

19 Aug 2024

These are unpolished notes I wrote for the first lecture in the ICTS program on Quantum Information, Quantum Field Theory and Gravity.

<https://www.icts.res.in/program/qftg>

Fuller accounts of the ideas we discuss can be found in standard texts on information theory, e.g.,

- (a) Thomas M Cover and Joy A Thomas, Elements of Information Theory
- (b) Imre Csiszar and Janos Korner, Information Theory

Some of these topics along with their generalization to the quantum case are also covered in

- (c) Edward Witten, A Mini-Introduction To Information Theory, <https://arxiv.org/abs/1805.11965>

-- Jaikumar

We will discuss $H[X]$ and $H[Y|X]$

Contents

1. The one-shot source coding problem
2. Prefix-free codes
3. Shannon's source coding theorem for one-shot coding
4. Kraft's inequality, Gibbs' inequality
5. Types and typical sequences
6. Shannon's source coding theorem for block codes with small error
7. Conditional entropy

A source or a probabilistic scheme consists of a set of symbols and their corresponding probabilities.

$X = [(a_i, p_i): i = 1, 2, \dots, n]$

The set of symbols $\{a_1, a_2, \dots, a_n\}$ will be called the alphabet of the source. Often we will write $p(a)$ or p_a for the probability of the symbol a . Such source may be used to model an experiment in a lab where there are n distinguishable outcomes, whose occurrence we model using probabilities. The toss of a coin, the roll of a die, the weather, the outcome of an election, indeed, from where the photon would emerge when aimed

at an apparatus with two slits, or when aimed at a beam splitter.

Shannon's source coding problem

Imagine two parties, Alice and Bob (for our imagination is so limited that we cannot think of other names). Alice observes the outcome of the experiment, Bob would like Alice to inform him by means of a message consisting of a string of bits.

Given: A source $X = [(a_i, p_i): i = 1, 2, \dots, n]$

Task: Assign to each a_i a unique string of bits, w_i , called the codeword for a_i , such that

(i) The codewords are prefix-free (a technical condition that allows Bob to know when the message from Alice has ended.

(ii) $\sum_i p_i |w_i|$ is minimum, $|w_i|$ is the number of bits in w_i .

We refer to the above optimum value as the transmission cost for X and denote it by $T[X]$.

The idea is that Alice when she observes a_i will send w_i of length ℓ_i . The expected cost is the objective value of the optimization problem.

Before we study this constrained optimization problem in a little more detail, let us understand the constraint (i).

What is prefix free? No sequence should look like the initial part of another.

$w_1 = 0011, w_2 = 011, w_3 = 10101$ is a prefix-free assignment

by

$w_1 = 0011, w_2 = 011, w_3 = 011101$ is not prefix-free, because w_2 is a prefix of w_3 .

Why prefix-free?

When a string of symbols is to be communicated, we should be able to tell immediately where the codeword for one symbol has ended. In the above example, on receiving 001 , Bob would not be able to tell if the Alice meant to send a_2 or a_3 ; perhaps some bits are yet to arrive. It turns out that the prefix-free property is not such a central assumption; it just makes some of our derivations easier.

Kraft's inequality

If w_1, w_2, \dots, w_n are prefix-free, the $\sum_i 2^{-|w_i|} \leq 1$.

Proof. Draw a binary tree with left edge leaving each node marked 0 and the right edge marked 1. The words w_i appear as leaves in this tree. A random walk from the root passes through w_i with probability $2^{-|w_i|}$. Since the w_i are prefix-free, the events of different w_i are disjoint. The inequality follows from this. (end of proof)

Converse of Kraft's inequality.

If numbers $\ell_1, \ell_2, \dots, \ell_n$ satisfy $\sum_i 2^{-\ell_i} \leq 1$, then there are words w_1, w_2, \dots, w_n such that (i) the w_i are prefix free and (ii) $|w_i| = \ell_i$.

So, we have the following reformulation of the Source Coding problem.

--- Determining $T[X]$, the transmission cost -----

minimize $\sum_i p_i \ell_i$
subject to $\sum_i 2^{-\ell_i} \leq 1$, and ℓ_i are integers.

Consider the above optimization problem without the constraint that the ℓ_i be integers.

Claim: The optimum solution is $\ell_i = \log_2 1/p_i$. (Today, our logs will be to base 2.)

Observe, that if the claim is true, then the optimum value is

$$\sum_i p_i \log 1/p_i,$$

(We will have more to say about this later. :)

Proof: First, the clearly $\ell_i = \log 1/p_i$ satisfies the constraints. Second, consider any solution to the problem: ℓ_1, ℓ_2, \dots , and compare the value of the objective function for it against what we have. The difference is

$$\begin{aligned} & \sum_i p_i \ell_i - \sum_i p_i \log 1/p_i \\ &= \sum_i p_i \log p_i 2^{\ell_i} \\ &\geq - \sum_i p_i \log 2^{-\ell_i}/p_i \\ &\geq \log \sum_i 2^{-\ell_i} \\ &\geq 0. \end{aligned}$$

So, no ℓ_i that satisfies the constraint can do better than what we already have. (end of proof)

It is not hard to see that if we round $\log 1/p_i$ up, then we obtain a feasible integral solution to the problem. We thus have

Theorem I (Shannon): $H[X] \leq T[X] < H[X] + 1$

(Shannon's source coding theorem, with variable length codewords.)

So what is entropy? Why is the entropy of a fair coin toss 1, but the entropy of a $(1/4, 3/4)$ coin toss only 0.811, and a $(1/3, 2/3)$ coin toss 0.919.

In all cases, it is clear that $T[X] = 1$.

QUESTIONS: Then, why is $H[X]$ different? Is there a slightly different engineering problem for which perhaps $H[X]$ provides the correct answer?

ANSWER: Shannon's source coding theory with block encoding.

Typical sequences

Imagine sampling from the source repeatedly and independently k times. (The source is memoryless.)

We get a sequence $\bar{x} = x_1, x_2, \dots, x_k$. We call this source X^k and its distribution is the product distribution, that is,

$$\Pr[X^k = \bar{x}] = \prod_i \Pr[X = x_i]$$

The number of such sequences \bar{x} is clearly n^k . However, most of the probability resides on the typical sequences. Fix \bar{x} . For a in A , let

$$N(a | \bar{x}) = \text{number of occurrences of } a \text{ in } \bar{x} \\ = |\{i: x_i = a\}|$$

Then, $(1/k) N(a | \bar{x})$ is the empirical distribution of a in \bar{x} . We say that \bar{x} is ϵ -typical if this empirical distribution is close to the original distribution:

- $\sum_a |(1/k) N(a | \bar{x}) - p(a)| \leq \epsilon$;
- if $p(a) = 0$, then $N(a | \bar{x}) = 0$.

We refer to the the empirical distribution of \bar{x} as its type. E.g., the type of 0011110000 (6/10, 4/10)

Fact: For all ϵ in $(0,1)$, $\Pr[\bar{x} \text{ is not } \epsilon\text{-typical}]$ goes to 0 as k goes infinity.

Let $T^k(P, \epsilon)$ be the set of ϵ -typical sequences.

[Draw a picture on the board.]

* What is the probability of an ϵ -typical sequences (wrt the produce distribution)?

$$\prod_a p(a)^{N(a | \bar{x})}$$

which is approximately $2^{-n H[X]}$ with some correction for ϵ . More precisely,

$$\begin{aligned} \lim_{k \rightarrow \infty} \max_{\bar{x}: \epsilon\text{-typical}} \log 1/p(\bar{x}) \\ = \lim_{k \rightarrow \infty} \min_{\bar{x}: \epsilon\text{-typical}} \log 1/p(\bar{x}) \\ = H[X] \end{aligned}$$

* How many typical sequences are there?

We know that as k becomes large, essentially all the probability resides on typical sequences, and each typical sequence has probability about $2^{-n H[X]}$.

$$\lim_{k \rightarrow \infty} (1/k) \log |T^k(p, \epsilon)| = H[X].$$

[Alternatively, write this number as

$$n! / \prod_a (n p_a)!$$

and use Stirling's formula.]

Let $S(k, \delta)$ be the smallest cardinality set S of sequences \bar{x} such that

$$P^k(S) = \Pr[X^k \text{ in } S] \geq 1 - \delta$$

Theorem II (Shannon): For all ϵ in $(0, 1)$

$$\lim_k (1/k) \log S(k, \delta) = H[X].$$

BACK TO THE TRANSMISSION PROBLEM

If we had decided to stick with the original encoding of one symbol at a time, and transmitted the entire block of k symbols by concatenating the individual codewords, we would get a cost close to $kT[X]$, that is, about $T[X]$ per source symbol. Indeed, with high probability, the length of the entire message would be close to $k T[X]$. It is not hard to see that $k H[X]$ is a lower bound, thus this method would be a

$$T[X]/H[X] \leq 1 + 1/H[X]$$

worse than the optimum. If $H[X]$ is large, this would perhaps

be acceptable. If $H[X]$ is small (say 0.001), however, the factor $1/H[X]$ can be significant.

In fact, for $(1/4, 3/4)$ or $(1/3, 2/3)$, this method of concatenating the codewords designed for one-shot communication would give us a rate of 1 bit per symbol, whereas the entropy of these distributions is strictly smaller than 1.

Theorem II tells us that if we block together k symbols, and tolerate a vanishingly small amount of error, then we can get by focussing on the set of typical sequences alone. Since almost all the probability resides on typical sequences, the error in transmission can be made vanishingly small as k increases, e.g.,

$$\text{probability of error} \leq \exp(-k^{1/3}).$$

We then transmit only $k H[X]$ bits, that is, we are able to compress the source to $H[X]$ bits per symbol; this is much less than 1 bit per symbol for the examples of biased coin tosses we considered earlier.

(This is what I wanted to say about $H[X]$.)

CONDITIONAL ENTROPY

(X, Y) : random variables with some joint distribution.
 X takes values in A
 Y takes values in B

Notation:

$p(a, b) = \Pr[X = a \text{ and } Y = b]$, the corresponding distribution is P
 $p(a) = \Pr[X = a]$, the corresponding distribution is P_X
 $p(b|a) = \Pr[Y=b \mid X=a]$
 $q(b) = \Pr[Y=b]$

We may define $H[(X, Y)]$ as before. That is,

$$H[(X, Y)] = \sum_{(a, b)} p(a, b) \log 1/p(a, b)$$

Conditional typicality

Suppose all conditional probabilities are given $\{p(b|a)\}$ are given.

Fix a sequence \bar{x} , say its type is $P_X = (p(a): a \text{ in } A)$. Define the joint distribution P on $A \times B$ by

$$p(a, b) = p(a) p(b|a)$$

We would like to know how many \bar{y} are there so that the pair (\bar{x}, \bar{y}) is jointly (P, ϵ) typical. We may count this as follows.

In \bar{x} there are $N(a|\bar{x}) = p(a)^k$ positions where a appears. Let us ask what appears in these position in \bar{y} . If (\bar{x}, \bar{y}) is to be jointly typical (with some small tolerance ϵ), then among these $N(a|\bar{x})$ position the symbols b in B should appear about $p(b|a) N(a|\bar{x})$ times. By our discussion above, entropy tells us how many possibilities we have, except that we must now replace k by $N(a|\bar{x}) = p(a)^k$ and the source by the conditional source $Y | X=a$ (whose distribution is given by number $p(b|a)$). Thus, roughly speaking, the number of extension of \bar{x} to (\bar{x}, \bar{y}) so that the pair is jointly typical is

$$\prod_{a \in A} 2^{\{ p(a)^k H[Y | X=a] \}} \\ = 2^{\{ k \sum_a p(a) H[Y | X=a] \}}$$

So the possibilities grow exponentially with k , but the coefficient k is the quantity

$$\sum_a p(a) H[Y | X=a]$$

We call the conditional entropy of Y given X , and denote it by

$$H[Y|X]$$

Now, the number of jointly typical pairs is $2^{\{k H[(X,Y)]\}}$.

Intuitively, We may think of generating jointly typical pair (\bar{x}, \bar{y}) by first generating \bar{x} for which we have about $2^{\{k H[X]\}}$ choices and then extending it to a jointly typical pair, we have another way of counting jointly typical pairs. Thus we have

$$2^{\{k H[(X,Y)]\}} = 2^{\{k H[X]\}} \times 2^{\{k H[Y|X]\}}$$

or

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

Next, time we will see that this follows more directly from our formulas, but this intuition will be useful for our discussion of Shannon's channel coding theorem. We often write $H[XY]$ when we mean $H[(X,Y)]$, when one is not likely to confuse XY to mean X times Y .

(End of lecture 1.)