

Chapter 8

Source coding

We are now ready to derive the direct analog of Shannon's lossless source coding theorem, that was first analyzed by Schumacher.

In the classical case we are given a memoryless source which produces strings $x_1x_2\dots x_n$ where each letter $x \in \{0, 1\}$ and occurs with probability $\text{Prob}(X = x) = p_x$. One shows that for n sufficiently large, the outputs of the source can be faithfully described by nR bits as long as $R > H(X)$ and also that this is not possible if $R < H(X)$. Thus length n messages can be compressed to length nR messages with negligible error as $n \rightarrow +\infty$.

In the quantum case a memoryless source produces tensor product states $\rho_{x_1} \otimes \dots \otimes \rho_{x_n}$ each "letter" ρ_x belonging to a finite set \mathcal{A} of $d \times d$ density matrices (the quantum alphabet) and occurring with probability p_x (so $\sum_{x \in \mathcal{A}} p_x = 1$). The quantum state of the source is therefore

$$\sum_{x_1 \dots x_n} p_{x_1} \dots p_{x_n} \rho_{x_1} \dots \rho_{x_n} = \left(\sum_x p_x \rho_x \right)^{\otimes n} = \rho^{\otimes n} \quad (8.1)$$

This is a density matrix of dimension $d^n \times d^n$ (the Hilbert space of pure states has dimension d^n ; for example $d = 2$ for Qbits, $d = 3$ for "quantum trits" etc). We want to compress the source: this means that we want to represent it faithfully by states (or density matrices) of a Hilbert space of dimension d^{nR} .

In general this problem is open. It is known that it is not possible to achieve a compression rate $R < \chi(\{p_x, \rho_x\})$, but it is not known that any rate higher than the Holevo quantity is achievable. However, Schumacher solved the special case where the alphabet letters $\rho_x \in \mathcal{A}$ are pure states $\rho_x = |\phi_x\rangle\langle\phi_x|$. Namely any rate $R > S(\rho)$ is achievable while it is not possible to faithfully compress at rates $R < S(\rho)$. Note that if the alphabet consists of orthonormal states (say $\{|0\rangle, |1\rangle\}$) $S(\rho) = H(X)$ so one recovers the classical Shannon theorem. This should be so, since orthonormal states

are perfectly distinguishable, so that the problem becomes equivalent to the classical one. In the sequel we concentrate on sources of pure states that live in \mathbb{C}^d .

8.1 Notion of typical subspace

In the classical case, the space of length n strings is partitioned into $\mathcal{T}_{n,\epsilon} \cup \mathcal{T}_{n,\epsilon}^c$ where

$$\mathcal{T}_{n,\epsilon} = \left\{ \{x_1 \dots x_n\} \mid \left| \frac{1}{n} \sum_{i=1}^n \log_d \frac{1}{p(x_i)} - H(X) \right| \leq \epsilon \right\} \quad (8.2)$$

is called the space of (weakly) typical sequences. This definition implies that all typical sequences have approximately the same probability, namely

$$d^{-n(H(X)+\epsilon)} \leq p_{x_1} \dots p_{x_n} \leq d^{-n(H(X)-\epsilon)} \quad (8.3)$$

By the law of large numbers, for any ϵ and δ small positive we can find n large enough such that

$$1 - \delta \leq \text{Prob}(\mathcal{T}_{n,\epsilon}) \leq 1 \quad (8.4)$$

Summing (8.3) over typical sequences and using (8.4) we also deduce an estimate on the number of typical sequences

$$(1 - \delta) d^{n(H(X)-\epsilon)} \leq |\mathcal{T}_{n,\epsilon}| \leq d^{n(H(X)+\epsilon)} \quad (8.5)$$

Finally, any set $S_{n,\epsilon}$ of sequences that is too small in the sense that $|S_{n,\epsilon}| \leq d^{nR}$ with $R \leq H(X) - \epsilon$ has negligible probability,

$$\text{Prob}(S_{n,\epsilon}) \leq \delta + d^{-n(H(X)-\epsilon-R)} \quad (8.6)$$

To see this write

$$S_{n,\epsilon} = (S_{n,\epsilon} \cap \mathcal{T}_{n,\epsilon}) \cup (S_{n,\epsilon} \cap \mathcal{T}_{n,\epsilon}^c) \quad (8.7)$$

and use (8.3) with the union bound.

These properties immediately suggest to encode only the typical sequences and to throw away or code non-typical ones into a junk state. Because of (8.4) this scheme will incur a decoding error with probability at most δ . Because of (8.5) it is enough to use $n(H(X) + \epsilon)$ nats for the encoding. Moreover because of (8.6) using less than $n(H(X) - \epsilon)$ nats will incur a finite probability of error¹.

¹See Cover and Thomas for more details

In the quantum case one defines a similar notion of typicality. Consider a memoryless source that outputs with probability p_x letters $|\phi\rangle_x \in \mathcal{A}$ which belong to the Hilbert space $\mathcal{H} = \mathbf{C}^d$. The density matrix for the source is

$$\sum_{x_1 \dots x_n} p_{x_1} \dots p_{x_n} |\phi_{x_1}\rangle \langle \phi_{x_1}| \otimes \dots \otimes |\phi_{x_n}\rangle \langle \phi_{x_n}| = \left(\sum_x p_x |\phi_x\rangle \langle \phi_x| \right)^{\otimes n} = \rho^{\otimes n} \quad (8.8)$$

One can find the spectral decomposition of this density matrix. Indeed let

$$\rho = \sum_a \lambda_a P_a \quad (8.9)$$

be the spectral decomposition for the length one case (here we assume for simplicity non-degeneracy of the eigenvalues, so $P_a = |a\rangle \langle a|$). Then

$$\rho^{\otimes n} = \sum_{a_1 \dots a_n} \lambda_{a_1} \dots \lambda_{a_n} P_{a_1} \otimes \dots \otimes P_{a_n} \quad (8.10)$$

The eigenvalues λ_a are positive and sum to one, thus define a probability distribution. Moreover the projectors P_a are mutually orthogonal, thus distinguishable. Therefore the density matrix $\rho^{\otimes n}$ is also the density matrix of a "mathematical" *memoryless classical source* that outputs letters a (or P_a or $|a\rangle$) with probabilities p_a . We stress that this is not the physical preparation of the state $\rho^{\otimes n}$. We can define a set of typical sequences of eigenvalues and/or eigenstates

$$\mathcal{T}_{n,\epsilon} = \left\{ a_1 \dots a_n \mid \left| \frac{1}{n} \sum_{i=1}^n \log_d \frac{1}{\lambda_{a_i}} - S(\rho) \right| \leq \epsilon \right\} \quad (8.11)$$

Definition: typical subspace. Consider the projector

$$P_{n,\epsilon} = \sum_{a_1 \dots a_n \in \mathcal{T}_{n,\epsilon}} P_{a_1} \otimes \dots \otimes P_{a_n} \quad (8.12)$$

The subspace $P_{n,\epsilon} \mathcal{H}^{\otimes n}$ is called the typical subspace. We have

$$\rho^{\otimes n} = P_{n,\epsilon} \rho^{\otimes n} P_{n,\epsilon} + (I - P_{n,\epsilon}) \rho^{\otimes n} (I - P_{n,\epsilon}) \quad (8.13)$$

The source coding scheme described in the next section is based on the following theorem, which is the quantum analog of (8.3), (8.5) and (8.6).

Theorem 1. [typical subspace theorem] Fix ϵ and δ positive, small. For n sufficiently large,

- the density matrix has almost all its support on the typical subspace

$$1 - \delta \leq \text{Tr} P_{n,\epsilon} \rho^{\otimes n} \leq 1, \quad (8.14)$$

- the dimension of the typical subspace is approximately $d^{nS(\rho)}$

$$(1 - \delta) d^{n(S(\rho) - \epsilon)} \leq \text{Tr} P_{n,\epsilon} \leq d^{n(S(\rho) + \epsilon)}, \quad (8.15)$$

- let $S_{n,\epsilon}$ be a projector on a subspace of dimension less than d^{nR} with $R \leq S(\rho) - \epsilon$. In other words $\text{Tr} S_{n,\epsilon} \leq d^{nR}$ with $R \leq S(\rho) - \epsilon$. For such a projector we have

$$\text{Tr} S_{n,\epsilon} \rho^{\otimes n} \leq \delta + d^{-n(S(\rho) - \epsilon - R)}. \quad (8.16)$$

Proof. The basic difference with the classical case is that one has to deal a bit more carefully with operator inequalities for the third statement²

First statement. Observe that

$$P_{n,\epsilon} \rho^{\otimes n} P_{n,\epsilon} = \sum_{a_1 \dots a_n \in \mathcal{T}_{n,\epsilon}} \lambda_{a_1} \dots \lambda_{a_n} P_{a_1} \otimes \dots \otimes P_{a_n} \quad (8.17)$$

So

$$\text{Tr} P_{n,\epsilon} \rho^{\otimes n} = \sum_{a_1 \dots a_n \in \mathcal{T}_{n,\epsilon}} \lambda_{a_1} \dots \lambda_{a_n} \quad (8.18)$$

which is the probability of the set $\mathcal{T}_{n,\epsilon}$. The statement follows by the law of large numbers (as in the classical case).

Second statement. Observe that

$$\text{Tr} P_{n,\epsilon} = \sum_{a_1 \dots a_n \in \mathcal{T}_{n,\epsilon}} 1 = |\mathcal{T}_{n,\epsilon}| \quad (8.19)$$

so the statement again follows like in the classical case. Note that here we have assumed that the eigenvalues are not degenerate (if $\text{Tr} P_a = g_a$ we have to modify the definition of typical sequences according to $\frac{1}{\lambda_a} \rightarrow \frac{g_a}{\lambda_a}$).

Third statement. We use the decomposition (8.13) to write $\text{Tr} S_{n,\epsilon} \rho^{\otimes n}$ as a sum of two contributions.

²We recall: a hermitian matrix $A = A^\dagger$ is said to be (semi-definite) positive iff $\langle \phi | A | \phi \rangle \geq 0$ for any $|\phi\rangle$; $A \geq B$ iff $(A - B) \geq 0$; and $A \geq 0$ implies $C^\dagger A C \geq 0$.

For the first one we have,

$$\begin{aligned}
\text{Tr} S_{n,\epsilon} P_{n,\epsilon} \rho^{\otimes n} P_{n,\epsilon} &= \text{Tr} S_{n,\epsilon} P_{n,\epsilon} \rho^{\otimes n} P_{n,\epsilon} S_{n,\epsilon} \\
&\leq d^{-n(S(\rho)-\epsilon)} \text{Tr} S_{n,\epsilon} \\
&\leq d^{-n(S(\rho)-\epsilon-R)}
\end{aligned} \tag{8.20}$$

In the first equality we use the cyclicity of the trace and for the first inequality we use the operator inequality $P_{n,\epsilon} \rho^{\otimes n} P_{n,\epsilon} \leq d^{-n(S(\rho)-\epsilon)} I$.

For the second contribution we observe that $M = (I - P_{n,\epsilon}) \rho^{\otimes n} (I - P_{n,\epsilon})$ is a positive operator so by the cyclicity of the trace

$$\begin{aligned}
\text{Tr} S_{n,\epsilon} (I - P_{n,\epsilon}) \rho^{\otimes n} (I - P_{n,\epsilon}) &= \text{Tr} \sqrt{M} S_{n,\epsilon} \sqrt{M} \\
&\leq \text{Tr} \sqrt{M} I \sqrt{M} = \text{Tr} M \\
&= \text{Tr} \rho^{\otimes n} (I - P_{n,\epsilon}) \\
&\leq \delta
\end{aligned} \tag{8.21}$$

In the inequality we used that $S_{n,\epsilon} \leq I$ (true for any projector). \square

8.2 Source coding scheme

The source outputs words of length n ,

$$|\phi_{x_1}\rangle \otimes \dots \otimes |\phi_{x_n}\rangle \tag{8.22}$$

with probability $p_{x_1} \dots p_{x_n}$. We specify a block coding scheme for these words: we would like to encode these words which belong to the Hilbert space $\mathcal{H}^{\otimes n}$ by states in a Hilbert space $\mathcal{H}^{\otimes nR}$ where $R < 1$. This encoding should be faithful in the sense that that it should be possible to recover, most of the time, the original words by some decoding procedure.

Encoding procedure. We have to rely on a slightly more general encoding process that encodes source states into density matrices. An encoding map is a map from states of dimension d^n to density matrices of dimension $d^{nR} \times d^{nR}$

$$\begin{aligned}
\mathcal{E}_n : \mathcal{H}^{\otimes n} &\rightarrow \mathcal{DM}(\mathcal{H}^{\otimes nR}) \\
|\phi_{x_1}\rangle \otimes \dots \otimes |\phi_{x_n}\rangle &\rightarrow \mathcal{E}(|\phi_{x_1}\rangle \otimes \dots \otimes |\phi_{x_n}\rangle)
\end{aligned}$$

Here $\mathcal{DM}(\mathcal{H}^{\otimes nR})$ is the space of density matrices of dimension $d^{nR} \times d^{nR}$. The compression rate per letter is $R = \frac{nR}{n}$.

Decoding procedure. Ideally we should map back the density matrix $\mathcal{E}(|\phi_{x_1}\rangle \otimes \dots \otimes |\phi_{x_n}\rangle)$ to the input word $|\phi_{x_1}\rangle \otimes \dots \otimes |\phi_{x_n}\rangle$. This cannot be done exactly, so we allow for a slightly more general definition,

$$\begin{aligned} \mathcal{D}_n : \mathcal{DM}(\mathcal{H}^{\otimes nR}) &\rightarrow \mathcal{DM}(\mathcal{H}^{\otimes n}) \\ \sigma &\rightarrow \mathcal{D}(\sigma) \end{aligned}$$

Reliability criterion. The scheme $(\mathcal{E}_n, \mathcal{D}_n)$ should be faithful. Let

$$\rho_{\text{output}} = \mathcal{D}(\mathcal{E}(|\phi_{x_1}\rangle \otimes \dots \otimes |\phi_{x_n}\rangle)) \quad (8.23)$$

We define a *fidelity* as the overlap between the input and output states,

$$F(|\phi_{x_1}\rangle \otimes \dots \otimes |\phi_{x_n}\rangle) = \langle \phi_{x_1}, \dots, \phi_{x_n} | \rho_{\text{output}} | \phi_{x_1}, \dots, \phi_{x_n} \rangle \quad (8.24)$$

The average fidelity is

$$\bar{F}_n = \sum_{x_1, \dots, x_n} p_{x_1} \dots p_{x_n} \langle \phi_{x_1}, \dots, \phi_{x_n} | \rho_{\text{output}} | \phi_{x_1}, \dots, \phi_{x_n} \rangle \quad (8.25)$$

One has to be careful with the notation here: in each term of the sum ρ_{output} depends on the input state $|\phi_{x_1}, \dots, \phi_{x_n}\rangle$.

The intuitive meaning of the fidelity can be better understood by looking at the classical case. If the letters $|\phi_x\rangle$ are orthonormal then we are reduced to a classical situation and the encoding-decoding operations can be done by looking at a “look-up table”. In other words for a typical source word we have perfect recovery so $\rho_{\text{output}} = |\phi_{x_1}, \dots, \phi_{x_n}\rangle \langle \phi_{x_1}, \dots, \phi_{x_n}|$ and $F = 1$; while for a non typical source word we have $\rho_{\text{output}} = |\text{junk}\rangle \langle \text{junk}|$ and the decoder simply declares an error and sets $F = 0$ (simply assume that $|\text{junk}\rangle$ is orthogonal to all source words). We see that in the classical case the average fidelity is precisely equal to $1 - \text{Prob}(\text{error})$.

Theorem 2. [Schumacher’s theorem] Fix $\epsilon > 0$, $\delta > 0$ small.

- Fix $R > S(\rho) + \epsilon$. Then one can find encoding-decoding schemes $(\mathcal{E}_n, \mathcal{D}_n)$ such that for n large enough $\bar{F}_n \geq 1 - 2\delta$. So asymptotically loss-less compression is possible.
- Fix $R < S(\rho) - \epsilon$. Then for any encoding-decoding scheme $(\mathcal{E}_n, \mathcal{D}_n)$ we have $\bar{F}_n \leq \delta + d^{-n(S(\rho) - \epsilon - R)}$. Loss-less compression is not possible.

This theorem says that compression rates above $S(\rho)$ are (faithfully) achievable, while this is not the case for compression rates below $S(\rho)$. Note

that in general $S(\rho) \leq H(X)$. The fact that a quantum source is more compressible than a classical one should not surprise the reader: this is an expression of the fact that non-orthogonal alphabet letters cannot be perfectly distinguished so that a quantum source word is more redundant than its classical counterpart.

8.3 Proof of theorem 2.

First we prove the achievability part and then proceed to the converse.

Achievability part. We specify the encoding map \mathcal{E} . Take the measurement apparatus defined by the two orthogonal projectors $\{P_{n,\epsilon}, I - P_{n,\epsilon}\}$ on the typical subspace and its orthogonal complement. Given a source word $|\phi_{x_1}, \dots, \phi_{x_n}\rangle$ perform a measurement. According to the measurement postulate the outcome is

$$\frac{P_{n,\epsilon}|\phi_{x_1}, \dots, \phi_{x_n}\rangle}{\langle\phi_{x_1}, \dots, \phi_{x_n}|\phi_{x_1}, \dots, \phi_{x_n}\rangle^{1/2}}, \text{ with prob } \langle\phi_{x_1}, \dots, \phi_{x_n}|P_{n,\epsilon}|\phi_{x_1}, \dots, \phi_{x_n}\rangle \quad (8.26)$$

or

$$\frac{(I - P_{n,\epsilon})|\phi_{x_1}, \dots, \phi_{x_n}\rangle}{\langle\phi_{x_1}, \dots, \phi_{x_n}|I - P_{n,\epsilon}|\phi_{x_1}, \dots, \phi_{x_n}\rangle^{1/2}}, \text{ with prob } \langle\phi_{x_1}, \dots, \phi_{x_n}|I - P_{n,\epsilon}|\phi_{x_1}, \dots, \phi_{x_n}\rangle \quad (8.27)$$

Now the first state is in the typical subspace $P_{n,\epsilon}\mathcal{H}^{\otimes n}$ so it can be described by $nS(\rho)$ quantum nats (because of theorem 1 the dimension of the typical subspace is $d^{nS(\rho)}$). One can find a basis of $\mathcal{H}^{\otimes n}$ such that this typical subspace is described by the first $nS(\rho)$ terms of the tensor product. In other words we can find a unitary operation U that transforms the state (8.26) to the form (this unitary depends only on the original basis and the typical space, not on the particular input state)

$$\sum_{b_1 \dots b_m} c_{x_1 \dots x_n}^{b_1 \dots b_n} | \underbrace{b_1 \dots b_m}_{nR \text{ terms}}, \underbrace{0, 0, \dots, 0}_{n(1-R) \text{ terms}} \rangle = |\psi_{\text{compressed}}\rangle \otimes | \underbrace{0, 0, \dots, 0}_{n(1-R) \text{ terms}} \rangle \quad (8.28)$$

The state $|0_{n(1-R)}\rangle$ is then discarded. The second possible outcome is not coded since it lies in the non typical subspace. More precisely we describe all such states as $|\text{junk}\rangle$ a single specified quantum state (in the typical subspace, say). We assume that the outcome is not observed during the compression stage so its state is described by the mixture

$$\mathcal{E}(|\phi_{x_1} \dots \phi_{x_n}\rangle) = \langle\phi_{x_1}, \dots, \phi_{x_n}|P_{n,\epsilon}|\phi_{x_1}, \dots, \phi_{x_n}\rangle |\psi_{\text{compressed}}\rangle \langle\psi_{\text{compressed}}| + \langle\phi_{x_1}, \dots, \phi_{x_n}|I - P_{n,\epsilon}|\phi_{x_1}, \dots, \phi_{x_n}\rangle |\text{junk}\rangle \langle\text{junk}| \quad (8.29)$$

For the decoding operation one first appends $n(1-R)$ quantum letters in the $|0_{n(1-R)}\rangle$ state that was discarded. then one performs the inverse unitary operation U^\dagger . So the decoder map is given by

$$\begin{aligned} \mathcal{D}(\mathcal{E}(|\phi_{x_1}\dots\phi_{x_n}\rangle)) & \quad (8.30) \\ &= P_{n,\epsilon}|\phi_{x_1}\dots\phi_{x_n}\rangle\langle\phi_{x_1},\dots,\phi_{x_n}|P_{n,\epsilon} \\ &+ \langle\phi_{x_1},\dots,\phi_{x_n}|I - P_{n,\epsilon}|\phi_{x_1},\dots,\phi_{x_n}\rangle U^\dagger |\text{junk}, 0_{n(1-R)}\rangle\langle\text{junk}, 0_{n(1-R)}| U \end{aligned}$$

We now estimate the fidelity associated to this scheme $(\mathcal{E}_n, \mathcal{D}_n)$. We replace ρ_{output} given by (8.30) in the definition of the average fidelity. The contribution from the first term is

$$\begin{aligned} & \sum_{x_1\dots x_n} p_{x_1}\dots p_{x_n} \langle\phi_{x_1},\dots,\phi_{x_n}|P_{n,\epsilon}|\phi_{x_1}\dots\phi_{x_n}\rangle^2 \quad (8.31) \\ & \geq \left\{ \sum_{x_1\dots x_n} p_{x_1}\dots p_{x_n} \langle\phi_{x_1},\dots,\phi_{x_n}|P_{n,\epsilon}|\phi_{x_1}\dots\phi_{x_n}\rangle \right\}^2 \\ & = (\text{Tr} \rho^{\otimes n} P_{n,\epsilon})^2 \\ & \geq (1 - \delta)^2 \end{aligned}$$

The first inequality is Cauchy-Schwartz, and the second comes from theorem 1. Finally the contribution from the second term is trivially positive (write it down and see!). Thus we conclude that

$$\bar{F} \geq (1 - \delta)^2 \geq 1 - 2\delta \quad (8.32)$$

Converse part. Let

$$\mathcal{E}_N : |\phi_{x_1}\dots\phi_{x_n}\rangle\langle\phi_{x_1}\dots\phi_{x_n}| \rightarrow \sigma \quad (8.33)$$

be a completely general encoding scheme (so σ is any $d^{nR} \times d^{nR}$ density matrix). The first step of the decoder is to append $|0_{n(1-R)}\rangle\langle 0_{n(1-R)}|$ to get a state

$$\sigma \otimes |0_{n(1-R)}\rangle\langle 0_{n(1-R)}| \quad (8.34)$$

in the original Hilbert space. Here we restrict the proof to the special case of *unitary decoders*³. So let

$$\mathcal{D} : \sigma \otimes |0_{nR}\rangle\langle 0_{nR}| \rightarrow U\sigma \otimes |0_{nR}\rangle\langle 0_{nR}| U^\dagger \quad (8.35)$$

³More general ones would correspond to a mappings between density matrices and would require a more complicated proof.

The density matrix (8.34) is constructed out of states of a d^{nR} dimensional subspace of $\mathcal{H}^{\otimes n}$. Let S_n be the projector on that subspace, and note that

$$U\sigma \otimes |0_{nR}\rangle\langle 0_{nR}|U^\dagger = US_n \left(\sigma \otimes |0_{nR}\rangle\langle 0_{nR}| \right) S_n U^\dagger \quad (8.36)$$

Now, the average fidelity is

$$\begin{aligned} \bar{F} &= \sum_{x_1 \dots x_n} p_{x_1} \dots p_{x_n} \langle \phi_{x_1} \dots \phi_{x_n} | US_n \left(\sigma \otimes |0_{nR}\rangle\langle 0_{nR}| \right) S_n U^\dagger | \phi_{x_1} \dots \phi_{x_n} \rangle \quad (8.37) \\ &\leq \sum_{x_1 \dots x_n} p_{x_1} \dots p_{x_n} \langle \phi_{x_1} \dots \phi_{x_n} | US_n U^\dagger | \phi_{x_1} \dots \phi_{x_n} \rangle \\ &= \text{Tr}(\rho^{\otimes n} US_n U^\dagger) \end{aligned}$$

We first used that any density matrix is smaller than the identity matrix, so $\sigma \otimes |0_{nR}\rangle\langle 0_{nR}| \leq I$, and then the cyclicity of the trace. Clearly $US_n U^\dagger$ is a projector on some d^{nR} dimensional subspace of $\mathcal{H}^{\otimes n}$ with $R < S(\rho) - \epsilon$. Then, the third statement of theorem 1 implies

$$\bar{F} \leq \delta + d^{-n(S(\rho) - \epsilon - R)} \quad (8.38)$$

This achieves the proof of the converse part for the class of unitary decoders.