

Shannon's Source Coding Theorem

Kim Boström

Institut für Physik, Universität Potsdam, 14469 Potsdam, Germany *

The idea of Shannon's famous source coding theorem [1] is to encode only *typical* messages. Since the typical messages form a *tiny* subset of all possible messages, we need less resources to encode them. We will show that the probability for the occurrence of non-typical strings tends to zero in the limit of large message lengths. Thus we have the paradoxical situation that although we “forget” to encode most messages, we lose no information in the limit of very long strings. In fact, we make use of *redundancy*, i.e. we do not encode “unnecessary” information represented by strings which almost never occur. Recall that a *random message* of length N is a string

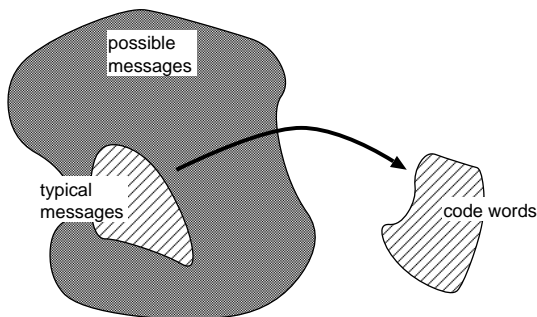


FIG. 1: Lossy coding.

$\mathbf{x} \equiv x_1 \cdots x_N$ of letters, which are independently drawn from an alphabet $\mathcal{A} = \{a_1, \dots, a_K\}$ with *a priori* probabilities

$$p(a_k) = p_k \in (0, 1], \quad k = 1, \dots, K \quad (1)$$

where $\sum_k p_k = 1$. Each given string \mathbf{x} of a random message is an *instance* or *realization* of the message ensemble $\mathbf{X} \equiv X_1 \cdots X_N$, where each random letter X_n is identical to a fixed *letter ensemble* X ,

$$X_n = X, \quad n = 1, \dots, N. \quad (2)$$

A particular message $\mathbf{x} = x_1 \cdots x_N$ appears with the probability

$$p(x_1 \cdots x_n) = p(x_1) \cdots p(x_n), \quad (3)$$

which expresses the fact that the letters are *statistically independent* from each other.

Now consider a very long message \mathbf{x} . Typically, the letter a_k will appear with the frequency $N_k \approx Np_k$. Hence, the probability of such *typical message* is roughly

$$p(\mathbf{x}) \approx p_{typ} \equiv p_1^{N_1} \cdots p_K^{N_K} = \prod_{k=1}^K p_k^{Np_k}. \quad (4)$$

We see that typical messages are *uniformly distributed* by p_{typ} . This indicates that the set T of typical messages has the size

$$|T| \approx \frac{1}{p_{typ}}. \quad (5)$$

If we encode each member of T by a binary string we need

$$I_N = \log |T| = -N \sum_{k=1}^K p_k \log p_k \equiv NH(X), \quad (6)$$

bits, where $H(X)$ is the Shannon entropy of the letter ensemble. Thus for very long messages the average number of bits per letter reads

$$I \equiv \frac{1}{N} I_N = H(X). \quad (7)$$

This is Shannon's source coding theorem in a nutshell. Now let us get a bit more into detail. In order to rigorously prove the theorem we need the concept of a random variable and the law of large numbers. Given the letter ensemble X , the function $f : \mathcal{A} \rightarrow \mathbb{R}$ defines a discrete, real *random variable*. The realizations of $f(X)$ are the real numbers $f(x), x \in \mathcal{A}$. The *average* of $f(X)$ is defined as

$$\langle f(X) \rangle := \sum_{x \in \mathcal{A}} p(x) f(x) = \sum_{k=1}^K p_k f(a_k), \quad (8)$$

and the *variance* is given by

$$\Delta^2 f(X) := \langle f^2(X) \rangle - \langle f(X) \rangle^2. \quad (9)$$

For the sequence $f(\mathbf{X}) \equiv f(X_1), \dots, f(X_N)$ we define its *arithmetic average* as

$$A := \frac{1}{N} \sum_{n=1}^N f(X_n), \quad (10)$$

which is also a random variable. Since the X_n are identical copies of the letter ensemble X , the average of A is equal to the average of $f(X)$,

$$\langle A \rangle = \frac{1}{N} \sum_{i=1}^N \langle f(X_n) \rangle = \langle f(X) \rangle, \quad (11)$$

*Electronic address: bostroem@qipic.org

and the variance of A reads

$$\Delta^2 A = \langle A^2 \rangle - \langle A \rangle^2 \quad (12)$$

$$= \frac{1}{N^2} \sum_{n,m} \langle f(X_n) f(X_m) \rangle - \frac{1}{N^2} \sum_{n,m} \langle f(X_n) \rangle \langle f(X_m) \rangle \quad (13)$$

$$= \frac{1}{N^2} \sum_n \{ \langle f^2(X_n) \rangle - \langle f(X_n) \rangle^2 \} \quad (14)$$

$$= \frac{1}{N} \Delta^2 f(X). \quad (15)$$

The relative standard deviation of A yields

$$\frac{\Delta A}{\langle A \rangle} = \frac{1}{\sqrt{N}} \left(\frac{\Delta f(X)}{\langle f(X) \rangle} \right). \quad (16)$$

Concluding, in the limit of large N the arithmetic average of the sequence $f(\mathbf{X})$ and the ensemble average of $f(X)$ coincide. This is the law of large numbers. It is responsible for the validity of statistical experiments. Without this law, we could never verify statistical properties of a system by performing many experiments. In particular, quantum mechanics would be free of any physical meaning.

Let us reformulate the law of large numbers in the ϵ, δ -language. For $\delta > 0$ we define the *typical set* T of a random sequence \mathbf{X} as the set of realizations $\mathbf{x} \equiv x_1 \cdots x_N$ such that

$$\langle f(X) \rangle - \delta \leq \frac{1}{N} \sum_{n=1}^N f(x_n) \leq \langle f(X) \rangle + \delta. \quad (17)$$

The law of large numbers implies that for every $\epsilon, \delta > 0$ there is a natural number N_0 , such that for all $N > N_0$ the total probability of all typical sequences fulfills

$$P_T \equiv \sum_{\mathbf{x} \in T} p(\mathbf{x}) \geq 1 - \epsilon. \quad (18)$$

The total probability P_T represents the probability for a randomly chosen sequence \mathbf{x} to lie in the typical set T . Now consider the special random variable

$$f(X) := -\log p(X). \quad (19)$$

The average of $f(X)$ equals the Shannon entropy of the ensemble X ,

$$\langle f(X) \rangle = - \sum_{x \in \mathcal{A}} p(x) \log p(x) = H(X). \quad (20)$$

The typical set now contains all messages \mathbf{x} whose probability fulfills

$$H - \delta \leq -\frac{1}{N} \sum_{n=1}^N \log p(x_n) \leq H + \delta, \quad (21)$$

or equivalently

$$2^{-N(H+\delta)} \leq p(\mathbf{x}) \leq 2^{-N(H-\delta)}, \quad (22)$$

where $H \equiv H(X)$. By the law of large numbers, the probability for a randomly drawn message \mathbf{x} to be a member of T reads

$$P_T \equiv \sum_{\mathbf{x} \in T} p(\mathbf{x}) \geq 1 - \epsilon. \quad (23)$$

If we encode only typical sequences, the probability of error

$$P_{err} := 1 - P_T \leq \epsilon \quad (24)$$

can be made arbitrarily small by choosing N large enough. Now let us determine how many typical sequences there are. The lefthand side of (22) gives

$$p(\mathbf{x}) \geq 2^{-N(H+\delta)} \quad (25)$$

$$\Leftrightarrow \sum_{\mathbf{x} \in T} p(\mathbf{x}) \geq |T| 2^{-N(H+\delta)}. \quad (26)$$

The righthand side of (22) gives

$$p(\mathbf{x}) \leq 2^{-N(H-\delta)} \quad (27)$$

$$\Leftrightarrow \sum_{\mathbf{x} \in T} p(\mathbf{x}) \leq |T| 2^{-N(H-\delta)}, \quad (28)$$

which yields together with (23)

$$|T| 2^{-N(H-\delta)} \geq 1 - \epsilon \quad (29)$$

$$\Leftrightarrow |T| \geq (1 - \epsilon) 2^{N(H-\delta)}. \quad (30)$$

Relations (28) and (30) can be combined into the crucial relation

$$(1 - \epsilon) 2^{N(H-\delta)} \leq |T| \leq 2^{N(H+\delta)}. \quad (31)$$

For $N \rightarrow \infty$ we can choose $\epsilon, \delta = 0$ and obtain the desired expression

$$|T| \rightarrow 2^{NH(X)}, \quad (32)$$

thus we need $I_N \rightarrow NH(X)$ bits to encode the message. Equivalently, the information content per letter reads $I = H(X)$ bits. Finally, let us investigate if we can further improve the compression. Relation (30) gives a lower bound for the size of the typical set. Let us compress below H bits per letter by fixing some $\epsilon' > 0$ and encode only sequences that lie in a “subtypical set” $T' \subset T$ whose size reads

$$|T'| \leq (1 - \epsilon) 2^{N(H-\delta-\epsilon')} < 2^{N(H-\delta-\epsilon')}. \quad (33)$$

The righthand side of (22) states that the probability of a typical sequence is bounded from above by

$$p(\mathbf{x}) \leq p_{max} \equiv 2^{-N(H-\delta)}. \quad (34)$$

If we encode only the typical sequences in the subtypical set T' , the probability that a sequence is in T' fulfills

$$P_{T'} = \sum_{\mathbf{x} \in T'} p(\mathbf{x}) \quad (35)$$

$$\leq |T'| \cdot p_{max} = 2^{N(H-\delta-\epsilon')} 2^{-N(H-\delta)} \quad (36)$$

$$= 2^{-N\epsilon'} \quad (37)$$

Because $\epsilon' > 0$, the probability of a successful encoding

goes to 0 for $N \rightarrow \infty$,

$$P_{T'} \rightarrow 0. \quad (38)$$

Concluding, if we compress the messages below $NH(X)$ bits, we are not able to encode all typical messages and for $N \rightarrow \infty$ we will lose all information. A good review on the issue can also be found in [2, 3].

[1] C. E. Shannon and W. Weaver. A mathematical Theory of communication, *The Bell System Technical Journal*, **27**, 379-423, 623-656, (1948).

[2] D.J.C. MacKay. Information theory, inference, and learning algorithms, <http://wol.ra.phy.cam.ac.uk/mackay/itprnn/book.html>,

(1995-2000).

[3] J. Preskill. Lecture notes. <http://www.theory.caltech.edu/people/preskill/ph219/>, (1997-1999).