, Shannon entropy quantifies, among other things, the absolute minimum amount of storage and transmission needed for succinctly capturing any information (as opposed to the raw data), and in typical cases that amount is less than what is required to store or transmit the raw data behind the information. Shannon entropy also suggests a way of representing the information in the calculated fewer # of bits.

the raw data corresponding to the coin toss
- whether the coin toss came up heads or not.

This Boolean variable can be represented in a single bit. This bit directly stores the value of the variable, i.e., the raw data corresponding to the outcome of the coin toss. It does not sufficiently capture the information in the coin toss, e.g., whether the coin is biased or unbiased and, if biased, how biased.

Whereas, Shannon entropy quantifies, among other things, the absolute minimum amount of storage and transmission needed for succinctly capturing any information (as opposed to the raw data), and in typical cases that amount is less than what is required to store or transmit the raw data behind the information. Shannon entropy also suggests a way of representing the information in the calculated fewer # of bits.

At a conceptual level, Shannon entropy is simply the amount of information in a variable. More mundanely, that translates to the amount of storage (e.g. # of bits) required to store the variable, which can intuitively be understood to correspond to the amount of information in that variable.

The calculation of the shannon entropy, and therefore the amount of information in a variable, is not simply the # of bits required to represent all the different values a variable might take on, which is just the raw data.

Ex:-

a variable may take on any of 4 different values. In digital storage, 2 bits would be sufficient to uniquely represent the 4 different values, and thus the variable can be stored in 2 bits. However, this is an upper limit on the required storage; it is the amount of storage required to store the raw data of the variable, not the information in that data.

Less storage might be sufficient to store the information, depending on the process by which the variable takes on different values.

Shannon entropy helps identify that amount of storage needed for the information.

→ Entropy can be looked as a measure of "compressibility" of the data.

i.e.,

a compression metric: how much can the raw data of a variable be compressed without losing the information in the variable?

Concept of Amount of information

Intuitively, one way to understand the concept of "amount of information" in a variable is to tie it to how difficult or easy it is to guess that information without having to look at the variable — the easier it is to guess the value of the variable, the less surprise is the variable and so the less information the variable has.

Another way to view information, is to contrast it with the amount of data. For example, two different Boolean variables could be stored in 1 bit each, but the amount of information in the two may be quite different:

Ex:-

- ① Suppose a coin is completely biased and always comes up heads when tossed.

The random variable representing the coin toss's outcome has probability 1 of coming up heads (it is constant), and thus there is no need to store or transmit that variable as it can be trivially guessed at any time.
∴ The amount of information in that variable is zero.

- ② If we had a perfect coin with 50:50 chances of coming up heads or tails upon a coin toss, then we can guess the outcome of a toss with only 50% accuracy

$$(\text{Probability} = \frac{1}{2})$$

∴ it is necessary to store/transmit the actual value of that coin toss outcome's random variable in order to know its value with better than 50% accuracy.

- ∴ The amount of information in this 2nd random variable is much higher than in the 1st case.

Contrast the fact that, for both of the above coins, the raw data regarding their toss outcomes need 1 bit each to be stored.

③ Suppose we had a perfect die with 6 possible outcomes on the toss (roll) of the die. The amount of information in the corresponding random variable is even higher, as it is even harder to guess the outcome of the die roll; we have only $\frac{1}{6}$ chance of guessing the outcome correctly.

Quantifying the Amount of Information

One way to represent the "amount of information" is the # of Bits it takes to represent/express the variable.

If a variable can take only 2 values, it can be represented with just 1 bit ; if it can take on any of 4 values, 2 bits are needed ; if 8 values then 3 bits are needed. etc.
But, this is the storage required for the raw data, not for the information content of the data.

If a variable is easier to guess, then we can leverage that fact to reduce the # of bits needed to store/transmit that information.

Ex:-

If a die is 80% likely to come up with the # "3", then

store/transmit only single bit with the value 0 whenever the die actually comes up with "3", and store/transmit more bits, starting with a 1st bit of value 1, when the die comes up with other values.

When the 1st bit is a 0, the receiver knows not to look for more bits associated with this particular storage/transmission.

i.e.,

While in a naive representation we would need 3 bits to represent each outcome of a die having 6 faces, but by leveraging probabilities we use only 1 bit whenever the die comes up with "3" - a more frequently occurring event, thus reducing the average # of bits needed (over multiple die throws) to store the die throw outcome.

One can leverage the probabilities of the different values to reduce the # of bits needed iff the variable has a non-uniform distribution. If not, there will be no reduction in the # of bits needed.

⇒ If a variable is more likely to take on one value than another, then it is easier to guess the value of the variable without looking at it, thus the variable has less information in it, and thus it takes fewer bits to store/transmit the value.

Example of Information Quantity

Suppose a variable can take on 3 different values a, b, c , but half the time it takes on the value 'a', and quarter of the time each the values 'b' and 'c'.

$$P(a) = 0.5 = \frac{1}{2}$$

$$P(b) = 0.25 = \frac{1}{4}$$

$$P(c) = 0.25 = \frac{1}{4}$$

We can represent 'a' with just 1 bit having the value "0"; 'b' with 2 bits "10", and "c" with 2 bits "11". This is Huffman coding.

When the 1st bit is a "0", the receiver knows to stop reading that "word" right there; when the 1st bit is a "1", the reader knows to also read the next bit to complete the "word".

With this, what's the # of bits needed to represent this variable?

We need just 1 bit half the time (when the value is 'a'), and 2 bits each the other 2 times (when the value is either 'b' or 'c')

∴ The average # of bits needed is,

$$\begin{aligned} &= P(a) \times 1 + P(b) \times 2 + P(c) \times 2 \\ &= \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{4} \times 2 = \frac{3}{2} = \underline{\underline{1.5}} \end{aligned}$$

⇒ The entropy of the above variable having these specified probabilities of taking on different values is 1.5

The Entropy Formula

Rationale Behind using Logarithms for Entropy

Defining information

If we encode an outcome of higher probability in such a way that a smaller codeword (ie., a smaller piece of information) can represent it, we are better off.

Can we say that it is the length of the codeword that is its information?

But in defining the length of the codeword that represents the outcome as its information we have not made use of the probability of the outcome at all which was seen to have a direct reciprocal effect on information that we needed to send.

Let's denote information with i and probability with p .

What if we say $i = \lceil p \rceil$?

There are at least 2 problems now!

1st: if the probability of an outcome = 1, we don't need to send any information, and so i should be 0, which is not happening with $i = \lceil p \rceil$.

Let's say we use $i = (\lceil p \rceil) - 1$

Say, the probability of Norway, Canada, Russia & Germany winning are all 0.25 each. We can encode all 4 of these outcomes using only 2 bits.

If we go by $i = \lceil p \rceil$ or $i = \lceil p \rceil - 1$, then i_{Norway} itself is 4 or 3.

Let's denote information with i and probability with p .

What if we say $i = \frac{1}{p}$?

There are at least 2 problems now!

1st: if the probability of an outcome = 1, we don't need to send any information, and so i should be 0, which is not happening with $i = \frac{1}{p}$.

Let's say we use $i = \left(\frac{1}{p}\right) - 1$

Say, the probability of Norway, Canada, Russia & Germany winning are all 0.25 each. We can encode all 4 of these outcomes using only 2 bits. If we go by $i = \frac{1}{p}$ (or) $i = \frac{1}{p} - 1$, then i_{Norway} itself is 4 or 3.

At this point of time, we do have a relationship that uses the probability of the outcome, but the "information" we defined this way seem to have no connection with the length of the codeword. It depends only on probability now

- We are looking for a formula for i that should be related to the length of the codeword that represents the outcome, and that length be somehow related to the probability of the outcome.

For finding the length of the codeword in order to describe an outcome is an event, we need to know the total # of outcomes. For 2 outcomes, we need 1 bit. For 4, we need 2 bits, and so on.

How can we find the total # of outcomes using the probability of each outcome?

Say, we have an event that can have n possible outcomes with each outcome denoted as O_i having a probability p_i . With $i = \log_2 n$ bits we'll be able to encode all the possible outcomes of this event, where $n = 2^i$.

What can be the upper limit of n if we know all the p_i 's?

- The smallest p solely limits the upper bound of n

Proof

Let p_j is the smallest. Then,

$p_i = p_j + a_i$ for every i , where $a_i = 0$ or a small positive number.

$$\begin{aligned} 1 &= p_1 + p_2 + \dots + p_n \\ &= (p_j + a_1) + (p_j + a_2) + \dots + (p_j + a_n) \\ &= n p_j + (a_1 + a_2 + \dots + a_n) \end{aligned}$$

$$1 \geq n p_j$$

$$\Rightarrow n \leq \frac{1}{p_j}$$

The total # of outcomes of such an event cannot be more than γ_p .

(γ_p is the smallest probability outcome having a probability p)

$\Rightarrow i = \log_2 \gamma_p$ bits can be used to encode all the outcomes of such an event, including γ_p .

$$n \leq \frac{1}{p} \Rightarrow \log_2 n \leq \log_2 \left(\frac{1}{p}\right) = i$$

We'll be able to encode γ_p and every other outcome in such an event using $\log_2(\gamma_p)$ bits.

So, can we define information for an outcome that has a probability p , as $\log_2(\gamma_p)$, in an event where this p is the smallest?

For an outcome with a probability p ,
 $i = \log_2(\gamma_p)$ bits are sufficient to encode it in an event in which this p is the smallest (among the probabilities of all other outcomes).

What if this p is not the smallest.

Do these many bits still suffice for encoding this outcome? ($\log(\gamma_p)$ bits).

For the # of outcomes n_p in an event in which p is the lowest probability outcome

$$n_p \leq \gamma_p \quad \text{where we chose } i = \log_2(\gamma_p)$$

$$\Rightarrow n_p \leq 2^i$$

$$\Rightarrow n_p = 2^i - \epsilon, \text{ where } \epsilon \text{ is zero or a +ve integer.}$$

Here, ϵ number of encodings are free
ie, not being used for representing any outcome.

Let's say we have added more outcomes to the event E with p' being the new lowest probability outcome.

⇒ The max. number of outcomes now becomes $\gamma_{p'}$

Let, i' : new # of bits needed to code any of these $\gamma_{p'}$ outcomes

Now, the max. # of additional outcomes will be given by,

$$\Delta n = n_{p'} - n_p$$

$$n_p \leq \gamma_{p'}$$

$$\leq \frac{1}{p'} \cdot \gamma_{p'}$$

$$= \frac{1}{p'} - \left(\frac{i}{2} - \epsilon \right), \text{ where } i = \log_2(\gamma_{p'})$$

$$= \frac{1}{p'} - \frac{i}{2} + \epsilon$$

$$\Delta n \leq \frac{i}{2} - \frac{i}{2} + \epsilon, \text{ where } i = \log_2(\gamma_{p'})$$

2^i : # of outcomes that can be encoded using i bits

$2^i - 2^j$: # of outcomes that can be encoded strictly using the additional $(i-j)$ bits

ϵ : # of outcomes that are anyway encodable using the unused arrangement of i bits.

Ex:-

We added 5 new outcomes to 3 outcomes (encodable by 2 bits), and now we have

8 outcomes (encodable by 3 bits)

Among the 5 newly added outcomes $2^3 - 2^2 = 4$.
is the max. # of outcomes that will strictly need the additional $3-2=1$ bit.

And the remaining 1 outcome was already encodable using 2 bits itself.

Thus,

the addition of 5 new outcomes (even if lower probability than any present before) doesn't affect the codability of the original 3 outcomes.
Those are still encodable using $i=2$ bits only.

All the additional outcomes can be encoded using only the additional bits (and the unused arrangement of i bits if needed), and the original i bits are free to be used for encoding all the original outcomes.

→ For an outcome that has a probability P , $i = \log_2(\frac{1}{P})$ bits can encode it in any event.

We have the length of the codeword that can represent an outcome based upon its probability P in any event - it is fair to call this length the information of the outcome.

$$i = \log_2(\frac{1}{P})$$

⇒ The information i associated with any outcome of probability p is the minimum # of bits that can be used for encoding this outcome in any event having any # of outcomes.

And this # is solely given by the probability of the outcome as, $i = \underline{\log_2(Y_p)}$.

- This definition depends upon the unit we are using, but given a unit, we have a well defined information.

- The information of an outcome (ie, the amount of data needed to represent this outcome) is independent of the total # of outcomes of the event.

Ex:-
 $P = 0.25 \Rightarrow i = \log_2\left(\frac{1}{0.25}\right) = \log_2(4) = 2$: we need 2 bits for encoding this outcome in any event (that can have 2, 4 or 16 outcomes).

- It is more correct to name the specific units
(along with information)

Ex: a given state has i bits of information associated to it.

Information or Entropy of an Event or System

What would be the information content of any event, i.e., the information of all its outcomes combined?

What's the information content of a system that can have different states with different probabilities?

- The information content of the event or the system is called Entropy.
- In statistics, such a system (event) that can have different states (outcomes) with different probabilities is called a probability distribution.

For example, if there are 3 states (outcomes) with information i_1, i_2 and i_3 , what is the entropy of this system (event)?

From the system's point of view these states are not equal, since these states or outcomes are not occurring with equal probability.

∴ We take the weighted average of the information (or entropy) of individual states of a system.

$$\begin{aligned} H &= \sum_i p_i^i \\ &= \sum_i p_i \log_2 (1/p_i) \\ &= - \sum_i p_i \log_2 (p_i) \end{aligned}$$

⇒ The entropy of a system describes the minimum number of bits that can be used to describe any state of the system on average.

$$S = H$$

Shannon entropy

Suppose we learn the value of a random variable X . The Shannon entropy of X quantifies how much information we gain, on average, when we learn the value of X .

Alternatively, the entropy of X measures the amount of uncertainty about X before we learn its value.

Intuitively, the information content of a random variable should not depend on the labels attached to the different values that may be taken by the random variable.

Ex:-

We expect that a random variable taking the values 'heads' and 'tails' with respective probabilities $\frac{1}{4}$ and $\frac{3}{4}$ contains the same amount of information as a random variable that takes the values 0 and 1 with respective probabilities $\frac{1}{4}$ and $\frac{3}{4}$.

∴ The entropy of a random variable is defined to be a function of the probabilities of the different possible values of the random variable takes, and is not influenced by the labels used for these values.

We often write the entropy as a function of a probability distribution P_1, P_2, \dots, P_n .

The Shannon entropy associated with the probability distribution P_1, P_2, \dots, P_n is defined by,

$$H(X) = H(P_1, P_2, \dots, P_n) = - \sum_x P_x \log(P_x)$$

When $P_x = 0$,

$\log 0$ is undefined

Intuitively, an event which can never occur should not contribute to the entropy, so by convention we agree that $0 \log(0) \equiv 0$.

More formally, $\lim_{x \rightarrow 0} x \log(x) = 0$

- Uniform distributions have maximum uncertainty
- Uniform distributions with more outcomes have more uncertainty.

- The best reason for this definition of entropy is that it can be used to quantify the resources needed to store information.

More concretely,

Suppose there is some source (perhaps a radio antenna) which is producing information of some sort, say in the form of a bit string.

Let's consider a very simple model for a source — we model it as producing a string x_1, x_2, \dots of independent, identically distributed random variables.

Most information sources don't behave quite this way, but it's often a good approximation to reality.

Shannon asked what minimal physical resources are required to store the information being produced by the source, in such a way that at a later time the information can be reconstructed?

ANS : Entropy

i.e., $H(X)$ bits are required per source symbol.
where $H(X) = H(X_1) = H(X_2) = \dots$ is the
entropy of each random variable modeling the
source. This result is known as Shannon's
noiseless coding theorem.

Example

Suppose an information source produces one
of the 4 symbols, 1, 2, 3 or 4. Without compression
2 bits of storage space corresponding to the 4
possible outputs are consumed for each use of
the source.

Suppose that the symbol 1 is produced by
the source with probability $\frac{1}{2}$, the symbol 2
with probability $\frac{1}{4}$, and the symbols 3 and 4
both with probability $\frac{1}{8}$.

We can make use of the bias b/w the source
outputs to compress the source, using fewer bits
to store commonly occurring symbols such as 1,
and more bits to store rarely occurring symbols
like 3 and 4.

One possible compression scheme is to encode 1 as the bit string 0, 2 as the bit string 10, 3 as the bit string 110, and 4 as the bit string 111.

The avg. length of the compressed string is,

$$\frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{8} \times 3 = \frac{7}{4} \text{ bits of}$$

information per use of the source.

This is less than is required in the naive approach to store this source.

This matches the entropy of the source,

$$\begin{aligned} H(X) &= -\sum P_i \log(P_i) \\ &= -\frac{1}{2} \log\left(\frac{1}{2}\right) - \frac{1}{4} \log\left(\frac{1}{4}\right) - \frac{1}{8} \log\left(\frac{1}{8}\right) - \frac{1}{8} \log\left(\frac{1}{8}\right) \\ &= \underline{\underline{\frac{7}{4}}} \end{aligned}$$

Moreover, it turns out that any attempt to compress the source further results in data being irretrievably lost; the entropy quantifies the optimal compression that may be achieved.

Ex: 11.1

What's the entropy associated with the toss of a fair coin? With the roll of a fair die?
How would the entropy behave if the coin or die were unfair?

Aus: $H(X) = -\sum_x P_x \log(P_x)$

② $H(X) = -\frac{1}{2} \log(\frac{1}{2}) - \frac{1}{2} \log(\frac{1}{2}) = 2 \times \frac{1}{2} \log 2 = 1$

③ $H(X) = -\frac{1}{6} \log \frac{1}{6} - \frac{1}{6} \log \frac{1}{6} - \dots$
 $= \log(6) \approx 2.58496 = 2.585 \text{ bits.}$

④ Unfair coin:
 $P(X = \text{heads}) = \frac{1}{4}$
 $P(X = \text{tails}) = \frac{3}{4}$

$$H(X) = -\frac{1}{4} \log \frac{1}{4} - \frac{3}{4} \log \frac{3}{4}$$

$$= \frac{1}{2} + \frac{3}{4} \log \frac{4}{3} = \frac{1}{2} + \frac{3}{4} [\log 4 - \log 3]$$

$$= \frac{1}{2} + \frac{3}{2} - \frac{3}{4} \log(3) = \frac{1}{2} + \frac{3}{2} - \frac{3}{4}(1.58496)$$

$$= 0.81128$$

d)

Unfair die: $P(1) = \frac{1}{2}$,

$$P(2) = \frac{1}{30}, P(3) = \frac{2}{30}, \dots, P(6) = \frac{5}{30}$$



1

$$H(X) = \frac{1}{2} \log(2) + \frac{1}{30} \log 30 + \frac{1}{15} \log 15 + \frac{1}{10} \log 10 + \frac{2}{15} \log \frac{15}{2} + \frac{1}{6} \log 6$$

=

2

$$P(1) = 8\%, P(2) = 17\%, P(3) = P(4) = 25\%, P(5) = 17\%, P(6) = 8\%$$

$$H(X) = \underline{\underline{2.454 \text{ bits}}}$$

⇒ Entropy has the value based on probability and is maximum when the coin is fair. But when the coin is fair we have the highest uncertainty about the outcome of the experiment. Hence, we could state that entropy is the amount of uncertainty in a random variable.

Basic properties of Entropy

A The binary entropy

The entropy of a two-outcome random variable is called the binary entropy, which is defined as,

$$H_{\text{bin}}(p) \equiv -p \log p - (1-p) \log(1-p)$$

where p and $1-p$ are the probabilities of the two outcomes.

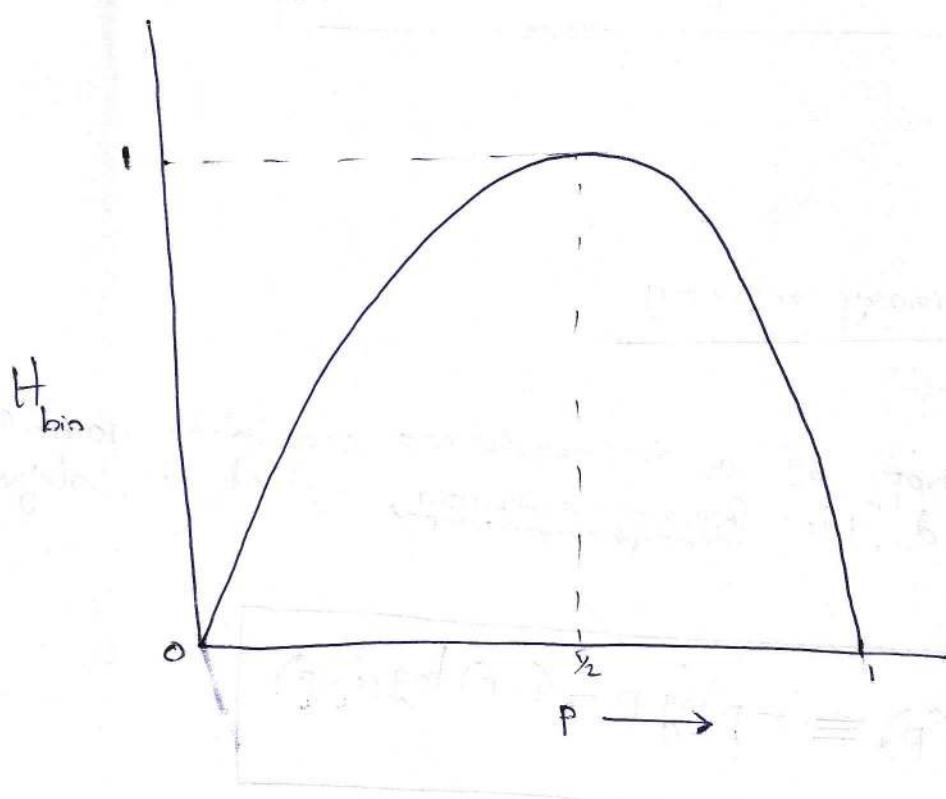
- $H(p) = H(1-p)$
- $H_{\text{bin}}^{\max}(p) = 1$ at $p = \frac{1}{2}$

Proof

$$H_{\text{bin}}(p) = -p \log p + (1-p) \log(1-p) = \log \left(\frac{1-p}{p} \right) = \alpha$$

$$\Rightarrow \frac{1-p}{p} = 1 \Leftrightarrow 1-p = p \Rightarrow 1 = 2p \Rightarrow p = \frac{1}{2}$$

* Binary entropy function, $H(p)$



How the entropy behaves when we mix 2 or more probability distributions ?

Example

Alice has in her possession 2 coins, one a quarter from the US, the other a dollar coin from Australia. Both coins have been altered to exhibit bias with

p_U : probability of heads on the US coin

p_A : probability of heads on the Australian coin.

Suppose,

Alice flips the US coin with probability $\frac{1}{2}$ and the Australian coin with probability $1 - \frac{1}{2}$, telling Bob whether the result was heads or tails.

How much information does Bob gain on average ?

Intuitively, it is clear that Bob should gain at least as much information as the average of the information he should have gained from a US coin flip or an Australian coin flip.

This can be expressed as an equation,

$$H(qP_U + (1-q)P_A) \geq qH(P_U) + (1-q)H(P_A)$$

where,

P_U/P_A : probability of heads on the US/Australian coin.

$q(1-q)$: probability of flipping the US/Australian coin.

Note: Sometimes the inequality can be strict, because Bob gains information not only about the value (heads or tails) of the coin, but also some additional information about the identity of the coin.

Consider the probability vector,

$$P(i) = q P_U(i) + (1-q) P_A(i)$$

This corresponds to a process involving two sequential choices: 1st a coin is tossed to decide b/w P_U and P_A , with the former picked up with probability q , and then another or the later with probability $(1-q)$, to decide b/w the possible outcome i according to the probabilities given by P_U or P_A .

① What's the information content corresponding to the possible outcomes?

In other words, completely disregarding the internal structure of P (meaning it is being actually composed of 2 sequential choices), what is its entropy?

$$\Rightarrow H(P) = H(q P_U + (1-q) P_A)$$

② If you already know the outcome of the 1st choice (ie., you know whether you are in the P_U or in the P_A branch). What's the information content corresponding to the possible outcomes?

This is $H(P_U)$ if we are in the P_U branch or $H(P_A)$ if we are in the other. If we don't want to assume we already know the outcome of the 1st coin toss, we ought to consider the (weighted) average of such entropy, which is

$$q H(P_U) + (1-q) H(P_A)$$

* The binary entropy is a concave function.

$$H_{\text{bin}}(px_1 + (1-p)x_2) \geq p H_{\text{bin}}(x_1) + (1-p) H_{\text{bin}}(x_2)$$

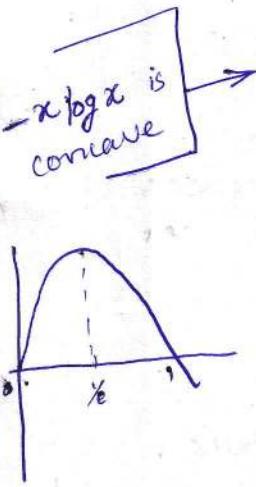
where $0 \leq p, x_1, x_2 \leq 1$.

→ Mixing 2 probability distributions increases the entropy.

The binary entropy is strictly concave,
i.e., the above inequality is an equality for the trivial cases

Proof

$$\begin{aligned}
 H_{\text{bin}}(p\alpha_1 + (1-p)\alpha_2) &= - \sum_i (p\alpha_1(i) + (1-p)\alpha_2(i)) \log(p\alpha_1(i) + (1-p)\alpha_2(i)) \\
 &\geq \sum_i \left(-p\alpha_1(i) \log(p\alpha_1(i)) - (1-p)\alpha_2(i) \log((1-p)\alpha_2(i)) \right) \\
 &\geq \sum_i \left(-p\alpha_1(i) \log(\alpha_1(i)) - (1-p)\alpha_2(i) \log(\alpha_2(i)) \right) \\
 &= -p \sum_i \alpha_1(i) \log(\alpha_1(i)) - (1-p) \sum_i \alpha_2(i) \log(\alpha_2(i)) \\
 &= p H_{\text{bin}}(\alpha_1) + (1-p) H_{\text{bin}}(\alpha_2)
 \end{aligned}$$



$f(\alpha_2, (1-p))$

$\geq f(\alpha_1) + (1-p)$

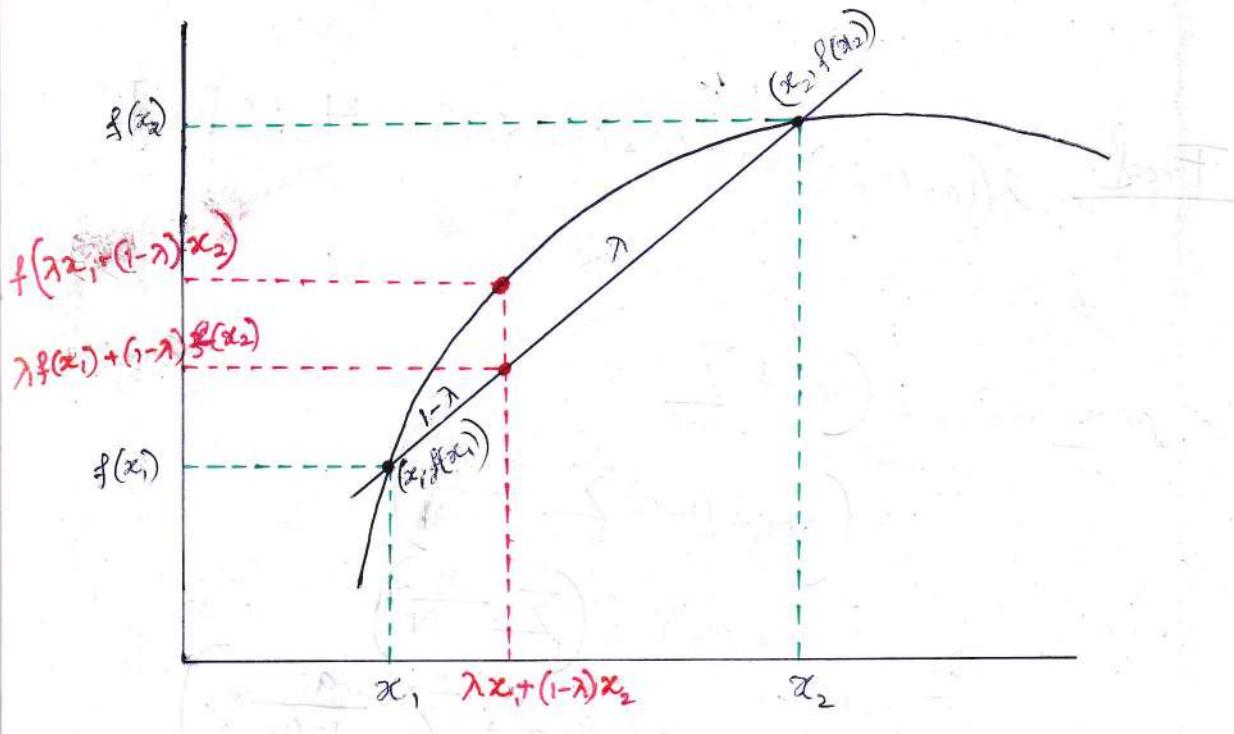
$f(x_1)$

Concave function

A real valued function f on an interval (a, b) is called concave, provided for each pair of points $x_1, x_2 \in (a, b)$ and each λ with $0 \leq \lambda \leq 1$, we have

$$f(\lambda x_1 + (1-\lambda)x_2) \geq \lambda f(x_1) + (1-\lambda)f(x_2) \quad \forall \lambda \in [0, 1]$$

i.e., a line connecting any 2 points on the function lies below or on the function.



* If f is a concave function of one real variable. Let $x_1, \dots, x_n \in \mathbb{R}$ and $a_1, \dots, a_n \geq 0$ satisfy $a_1 + \dots + a_n = 1$. Then,

$$f(a_1x_1 + \dots + a_nx_n) \geq a_1f(x_1) + \dots + a_nf(x_n)$$



$$f\left(\sum_{i=1}^n a_i x_i\right) \geq \sum_{i=1}^n a_i f(x_i)$$

Proof $f(ta + (1-t)b) \geq t f(a) + (1-t) f(b) \quad \forall t \in [0, 1]$

Start
1/1/2023

$$\begin{aligned} f\left(\sum_{i=1}^n a_i x_i\right) &= f\left(a_1 x_1 + \sum_{i=2}^n a_i x_i\right) \\ &= f\left(a_1 x_1 + (1-a_1) \sum_{i=2}^n \frac{a_i x_i}{1-a_1}\right) \\ &\geq a_1 f(x_1) + (1-a_1) f\left(\sum_{i=2}^n \frac{a_i x_i}{1-a_1}\right) \\ &= a_1 f(x_1) + (1-a_1) f\left(\frac{a_2}{1-a_1} x_2 + \sum_{i=3}^n \frac{a_i}{1-a_1} x_i\right) \\ &= a_1 f(x_1) + (1-a_1) f\left(\frac{a_2}{1-a_1} x_2 + \frac{1-a_1-a_2}{1-a_1-a_2} \sum_{i=3}^n \frac{1-a_1}{1-a_1-a_2} \frac{a_i}{1-a_1} x_i\right) \\ &\geq a_1 f(x_1) + (1-a_1) \left[\frac{a_2}{1-a_1} f(x_2) + \frac{1-a_1-a_2}{1-a_1} f\left(\sum_{i=3}^n \frac{a_i}{1-a_1-a_2} x_i\right) \right] \end{aligned}$$

$$= a_1 f(x_1) + a_2 f(x_2) + (1-a_1-a_2) f\left(\sum_{i=3}^n \frac{a_i}{1-a_1-a_2} x_i\right)$$

$$\geq a_1 f(x_1) + \dots + a_{n-1} f(x_{n-1}) + \\ (1-a_1-\dots-a_{n-1}) f\left(\sum_{i=n}^n \frac{a_i}{1-a_1-\dots-a_{n-1}} x_i\right) \\ = a_1 f(x_1) + \dots + a_{n-1} f(x_{n-1}) + a_n f(x_n) = \sum_{i=1}^n a_i f(x_i)$$

* Concavity of f on (a, b) implies that the slope of the chord from $(x_1, f(x_1))$ to $(x_2, f(x_2))$ is greater than or equal to the slope of the chord from $(x_1, f(x_1))$ to $(x_n, f(x_n))$, i.e.,

$$\frac{f(x) - f(x_1)}{x - x_1} \geq \frac{f(x_2) - f(x_1)}{x_2 - x_1} \quad \forall x_1 < x < x_2 \text{ in } (a, b)$$

Proof

With $\alpha_1 < \alpha_2$ in (a, b) and $\alpha \in (\alpha_1, \alpha_2) \Rightarrow$ we have $\alpha = \lambda \alpha_1 + (1-\lambda) \alpha_2$ for $\lambda = \frac{\alpha_2 - \alpha}{\alpha_2 - \alpha_1}$

$$f(\lambda \alpha_1 + (1-\lambda) \alpha_2) \geq \lambda f(\alpha_1) + (1-\lambda) f(\alpha_2)$$

$$f(\alpha) \geq \left(\frac{\alpha_2 - \alpha}{\alpha_2 - \alpha_1} \right) f(\alpha_1) + \left(\frac{\alpha - \alpha_1}{\alpha_2 - \alpha_1} \right) f(\alpha_2)$$

$$(\alpha_2 - \alpha) f(\alpha) \geq (\alpha_2 - \alpha) f(\alpha_1) + (\alpha - \alpha_1) f(\alpha_2)$$

$$(\alpha_2 - \alpha + \alpha - \alpha_1) f(\alpha) \geq (\alpha_2 - \alpha_1) f(\alpha_1) + (\alpha - \alpha_1) f(\alpha_2)$$

$$(\alpha_2 - \alpha) f(\alpha) + (\alpha - \alpha_1) f(\alpha) \geq (\alpha_2 - \alpha_1) f(\alpha_1) + (\alpha - \alpha_1) f(\alpha_2)$$

$$(\alpha_2 - \alpha)(f(\alpha) - f(\alpha_1)) \geq (\alpha - \alpha_1)(f(\alpha_2) - f(\alpha_1))$$

$$\frac{f(\alpha) - f(\alpha_1)}{\alpha - \alpha_1} \geq \frac{f(\alpha_2) - f(\alpha_1)}{\alpha_2 - \alpha_1} \quad \forall \alpha_1 < \alpha < \alpha_2 \text{ in } (a, b).$$

Proof

fC

(x2)

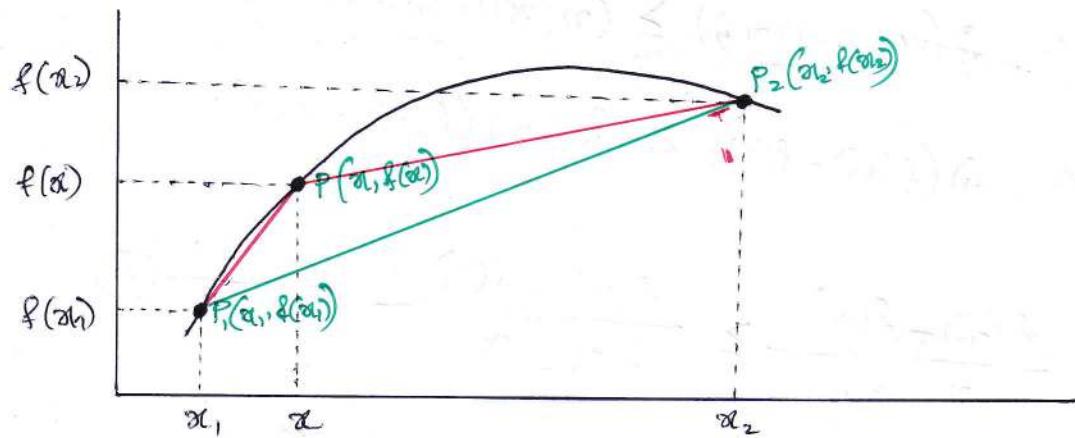
(x1)

* The Chordal Slope Lemma

Let f be concave on (a, b) . If $x_1 < x < x_2$ are in (a, b) then for points $P_1 = (x_1, f(x_1))$, $P = (x, f(x))$ and $P_2 = (x_2, f(x_2))$ we have

$$\frac{f(x) - f(x_1)}{x - x_1} \geq \frac{f(x_2) - f(x_1)}{x_2 - x_1} \geq \frac{f(x_2) - f(x)}{x_2 - x} \quad \forall x_1 < x < x_2 \text{ in } (a, b)$$

\Rightarrow Slope of $\overline{P_1 P} \geq$ Slope of $\overline{P_1 P_2} \geq$ Slope of $\overline{P P_2}$



Proof

$$f(x) \geq \left(\frac{x_2 - x}{x_2 - x_1} \right) f(x_1) + \left(\frac{x - x_1}{x_2 - x_1} \right) f(x_2)$$

$$(x_2 - x) f(x) \geq (x_2 - x_1) f(x_1) + (x - x_1) f(x_2)$$

$$\begin{aligned} (x_2 - x_1) f(x) &\geq (x_2 - x_1 + x_1 - x) f(x_1) + (x - x_1) f(x_2) \\ &= (x_2 - x_1) f(x_1) + (x - x_1) (f(x_2) - f(x_1)) \end{aligned}$$

$$(\alpha_2 - \alpha_1)(f(\alpha) - f(\alpha_1)) \geq (\alpha - \alpha_1)(f(\alpha_2) - f(\alpha))$$

$$\frac{f(\alpha) - f(\alpha_1)}{\alpha - \alpha_1} \geq \frac{f(\alpha_2) - f(\alpha)}{\alpha_2 - \alpha_1} \quad \text{--- } ①$$

Similarly,

$$\begin{aligned} (\alpha_2 - \alpha_1)f(\alpha) &\geq (\alpha_2 - \alpha)f(\alpha_1) + (\alpha - \alpha_1)f(\alpha_2) \\ &= (\alpha_2 - \alpha)f(\alpha_1) + (\alpha - \alpha_2 + \alpha_2 - \alpha_1)f(\alpha_2) \\ &= (\alpha_2 - \alpha)f(\alpha_1) + (\alpha - \alpha_2)f(\alpha_2) + (\alpha_2 - \alpha_1)f(\alpha_2) \end{aligned}$$

$$(\alpha_2 - \alpha_1)\underline{(f(\alpha) - f(\alpha_2))} \geq (\alpha_2 - \alpha)\underline{(f(\alpha_1) - f(\alpha_2))}$$

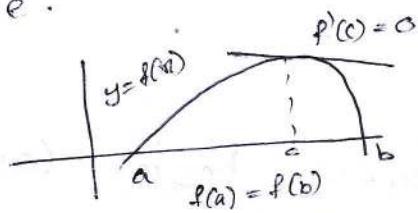
$$(\alpha_2 - \alpha_1)(f(\alpha_2) - f(\alpha)) \leq (\alpha_2 - \alpha)(f(\alpha_2) - f(\alpha_1))$$

$$\frac{f(\alpha_2) - f(\alpha_1)}{\alpha_2 - \alpha_1} \geq \frac{f(\alpha_2) - f(\alpha)}{\alpha_2 - \alpha} \quad \text{--- } ②$$

$$① \& ② \Rightarrow$$

$$\frac{f(\alpha) - f(\alpha_1)}{\alpha - \alpha_1} \geq \frac{f(\alpha_2) - f(\alpha_1)}{\alpha_2 - \alpha_1} \geq \frac{f(\alpha_2) - f(\alpha)}{\alpha_2 - \alpha}$$

- * If f is differentiable on (a, b) and its derivative f' is decreasing, then f is concave.
- If f'' exists on (a, b) and $f'' \leq 0$ on (a, b)
then f is concave.

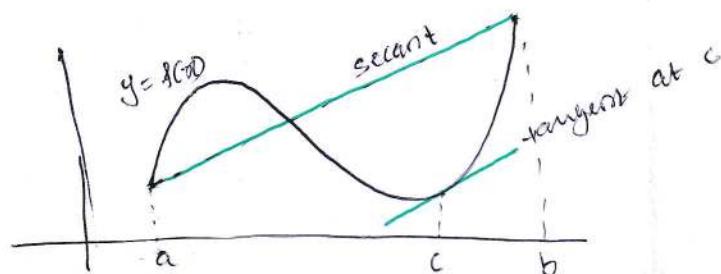


Proof

Rolle's theorem: If a real valued function f is continuous on a proper closed interval $[a, b]$, differentiable on the open interval (a, b) and $f(a) = f(b)$, then there exists at least one c in (a, b) such that $f'(c) = 0$.

i.e., any real valued differentiable function that attains equal values at 2 distinct points must have at least one stationary point somewhere between them.

Mean value theorem: If f is a continuous function on the closed interval $[a, b]$ and differentiable on the open interval (a, b) , then there exists a point $c \in (a, b)$ such that the tangent at c is parallel to the secant line thro' the end points $(a, f(a))$ and $(b, f(b))$, i.e., $f'(c) = \frac{f(b) - f(a)}{b - a}$



Let $x_1 < x_2$ be in (a, b) and let $x \in (x_1, x_2)$.

f is differentiable \Rightarrow Mean value theorem applies to f on intervals $[x_1, x]$ and $[x_1, x_2]$.

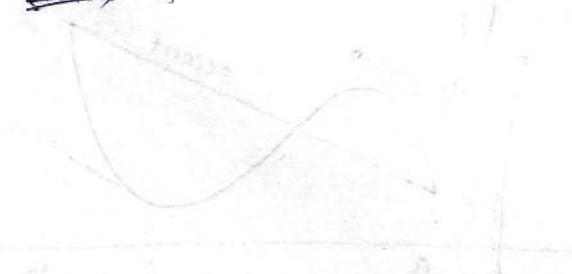
Choose $c_1 \in (x_1, x)$ and $c_2 \in (x_1, x_2)$ for which

$$f'(c_1) = \frac{f(x) - f(x_1)}{x - x_1} \text{ and } f'(c_2) = \frac{f(x_2) - f(x)}{x_2 - x}$$

Since f' is decreasing then,

$$\frac{f(x) - f(x_1)}{x - x_1} = f'(c_1) \geq f'(c_2) = \frac{f(x_2) - f(x)}{x_2 - x}$$

$\Rightarrow f$ is concave.



(B) The relative entropy (Kullback-Leibler divergence)

Suppose $p(x)$ and $q(x)$ are probability distributions on the same index set, \mathcal{X} . The relative entropy of $p(x)$ to $q(x)$ is defined as,

$$\begin{aligned}
 H(p(x) \| q(x)) &\equiv \sum_x p(x) \log \frac{p(x)}{q(x)} \\
 &\equiv \sum_x p(x) \left(\log \frac{1}{q(x)} - \log \frac{1}{p(x)} \right) \\
 &\equiv -\sum_x p(x) \log q(x) + \sum_x p(x) \log p(x) \\
 &\equiv -\sum_x p(x) \log q(x) - H(X)
 \end{aligned}$$

- * An outcome x that has a probability $p(x)$, we can encode it in any event using $i = \log \frac{1}{p(x)}$ bits.

- We use the convention that $-\log 0 = 0$ & $-p(x) \log 0 = +\infty$ if $p(x) > 0$

$\log \frac{0}{0} = 0$ and based on continuity arguments
that: $\log \frac{0}{0} = 0$ & $p \log \frac{P}{0} = \infty$

$\Rightarrow H(P(x) || Q(x)) = \infty$ if there exists any $x \in X$ such
that $P(x) > 0$ and $Q(x) = 0$

Since,

$$H(P(x) || Q(x)) = \sum_x P(x) \log \frac{P(x)}{Q(x)} = -\sum_x P(x) \log \frac{Q(x)}{P(x)}$$

and $-P(x) \log 0 = \infty$ if $P(x) > 0$.

$$* H(P(x) \parallel Q(x)) = \sum_x P(x) \log \frac{P(x)}{Q(x)} = \sum_x P(x) \left(\log \frac{1}{Q(x)} - \log \frac{1}{P(x)} \right)$$

An outcome x that has a probability $P(x)$, we can encode it in any event using $i = \log \frac{1}{P(x)}$ bits.

\Rightarrow Relative entropy $H(P(x) \parallel Q(x))$ is the # of extra bits required per letter on average to encode a source with a distribution $Q(x)$ when the true underlying distribution is $P(x)$.

ie., expected extra-message length per datum that must be communicated if a code that is optimal for a given (wrong) distribution $Q(x)$ is used, compared to using a code based on the true distribution $P(x)$. : it is the excess entropy.

ie., If we knew the true distribution of the random variable, then we could construct a code with average description length $H(P)$.

If instead, we used the code for a distribution Q , we would need $\sum_x P(x) \log \frac{1}{Q(x)} = H(P) + H(P \parallel Q)$ bits on the average to describe the random variable.

- The relative entropy / Kullback - Leibler divergence is a very useful entropy-like measure of the closeness of 2 probability distributions, $p(x)$ and $q(x)$, over the same index set \mathcal{X} .

Relative entropy is always non-negative and is zero if and only if $P = Q$.

However, it is not a true distance b/w distributions since it is not symmetric and does not satisfy the triangle inequality — it is not a metric.

* The relative entropy is non-negative

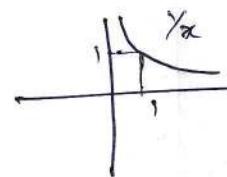
$H(P(x) \parallel Q(x)) \geq 0$, with equality iff $P(x) = Q(x)$

for all x .

Proof

Lemma: $\ln(x) \leq x-1$ for all $x > 0$

with $\ln(x) = x-1$ iff $x=1$



Proof

case 1: $0 < x < 1 \rightarrow \frac{1}{x} > 1$

$$\Rightarrow \int_1^x \frac{dx}{x} \geq \int_1^x dx \Rightarrow [\ln(x)]_1^x = [x]_1^x$$

$$\Rightarrow \ln(1) - \ln(x) \geq 1 - x$$

$$\Rightarrow -\ln(x) \geq 1 - x \rightarrow \underline{\ln(x) \leq x-1}$$

case 2: $x \geq 1 \rightarrow \frac{1}{x} \leq 1$

$$\Rightarrow \int_1^n \frac{dx}{x} \leq \int_1^n dx \Rightarrow [\ln(x)]_1^n \leq [x]_1^n$$

$$\Rightarrow \ln(n) - \ln(1) = \underline{\ln(n) \leq n-1}$$

$$\therefore \ln(x) \leq x-1 \quad \forall x > 0$$

② For $x > 0$,

$$f(x) = \ln(x) - x + 1 \Rightarrow f'(x) = \frac{1}{x} - 1$$

$$f'(x) = 0 \Rightarrow x = 1$$

$$f''(x) = -\frac{1}{x^2} < 0 \quad \forall x > 0$$

$\therefore x = 1$ is a maximum.

$$\lim_{n \rightarrow \infty} f(n) = -\infty = \lim_{x \rightarrow \infty} f(x)$$

$x = 1$ is global maximum

$$\forall x > 0, f(x) \leq f(1) = 0$$

$$\therefore \ln(x) - x + 1 \leq 0 \Rightarrow \ln(x) \leq x - 1 \quad \forall x > 0.$$

$$\textcircled{3}. \quad \ln(x) \leq x - 1 \Rightarrow x \leq e^x - 1 \Rightarrow e^x \geq x + 1 \quad \text{by putting } e^x \text{ for } x.$$

$$e^x = \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n \quad \text{Bernoulli's inequality}$$

$$(1+t)^n \geq 1+nt \quad \text{if } t > -1 \text{ and } n > 0.$$

$$\frac{x}{n} > 0 \Rightarrow \left(1 + \frac{x}{n}\right)^n \geq 1 + x \Rightarrow e^x = \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n \geq 1 + x.$$

$$\Rightarrow \ln(x) \leq x - 1 \quad \forall x > 0$$

$$1 - \frac{1}{n} < \ln(2) \Leftrightarrow \ln(n) > n - 1$$

$$1 = \ln(2) \Leftrightarrow 2 = e^1$$

$$\frac{\ln(n)}{\ln(2)} = \log_2(n) \leq n-1 \quad \forall n > 0$$

$$\Rightarrow \log(n) \ln(2) = \ln(n) \leq n-1$$

with equality iff $n=1$.

$$-\log(\ln(n)) \ln(2) \geq 1-n \Rightarrow -\log(n) \geq \frac{1-n}{\ln(2)}$$

$$\begin{aligned} H(p(x) \| q(x)) &= \sum_x p(x) \log \frac{p(x)}{q(x)} \\ &= -\sum_x p(x) \log \left(\frac{q(x)}{p(x)} \right) \\ &\geq \frac{1}{\ln(2)} \sum_x p(x) \left(1 - \frac{q(x)}{p(x)} \right) \\ &\leq \frac{1}{\ln(2)} \sum_x (p(x) - q(x)) \end{aligned}$$

$$H(p(x) \parallel q(x)) \geq \frac{1}{\ln(2)} \left[\sum_x p(x) - \sum_x q(x) \right]$$

$$= \frac{1}{\ln(2)} (1 - 1) = 0$$

$$\Rightarrow H(p(x) \parallel q(x)) \geq 0$$

$$\Rightarrow H(p(x) \parallel q(x)) = 0 \quad \text{iff} \quad \begin{array}{l} q(x)/p(x) = 1 \quad \forall x \\ p(x) = q(x) \quad \forall x \end{array}$$

- * χ^2 divergence is a function (binary) that takes 2 probability distributions as input, and returns a # that measures how much they differ. The # returned must be non-negative, and equal to zero iff the 2 distributions are identical.
Bigger numbers indicate greater dissimilarity.

In addition to the requirements above for the divergence, a distance metric must also be symmetric : $D(a|b) = D(b|a)$. And, it must satisfy the triangle inequality : $D(a|c) \leq D(a|b) + D(b|c)$

Divergences are defined specifically on probability distributions, whereas distance metric can be defined on other types of objects too.

All distance metrics b/w probability distributions are also divergences, but the converse is not true.

→ divergences measure the dissimilarity b/w distributions.

The simplest divergence is squared Euclidean distance (SED) and divergences can be viewed as generalizations of SED.

The relative entropy is often useful, not in itself but because other entropic quantities can be regarded as special cases of the relative entropy.

* Suppose X is a random variable with d outcomes.

$$\text{Then, } H(X) \leq \log(d)$$

with equality iff X is uniformly distributed over those d outcomes.

Proof

Suppose $p(\alpha)$ is a probability distribution for X , over d outcomes.

Let $q(\alpha) = \frac{1}{d}$ be the uniform probability distribution over these outcomes. Then,

$$H(p(\alpha) \parallel q(\alpha)) = H(p(\alpha) \parallel \frac{1}{d}) = - \sum_{\alpha} p(\alpha) \log q(\alpha) - H(X)$$

$$= - \sum_{\alpha} p(\alpha) \log \left(\frac{1}{d} \right) - H(X)$$

$$= - \log \left(\frac{1}{d} \right) \sum_{\alpha} p(\alpha) - H(X)$$

$$= \log d - H(X) \geq 0$$

$$\Rightarrow \underline{H(X) \leq \log d}$$

Ex: 11.5

Subadditivity of the Shannon entropy

$$H(X,Y) \leq H(X) + H(Y)$$

with equality iff X and Y are independent random variables.

Proof

$$\begin{aligned} H(X) + H(Y) - H(X,Y) &= \sum_x p(x) \log p(x) - \sum_y p(y) \log p(y) \\ &\quad + \sum_{x,y} p(x,y) \log p(x,y) \\ &= - \sum_{x,y} p(x,y) \log p(x) - \sum_{x,y} p(x,y) \log p(y) + \sum_{x,y} p(x,y) \log p(x,y) \\ &= \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} = H(p(x,y) || p(x)p(y)) \geq 0 \end{aligned}$$

Equality is achieved iff $p(x,y) = p(x)p(y)$ if and only if X and Y are independent.

© Conditional entropy & mutual information

Suppose X and Y are random variables.

— How is the information content of X related to the information content of Y ?

— The joint entropy $H(X,Y)$ of a pair of discrete random variables (X,Y) with a joint distribution $p(x,y)$ is defined as

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

The joint entropy measures over total uncertainty about the pair (X,Y) .

There is nothing really new in this definition because (X, Y) can be considered to be a single vector-valued random variable.

$$(X, Y) \sim f_{(X, Y)}(x, y) = f_X(x)f_Y(y)$$

Suppose we know the value of y , so we have acquired $H(y)$ bits of information about the pair, (x,y) . The remaining uncertainty about the pair (x,y) , is associated with our remaining lack of knowledge about x , given that we know y .

The entropy of x conditional on knowing y is defined as,

$$\begin{aligned} H(x|y) &= \sum_y p(y) H(x|y) \\ &= \sum_y p(y) - \sum_x p(x|y) \log p(x|y) \\ &= - \sum_{x,y} p(y) p(x|y) \log p(x|y) \\ &= - \sum_{x,y} p(x,y) \log p(x|y) \end{aligned}$$

→ The conditional entropy $H(X|Y)$ is a measure of how uncertain we are, on average, about the value of X , given that we know the value of Y .

→ If you know the value of Y , then you know the value of X .

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

$$\text{Entropy of } X = - \sum_x p(x) \log p(x)$$

$$\text{Entropy of } Y = - \sum_y p(y) \log p(y)$$

$$\text{Mutual Information } I(X;Y) = - \sum_{x,y} p(x,y) \log p(x,y)$$

*

$$\begin{aligned} H(x, y) &= H(x) + H(y|x) \\ &= H(y) + H(x|y) \end{aligned}$$

Proof

$$\begin{aligned} H(x, y) &= - \sum_{x,y} p(x,y) \log p(x,y) \\ &= - \sum_{x,y} p(x,y) \log p(y) p(x|y) \\ &= - \sum_{x,y} p(x,y) \log p(y) - \sum_{x,y} p(x,y) \log p(x|y) \\ &= - \sum_y \left(\sum_x p(x,y) \right) \log p(y) - \sum_{x,y} p(x,y) \log p(x|y) \\ &= - \sum_y p(y) \log p(y) - \sum_{x,y} p(x,y) \log p(x|y) \\ &= \underline{H(y) + H(x|y)} \end{aligned}$$

The mutual information $H(X:Y)$ content of X and Y , measures how much information X and Y have in common.

i.e.,

is a measure of the amount of information that one random variable contains about another random variable.

Suppose we add the information content of X , $H(X)$, to the information content of Y . Information which is common to X and Y will have been counted twice in this sum, while information which is not common will have exactly once.

Subtracting off the joint information of (X,Y) , $H(X,Y)$, we obtain the common or mutual information of X and Y :

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

$$\begin{aligned}
 H(X,Y) &= H(X) + H(Y|X) \\
 &= H(Y) + H(X|Y)
 \end{aligned}$$

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

∴

$$\begin{aligned}
 H(X:Y) &= H(X) - H(X|Y) \\
 &= H(Y) - H(Y|X)
 \end{aligned}$$

(Ex: 11.5)

→ The mutual information is the reduction in the uncertainty of one random variable due to the knowledge of the other.

Ex: 11.7

$I(X;Y)$ is called mutual information. It is defined as $I(X;Y) = H(Y) - H(Y|X)$

Ans:

QESTAK
23/12/2022

$$(P(X)H - P(X)H + (X)H) = (P(X)H)$$

* Consider 2 random variables X and Y with a joint probability mass function $p(x,y)$ and marginal probability mass functions $p(x)$ and $q(y)$. The mutual information $H(X:Y)$ can be defined as the relative entropy w.r.t. the joint distribution and the product distribution $p(x)P(y)$,

$$H(X:Y) = H(p(x,y) \parallel p(x)p(y))$$

$$= \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

(Ex: 11.5)

Ex: 11.7. Find an expression for the conditional entropy $H(Y|X)$ as a relative entropy w.r.t. a probability distributions. Use this to deduce that $H(Y|X) \geq 0$, and find the equality conditions.

Ans:

$$H(Y|X) = \sum_x p(x) H(Y|X=x)$$

$$= \sum_x p(x) - \sum_y p(y|x) \log p(y|x)$$

$$= - \sum_{x,y} p(x) p(y|x) \log p(y|x)$$

$$= - \sum_{x,y} p(x,y) \log p(y|x)$$

Q.E. Stark
23/12/2022

$$P(Y|X) = \frac{P(X,Y)}{P(X)} \quad \iff P(X,Y) = P(X)P(Y|X)$$

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

$$= - \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x) \times \frac{1}{d_y}} = - \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)} + \sum_{x,y} p(x,y) \log \frac{1}{d_y}$$

$$H(P(x)/\#q(x)) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

$p(\cdot)$ & $q(\cdot)$ must be defined over the same set.

where d_y : # of outcome for the random variable y .

$$= - \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)} + \sum_{x,y} p(x,y) \log \left(\frac{1}{d_y} \right)$$

We have the joint distribution $p(x,y)$ b/w 2 random variables X and Y , and consider a distribution q which is defined as

$$q(y) := p(x) \times \frac{1}{d_y} = \frac{p(x)}{d_y}$$

i.e., q is a product of distribution of the marginal distribution of X together with a uniform distribution over the outcome space of Y .

$$H(Y|X) = \underline{-H(P(X,Y) \parallel P(X)/dY) + \log(dY)}$$

□ Basic properties of Shannon entropy

$$1. H(x,y) = H(y,x), \quad H(x:y) = H(Y:X)$$

$$2. H(Y|X) \geq 0 \Rightarrow H(X:Y) \leq H(Y)$$

with equality iff Y is a function of X ,

$$Y = f(X)$$

$$3. H(X) \leq H(X,Y)$$

with equality iff Y is a function of X

$$4. \text{Subadditivity : } H(X,Y) \leq H(X) + H(Y)$$

Ex: 11.5 with equality iff X and Y are independent random variables.

$$5. H(Y|X) \leq H(Y) \text{ and thus } H(X:Y) \geq 0.$$

with equality in each iff X and Y are independent random variables.

$$6. \text{Strong subadditivity :}$$

$$H(X,Y,Z) + H(Y) \leq H(X,Y) + H(Y,Z)$$

with equality iff $Z \rightarrow Y \rightarrow X$ forms a markov chain.

$$H(X|Y, Z) \leq H(X|Y)$$

$$(X, Y, Z) \in \Omega \text{ mit } H(X, Y|Z) = H(X|Y)$$

$$H(Y|H) \geq H(X|Y) \Leftrightarrow H(X|Y) \leq H(X|H)$$

X ist unabhängig von Y \Rightarrow gemeinsame Bedingung

$$H(Y|H) \geq H(X|H) \Leftrightarrow H(X|H) \leq H(X|Y)$$

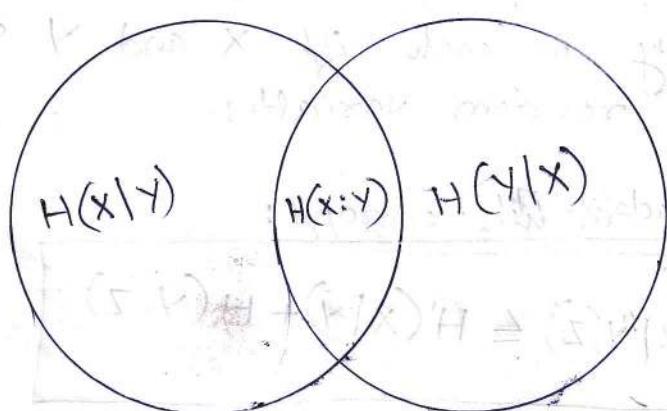
$$H(Y|H) \geq H(X|H) \Leftrightarrow H(X|H) \leq H(X|Y)$$

X ist unabhängig von Y \Rightarrow gemeinsame Bedingung

$$(V) H(X|Y) \geq H(X|H) \Leftrightarrow H(X|H) \leq H(X|Y)$$

Bedingung von Y hat X kein gemeinsame Bedingung

$H(X|Y) \geq H(X|H) \Leftrightarrow H(X|H) \leq H(X|Y)$



$$(S) H(X|Y) + H(Y|X) \geq H(X:Y) \Leftrightarrow H(X:Y) \leq H(X|Y) + H(Y|X)$$

7. Conditioning reduces entropy

$$H(X|Y, Z) \leq H(X|Y)$$

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

$$(X|Y)H \geq (Y|X)H \leftarrow \text{obs } (X|Y)H$$

X is random in $(Y|X)H$ entropy form.

$$(X|Y)H$$

$$(Y|X)H \geq (X|Y)H$$

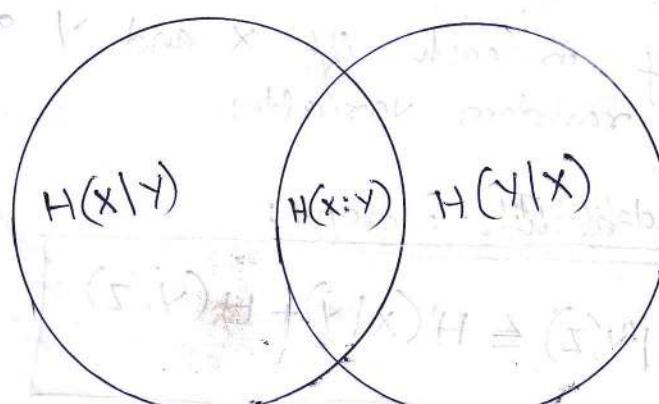
X is random in $(X|Y)H$ entropy form.

$$(Y|X)H + (X|Y)H \geq (Y|X)H$$

conditioning on Y less X bits added.

so $H(Y|X) \leq H(Y)$

so $H(X|Y) \leq H(X)$ with loss $(Y|X)H$ and $(X|Y)H$.



$$(X|Y)H + (Y|X)H \geq (X|Y)H + (Y|X)H$$

$$H(X) \quad X \leftarrow Y \leftarrow Z \quad H(Y)$$

Proof

$$H(x,y) = H(x) + H(y|x) = H(y) + H(x|y)$$

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

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

$$= H(x) - H(x|y) = H(y) - H(y|x)$$

2. $H(x:y) = H(y) - H(y|x)$

$$H(y|x) = \sum_x p(x) H(y|x=x)$$

$$= \sum_x p(x) - \sum_y p(y|x) \log p(y|x)$$

$$= - \sum_{x,y} p(x,y) \log p(y|x)$$

$$0 \leq p(y|x) \leq 1 \Rightarrow \log p(y|x) \leq 0 \Rightarrow -\log p(y|x) \geq 0$$

$$H(y|x) \geq 0$$

$$H(x:y) = H(y) - H(y|x) \Rightarrow H(y) - H(x:y) \geq 0$$

$$\Rightarrow H(x:y) \leq H(y)$$

$H(Y|X)=0$ Equality iff Y is a deterministic function of X
 i.e., knowing X determines the value of Y
 and vice versa.
 $H(Y|X) = -\sum_{x,y} p(x,y) \log p(y|x) = 0 \Rightarrow$ whenever $p(X=x, Y=y) > 0$
 we have $p(Y=y | X=x) = 1$

$$H(Y|X)=0 \implies H(X:Y)=H(Y)$$

\therefore The mutual information is the same as the uncertainty contained in Y alone, namely the entropy of Y .

- $H(X) \geq 0$

$H(X)=0$ iff X is deterministic.

$$\begin{aligned}
 3. \quad H(X, Y) &= - \sum_{x,y} P(x,y) \log P(x,y) \\
 &= - \sum_{x,y} P(x,y) \log P(x) P(y|x) \\
 &= - \sum_{x,y} P(x,y) \log P(x) - \sum_{x,y} P(x,y) \log P(y|x) \\
 &= - \sum_x \left(\sum_y P(x,y) \right) \log P(x) - \sum_{x,y} P(x,y) \log P(y|x) \\
 &= - \sum_x P(x) \log P(x) - \sum_{x,y} P(x,y) \log P(y|x) \\
 &= H(X) - \sum_{x,y} P(x,y) \log P(y|x).
 \end{aligned}$$

\downarrow

$$\begin{aligned}
 P(y|x) \leq 1 \Rightarrow \log P(y|x) \leq 0 \Rightarrow -\log P(y|x) \geq 0 \\
 \therefore H(Y|X) = - \sum_{x,y} P(x,y) \log P(y|x) \geq 0
 \end{aligned}$$

\downarrow

$$\begin{aligned}
 \therefore H(X, Y) - H(X) \geq 0 \Rightarrow \underline{H(X) \leq H(X, Y)}
 \end{aligned}$$

$$5. H(x,y) = H(x) + H(y) - H(x:y)$$

$$H(x,y) \leq H(x) + H(y)$$

$$\Rightarrow H(x:y) \geq 0$$

$$H(x:y) = H(y) - H(y|x) \geq 0$$

$$\Rightarrow H(y|x) \leq H(y)$$

$$H(x:y) = H(p(x|y) || p(y) p(y)) = 0$$

$$\Rightarrow p(x|y) = p(x) p(y)$$

\therefore equality iff x & y are independent random variables.

$$\textcircled{a} \quad H(x,y,z) + H(y) \leq H(x,y) + H(y,z)$$

Proof

$$H(x,y,z) = - \sum_{x,y,z} p(x,y,z) \log p(x,y,z)$$

Chain rule: $p(x,y) = p(y) p(x|y)$

$$P(x_1, y, z) = P(y, z) P(x_1 | y, z)$$

$$\begin{aligned}
 H(x_1, y, z) &= - \sum_{x_1, y, z} P(x_1, y, z) \log P(x_1 | y, z) P(y, z) \\
 &= - \sum_{x_1, y, z} P(x_1, y, z) \left\{ \log P(x_1 | y, z) + \log P(y, z) \right\} \\
 &= - \sum_{x_1, y, z} P(x_1, y, z) \log P(x_1 | y, z) - \sum_{y, z} \sum_x P(x, y, z) \times \log P(y, z) \\
 &= - \sum_{x_1, y, z} P(x_1, y, z) \log P(x_1 | y, z) - \sum_{y, z} P(y, z) \log P(y, z) \\
 &= - \sum_{x_1, y, z} P(x_1, y, z) \log P(x_1 | y, z) + H(y, z)
 \end{aligned}$$

$$\begin{aligned}
 H(x, y) &= - \sum_{x, y} P(x, y) \log P(x, y) \\
 &= - \sum_{x, y} P(x, y) \log P(y) P(x | y) \\
 &= - \sum_{x, y} \left(\sum_z P(x, y, z) \right) \log P(x | y) - \sum_y \left(\sum_x P(x, y) \right) \log P(y) \\
 &= - \sum_{x, y, z} P(x, y, z) \log P(x | y) - \sum_y P(y) \log P(y) \\
 &= - \sum_{x, y, z} P(x, y, z) \log P(x | y) + H(y)
 \end{aligned}$$

$$\begin{aligned} & \left\{ H(x,y,z) - H(y,z) \right\} + \left\{ H(y) - H(x,y) \right\} = \\ & = - \sum_{x,y,z} p(x,y,z) \log p(x|y,z) + \sum_{x,y,z} p(x,y,z) \log p(x|y) \end{aligned}$$

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

$x > 0$

$$\log(x) \ln(2) = \ln(x) \leq x-1$$

*check
for proof*

$$\leq \frac{1}{\ln(2)} \sum_{x,y,z} p(x,y,z) \left\{ \frac{p(x|y)}{p(x|y,z)} - 1 \right\}.$$

$$= \frac{1}{\ln(2)} \sum_{x,y,z} \left\{ p(x|y,z) p(y,z) \frac{p(x|y)}{p(x|y,z)} - p(x|y,z) \right\}$$

$$= \frac{1}{\ln(2)} \sum_{x,y} p(x|y) \sum_z p(y,z) - \frac{1}{\ln(2)} \sum_{x,y,z} p(x|y,z)$$

$$= \frac{1}{\ln(2)} \sum_{x,y} p(x|y) p(y) - \frac{1}{\ln(2)} * 1$$

$$= \frac{1}{\ln(2)} \left\{ \sum_{x,y} p(x|y) - 1 \right\} = 0.$$

$$\therefore H(x, y, z) + H(y) \leq H(x, y) + H(y, z)$$

$$I_n(x) = n - 1 \quad \text{when } n=1. \quad \left\{ \begin{array}{l} P(x|y) \\ P(x|y, z) \end{array} \right\} = 1$$

Equality is achieved iff $P(x|y) = P(x|y, z)$.

i.e., $Z \rightarrow Y \rightarrow X$ is a Markov chain.

$$\textcircled{7} \quad H(X|Y,Z) \leq H(X|Y)$$

Proof

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

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

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

$$\underbrace{H(X,Y|Z)}_{H(X|Y,Z)} - H(Y|Z) \leq H(X|Y) - H(Y)$$

$$\underline{\underline{H(X|Y,Z) \leq H(X|Y)}}$$

Intuitively, we expect that the uncertainty about X , given that we know the value of Y and Z is less than our uncertainty about X , given that we only know Y .

Chaining rule for conditional entropies

Theorem 11.4: Let x_1, \dots, x_n and Y be any set of random variables. Then,

$$H(x_1, \dots, x_n | Y) = \sum_{i=1}^n H(x_i | Y, x_1, \dots, x_{i-1})$$

$$H(x_1, x_2 | Y) = H(x_1 | Y) + H(x_2 | Y, x_1)$$

$$H(x_1, x_2, x_3 | Y) = H(x_1 | Y) + H(x_2 | Y, x_1) + H(x_3 | Y, x_1, x_2)$$

$$\begin{aligned} H(X, Y) &= H(X) + H(Y|X) \\ &= H(Y) + H(X|Y) \end{aligned}$$

Proof.

For $n=2$,

$$\begin{aligned} H(x_1, x_2 | Y) &= H(x_1, x_2, Y) - H(Y) \\ &= H(x_1, x_2, Y) - H(x_1, Y) + H(x_1, Y) - H(Y) \\ &= H(x_2 | Y, x_1) + H(x_1 | Y) \end{aligned}$$

which is the required result for $n=2$.

Assume that the result holds for general n ,

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

Using the $n=2$ case,

$$\begin{aligned} H(x_1, \dots, x_{n+1} | y) &= H(\underline{x_1, x_2, \dots, x_n} | y) \\ &= H(x_1 | y) + H(x_2, \dots, x_{n+1} | y, x_1) \end{aligned}$$

Applying the inductive hypothesis to the 1st term on the right hand side gives,

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

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

Ex: 11.8

Mutual information is not always subadditive

Let X & Y be independent identically distributed random variables taking the values 0 or 1, with probability $\frac{1}{2}$. Let $Z = X \oplus Y$, where \oplus denotes addition modulo 2. Show that the mutual information in this case is not subadditive.

$$H(X, Y : Z) \neq H(X : Z) + H(Y : Z)$$

Ans:

X	Y	Z	P	
0	0	0	$\frac{1}{4}$	$H(X : Y) = H(X) + H(Y) - H(X, Y)$
0	1	1	$\frac{1}{4}$	
1	0	1	$\frac{1}{4}$	
1	1	0	$\frac{1}{4}$	

$$H(X) = H(Y) = - \sum p(x) \log p(x) = 2 \times -\frac{1}{2} \log \frac{1}{2} = 1 = H(Z)$$

$$H(X, Y, Z) = - \sum_{x,y,z} p(x,y,z) \log p(x|y,z) = 4 \times \frac{1}{4} \log \frac{1}{4} = 4 \times \frac{1}{2} = 2$$

$$H(X, Y) = - \sum_{x,y} p(x,y) \log p(x,y) = 4 \times \frac{1}{4} \log \frac{1}{4} = 2 = H(Y, Z) = H(X, Z)$$

$$H(X : Z) = H(X) + H(Z) - H(X, Z) = 1 + 1 - 2 = 0 = H(Y : Z)$$

$$H(X, Y : Z) = H(X, Y) + H(Z) - H(X, Y, Z) = 2 + 1 - 2 = 1$$

$$\therefore H(X, Y : Z) = 1 \neq 0 = H(X : Z) + H(Y : Z)$$

Definition of Mutual Information: If X and Y are two random variables, then the mutual information between them is defined as $I(X;Y) = H(X) + H(Y) - H(X,Y)$. It measures the reduction in uncertainty of one variable due to knowledge of the other.

(Ex. 11.9) Mutual Information

Ex. 11.9 Mutual information is not always superadditive

Let X_1 be a random variable taking values 0 or 1 with respective probabilities of $\frac{1}{2}$, and $X_2 \equiv Y_1 \equiv Y_2 \equiv X_1$. Show that the mutual information in this case is not superadditive.

$$(S) H(X_1, X_2, Y_1, Y_2) \neq H(X_1) + H(X_2; Y_1, Y_2)$$

Ans: 1

X_1	X_2	Y_1	Y_2	P
0	0	0	0	$\frac{1}{2}$
1	1	1	1	$\frac{1}{2}$

$$H(x_1:y_1) = H(x_1) + H(y_1) - H(x_1,y_1)$$

$$= 1+1 - \left(- \sum_{x_1,y_1} p(x_1,y_1) \log p(x_1,y_1) \right)$$

$$= 1+1 + 2 \times \frac{1}{4} \log \frac{1}{4} = 1+1 - 2 \times \frac{1}{4} \times 2$$

$$= 1+1-1 = \underline{\underline{1}}$$

$$H(x_2:y_2) = H(x_2) + H(y_2) - H(x_2,y_2)$$

$$= 1$$

$$H(x_1,x_2:y_1,y_2) = H(x_1,x_2) + H(y_1,y_2) - H(x_1,x_2,y_1,y_2)$$

$$= 1+1-1 = 1$$

$$\therefore H(x_1:y_1) + H(x_2:y_2) = 1+1 = 2 \neq 1 = H(x_1,x_2:y_1,y_2)$$

□ The data processing inequality

- In information theory, the data processing inequality states that information about the off of a source can only decrease with time: once information has been lost, it is gone forever.

The intuitive notion of information processing is captured in the idea of a Markov chain of random variables.

- A Markov chain is a sequence $X_1 \rightarrow X_2 \rightarrow \dots$ of random variables such that X_{n+1} is independent of X_1, \dots, X_{n-1} , given X_n .

$$P(X_{n+1} = x_{n+1} | X_n = x_n, \dots, X_1 = x_1) = P(X_{n+1} = x_{n+1} | X_n = x_n)$$

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

i.e.,

X, Y and Z form a Markov chain $X \rightarrow Y \rightarrow Z$ if the joint probability mass function can be written as.

$$\begin{aligned} p(x, y, z) &= P(x) P(y|x) P(z|xy) \\ &= P(x) P(y|x) P(z|y) \end{aligned}$$

Under what conditions does a Markov chain lose information about its early values, as time progresses?

- Data processing inequality

Suppose $X \rightarrow Y \rightarrow Z$ is a Markov chain.
Then,

$$H(X) \geq H(X:Y) \geq H(X:Z)$$

The 1st inequality is saturated ($H(X) = H(X:Y)$) iff, given Y , it is possible to reconstruct X .

⇒ If a random variable X is subject to noise, producing Y , then further actions on our part (data processing) cannot be used to increase the amount of mutual information bw the output of the process and the original information X .

Proof

$$\text{Part 1: } H(X) \geq H(X:Y)$$

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

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

$$\begin{aligned} 0 \leq p(x|y) \leq 1 &\Rightarrow \log p(x|y) \leq 0 \\ &\Rightarrow -\log p(x|y) \geq 0 \end{aligned}$$

$$\text{Therefore: } H(X|Y) \geq 0$$

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

$$H(X|Y) = H(X) - H(X:Y) \geq 0$$

$$\Rightarrow H(X) \geq H(X:Y)$$

$$\text{Part 2: } H(x:y) \geq H(x:z)$$

$$H(x) - H(x|y) \geq H(x) - H(x|z)$$

$$H(x|y) \leq H(x|z)$$

- If $x \rightarrow y \rightarrow z$ is a Markov chain then $z \rightarrow y \rightarrow x$ is also a Markov chain.

Proof

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

$$\iff \frac{P(x,y,z)}{P(y,z)} = \frac{P(x,y)}{P(y)}$$

$$\iff P(x|y,z) = P(x|y)$$

$$\therefore H(x|y) = H(x|y,z) = H(x,y,z) - H(y,z) \quad \begin{aligned} & H(x,y) = H(x) + H(y|x) \\ & = H(y) + H(x|y) \end{aligned}$$

$$H(x|z) = H(x,z) - H(z)$$

$$\therefore H(x|y) \leq H(x|z) \iff H(x,y,z) - H(y,z) \leq H(x,z) - H(z)$$

$$\iff H(x,y,z) + H(z) \leq H(x,z) + H(y,z)$$

which is the strong subadditivity inequality.

If $H(X:Y) < H(X)$,

only if Z is the attempted reconstruction based on knowledge of Y , then $X \rightarrow Y \rightarrow Z$ must be a Markov chain. \Leftarrow

Data processing inequality : $H(X) \geq H(X:Y) \geq H(X:Z)$

$$\Rightarrow H(X) > H(X:Z)$$

$$\Rightarrow Z \neq X$$

\therefore It is not possible to reconstruct X from Y .

If $H(X:Y) = H(X)$,

$$H(X:Y) = H(X) - H(X|Y) \Rightarrow H(X|Y) = 0$$

$$\begin{aligned} H(X|Y) &= \sum_y p(y) H(X|Y=y) = \sum_y p(y) * -\sum_{\alpha} p(\alpha|y) \log p(\alpha|y) \\ &= -\sum_{x,y} p(y) p(x|y) \log p(x|y) = -\sum_{\alpha,y} p(\alpha,y) \log p(\alpha|y) \end{aligned}$$

$H(X|Y) = 0 \Rightarrow X$ is a deterministic function of Y

i.e., knowing Y determines the value of X .

Whenever $P(X=x, Y=y) > 0$ we have

$$P(X=x | Y=y) = 1$$

\Rightarrow If $Y=y$, then we can infer with certainty that X was equal to x , allowing us to reconstruct X .

$$(X)_H \leq (x)_H \iff$$

$$X \neq h \iff$$

X and x are one of different types

$$(X)_H = \bigcap_{h \in H} X_h$$

$$\text{and } M_H = \bigcup_{h \in H} h = \bigcup_{h \in H} X_h$$

Because $\bigcap_{h \in H} X_h = \bigcap_{h \in H} h = \bigcap_{h \in H} X_h$

Because $\bigcap_{h \in H} h = \bigcap_{h \in H} X_h = \bigcap_{h \in H} X_h$

Because $\bigcap_{h \in H} X_h = \bigcap_{h \in H} h = \bigcap_{h \in H} X_h$

Because $\bigcap_{h \in H} h = \bigcap_{h \in H} X_h = \bigcap_{h \in H} X_h$

Because $\bigcap_{h \in H} X_h = \bigcap_{h \in H} h = \bigcap_{h \in H} X_h$

Because $\bigcap_{h \in H} h = \bigcap_{h \in H} X_h = \bigcap_{h \in H} X_h$

Because $\bigcap_{h \in H} X_h = \bigcap_{h \in H} h = \bigcap_{h \in H} X_h$

Because $\bigcap_{h \in H} h = \bigcap_{h \in H} X_h = \bigcap_{h \in H} X_h$

Because $\bigcap_{h \in H} X_h = \bigcap_{h \in H} h = \bigcap_{h \in H} X_h$

Because $\bigcap_{h \in H} h = \bigcap_{h \in H} X_h = \bigcap_{h \in H} X_h$

If $X \rightarrow Y \rightarrow Z$ is a Markov chain, then
so is $Z \rightarrow Y \rightarrow X$.

Data processing inequality: $H(X) \geq H(X:Y) \geq H(X:Z)$

Data pipelining inequality: $H(Z:Y) \geq H(Z:X)$

⇒ any information Z shares with X must be
information which Z also shares with Y ;
the information is 'pipelined' from X thro'
 Y to Z .

classmate

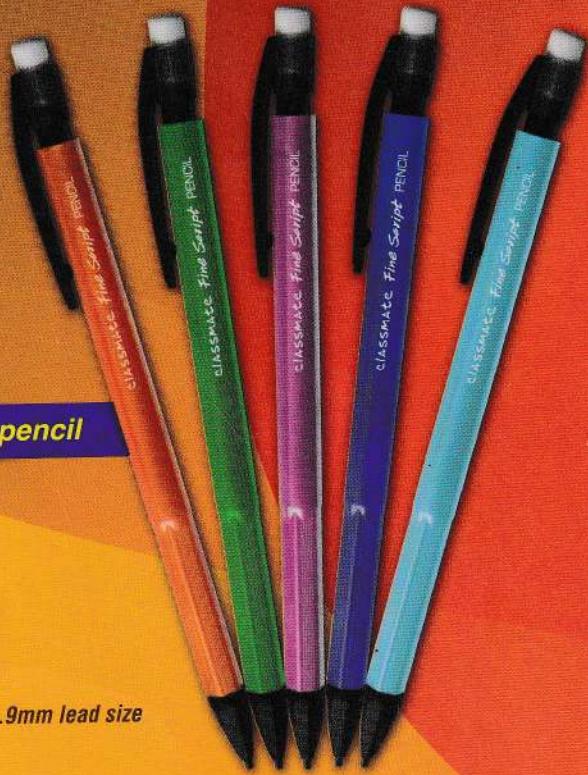
Fine Script

MECHANICAL PENCIL

0.7mm

MRP Rs. 5/- per pencil

Available in 0.7mm & 0.9mm lead size



0.7 mm lead for
precise writing



Strong lead for
uninterrupted writing



Built-in Eraser

At Classmate, INDIA'S NO. 1 NOTEBOOK BRAND*, we are committed to providing high quality stationery products that are a result of a deep understanding of our consumers, thoughtful ideation, innovative designs and superior craftsmanship, that have, in turn, helped us become one of India's leading stationery brands!

Our Classmate range of products include: NOTEBOOKS, Writing Instruments - PENS (ball, gel and roller), PENCILS (mechanical), MATHEMATICAL DRAWING INSTRUMENTS, ERASERS, SHARPENERS and ART STATIONERY (wax crayons, colour pencils, sketch pens and oil pastels).

*Survey conducted by IMRB in Aug, 2020

classmate



Forest Safe, is not just a badge or an icon that this exercise book sports. To us it is a label of merit that ensures this Classmate exercise book of ours makes the world and environment a better place. Wood harvested to make paper and board of this Classmate exercise book is not cut from natural forests, but is ethically and responsibly sourced from sustainably managed plantations, showing our commitment towards building an inclusive and secure future for our stakeholders, and the society at large. Join us in our efforts to create a better future, page by page.



Classmate uses eco-friendly and chlorine free paper



Scan with your smartphone.
Visit us at
www.classmateshop.com

FEEDBACK?
SUGGESTIONS?

Quality Manager, ITC Ltd.- ESPB,
ITC Centre, 5th Floor,
760, Anna Salai, Chennai. 600 002.
classmate@itc.in 18004253242
(Toll-Free from MTNL/BSNL lines)

A quality product marketed by



Exercise Book

172 Pages
(Total Pages Include Index
& Printed Information)

Size : 24 x 18 cm

MRP Rs. 48.00
Inclusive of all taxes

Batch: B/AE FT

02000222



8902519002228

©ITC Limited

Type of Ruling :

Unruled