

---

## TUTORIAL IV

---

### 1 Fixed-length almost lossless compressor: source coding theorem

Recall that a *fixed length compressor* for source  $Y \in \mathcal{Y}$  of length  $\ell$  is a function  $C : \mathcal{Y} \rightarrow \{0, 1\}^\ell$ . It has error probability at most  $\delta$  if there exists a decompressor  $D : \{0, 1\}^\ell \rightarrow \mathcal{Y}$  such that  $\mathbf{P}[D(C(Y)) = Y] \geq 1 - \delta$ . Let define

$$\ell^{opt}(Y, \delta) = \min\{\ell : \text{there exists a length } \ell \text{ compressor for } Y \text{ with error probability } \delta\}.$$

We will prove what is usually called *Shannon Source Coding Theorem*: Let  $X^n = X_1 \dots X_n$  be a sequence of independent and distributed as  $X \in \mathcal{X}$ . For any  $\delta \in (0, 1)$ ,

$$\lim_{n \rightarrow \infty} \frac{\ell^{opt}(X^n, \delta)}{n} = H(X). \quad (1)$$

We will first show the lower bound (also called converse) and then give another proof for the upper bound.

#### 1.1 Converse of Shannon source coding theorem

In this section, we will show that for any  $\varepsilon > 0$  and for large enough  $n$ , we have  $H(X) - 2\varepsilon \leq \ell^{opt}(X^n, \delta)/n$ .

Recall that we showed in class that  $\ell^{opt}(Y, \delta) = \min\{\lceil \log |S| \rceil : \mathbf{P}\{Y \in S\} \geq 1 - \delta\}$ .

1. Prove that we are done if we can show that: for every set  $S \subseteq \mathcal{X}^n$  such that  $\mathbf{P}[X^n \in S] \geq 1 - \delta$  we have  $|S| \geq 2^{n(H(X) - 2\varepsilon)}$ .
2. Now suppose that there is  $S \subseteq \mathcal{X}^n$  such that  $|S| \leq 2^{\ell n}$  for some  $\ell$  and  $\mathbf{P}[X^n \in S] \geq 1 - \delta$ . Prove that

$$\mathbf{P}[X^n \in S] \leq \mathbf{P}\left[-\sum_{i=1}^n \log P_X(X_i) \leq \ell n + \varepsilon n\right] + 2^{-\varepsilon n}.$$

3. Deduce that under the assumption of Question 2,  $\ell \geq H(X) - 2\varepsilon$ , and so the lower bound holds.

#### 1.2 Achievability using random coding

Recall that in order to prove achievability of the source coding theorem, we chose the set  $S$  of correctly encoded symbols to be the set of  $x^n \in \mathcal{X}^n$  such that  $P_{X^n}(x^n) \geq 2^{-n(H(X) - \varepsilon)}$ . We will now show a similar result by choosing the set  $S$  at random. In fact, we start by considering a general source (i.e., not necessarily iid) and derive an upper bound on the probability of error and will give us the desired result in the special case of an iid source.

Our objective is to show that for any source  $X$  and any integer  $l \geq 0$ , there exists a compressor with error probability

$$\delta \leq \mathbf{P}[-\log_2(P_X(X)) > l - \tau] + 2^{-\tau}, \quad \forall \tau > 0. \quad (2)$$

1. Let  $\tau > 0$ ,  $X$  be a random variable and  $C$  be a length  $l$  compressor. Let  $x_0$  be a fixed letter of  $\mathcal{X}$ . Define  $D = \{0, 1\}^l \rightarrow \mathcal{X}$  by

$$D(y) = \begin{cases} x, & \text{if } \exists! x \in \mathcal{X} \text{ s.t. } C(x) = y \text{ and } -\log_2(P_X(x)) \leq l - \tau \\ x_0, & \text{otherwise} \end{cases} \quad (3)$$

Define also

$$J(x, C) = \{x' \in \mathcal{X} : C(x) = C(x'), x \neq x', \text{ and } -\log_2(P_X(x')) \leq l - \tau\}$$

Show that

$$\mathbf{P}[D(C(X)) \neq X] \leq \mathbf{P}[-\log_2(P_X(X)) > l - \tau] + \mathbf{P}[J(X, C) \neq \emptyset]$$

2. Let  $C$  be a random length- $l$  compressor, that is for each  $x \in \mathcal{X}$ ,  $C(x)$  is a random bit string of length  $l$ , with each bit chosen independently and uniformly from  $\{0, 1\}$ . Show that

$$\mathbb{E}_C[\mathbf{P}[J(X, C) \neq \emptyset]] \leq 2^{-\tau}$$

where we compute the mean on the randomness of  $C$  but not on  $X$ .

3. Prove Eq. (2)  
 4. Can you build from the proof a length  $l$  compressor with error  $\delta \leq \mathbf{P}[-\log_2(P_X(X)) > l - \tau] + 2^{-\tau}$ ?  
 5. Use Eq. (2) to give a proof of the upper bound in Shannon source coding theorem (1).

## 2 Code for unknown distribution

Recall that we can build a code  $C$  that achieves an expected length within 1 bit of the lower bound, that is:

$$H(X) \leq \mathbb{E}(|C(X)|) < H(X) + 1$$

This is done either by using Huffman's algorithm or using the following choice of word lengths:  $l_x = \lceil \log \frac{1}{p(x)} \rceil$ , where  $p$  is the distribution of  $X$ . In some cases, we don't know the true distribution  $p$ , but only have an approximation  $q$ , and still want to find a code.

1. Show that if we use the same choice of word lengths:  $l_i = \lceil \log \frac{1}{q_i} \rceil$ , we have:

$$H(p) + D(p||q) \leq \mathbb{E}(|C(X)|) < H(p) + D(p||q) + 1$$

Extra for those who know Huffman's algorithm: What about Huffman's algorithm?

## 3 Entropy of Markov chains

A *Markov chain* is an indexed sequence  $\{X_i\}$  of random variables such that the variable  $X_{n+1}$  only depends on the value of  $X_n$ . In other terms:

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

In the following, we will always assume that the Markov chains are time-independent, i.e. the following holds:

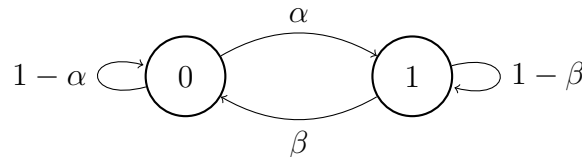
$$\mathbf{P}(X_{n+1} = a | X_n = b) = \mathbf{P}(X_1 = a | X_0 = b)$$

In this case, the evolution of the system depends only on the conditional distribution  $P(X_1 | X_0)$ , and we will usually describe this distribution using a *probability transition matrix*  $P = [P_{ij}]$ , where  $P_{ij} = \mathbf{P}(X_1 = j | X_0 = i)$ . If all the  $X_i$ 's can only take a finite number of values, we usually represent  $X_i$  by its distribution  $p_i = (\mathbf{P}(X_i = 0), \mathbf{P}(X_i = 1), \dots, \mathbf{P}(X_i = l))$ .

Those notations allow us to use the tools of linear algebra, since we can describe the dependency between  $X_{i+1}$  and  $X_i$  using the matrix product:  $p_{i+1} = p_i \cdot P = p_0 \cdot P^i$ . For instance, under reasonable assumptions, we know that  $P^i$  converges to a certain matrix  $P^\infty$ , and that the resulting limit distribution  $p_\infty = p_0 \cdot P^\infty$  is the only fixedpoint of  $P$  (i.e. the only  $p$  such that  $p = p \cdot P$ ).

1. Find the stationary/limit distribution of a two-states Markov chain with a probability transition matrix of the form:

$$\begin{pmatrix} 1 - \alpha & \alpha \\ \beta & 1 - \beta \end{pmatrix}$$



2. In the case of a system with memory, the basic notion of entropy don't capture the dependency between states. Thus, we define another notion of entropy: the *entropy rate* is defined as

$$H(\mathcal{X}) = \lim_{n \rightarrow +\infty} H(X_n | X_{n-1}, \dots, X_0) = \lim_{n \rightarrow +\infty} \frac{1}{n} H(X_1, \dots, X_n)$$

In the case of Markov chain, we thus have:  $H(\mathcal{X}) = \lim_{n \rightarrow +\infty} H(X_n | X_{n-1})$ . If we are in a convergent case, we have:  $H(\mathcal{X}) = H(X_1 | X_0)$ , where the conditional entropy is calculated using the stationary distribution, i.e. with  $X_0 \sim \mu$ .

Compute the entropy rate of the Markov chain of question 1.

3. What is the maximum value of  $H(\mathcal{X})$  in this example?
4. We now take the special case where  $\beta = 1$ . Give a simplified expression of the entropy rate.
5. Find the maximum value of  $H(\mathcal{X})$  in this case. Is it normal that this maximum is achieved for  $\alpha < 1/2$ ?
6. Let  $N(t)$  be the number of allowable state sequences of length  $t$  for the Markov chain (with  $\beta = 1$ ). Find  $N(t)$  and calculate:

$$H_0(\mathcal{X}) = \lim_{t \rightarrow +\infty} \frac{1}{t} H_0(X_0, \dots, X_{t-1}) = \lim_{t \rightarrow +\infty} \frac{1}{t} \log N(t)$$

Why is  $H_0$  an upper bound on the entropy rate of the Markov chain? Compare  $H_0$  with the maximum entropy found in the previous question.