

Theoretical physics of information: Lecture notes

Paolo Perinotti

November 21, 2022

Contents

Is information physical? Is physics informational?	vii
I Classical information theory	1
1 Lecture 1: Introduction	3
1.1 Basic ideas of classical information theory	3
1.1.1 Data compression	3
1.1.2 Error-correction	5
2 Lecture 2: Operational Probabilistic Theories	15
2.1 Operational theories	15
2.2 Coarse graining	18
2.3 Probabilistic theories	18
2.4 States and effects	18
2.5 Transformations	19
2.6 Classical information theory	20
3 Lecture 3: Classical systems and Random variables	23
3.1 Random variables	23
3.1.1 Inequalities and probability bounds	26
3.1.2 Convex functions	29
4 Lecture 4: Shannon entropy and typical sequences	33
4.1 Shannon entropy	33
4.1.1 Shannon information content and Shannon entropy	33
4.1.2 Relative entropy	36
4.2 Bound on Shannon entropy	38
4.2.1 Raw bit content	38
4.3 Lossy and lossless compressors	38
4.3.1 Smallest δ -sufficient set and essential bit content	39
4.4 Source coding	41

4.4.1	Typical sets and asymptotic equipartition	41
4.4.2	Shannon source coding theorem	43
5	Lecture 5: Symbol codes and stream codes	47
5.1	Symbol codes	47
5.2	Kraft-McMillan theorem	50
5.3	Source coding theorem for symbol codes	51
6	Lecture 6: Huffman coding and stream codes	55
6.1	Optimal code: Huffman's coding	55
6.2	Stream codes and redundancy	59
6.3	Arithmetic code	61
6.4	Lempel-Ziv algorithm	65
7	Lecture 7: Mutual information and Markov chains	69
7.1	Entropies for bi- and multi-variate random variables	70
7.1.1	Joint entropy and conditional entropy	70
7.1.2	Mutual information	73
7.2	Markov chains	75
7.3	Data Processing theorem	77
8	Lecture 8: Channel capacity and joint typicality	79
8.1	Fano's inequality	79
8.2	Discrete memoryless channels	81
8.3	Capacity of a channel	85
8.4	Block codes	89
8.5	Error probabilities	90
8.6	Optimal decoder	91
8.7	Joint typicality	92
9	Lecture 9: Shannon noisy channel coding theorem	97
9.1	Random coding and typical-set decoding	97
9.2	The second theorem of Shannon's	100
9.3	Communication above capacity [extra]	105
II	Quantum information theory	109
10	Lecture 10: Elements of linear algebra	111
10.1	Hilbert spaces	111
10.1.1	Subspaces	112
10.2	Linear operators on complex Hilbert spaces	112
10.2.1	Unitary operators	115
10.2.2	Selfadjoint operators	117
10.2.3	Normal operators and the spectral theorem	120
10.3	Isometric operators	126

10.4 Polar decomposition	127
10.4.1 Singular value decomposition	130
10.5 Tensor product	130
10.5.1 Linear operators	131
10.6 The partial trace	133
10.7 The Vec isomorphism: double-ket notation	135
10.8 Linear maps on $\mathcal{L}(\mathcal{H})$	138
10.8.1 Positive maps	139
10.9 Maps on tensor products	139
10.10 Complete positivity	141
10.11 Trace-preserving and unital maps	142
11 Lecture 11: Quantum Theory and the Choi isomorphism	145
11.1 Quantum theory	145
11.2 States	145
11.2.1 Quantum operations and channels	146
11.2.2 Effects	147
11.3 The qubit	147
11.4 Purification	148
11.5 The Choi correspondence	150
11.5.1 Effects	155
11.6 Unitary dilation of channels	155
12 Lecture 13: von Neumann entropy	157
12.1 Quantum transmission of classical information	157
12.2 Von Neumann entropy	158
12.2.1 Klein's inequality	161
12.2.2 Preliminary lemmas	163
12.3 Mathematical properties of the von Neumann entropy	164
13 Lecture 15: Properties of quantum relative entropy and Lieb's theorem	171
13.1 Mathematical properties of the quantum relative entropy	171
13.2 Preliminary results for Lieb's theorem	172
13.3 Lieb's theorem	174
14 Lecture 16: Monotonicity and Holevo bound	179
14.1 Joint convexity of quantum relative entropy	180
14.2 Concavity of quantum conditional entropy	181
14.3 Strong subadditivity of von Neumann entropy	182
14.4 Uhlmann monotonicity theorem	184
14.5 The Holevo bound	186
14.6 Holevo-Schumacher-Westmoreland theorem	189
15 Lecture 17: Quantum information, compression, and Uhlmann fidelity	193
15.1 Quantum compression	194

15.2 Fidelity	196
16 Lecture 19: Schumacher quantum source coding theorem	203
16.1 Entanglement fidelity	203
16.2 Refinement set	205
16.3 Equality upon input	208
16.4 Schumacher's quantum source coding theorem	210
17 Lecture 20: Entropy exchange and coherent information	215
17.1 Entropy exchange	215
17.1.1 Quantum data-processing theorem	215
17.1.2 Quantum Fano inequality	218
17.2 Reversibility upon input	219
17.3 Coherent information and quantum data processing theorem	221
18 Lecture 21: The theory of entanglement	225
18.1 Entanglement	225
18.2 What does LOCC mean, precisely?	226
18.2.1 Entanglement and LOCC	228
18.2.2 Resource theories	228
18.2.3 Maximally entangled states	229
18.3 Entanglement cost and distillable entanglement	232
18.4 Criteria for entanglement	233
18.4.1 The PPT criterion - Positive partial transpose	233
18.4.2 Entanglement witnesses	234
18.4.3 Bound entanglement	235
18.4.4 Entanglement of formation	236
Bibliography	239

Is information physical? Is physics informational?

Is information physical? We may be led to think so, based on the fact that every piece of information we receive about the outside world is carried by some physical system. If we start analysing the properties of information from this perspective we realise that the fundamental information processing tasks are ruled by fundamental physical laws, thus discovering fascinating results. An outstanding example is *Landauer's principle* stating that erasure of one bit of information always requires dissipation of an amount of energy equal to $KT \ln 2$, T being the temperature of the memory storing the bit.

However, if we think more deeply, we realise the way in which we become aware of physical processes is by gathering information contained in events: Physical systems and their evolution laws are a clever architecture of our mind for organising and explaining events, but they are a secondary notion, while primitive facts are events. The perception of their very reality is only a powerful belief that stands out of doubt only by virtue of its astonishing predictive power. Thus, we could reverse the starting question and ask: Is physics informational? In present-day physics, the number of scientists that started exploring the program summarised by the motto “it from bit”—coined by John A. Wheeler—is increasingly large, and this line of research is counting more and more successful results.

The aim of the course is not to answer this fundamental question, but to provide some operational and quantitative substance to the notion of information, which is one step in the long journey that takes to explore the relation between information and physics.

Messages, which are pieces of information, are encoded in some alphabet. However, as physicists we cannot accept this as the definition of a quantity, unless we provide an operational definition, which also requires a way to quantify it. The course will illustrate the basic operational contexts in which information can be quantified, along with the precise mathematical measures accomplishing this task.

The course is structured in three parts:

1. Classical information theory, where information is supposed to be carried by classical systems.
2. Quantum information theory, where information is supposed to be carried by quantum systems. This part in turn is split into two parts:

viii Is information physical? Is physics informational?

- a) Encoding of classical information over quantum carriers.
- b) Encoding of quantum information.

The main tasks that allow us to define precise notions of information are *compression* and *transmission over noisy channels*. While in the classical case, under rather general assumptions, both problems have been exhaustively analysed in the early work by Shannon [2], in the quantum case only compression has been thoroughly understood, while transmission is the really challenging problem, the one where the unintuitive feature of *entanglement* plays a crucial role, which is not yet completely tamed.

Part I

Classical information theory

Chapter 1

Lecture 1: Introduction

1.1 Basic ideas of classical information theory

1.1.1 Data compression

Our first question regards the amount of resources needed to *store* or *transmit* a given piece of information. This question is very vague at this stage, because we need to define what are the resources and what are the requirements of our storage. In the first place, let us set the constraint that information must be completely recoverable, without losses or errors. The question thus regards the minimal support that grants such a perfect storage or transmission. We imagine here to face the task of storage, and the resource will be memory cells.

Let us use a language L written in an alphabet Σ with four symbols, say $\Sigma = \{a, b, c, d\}$. Suppose we want to store or transmit a file written in the language L . For example, let us consider the string

$$bacadaba$$

Storing this message costs 8 memory cells, each with 4 levels. Suppose now that we use a different language L' whose alphabet X has only two symbols $X = \{0, 1\}$. If we choose the encoding

$$a \mapsto 00, \quad b \mapsto 01, \quad c \mapsto 10, \quad d \mapsto 11, \tag{1.1}$$

we can write the same string as

$$0100100011000100$$

which is longer: For the same message we now need 16 cells, each with 2 levels.

What alphabet do we choose, then? It is clear that the answer requires to establish one basic fact: whether the cost of a cell with $m \times n$ levels is higher or smaller than m times the cost of a cell with n levels. The basic assumption of information theory is that the two costs are *identical*. In this way, the issue of the choice of alphabet is solved by assuming that counting resources for storage is equivalent to counting the sheer number of levels used to store a given message. In other words, the choice of alphabet is not relevant for the purpose of counting storage resources. We can then choose one alphabet once for all, and it will be the most elementary one: $X = \{0, 1\}$. The cell for such an

alphabet is a *bit*, allocating one character which can be either 0 or 1. More formally, a bit is a variable that takes values in the alphabet X . Before analysing any given source using a different alphabet, we will then find a way to encode messages into bit strings.

Let us now come back to our original file, written in the alphabet $\Sigma = \{a, b, c, d\}$. Let us suppose that such symbols have probability distribution

$$p(a) = \frac{1}{2}, \quad p(b) = \frac{1}{4}, \quad p(c) = \frac{1}{8}, \quad p(d) = \frac{1}{8}.$$

This is the case, in particular, for our message *bacabadabba*. Can we exploit our prior information about probabilities of characters to produce a better encoding? For example, we know that the character *a* will occur more frequently than the others, thus we choose to encode it in a short string: $a \mapsto 0$. If we now encode the second most frequent character *b* into 1, we will end up with an un-decodable storage: Every encoded word will be a sequence of 0 and 1, and could then be mistaken for a word made only of *a*'s and *b*'s. We then choose to encode $b \mapsto 10$. What about the last two characters? If we do not want the encoding of *a* and *b* to appear as a prefix of the encodings for *c* and *d*, we need to consider strings starting with 11: $c \mapsto 110$ and $d \mapsto 111$. Summarizing,

$$a \mapsto 0, \quad b \mapsto 10, \quad c \mapsto 110, \quad d \mapsto 111.$$

In this way, with respect to the previous encoding of (1.1), we have one character with a shorter encoding, one with the same length, and two with a longer encoding. However, in order to make a reasonable comparison we need to consider the *average* amount of bits per character of the two encodings. In this case, then, we have

$$l = \sum_{x=a,b,c,d} p(x)l(x) = \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{8} \times 3 = \frac{7}{4} \text{ bits/character},$$

which is clearly smaller than 2 bits/character provided by the encoding of equation (1.1). For example, the encoding of our string *bacabadabba* is given by the following 14-bits string

10011001110100

In this mapping, every character is represented by a string of bits in such a way that no character is mapped to a string with a prefix that represents another character. This feature makes our encoding perfectly decodable. For example, the string

110011010110000

corresponds to *cacbcaaa*. We thus achieved a compression rate of $7/4$ bits/character as opposed to the 2 bits/character of the encoding in equation (1.1), exploiting our prior knowledge summarized by a non-uniform *prior* distribution.

The above example leads us to the first big question in information theory:

First Shannon question: *What is the smallest physical resource needed to store/transmit a given file without loss?* In other words, *What is the best compression rate that we can achieve in transmitting a given file faithfully?*

The first Shannon theorem—the so-called *source-coding theorem* will provide an answer to this question, along with a first operational way to define information.

1.1.2 Error-correction

In conjunction with the problem of information compression, the most fundamental problem in information theory is that of error correction. Dealing with compression, we assumed that our storage device is *noiseless*, i.e. an ideal storage device from which we can retrieve every bit of information untouched. Alternatively, we could have imagined that information was sent through an ideal channel, and tried to minimise the number of uses of the channel.

Suppose now that our physical resource that stores or transmits information is noisy, e.g. a disk with some bad blocks, or a transmission line that often outputs a flipped bit. In most practical cases, a more reliable resource is unaffordable, and communication can be improved only using a system solution, adding an encoder before the channel and a decoder after the channel. This kind of solution suitably exploits data redundancy and error correction techniques in the encoding.

Then, the following question naturally arises

Second Shannon question: *What is the best rate for storing or transmitting a given file reliably over a noisy resource?*

The answer to this question will be given by Shannon's second theorem—the so-called noisy-channel-coding theorem. To understand how the problem of error correction is foundational for information theory, consider a fictitious scenario in which somebody is claiming to have the ability of transmitting information telepathically, but with some probability of error p . Clearly, if it turns out that $p = 1/2$, then we immediately understand that they are actually transmitting no information at all, since $p = 1/2$ is just the probability of pure-guessing. On the other hand, from the second Shannon theorem we will see that any value of p greater than $1/2$ (actually, even any value *smaller* than $1/2$!) can be used to transmit information reliably—i.e. with arbitrarily small error—upon reducing the information transmission rate. An intuitive understanding of this fact can be achieved considering a redundant encoding, e.g. substituting 00000 for 0 and 11111 for 1. In this way, it is more likely that the majority of the sent bits will be unaffected rather than the contrary. Thus, a simple *majority-voting* decoding would provide a correct transmission. Clearly, the larger is the transmitted string, the smaller is the error probability. Thus, by taking arbitrarily large encoding it is possible to make the error probability arbitrarily small. The price to pay for this strategy, however, is that the number of physical bits used to transmit a single logical bit becomes arbitrarily large. In other words, the *transmission rate*, defined as the ratio between the number of transmitted logical bits and the number of transmitted physical bits, scales down to zero if the error probability has to become arbitrarily small.

What is surprising about the second Shannon theorem, however, is the crucial result that we do not need to make the transmission rate vanishingly small: there exists a non null transmission rate at which we can transmit reliably asymptotically for large transmitted files, provided that we exploit an encoding over the entire file. What is wrong about our very naive strategy is thus the attempt to encode every single bit independently. Therefore, telepathy would have been possible if the bit error probability were just slightly different from $1/2$!

The amount of information that can be actually transmitted along a channel is connected to both the possibility of compressing and of correcting it, and the optimal encoding will address both issues jointly.

Let us start introducing the simplest example of a noisy channel, which is the one we will use to illustrate the basic ideas behind error correcting codes.

Binary symmetric channel. A *binary symmetric channel* is a channel that accepts a bit b_{in} as an input and provides a bit b_{out} as an output, with probability of correct transmission given by $1 - f$, and probability of flipping the bit f , where $0 < f < 1$.

The behaviour of the binary symmetric channel is summarised by the following conditional probabilities

$$\begin{aligned} P(b_{\text{out}} = 0 | b_{\text{in}} = 0) &= P(b_{\text{out}} = 1 | b_{\text{in}} = 1) = 1 - f, \\ P(b_{\text{out}} = 1 | b_{\text{in}} = 0) &= P(b_{\text{out}} = 0 | b_{\text{in}} = 1) = f. \end{aligned}$$

Notice that it is not restrictive to consider $0 < f \leq 1/2$, because in the opposite case we can flip the output bit $b'_{\text{out}} = b_{\text{out}} \oplus 1$, thus reducing to the case under consideration. Let us now define some notation (\oplus denotes bit-wise sum modulo 2)

- s is the *source message*, a string of M bits;
- t is the *transmitted string*, a string of N bits;
- n is the *noise*, a string of N bits representing the errors introduced by the channel;
- $r = t \oplus n$ is the *received string*, a string of N bits;

Majority-voting error correction

The purpose of the encoding is to reduce the effect of the noise, at the expense of a reduced transmission rate. Roughly speaking, this means that we encode information redundantly. The most intuitive strategy is to encode every single bit in our source message s into three copies of the same bit. It will be then very unlikely that an error affects more than one bit in the three-bits block. The decoding strategy thus consists in decoding each block to the most frequent value in the block. As an example, consider the following instance of transmission

s	0	0	1	0	1	0	1
t	000	000	111	000	111	000	111
n	010	000	001	000	100	000	000
r	010	000	110	000	011	000	111

Consider now a transmitted three-bits block $\mathbf{t}_j = t_{3j}t_{3j+1}t_{3j+2}$, and analogously a received three-bits block $\mathbf{r}_j = r_{3j}r_{3j+1}r_{3j+2}$. Denote by $P(\mathbf{r}_j | s_j)$ the conditional probability that the block \mathbf{r}_j is received, given that the signal bit is s_j . Then, the

posterior probability that s_j was encoded given that \mathbf{r}_j has been received is given by Bayes' rule

$$P(s_j|\mathbf{r}_j) = \frac{P(\mathbf{r}_j|s_j)P(s_j)}{P(\mathbf{r}_j)}, \quad P(\mathbf{r}_j|s_j) = P(\mathbf{r}_j|\mathbf{t}_j) = \prod_{k=0}^2 P(r_{3j+k}|t_{3j+k}), \quad (1.2)$$

where by definition $\mathbf{t}_j = s_j s_j s_j$, namely $t_{3j+k} = s_j$ for $k = 0, 1, 2$. We remind that the noisy channel is completely specified by the conditional probability $P(r|t)$, whereas $P(s_j)$ is the *prior* probability of signal s_j , and $P(\mathbf{r}_j)$ is the probability of receiving \mathbf{r}_j , namely the marginal probability $P(\mathbf{r}_j) = \sum_{s_j=0}^1 P(\mathbf{r}_j|s_j)P(s_j)$. All the functions involved in the expression of equation (1.2) are thus known.

Upon receiving \mathbf{r}_j we can decode the value \bar{s}_j maximising the posterior probability of equation (1.2), that for equal prior probabilities $P(s_j = 0) = P(s_j = 1) = 1/2$, is equivalent to the s_j maximising the likelihood $P(\mathbf{r}_j|s_j)$, given by

$$P(\mathbf{r}_j|\mathbf{t}_j) = \prod_{k=0}^2 P(r_{3j+k}|t_{3j+k}).$$

We remind that for the binary symmetric channel we have

$$P(r_i|t_i) = \begin{cases} 1-f & r_i = t_i, \\ f & r_i \neq t_i. \end{cases}$$

The ratio between the likelihood for $s_j = 0$ and $s_j = 1$ is thus

$$\frac{P(s_j = 1|\mathbf{r}_j)}{P(s_j = 0|\mathbf{r}_j)} = \frac{P(\mathbf{r}_j|s_j = 1)}{P(\mathbf{r}_j|s_j = 0)} = \prod_{k=0}^2 \frac{P(r_{3j+k}|1)}{P(r_{3j+k}|0)}.$$

Defining $\gamma := (1-f)/f$, each factor in the product is equal to $\gamma > 1$ for $r_{3j+k} = 1$ and to $\gamma^{-1} < 1$ for $r_{3j+k} = 0$. Therefore, the value \bar{s}_j that maximises the likelihood corresponds to the one obtained by majority-voting: every vote for 1, namely $r_{3j+k} = 1$, contributes with a factor γ in the likelihood ratio, while every vote for 0, namely $r_{3j+k} = 0$, contributes with a factor γ^{-1} . Clearly this strategy will produce the wrong correction if two errors occur within the same block. However, assuming that errors are statistically independent, and occurring with small probability, the probability of double or triple flip within each block amounts to $3f^2(1-f) + f^3$, and is much smaller than the probability of no flip or single flip $(1-f)^3 + 3f(1-f)^2$. A practical example is the following instance of transmission

s	0	0	1	0	1	0	1
t	000	000	111	000	111	000	111
n	010	000	101	000	100	000	000
r	010	000	010	000	011	000	111
\bar{s}	0	0	0	0	1	0	1
successful correction	*				*		
wrong correction				*			

Using the majority-vote error correction the error probability is improved from $\sim f$ to $\sim f^2$, being now dominated by the two-flip errors. We want to emphasise that the majority-vote code has been just considered for didactic reasons, since it is very simple, but it is also very inefficient: it uses too many bits, and still has a high probability of error. Indeed, performances of different error-correction codes can be much better. For example, for a model of disk drive making errors as a binary symmetric channel with $f = 0.1$, in order to achieve a bit probability error 10^{-15} using a majority-vote code one would need *sixty* disks to store an amount of information equivalent to one disk. On the other hand, according to the second Shannon theorem we will see that only *two* disks are needed. Indeed, the astonishing fact asserted by the second Shannon theorem is that above some threshold of error probability on a single bit, it is always possible to encode information with a vanishingly small probability of error by keeping the rate of transmission *non null*. Unfortunately, a naïve majority-vote encoding would provide a vanishingly small rate versus error probability.

The (7, 4) Hamming block code

In order to improve the encoding, one can exploit redundancy more conveniently by encoding an entire block of data instead of encoding every bit independently, as in the majority-voting strategy. This feature is implemented in block codes. A (N, K) block code is a rule for converting a sequence of source bits s of length K into a transmitted bit string t of length N . Redundancy is introduced taking $N > K$. In a linear block code the additional $N - K$ encoded bits are obtained as a linear function of the original K bits, and are called *parity-check* bits. An example of block code is the (7, 4) Hamming code, in which one transmits $N = 7$ bits for each block of $K = 4$ source bits. The rule for the encoding is to set $t_j = s_j$ for $j = 1, 2, 3, 4$, whereas the parity-check bits $t_5 t_6 t_7$ are set so that the parity within each circle in Fig. 1.1 is even.

An instance of encoding is given in Fig. 1.2.

In order to treat the block code appropriately, let us introduce some notation here.

Definition 1.1 (Inner product). Let x, y be strings in \mathbb{F}_2^M . Then the *inner product* of x and y is

$$x^T y = x_1 y_1 \oplus x_2 y_2 \oplus \cdots \oplus x_M y_M,$$

where \oplus denotes sum modulo 2.

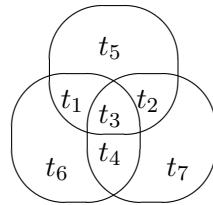


Figure 1.1 Representation of the $(7, 4)$ Hamming code.

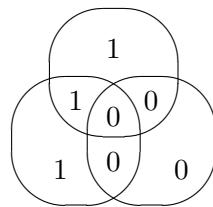


Figure 1.2 Example of $(7, 4)$ Hamming encoding.

Considering that a matrix $A \in \mathbb{F}_2^M \times \mathbb{F}_2^L$ can be written as

$$A = \begin{pmatrix} \alpha_r^{(1)T} \\ \alpha_r^{(2)T} \\ \vdots \\ \alpha_r^{(M)T} \end{pmatrix} = (\alpha_c^{(1)} \ \alpha_c^{(2)} \ \dots \ \alpha_c^{(L)}),$$

we can then define the matrix product as follows

Definition 1.2 (Matrix product). Let A be a matrix in $M_{M \times L}(\mathbb{F}_2)$ and $B \in M_{L \times J}(\mathbb{F}_2)$. Then the *matrix product* of A and B is the matrix C in $M_{M \times J}(\mathbb{F}_2)$ with elements $C_{i,j} = \alpha_r^{(i)T} \beta_c^{(j)}$

The fact that the Hamming code is linear allows one to write it compactly in form of

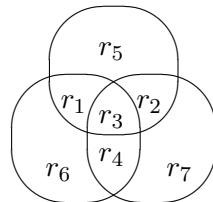
matrix multiplication as

$$t = Gs$$

where G is the *generator matrix* of the code, given by

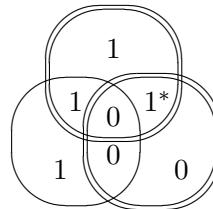
$$G = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{pmatrix}$$

The received vector r is represented as follows

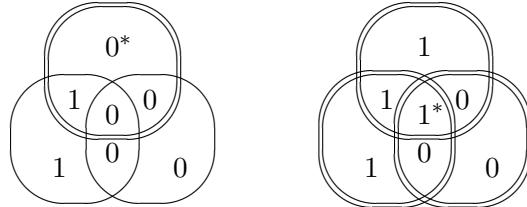


where $r = t \oplus n$. What makes it difficult to derive the decoding procedure is the fact that each bit—including the parity-check ones—may be flipped by the noise, thus we do not know where the error occurred. However, if the parity in one circle is violated, we know that one error must have occurred within that circle. We can then define three bits z_1 , z_2 and z_3 corresponding to the parities of the three circles, and since their pattern heralds the occurrence of errors, the string z is called *syndrome*. One can decode the string by figuring out what are the bits affected by noise, which are supposed to be those in the intersection of circles with syndrome 1 and complements of circles with syndrome 0.

For example, in the following instance the syndrome is $z = 101$, with flipped bit r_2



where the syndrome 1 circles have been represented with a double line. Other examples are the following ones, with syndrome $z = 100$ and $z = 111$, respectively.



Also the decoding can be described in matrix form, as is the case for the encoding. In the presence of noise the received vector r is given by

$$r = Gs \oplus n$$

The matrix G can be conveniently rewritten in block form as follows

$$G = \begin{pmatrix} I_4 \\ P \end{pmatrix}, \quad P = \begin{pmatrix} 1 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{pmatrix}$$

I_4 denoting the four-dimensional identity matrix. The computation of the syndrome vector z is linear, and can be obtained as $z = Hr$, where H is called *parity-check* matrix. In our case H is given by

$$H = \begin{pmatrix} 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 & 1 \end{pmatrix} = (P \quad I_3). \quad (1.3)$$

Notice that H is designed to satisfy $HG = 0$, and thus, since codewords are obtained as $t = Gs$, the matrix H satisfies $Ht = 0$ for every codeword t . Therefore, for the syndrome one has

$$z = H(Gs \oplus n) = Hn$$

Indeed, there may be many noise strings n leading to the same syndrome z , and the syndrome-decoding problem is to find the most probable noise vector n satisfying the above equation, so that we can reconstruct the message through a decoding function $t = t(r)$. A decoding algorithm that solves this problem is called a maximum-likelihood decoder. For example, for small probability of flips, we have sizeable probabilities only for a single flip, and such a flip in the j -th position is described by a noise vector e_j . Therefore, one concludes that the error occurred at the j -th place for j satisfying the equation $Hr = He_j$ (single-flip syndrome), and the decoded source block would be $t = r \oplus e_j$.

The question is now: how to build a general block code? Let us analyse the question thoroughly.

Definition 1.3 (Linear block code). Let $G : \mathbb{F}_2^K \rightarrow \mathbb{F}_2^N$ be a linear *encoding* function and $H : \mathbb{F}_2^N \rightarrow \mathbb{F}_2^{N-K}$ a linear *decoding* function. A (N, K) *linear block code* is a couple (G, H) . Given such a code, we say that $y \in \mathbb{F}_2^N$ is a *codeword* if $y \in G(\mathbb{F}_2^K)$.

In the analysis of block codes a useful notion is that of Hamming distance. The *weight* of a bit string $x \in \mathbb{F}_2^K$ is the function $w : Z_2^K \rightarrow \mathbb{N}$ defined as $w(x) := \sum_{i=1}^K x_i$.

Definition 1.4 (Hamming distance). Given $x, y \in \mathbb{F}_2^N$, their *Hamming distance* is given by the function $d : \mathbb{F}_2^N \times \mathbb{F}_2^N \rightarrow \mathbb{N}$ defined as $d(x, y) := w(x \oplus y)$.

The Hamming distance between two strings is then the number of sites at which the bits of the two strings differ. For general linear codes, the optimal decoding consists in associating the codeword t to the received string r such that $d(t, r)$ is minimised. Such a decoding is called *minimum distance decoding*. A code is characterised by its minimum distance, defined as the minimum distance between two codewords: $d_m(G, H) := \min_{x, y \in G(\mathbb{F}_2^K)} d(x, y)$. This quantity determines the performance of a code as follows.

Theorem 1.5. Let (G, H) be a code with minimum distance $d_m(G, H) = 2j + 1$. Then the two following statements hold

1. (G, H) can correct up to j errors in a block
2. (G, H) can detect up to $2j$ errors per block.

Proof. Suppose that the message is the string s , with codeword $t := Gs$ and the transmission is affected by the noise n with weight $w(n) \leq j$. Then, for any codeword $z \in G(\mathbb{F}_2^K)$ one has

$$2j + 1 \leq d(t, z) \leq d(t, r) + d(r, z) = w(n) + d(r, z) \leq j + d(r, z), \quad (1.4)$$

which implies $d(r, z) \geq j + 1$, while by hypothesis $d(r, t) \leq j$. Thus, the minimum-distance decoding allows one to correct up to j errors correctly. Let us now suppose that $1 \leq w(n) \leq 2j$. In this case one has

$$1 \leq d(t, r) \leq 2j, \quad (1.5)$$

and thus r cannot be a different codeword. This implies that the occurrence of less than $2j + 1$ errors is always detected. However, in case $w(n) > j$ correct decoding is not possible. \square

Let us now show the general structure of a block code. Most of the analysis is based on the decoding matrix H . For our purposes, it is useful to introduce the *kernel* of a matrix $F : \mathbb{F}_2^M \rightarrow \mathbb{F}_2^J$, which is defined as the set $\text{Ker}_F := \{y \in Z_2^M \mid Fy = 0\}$, where $0 \in \mathbb{F}_2^J$ is the string of J bits 0, i.e. 00...0.

Definition 1.6 (Associated parity check and generator matrices). A *canonical parity check matrix* in $M_{N-K \times N}(\mathbb{F}_2)$ is a matrix of the form

$$H = (A \quad I_{N-K}). \quad (1.6)$$

A *standard generator matrix* in $M_{N \times K}(\mathbb{F}_2)$ is a matrix of the form

$$G = \begin{pmatrix} I_K \\ A \end{pmatrix}. \quad (1.7)$$

Two such matrices with a common $(N - K) \times K$ block A are *associated* to each other. They define a (N, K) block code and set of codewords $\mathcal{C} := G(\mathbb{F}_2^N)$.

The main consequence of the above definition is the following result that we already exploited for the case of the $(7, 4)$ block code.

Theorem 1.7. *Let G, H be associated to each other. Then $HG = 0$.*

Proof. If G, H are associated to each other we have

$$HG = (A \quad I_{N-K}) \begin{pmatrix} I_K \\ A \end{pmatrix} = A \oplus A = 0, \quad (1.8)$$

where we used the property of element-wise sum modulo two that $X \oplus X = 0$. \square

Actually, we can prove a slightly stronger statement.

Theorem 1.8. *Let G, H be associated. Then $\mathcal{C} = \text{Ker}_H$.*

Proof. By theorem 1.7 we clearly have $\mathcal{C} \subseteq \text{Ker}_H$. We then need to prove that $\text{Ker}_H \subseteq \mathcal{C}$. Indeed, let $y \in \text{Ker}_H$. Let us decompose Hy as follows

$$Hy = (A \quad I_{N-K}) \begin{pmatrix} y_u \\ y_d \end{pmatrix} = Ay_u \oplus y_d = 0, \quad (1.9)$$

where $y_u \in \mathbb{F}_2^N$ and $y_d \in \mathbb{F}_2^{N-K}$. Then, $Ay_u = y_d$, namely one can write

$$y = \begin{pmatrix} y_u \\ Ay_u \end{pmatrix} = \begin{pmatrix} I_N \\ A \end{pmatrix} y_u = Gy_u. \quad (1.10)$$

Thus, we proved that $y \in G(\mathbb{F}_2^N) = \mathcal{C}$. \square

Finally, we can prove the following theorem.

Theorem 1.9. *Let G, H define a (K, N) block-code. The minimum distance $d_m(G, H)$ is j iff the parity check matrix H is such that every set of $l < j$ columns $\{h_c^{(i_1)}, h_c^{(i_2)}, \dots, h_c^{(i_l)}\}$ satisfies $h_c^{(i_1)} \oplus h_c^{(i_2)} \oplus \dots \oplus h_c^{(i_l)} \neq 0$, while there is a set of j columns $\{h_c^{(i_1)}, h_c^{(i_2)}, \dots, h_c^{(i_j)}\}$ such that $h_c^{(i_1)} \oplus h_c^{(i_2)} \oplus \dots \oplus h_c^{(i_j)} = 0$.*

Proof. Let e_i denote the canonical basis vector in \mathbb{F}_2^N , then $h_c^{(i)} = He_i$. Now, by theorem 1.8, two strings $y, z \in \mathbb{F}_2^{N-K}$ are codewords, i.e. elements of C , if and only if $y, z \in \text{Ker}_H$. The minimum distance $d_m(G, H)$ is then larger than $j - 1$ if and only if $y \oplus e_{i_1} \oplus e_{i_2} \oplus \dots \oplus e_{i_l} \notin \text{Ker}_H$ for all $y \in \text{Ker}_H$ and for any choice of the l vectors $\{e_{i_k}\}_{k=1}^l$ with $l < j$. Thus, the Hamming distance of any two codewords y, z is larger than $j - 1$ if and only if $H(e_{i_1} \oplus e_{i_2} \oplus \dots \oplus e_{i_l}) \neq 0$ for any choice of the l vectors $\{e_{i_k}\}_{k=1}^l$, with $l < j$. Now, since the vectors He_i are columns $h_c^{(i)}$ of H , the last condition is equivalent to the first statement of the thesis, namely every set of $l < j$ columns $\{h_c^{(i_1)}, h_c^{(i_2)}, \dots, h_c^{(i_l)}\}$ satisfies $h_c^{(i_1)} \oplus h_c^{(i_2)} \oplus \dots \oplus h_c^{(i_l)} \neq 0$. Finally, the minimum distance is precisely j if there are two codewords y, z such that $y = z \oplus e_{i_1} \oplus e_{i_2} \oplus \dots \oplus e_{i_j}$, and then $H(e_{i_1} \oplus e_{i_2} \oplus \dots \oplus e_{i_j}) = h_c^{(i_1)} \oplus h_c^{(i_2)} \oplus \dots \oplus h_c^{(i_j)} = 0$. \square

Exercise 1.1

Prove that the $(7, 4)$ block code can detect up to 2 errors, and correct up to 1.

Answer of exercise 1.1

Let us consider the matrix H in equation (1.3). Every couple of columns is different, but if we consider e.g. the columns h_3, h_4, h_5 we have $h_3 \oplus h_4 \oplus h_5 = 0$. Thus by theorem 1.9 we have $d_m(G, H) = 3$. Finally, by theorem 1.5 the $(7, 4)$ block code can detect 2 errors and correct one.

Chapter 2

Lecture 2: Operational Probabilistic Theories

In this lecture we introduce Operational Probabilistic Theories (OPTs). Every OPT represents a possible theory of information. What do we expect a theory of information to look like? In order to answer this question, let us think about our experience. One situation where we deal with information and its processing in our everyday life is when handling a smartphone. A smartphone is a device that manipulates the electromagnetic field in its surroundings and within its silicon circuits, as well as the field of phonons (i.e. vibration modes of the surrounding air, whose excitations are acoustic waves). Information in this case is encoded in the carrier systems—modes of the electromagnetic field and of acoustic vibrations—and processed by suitable operations. A theory of information is a mathematical model aimed at describing such kind of situation. In the first place, we require our theory to provide a way to represent information carrying systems, as well as operations that such systems can undergo.

2.1 Operational theories

In our example, every operation reads the state of a set of *input* systems (e.g. electromagnetic pulses that are absorbed), produces a new state of a set of *output* systems (e.g. re-writes some blocks of the memory), and displays some *outcome* (e.g. a message on the display). In the context of OPTs, we call such an operation a *test*. In classical information processing the state of the output system and the outcome may often coincide, but this is not a necessary condition. A better understanding of what a test is can be given by a quantum measurement: one has an input state, a process that produces an outcome—e.g. a string of digits on a display—and an output state, that possibly depends on the outcome—e.g. the state of the quantum system after a von Neumann-Lüders measurement. An operational theory Θ consists then in i) a collection $T(\Theta)$ of tests $T_X^{A \rightarrow B}$, each labelled by input and output letters from a collection $Sys(\Theta)$ denoting system types, e.g. $A \rightarrow B$, and by a finite set of outcomes X ; for every pair of types $A, B \in Sys(\Theta)$ the set of tests of type $A \rightarrow B$ is denoted $Test(A \rightarrow B)$. A test A_X of type $A \rightarrow B$ can be represented by a diagram as the following

$$\xrightarrow{A} \boxed{A_X} \xrightarrow{B} .$$

In general, an information processing algorithm requires more than a single operation, and is articulated in subsequent steps. Thus, an operational theory is also specified by ii) a rule for sequential composition: the test $T_X \in \text{Test}(A \rightarrow B)$ can be followed by the test $R_Y \in \text{Test}(A' \rightarrow B')$ if $A' \equiv B$, thus obtaining the sequential composition $R_Y T_X := R T_{X \times Y} \in \text{Test}(A \rightarrow B')$. Sequential composition must have the following properties.

1. Associativity: given $R_X \in \text{Test}(A \rightarrow B)$, $T_Y \in \text{Test}(B \rightarrow C)$, and $W_Z \in \text{Test}(C \rightarrow D)$, one has

$$W_Z(T_Y R_X) = (W_Z T_Y) R_X \in \text{Test}(A \rightarrow D).$$

2. Identity: for every $A \in \text{Sys}(\Theta)$, there exists a test with $X = \{\ast\}$, denoted by $I^A \in \text{Test}(A \rightarrow A)$ such that $I^B R_X = R_X I^A = S_X$, for every $R_X \in \text{Test}(A \rightarrow B)$.

In terms of diagrams, sequential composition of $A_X \in \text{Test}(A \rightarrow B)$ and $B_Y \in \text{Test}(B \rightarrow C)$ is represented as follows

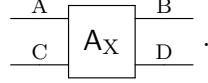
$$\xrightarrow{A} [BA_{X \times Y}] \xrightarrow{C} = \xrightarrow{A} [A_X] \xrightarrow{B} [B_Y] \xrightarrow{C} .$$

Now, let us consider e.g. the blocks of a hard disk. Every block is a system, but an array of blocks can be treated as a composite system itself, carrying information that is stored in a delocalised way. For example, if we want to encode four symbols, we need two bits. If we want to process information encoded in four symbols, we need to operate jointly on the two cells containing the two bits. More generally, every pair of systems can be composed *in parallel* to constitute composite systems. When operating on composite systems, one can run a test that involves all component systems at once. However, a distinctive feature of composite systems is that one can carry independent tests on each of their components simultaneously. For example, when operating on an array of two bits, one can “reset the first bit to the value 0, and flip the second bit”. The result of such an operation is independent of the order of the two operations on the single bits, and these can be also run simultaneously. Such a test is the *parallel composition* of two tests: “reset the first bit to the value 0”, and “flip the second bit”. To complete the picture of an operational theory we then require iii) a rule $\otimes : (A, B) \mapsto AB$ for composing systems in parallel, and a corresponding rule for tests $\otimes : (S_X, T_Y) \mapsto (S \otimes T)_{X \times Y}$, with the following properties.

1. Associativity: $(AB)C = A(BC)$.
2. Unit: there is a unique system I such that $IA = AI = A$ for every $A \in \text{Sys}(\Theta)$.
3. For every $R_X \in \text{Test}(A \rightarrow B)$ and $T_Y \in \text{Test}(C \rightarrow D)$, one has $R_X \otimes T_Y \in \text{Test}(AC \rightarrow BD)$. Associativity of \otimes holds:

$$(R_X \otimes T_Y) \otimes W_Z = R_X \otimes (T_Y \otimes W_Z).$$

For composite systems we use diagrams with multiple wires, e.g.



In particular, for tests that are parallel compositions we will draw

$$\begin{array}{c} \text{A} \\ \text{C} \end{array} \xrightarrow{\quad} \boxed{A_X \otimes B_Y} \begin{array}{c} \text{B} \\ \text{D} \end{array} = \begin{array}{c} \text{A} \\ \text{C} \end{array} \xrightarrow{\quad} \boxed{A_X} \begin{array}{c} \text{B} \\ \text{D} \end{array} + \begin{array}{c} \text{A} \\ \text{C} \end{array} \xrightarrow{\quad} \boxed{B_Y} \begin{array}{c} \text{B} \\ \text{D} \end{array} .$$

Altogether, sequential and parallel composition must satisfy the following properties

1. For every $A_X \in \text{Test}(A \rightarrow B)$, $B_Y \in \text{Test}(B \rightarrow C)$, $D_Z \in \text{Test}(D \rightarrow E)$, $E_W \in \text{Test}(E \rightarrow F)$, one has

$$(B_Y \otimes E_W)(A_X \otimes D_Z) = (B_Y A_X) \otimes (E_W D_Z). \quad (2.1)$$

2. Swap: for every pair of system types A, B , there exists a test $S^{AB} \in \text{Test}(AB \rightarrow BA)$ such that $S^{BA} S^{AB} = I^{AB}$, and $S^{AB}(A_X \otimes B_Y) = (B_Y \otimes A_X)S^{AB}$. Moreover,

$$\begin{aligned} S^{(AB)C} &= (I^A \otimes S^{BC})(S^{AC} \otimes I^B), \\ S^{A(BC)} &= (S^{AB} \otimes I^C)(I^B \otimes S^{AC}). \end{aligned}$$

The identity test will be often omitted: $\frac{A}{I} \frac{A}{A} = \frac{A}{A}$.

All tests of an operational theory are (finite) collections of *events*: $\text{Test}(A \rightarrow B) \ni R_X = \{\mathcal{R}_i\}_{i \in X}$. Now, if $\text{Test}(A \rightarrow B) \ni R_X = \{\mathcal{R}_i\}_{i \in X}$ and $\text{Test}(B \rightarrow C) \ni T_Y = \{\mathcal{T}_j\}_{j \in Y}$, then

$$\text{Test}(A \rightarrow C) \ni (TR)_{X \times Y} := \{\mathcal{T}_j \mathcal{R}_i\}_{(i,j) \in X \times Y}.$$

Similarly, for $R_X \in \text{Test}(A \rightarrow B)$ and $T_Y \in \text{Test}(C \rightarrow D)$,

$$(R \otimes T)_{X \times Y} := \{\mathcal{R}_i \otimes \mathcal{T}_j\}_{(i,j) \in X \times Y}.$$

The set of events of tests in $\text{Test}(A \rightarrow B)$ is denoted by $\text{Transf}(A \rightarrow B)$. By the properties of sequential and parallel composition of tests, one can easily derive associativity of sequential and parallel composition of events, as well as the analogue of Eq. (2.1). We will use for events the same diagrammatic notation as for tests:

$$\begin{array}{c} \text{A} \\ \text{C} \end{array} \xrightarrow{\quad} \boxed{\mathcal{A}_i} \begin{array}{c} \text{B} \\ \text{D} \end{array}, \quad \begin{array}{c} \text{A} \\ \text{C} \end{array} \xrightarrow{\quad} \boxed{\mathcal{A}_i} \begin{array}{c} \text{B} \\ \text{D} \end{array} \xrightarrow{\quad} \boxed{\mathcal{B}_j} \begin{array}{c} \text{C} \\ \text{D} \end{array}, \quad \begin{array}{c} \text{A} \\ \text{C} \end{array} \xrightarrow{\quad} \boxed{\mathcal{A}_i} \begin{array}{c} \text{B} \\ \text{D} \end{array}, \quad \begin{array}{c} \text{A} \\ \text{C} \end{array} \xrightarrow{\quad} \boxed{\mathcal{B}_j} \begin{array}{c} \text{B} \\ \text{D} \end{array} .$$

As for the identity test, also the identity event will be often omitted: $\frac{A}{J} \frac{A}{A} = \frac{A}{A}$.

2.2 Coarse graining

For every test $T_X \in \text{Test}(A \rightarrow B)$ with $T_X = \{\mathcal{T}_i\}_{i \in X}$, and every disjoint partition $\{X_j\}_{j \in Y}$ of $X = \bigcup_{j \in Y} X_j$, one can define the *coarse graining* that maps T_X to $T'_Y \in \text{Test}(A \rightarrow B)$, with $T'_Y = \{\mathcal{T}'_j\}_{j \in Y}$. We define $\mathcal{T}_{X_j} := \mathcal{T}'_j$. This mathematical map represents a test where some outcome is read, and part of the information is forgotten. Consider, e.g. the roll of a die. The system is the die, the state is the upper face. A roll will change the state and produce an outcome which is the reading of the upper face. A coarse graining is e.g. the test corresponding to a roll of the die with a device only announcing whether the upper face is even or odd.

Parallel and sequential composition distribute over coarse graining:

$$\begin{aligned}\mathcal{T}_{X_j} \otimes \mathcal{R}_k &= (\mathcal{T} \otimes \mathcal{R})_{X_j \times \{k\}}, \\ \mathcal{A}_l \mathcal{T}_{X_j} \mathcal{B}_k &= (\mathcal{A} \mathcal{T} \mathcal{B})_{\{l\} \times X_j \times \{k\}}.\end{aligned}$$

Notice that for every test $T_X \in \text{Test}(A \rightarrow B)$ there exists the singleton test $T'_* := \{\mathcal{T}_X\}$. One can easily prove that the identity event \mathcal{I}_A , such that $I^A = \{\mathcal{I}_A\}$, satisfies $\mathcal{I}_B \mathcal{T} = \mathcal{T} \mathcal{I}_A$ for every event $\mathcal{T} \in \text{Transf}(A \rightarrow B)$. Similarly, for $S_{AB} = \{\mathcal{S}_{AB}\}$ we have $\mathcal{S}_{BA} \mathcal{S}_{AB} = \mathcal{I}_{AB}$. The collection of events of a theory Θ will be denoted by $\text{Ev}(\Theta)$.

2.3 Probabilistic theories

Operational theories as we defined them allow one to *describe* a cluster of processes in parallel and in a sequence, but at the present stage they are devoid of any predictive power. In the present section we introduce new structures that make an operational theory predictive. These come in the form of rules for assessing *probabilities* of events. An operational theory is an *operational probabilistic theory* (OPT) if the tests $\text{Test}(I \rightarrow I)$ are probability distributions: $1 \geq \mathcal{T}_i = p_i \geq 0$, so that $\sum_{i \in X} p_i = 1$, and given two tests $S_X, T_Y \in \text{Test}(I \rightarrow I)$ with $\mathcal{S}_i = p_i$ and $\mathcal{T}_i = q_i$, the following identities hold

$$\begin{aligned}\mathcal{S}_i \otimes \mathcal{T}_j &= \mathcal{S}_i \mathcal{T}_j := p_i q_j, \\ \mathcal{T}_{X_j} &:= \sum_{i \in X_j} p_i,\end{aligned}$$

meaning that events in the same test are mutually exclusive and events in different tests of system I are independent. While it is immediate that $1 \in \text{Transf}(I \rightarrow I)$, since 1 is the only singleton test, we will assume that $0 \in \text{Transf}(I \rightarrow I)$.

An OPT Θ is specified by the collections $(\text{Sys}(\Theta), \text{T}(\Theta))$, along with rules to calculate probabilities for every test $\text{Test}(I \rightarrow I)$

2.4 States and effects

Two special families of events are those of type $\text{Transf}(I \rightarrow A)$ —from a trivial system to a non-trivial one—and $\text{Transf}(A \rightarrow I)$ —transforming a non-trivial system to a trivial one. The first type can be interpreted as a process where a system is initialised, i.e. a procedure

of state preparation. Notice that the preparation might not prepare a unique state, but could result in a collection of different preparations $\{\rho_i\}_{i \in X}$, heralded by some outcome $i \in X$. We call these tests *preparation-tests*. The second type of test can be interpreted as a measurement whose aftermath is completely disregarded. This kind of test is called an *observation-test*. We will denote the sets $\text{Transf}(I \rightarrow A)$ and $\text{Transf}(A \rightarrow I)$ by the symbols $\text{St}(A)$ and $\text{Eff}(A)$, respectively. Events in $\text{St}(A)$ are called *states*, and denoted by lower-case greek letters, e.g. ρ , while events in $\text{Eff}(A)$ are called *effects*, and denoted by lower-case latin letters, e.g. a . The following diagrammatic notation, where we denote states and effects by the symbols

$$\langle \rho \rangle^A , \xrightarrow{A} [a] ,$$

respectively, turns out to be very useful and will be used extensively in the course.

As a consequence of the above definitions, we have that a sequence made of a preparation-test and an observation test is a joint probability distribution:

$$\langle \rho_i \rangle^A [a_j] = P(i, j | \rho, a).$$

If we consider a single event $\rho_0 \in \text{St}(A)$ in a preparation test, we then have a *functional* that maps every effect $a_j \in \text{Eff}(A)$ to a real number

$$0 \leq \langle \rho_0 \rangle^A [a_j] \leq 1 .$$

Given two states $\rho, \sigma \in \text{St}(A)$ one can then formally define their linear combination $x\rho + y\sigma$ as the functional on effects $a \in \text{Eff}(A)$ given by

$$\langle x\rho + y\sigma \rangle^A [a] := x \langle \rho \rangle^A [a] + y \langle \sigma \rangle^A [a] .$$

The set $\text{St}(A) := \text{Transf}(I \rightarrow A)$ is then a spanning subset of a real vector space $\text{St}(A)_{\mathbb{R}}$ of real functionals on effects. On the other hand, $\text{Eff}(A)$ is a *separating* set of positive linear functionals on $\text{St}(A)$, namely for every pair of different states $\rho, \sigma \in \text{St}(A)$ there exists an effect $a \in \text{Eff}(A)$ such that

$$\langle \rho \rangle^A [a] \neq \langle \sigma \rangle^A [a] .$$

One can prove that, by the above property, $\text{Eff}(A)$ spans the dual space $\text{St}(A)_{\mathbb{R}}^* =: \text{Eff}(A)_{\mathbb{R}}$. The dimension D_A of $\text{St}(A)_{\mathbb{R}}$ (which is the same as that of $\text{Eff}(A)_{\mathbb{R}}$) is called *size* of system A. Using the properties of parallel composition, one can easily prove that $D_{AB} \geq D_A D_B$, and that in any OPT Θ , I is the unique system with unit size $D_I = 1$.

2.5 Transformations

Let us now consider a general event $\mathcal{A} \in \text{Transf}(A \rightarrow B)$ of an OPT Θ . By definition of sequential composition, given a state $\rho \in \text{St}(A) = \text{Transf}(I \rightarrow A)$, we have that

$$\langle \rho \rangle^A \xrightarrow{A} [\mathcal{A}] \xrightarrow{B} \in \text{Transf}(I \rightarrow B) = \text{St}(B).$$

This means that an event $\mathcal{A} \in \text{Transf}(A \rightarrow B)$ transforms input states of system A to output states of system B. It is possible to prove that an event defines a linear transformation from the real vector space $\text{St}(A)_{\mathbb{R}}$ to the real vector space $\text{St}(B)_{\mathbb{R}}$. Events will then be often also called *transformations*¹. A test is then a collection of possible transformations, labelled by outcomes that can be thought of as the reading of a pointer, heralding the occurrence of a specific event within the test. Coarse graining in an OPT corresponds to the *sum* of linear maps:

$$\mathcal{T}_{X_j} = \sum_{i \in X_j} \mathcal{T}_i.$$

A *singleton test*, i.e. a test with a unique outcome (such as the full coarse graining of any test) is called *deterministic*, and its (unique) event is called a *channel*. A channel $\mathcal{C} \in \text{Transf}(A \rightarrow B)$ is *reversible* if there exists a channel $\mathcal{D} \in \text{Transf}(B \rightarrow A)$ such that $\mathcal{DC} = \mathcal{I}_A$ and $\mathcal{CD} = \mathcal{I}_B$.

2.6 Classical information theory

Classical information theory is a special Operational Probabilistic Theory. As such, it is specified by a class of systems and their possible tests. In practice, classical information carrying systems can be thought of as physical memory cells with n internal levels, each one encoding one symbol from an alphabet with n symbols. The theoretical description of classical systems X is thus given by alphabets, such as $X \leftrightarrow \{x_1, x_2, \dots, x_n\}$ having an integer number of characters n , every character labelling one level. The *type* of a system X is the cardinality n of the corresponding alphabet. For example $X = \{0, 1, 2\}$ corresponds to a system of type $n = 3$. For the purpose of information encoding, all systems of the same type are equivalent, independently of the specific alphabet. For example, encoding in the alphabet $\{A, B\}$ is completely equivalent to encoding in $\{0, 1\}$.

The encoding of a symbol in a system corresponds to the preparation of a special *state* of the system. The preparation may occur according to a probabilistic algorithm. Consequently, it is convenient to consider states of a system X of type n as (generally sub-normalised) probability distributions $\mathbf{p} = \{p_1, p_2, \dots, p_n\}$ over n values: $\sum_{i=1}^n p_i \leq 1$. A sub-normalised distribution \mathbf{p} corresponds to a preparation event that might occur within a preparation test, with probability $\mathbb{P}_X[\mathbf{p}] = \sum_{i=1}^n p_i$. The state \mathbf{p} is *deterministic* if and only if it is *normalised*, namely $\sum_{i=1}^n p_i = 1$, i.e. it occurs with certainty when the corresponding preparation test is performed. In the course of these lectures, we will only deal with normalised states.

A state \mathbf{p} of system X can be represented by the following diagram

$$\boxed{\mathbf{p}} \xrightarrow{\quad} X \tag{2.2}$$

Notice that not only the particular alphabet is irrelevant for evaluation of the information content of a system—what matters being only the cardinality of the alphabet—but even

¹The notion of transformation is actually more complicate than the one reported here, but this will be sufficient for the theories of interest in the present course, i.e. classical and quantum theory.

the details of the encoding procedure are irrelevant. For example, is not important if a symbol is produced by flipping a coin or by rolling a die and recording the parity of the upper face. Indeed, the information content of a symbol is fully captured by the probability distribution \mathbf{p} alone.

Every state \mathbf{p} of a system of type n is then a probability distribution over a finite alphabet $X = \{x_1, x_2, \dots, x_n\}$, namely a vector with n non-negative real entries bounded by 1. The set of states of system X is then given by

$$\text{St}(X) = \{\mathbf{p} \in \mathbb{R}^n \mid \sum_{i=1}^n p_i \leq 1, \forall i p_i \geq 0\}, \quad (2.3)$$

while the set of deterministic states is $\text{St}(X)_1 \subseteq \text{St}(X)$, defined as

$$\text{St}(X)_1 = \{\mathbf{p} \in \mathbb{R}^n \mid \sum_{i=1}^n p_i = 1, \forall i p_i \geq 0\}. \quad (2.4)$$

The *state space* of a system X , instead, is $\text{St}(X)_{\mathbb{R}} = \mathbb{R}^n$. A basis of such space is given by the vectors \mathbf{e}_i with components $(\mathbf{e}_i)_j = \delta_{ij}$. Notice that for every $1 \leq i \leq n$ one has

$$\mathbf{e}_i \in \text{St}(X)_1 \subseteq \text{St}(X). \quad (2.5)$$

Summarising what we observed so far, from the point of view of classical information theory, the encoding of a symbol is completely described by a random variable.

The first purpose of classical information theory is to quantify the amount of information that a given preparation encodes in a system. The quantification of information typically pertains a file or string. For this reason, it is important to specify the way in which systems compose into composite systems, carrying strings of symbols. Notice that strings of symbols from an alphabet X of type n of a given length L can be thought of as symbols from a new alphabet X^L of type n^L . More generally, composing two systems, say X of type n and Y of type m , *in parallel* allows one to encode messages in a new system Z whose alphabet has mn characters, namely $XY = Z$ of type mn . The alphabet of a composite system is made of pairs of symbols (or more generally by strings of symbols, if the system is made of more than two components), each one from the alphabet of the corresponding component system. Their probability distribution can in general exhibit correlations, and will be diagrammatically represented as follows

$$\boxed{\mathbf{p}} \xrightarrow{Z} = \boxed{\mathbf{p}} \begin{array}{c} \xrightarrow{X} \\ \xrightarrow{Y} \end{array} \quad (2.6)$$

Factorised distributions, such as $p_{mn} = p'_m p''_n$, do not have correlations, and represent *independent* preparations of the subsystems. We will represent such special preparations as follows

$$\boxed{\mathbf{p}} \begin{array}{c} \xrightarrow{X} \\ \xrightarrow{Y} \end{array} = \boxed{\mathbf{p}'} \xrightarrow{X} \boxed{\mathbf{p}''} \xrightarrow{Y} .$$

Notice that Each symbol $x_{ij} \in X_j$, in a string of a composite system made of equivalent components $X_1 X_2 \dots X_L = X^L$, is distributed according to a probability distribution $\mathbf{p}^{(j)}$ —the marginal distribution—that can be the same for every symbol, for example the sheer statistics of symbols in a long string (this however is too simplistic a model in most cases, because it completely neglects correlations), or may represent some prior knowledge, e.g. knowing that “the symbol x_j occurs at position j within a text that is written in Greek”.

Let us now consider transformations in classical information theory. From the framework of OPTs we know that transformations of type $X \rightarrow Y$ with X of type n and Y of type m are in correspondence with linear maps from $\text{St}(X)_{\mathbb{R}} = \mathbb{R}^n$ to $\text{St}(Y)_{\mathbb{R}} = \mathbb{R}^m$. A transformation in $\text{Transf}(X \rightarrow Y)$ is then represented by a $m \times n$ real matrix (C_{ij}) . But what matrices do actually represent an event? Let us start considering deterministic transformations, namely those that belong to singleton tests. Let C be such a transformation. Then, reminding Eq. (2.5) we have $C\mathbf{e}_i \in \text{St}(Y)_1$. Then $\mathbf{p}_i := C\mathbf{e}_i$, that represents the i -th column of the matrix (C_{ji}) , is a probability distribution, namely $p_j = (\mathbf{p}_i)_j = (C\mathbf{e}_i)_j = C_{ji} \geq 0$. Moreover, summing over the row index j one has $\sum_{j=1}^m C_{ji} = 1$. A matrix C with these properties, i.e. $C_{ji} \geq 0$ and $\sum_{j=1}^m C_{ji} = 1$ for all i , is called a *stochastic matrix*.

The columns of a stochastic matrix can be interpreted as conditional probability distributions $C_{ji} = p(j|i)$: they represent the probabilities that the channel provides output symbols j given that the input symbol was i .

Finally, what are the most general events E of the theory, i.e. transformations that occur within a test? First of all, events must map states to states. This implies that the corresponding matrices (E_{ji}) must have positive entries, since E_{ji} is the j -th component of the vector $E\mathbf{e}_i$, which must be a state. Moreover, the general (sub-normalized) states $E\mathbf{e}_i$ must belong to the set $\text{St}(Y)$ characterised by Eq. (2.3), namely

$$\sum_{j=1}^m E_{ji} \leq 1.$$

General transformations then correspond to *sub-stochastic matrices*, and a test is a collection of transformations whose sum is a stochastic matrix.

Chapter 3

Lecture 3: Classical systems and Random variables

3.1 Random variables

Since we defined systems of classical information theory as alphabets and their states as random variables, we now take a closer look at the notion of random variable. A random variable X can take values in a set called *range* of X and denoted as $\text{Rng}(X) = \{x_1, x_2, \dots, x_n\}$. The random variable is also defined by the probability distribution $\mathbf{p} = \{p(x_1), p(x_2), \dots, p(x_n)\}$. More formally, in order to define a random variable we need to introduce the notion of *probability space*, i.e. a triple $(\Omega, \mathcal{A}, \mu)$ where \mathcal{A} is a σ -algebra of subsets of the set Ω , namely a family of subsets such that

1. $\emptyset, \Omega \in \mathcal{A}$,
2. $A \in \mathcal{A} \Rightarrow \Omega \setminus A \in \mathcal{A}$,
3. $\{A_i\}_{i \in I} \subseteq \mathcal{A} \Rightarrow \bigcup_{i \in I} A_i \in \mathcal{A}, \forall I \subseteq \mathbb{N}$.

and μ is a *probability measure*, namely a function $\mu : \mathcal{A} \rightarrow \mathbb{R}$ that satisfies

1. $\forall A \in \mathcal{A}, \mu(A) \geq 0$,
2. $\mu(\Omega) = 1$,
3. $\{A_i\}_{i \in I} \subseteq \mathcal{A}, A_i \cap A_j = \emptyset \Rightarrow \mu \left(\bigcup_{i \in I} A_i \right) = \sum_{i \in I} \mu(A_i), \forall I \subseteq \mathbb{N}$.

Then $\mu(A)$ is called *probability* of A . The probability space (as any measurable space) is called *complete* if for all $A \in \mathcal{A}$ with $\mu(A) = 0$ and all $B \subseteq A$ one has $B \in \mathcal{A}$. Every probability space can be completed by including all subsets B of null-probability sets A in the extended σ -algebra $\tilde{\mathcal{A}}$ and extending the probability measure μ to μ^* which amounts to 0 on the new sets B .

A random variable is a measurable function X from a complete probability space $(\Omega, \mathcal{A}, \mu)$ to the measurable space $(\mathbb{R}^n, \mathcal{B})$, where \mathcal{B} is the σ -algebra of Borel sets. We recall that a function $X : \Omega \rightarrow \mathbb{R}^n$ is measurable if

$$X^{-1}(B) \in \mathcal{A}, \quad \forall B \in \mathcal{B}.$$

Notice that a random variable *pushes forward* the probability measure to the Borel sets of \mathbb{R}^n , by the probability measure \mathbb{P}_X defined as

$$\mathbb{P}_X(B) := \mu[X^{-1}(B)]$$

More precisely, we construct \mathbb{P}_X as follows. Let us first define $\mathbf{a} < \mathbf{b}$ if $a_i < b_i$ for all $1 \leq i \leq n$, and then introduce the measurable sets $\mathcal{A} \ni A_{\mathbf{r}} := \{\omega \in \Omega | X(\omega) < \mathbf{r}\} = X^{-1}(\{\mathbf{a} \in \mathbb{R}^n | \mathbf{a} < \mathbf{r}\})$. The *cumulative* function $F_X(\mathbf{r}) := \mu(A_{\mathbf{r}}) = \mathbb{P}_X(X < \mathbf{r})$ defines a measure on \mathbb{R}^n , whose completion is the Lebesgue-Stieltjes measure induced by F_X .

The case of a discrete probability space is a special case, corresponding to a cumulative function F_X that is piecewise constant. In the case of discrete random variables it is sufficient to consider the case where $\text{Rng}(X) \subseteq \mathbb{N}$, as we will do in the following—unless otherwise stated.

Example 3.1. Let $\Omega = \{\square, \square, \square, \square, \square, \square\}$ and $\mu(\{f\}) = \frac{1}{6}$ for every face f of the dice. Then we can define $X(\square) = 1$, $X(\square) = 2$, $X(\square) = 3$, $X(\square) = 4$, $X(\square) = 5$, $X(\square) = 6$, and the corresponding cumulative function is

$$F_X(r) = \begin{cases} 0 & -\infty < r < 1, \\ \frac{1}{6} & 1 \leq r < 2, \\ \frac{1}{3} & 2 \leq r < 3, \\ \frac{1}{2} & 3 \leq r < 4, \\ \frac{2}{3} & 4 \leq r < 5, \\ \frac{5}{6} & 5 \leq r < 6, \\ 1 & 6 \leq r < +\infty. \end{cases}$$

Example 3.2. Let $\Omega = \{H, T\}$ with $\mu(\{H\}) = p$ and $\mu(\{T\}) = 1 - p$. Then we can define $X(H) = 0$, $X(T) = 1$, and the corresponding cumulative function is

$$F_X(r) = \begin{cases} 0 & -\infty < r < 0, \\ p & 0 \leq r < 1, \\ 1 & 1 \leq r < +\infty, \end{cases}$$

and the σ -algebra of Ω is $\mathcal{A} = \{\emptyset, \{H\}, \{T\}, \{H, T\}\}$.

Let us now consider a measurable function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$. Then, $f(X)$ is another random variable: $Y := f \circ X : \Omega \rightarrow \mathbb{R}^m$. The measurable sets defining the cumulative function in this case are $B_{\mathbf{r}} = \{\omega \in \Omega | f(X(\omega)) < \mathbf{r}\} = X^{-1}[f^{-1}(\{\mathbf{a} \in \mathbb{R}^m | \mathbf{a} < \mathbf{r}\})]$. The cumulative function $F_{f(X)}$ is then $F_{f(X)}(\mathbf{r}) = \mu(B_{\mathbf{r}}) = \mathbb{P}_X[f(X) < \mathbf{r}] = \mathbb{P}_{f(X)}[Y < \mathbf{r}]$.

Example 3.3. Let $f(x) = x^2$. Then $F_{f(X)}(y) = \mathbb{P}_X[X^2 < y]$. In this case we have

$$F_{f(X)}(y) = \begin{cases} 0 & y < 0, \\ \mathbb{P}_X[X^2 < y] & y \geq 0. \end{cases}$$

Since $\mathbb{P}_X[X^2 < y] = \mathbb{P}_X[|X| < \sqrt{y}] = \mathbb{P}_X[-\sqrt{y} < |X| < \sqrt{y}]$, we conclude that

$$F_{f(X)}(y) = F_X(\sqrt{y}) - F_X(-\sqrt{y}).$$

For a discrete random variable X the probability distribution \mathbf{p} is defined by $p_i := \mathbb{P}_X(X = x_i)$. With a slight abuse of notation, we will write $\mathbb{P}_X(x)$ instead of $\mathbb{P}_X(\{x\})$, thus we have $p_i = \mathbb{P}_X(X = x_i) = \mathbb{P}_X(x_i)$. In the following we will interchangeably use the notation \mathbf{p} or \mathbb{P}_X , depending on the context.

Thus, summarising, we have

Alphabet	Σ	a	b	...	z
Range	$\text{Rng}(X)$	1	2	...	26
Probability function	\mathbb{P}_X	$\mathbb{P}_X(X = 1)$	$\mathbb{P}_X(X = 2)$...	$\mathbb{P}_X(X = 26)$
Probability distribution	\mathbf{p}	p_1	p_2	...	p_{26}

Expectations

Given a random variable X with probability space $(\Omega, \mathcal{A}, \mu)$ and measure space $(\mathbb{R}^n, \mathcal{B})$, we define the *expectation value*, or simply *expectation* (sometimes called *mean* or *average*) of $f(X)$ for a measurable $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ as follows

$$\mathbb{E}[f(X)] = \int_{\Omega} d\mu(\omega) f[X(\omega)] := \int_{\mathbb{R}^n} dF_X(\mathbf{r}) f(\mathbf{r}), \quad (3.1)$$

where the latter integral is defined in the Lebesgue-Stieltjes sense, as F_X is monotonic and right continuous. Notice that

$$\mathbb{P}_X[A] = \mathbb{E}[\chi_A(X)], \quad (3.2)$$

where $\chi_A(\mathbf{r})$ is the *characteristic function* of the set A , i.e. $\chi_A(\mathbf{r}) = 1$ for $\mathbf{r} \in A$ and $\chi_A(\mathbf{r}) = 0$ for $\mathbf{r} \notin A$. If $f(x) \geq 0$ for every $x \in \text{Rng}(X)$, then clearly $\mathbb{E}[f(X)] \geq 0$. Notice that, given two functions f, g such that $f(\mathbf{r}) \leq g(\mathbf{r})$, one clearly has $\mathbb{E}[f(X)] \leq \mathbb{E}[g(X)]$. Indeed, $g(\mathbf{r}) - f(\mathbf{r}) \geq 0$, and $0 \leq \mathbb{E}[(g - f)(X)] = \mathbb{E}[g(X)] - \mathbb{E}[f(X)]$. An expectation playing a crucial role is the following

$$\sigma_X^2 := \mathbb{E}[(X - \mathbb{E}(X))^2] = \mathbb{E}(X^2) - \mathbb{E}(X)^2. \quad (3.3)$$

The quantity σ_X is called *root mean square* of X , and σ_X^2 is the *variance* of X .

In the case of a discrete random variable we have

$$\mathbb{E}[f(X)] = \sum_{\omega \in \Omega} \{\mu(\omega) f[X(\omega)]\} = \sum_{x_i \in \text{Rng}(X)} \{\mathbb{P}_X(x_i) f(x_i)\} = \sum_{x_i \in \text{Rng}(X)} \{p_i f(x_i)\}.$$

Independent random variables

Let \mathcal{A}_1 and \mathcal{A}_2 be two σ -algebras $\mathcal{A}_i \subseteq \mathcal{A}$ in a measurable space $(\Omega, \mathcal{A}, \mu)$. We say that \mathcal{A}_1 and \mathcal{A}_2 are *independent* if $\mu(A_1 \cap A_2) = \mu(A_1)\mu(A_2)$ for all $A_1 \in \mathcal{A}_1$ and $A_2 \in \mathcal{A}_2$. Let now Z be a random variable with range $\text{Rng}(Z) = \text{Rng}(X) \times \text{Rng}(Y) \subseteq \mathbb{R}^m \times \mathbb{R}^n$. We can define the two projections $\Pi_X : \text{Rng}(Z) \rightarrow \text{Rng}(X)$ and $\Pi_Y : \text{Rng}(Z) \rightarrow \text{Rng}(Y)$ defined by $\Pi_X(x, y) = x$ and $\Pi_Y(x, y) = y$. We say that X and Y are *independent random*

variables if both projections are measurable and $Z^{-1}[\Pi_X^{-1}(\mathcal{B}_X)]$ and $Z^{-1}[\Pi_Y^{-1}(\mathcal{B}_Y)]$ are independent σ -algebras with respect to the measure μ_Z , where \mathcal{B}_X and \mathcal{B}_Y are the Borel σ -algebras in $\text{Rng}(X)$ and $\text{Rng}(Y)$, respectively. The measures \mathbb{P}_X and \mathbb{P}_Y are given by $\mathbb{P}_X(A) := \mathbb{P}_Z[\Pi_X^{-1}(A)]$, and $\mathbb{P}_Y(B) := \mathbb{P}_Z[\Pi_Y^{-1}(B)]$. For discrete random variables this definition reduces to the intuitive notion given in the following.

Definition 3.4 (Marginal distribution). Let Z be a discrete random variable with $\text{Rng}(Z) = \text{Rng}(X) \times \text{Rng}(Y)$. The *marginal distributions* of X and Y are the probability distributions $\mathbb{P}_X(x) := \sum_{y \in \text{Rng}(Y)} \mathbb{P}_Z(x, y)$ and $\mathbb{P}_Y(y) := \sum_{x \in \text{Rng}(X)} \mathbb{P}_Z(x, y)$, respectively.

The definition is generalised to $\text{Rng}(Z) = \text{Rng}(X^{(1)}) \times \text{Rng}(X^{(2)}) \times \cdots \times \text{Rng}(X^{(N)})$ for finite N , as an easy exercise.

Definition 3.5 (Independent random variables). Let Z be a discrete random variable with $\text{Rng}(Z) = \text{Rng}(X) \times \text{Rng}(Y)$. We say that X and Y are *independent* if $\mathbb{P}_Z(x, y) = \mathbb{P}_X(x)\mathbb{P}_Y(y)$, where \mathbb{P}_X and \mathbb{P}_Y are the marginal distributions of \mathbb{P}_Z .

Also this definition can be easily generalised to the case of N independent random variables for arbitrary finite N .

We remind a remark that we already made in the previous lecture, i.e. that random variables Z with $\text{Rng}(Z) = \text{Rng}(X) \times \text{Rng}(Y)$ represent states of composite systems in our classical information theory. Independent random variables play a special role in the correspondence between classical states and random variables, since they represent *independent classical preparations*, represented diagrammatically as follows

$$\boxed{\mathbf{p}} \xrightarrow{Z} = \boxed{\mathbf{p}} \xrightarrow{X} \boxed{\mathbf{p}'} \xrightarrow{Y} \boxed{\mathbf{p}''}, \quad (3.4)$$

where \mathbf{p} , \mathbf{p}' , and \mathbf{p}'' denote the probability distributions \mathbb{P}_Z , \mathbb{P}_X , and \mathbb{P}_Y , respectively. On the other hand, a state of the composite system $Z = XY$ generally corresponds to a joint probability distribution $\mathbb{P}_Z(x, y) \neq \mathbb{P}_X(x)\mathbb{P}_Y(y)$ exhibiting *correlations*.

3.1.1 Inequalities and probability bounds

In this subsection we introduce some bounds on probabilities that will appear ubiquitously in the proofs of theorems. Some of the bounds regard *i.i.d.* variables, namely a finite number N of variables with the same range and with the same distribution, and with no correlation. These variables correspond to N copies of the same system, all prepared in the same state. Let us give a formal definition.

Definition 3.6 (I.i.d. random variables). Let X be a random variable with range $\text{Rng}(X) = \{x_1, x_2, \dots, x_n\} \subseteq \mathbb{R}$ and probability distribution $\mathbb{P}_X(X = x_i) = p_i$. We define the *identically independently distributed* (i.i.d.) random variables $X^{(1)}X^{(2)}\dots X^{(N)}$, as a single random variable X^N , with range $\text{Rng}(X^N) = \text{Rng}(X)^{\times N}$ and

$$\mathbb{P}_{X^N}(x_{i_1}, x_{i_2}, \dots, x_{i_N}) := \mathbb{P}_X(x_{i_1})\mathbb{P}_X(x_{i_2})\dots\mathbb{P}_X(x_{i_N}). \quad (3.5)$$

In the following for every $\mathbf{i} := i_1 i_2 \dots i_N$ we define $x_{\mathbf{i}} := x_{i_1} x_{i_2} \dots x_{i_n}$. Thus $\mathbb{P}_{X^N}(x_{\mathbf{i}}) = \mathbb{P}_{X^N}(x_{i_1}, x_{i_2}, \dots, x_{i_n})$.

Lemma 3.7 (Chebyshev inequality 1 - Markov Inequality). *For a non-negative random variable X and any positive real number a , one has*

$$\mathbb{P}_X[X \geq a] \leq \frac{\mathbb{E}(X)}{a}. \quad (3.6)$$

Proof. We remind (3.2), which implies that $\mathbb{P}_X[X \geq a] = \mathbb{E}[\theta(X - a)]$, where $\theta(x)$ is the Heavyside step function. By definition, the function satisfies the bound

$$\theta(x - a) \leq \frac{x}{a}, \quad (3.7)$$

and thus, taking the expectation of both sides and remembering the properties of expectation values, we obtain the thesis. \square

Lemma 3.8 (Chebyshev inequality 2). *For a random variable X and any positive real number a , one has*

$$\mathbb{P}_X[(X - \mathbb{E}[X])^2 \geq a] \leq \frac{\sigma_X^2}{a}, \quad \sigma_X = \mathbb{E}[(X - \mathbb{E}(X))^2]^{\frac{1}{2}}. \quad (3.8)$$

Proof. This is just a special case of Lemma 3.7 for the random variable $(X - \mathbb{E}[X])^2$. A common way of stating the Chebyshev bound is also

$$\mathbb{P}_X[|X - \mathbb{E}[X]| \geq k\sigma_X] \leq \frac{1}{k^2}. \quad \square$$

Lemma 3.9 (Weak law of large numbers). *Let $X^{(1)}, X^{(2)}, \dots, X^{(N)}$ be N i.i.d. random variables, with mean $\mathbb{E}(X) = \bar{x}$ and variance $\sigma_X^2 = \sigma^2$. Then for the random variable $Z := \frac{1}{N} \sum_{i=1}^N X^{(i)}$ and any positive real number a one has*

$$\mathbb{P}_Z[(Z - \bar{x})^2 \geq a] \leq \frac{\sigma^2}{aN}. \quad (3.9)$$

Proof. First, notice that using the definition of i.i.d. random variables one can easily find $\mathbb{E}(Z) = \bar{x}$. Thus, applying the Chebyshev bound of lemma 3.8 one has

$$\mathbb{P}_Z[(Z - \bar{x})^2 \geq a] \leq \frac{\sigma_Z^2}{a} \quad (3.10)$$

We then evaluate σ_Z^2 as follows.

$$\begin{aligned}\sigma_Z^2 &= \mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N X^{(i)} - \bar{x} \right)^2 \right] = \mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N (X^{(i)} - \bar{x}) \right)^2 \right] = \frac{1}{N^2} \sum_{i,j=1}^N \mathbb{E}[(X^{(i)} - \bar{x})(X^{(j)} - \bar{x})] \\ &= \frac{1}{N^2} \left\{ \sum_{i=1}^N \mathbb{E}[(X^{(i)} - \bar{x})^2] + \sum_{i=1}^N \sum_{j \neq i} \mathbb{E}[(X^{(i)} - \bar{x})(X^{(j)} - \bar{x})] \right\} \\ &= \frac{\sigma^2}{N}.\end{aligned}$$

Finally, substituting the last result in equation (3.10) one has the thesis. \square

Lemma 3.10 (Chernoff bound). *Suppose that $X^{(1)}, X^{(2)}, \dots, X^{(N)}$ are binary i.i.d. random variables, with range $\text{Rng}(X^{(i)}) = \{0, 1\}$ and $\mathbb{P}_X(0) = \frac{1}{2} + \varepsilon$, and $\mathbb{P}_X(1) = \frac{1}{2} - \varepsilon$, with $0 \leq \varepsilon \leq \frac{1}{2}$. Then for the random variable $Z := \sum_{i=1}^N X^{(i)}$ one has*

$$\mathbb{P}_Z \left[Z \geq \frac{N}{2} \right] \leq e^{-2\varepsilon^2 N}. \quad (3.11)$$

Proof. By the Markov inequality (3.6) we have

$$\mathbb{P}_Z[Z \geq c] = \mathbb{P}_Z[e^{tZ} \geq e^{tc}] \leq \frac{\mathbb{E}(e^{tZ})}{e^{tc}}, \quad t > 0.$$

Since the variables $X^{(i)}$ are i.i.d., we have

$$\mathbb{E}(e^{tZ}) = \mathbb{E} \left(e^{t \sum_{i=1}^N X^{(i)}} \right) = \mathbb{E} \left(\prod_{i=1}^N e^{tX^{(i)}} \right) = \prod_{i=1}^N \mathbb{E}(e^{tX^{(i)}}),$$

and for every variable $X^{(i)}$ one has

$$\mathbb{E}(e^{tX^{(i)}}) = \left(\frac{1}{2} + \varepsilon \right) + \left(\frac{1}{2} - \varepsilon \right) e^t. \quad (3.12)$$

This implies that

$$\mathbb{P}_Z[Z \geq c] \leq f(t), \quad f(t) := \frac{[(\frac{1}{2} + \varepsilon) + (\frac{1}{2} - \varepsilon)e^t]^N}{e^{tc}}. \quad (3.13)$$

For $c < N$, it is a simple exercise to find the minimum of f , which for $c = N/2$ corresponds to $t_0 = \ln \frac{\frac{1}{2} + \varepsilon}{\frac{1}{2} - \varepsilon}$. Since the bound (3.13) holds for every $t > 0$, choosing $t = t_0$ we have

$$\mathbb{P}_Z[Z \geq \frac{N}{2}] \leq \frac{[1 + 2\varepsilon]^N}{\left(\frac{1}{2} + \varepsilon\right)^{\frac{N}{2}}} = [1 - 4\varepsilon^2]^{\frac{N}{2}}.$$

We now apply a bound that will be recurrent, and is known as the *fundamental inequality of information theory*, namely $1 - x \leq e^{-x}$, which for $x = 4\varepsilon^2$ gives $1 - 4\varepsilon^2 \leq e^{-4\varepsilon^2}$, and

$$\mathbb{P}_Z[Z \geq \frac{N}{2}] \leq e^{-2\varepsilon^2 N}. \quad \square$$

3.1.2 Convex functions

Let us now recapitulate some elementary properties of convex functions. First of all, we remind the definition of convex set and convex function.

Definition 3.11 (Convex set). Let $C \subseteq \mathbb{R}^n$. The set C is *convex* if for any couple of points $\mathbf{x}_1, \mathbf{x}_2 \in C$ and every $0 \leq \lambda \leq 1$ the point $\lambda\mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2$ belongs to C .

Definition 3.12 (Epigraph). Let $f : D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$. The *epigraph* of f is the set $\text{Epi}(f) \subseteq \mathbb{R}^{n+1}$ defined as

$$\text{Epi}(f) := \{(\mathbf{r}, y) \in D \times \mathbb{R} \mid y \geq f(\mathbf{r})\}$$

Definition 3.13 (Convex function). Let $f : D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$. We say that f is *convex* if $\text{Epi}(f)$ is convex.

Notice that if the epigraph of f is convex, given $\mathbf{x}_1, \mathbf{x}_2 \in D$ $\lambda(\mathbf{x}_1, f(\mathbf{x}_1)) + (1 - \lambda)(\mathbf{x}_2, f(\mathbf{x}_2)) = (\lambda\mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2, \lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2)) \in \text{Epi}(f)$, and thus $\lambda\mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2 \in D$, namely the domain D of f must be convex.

Proposition 3.14. A function $f : D \rightarrow \mathbb{R}$ is convex if and only if D is convex and for every $0 \leq \lambda \leq 1$ and every $\mathbf{x}_1, \mathbf{x}_2 \in D$,

$$f[\lambda\mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2] \leq \lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2).$$

Proof. Let us prove necessity. Notice that for every $\mathbf{x}_1, \mathbf{x}_2 \in D$ we have $(\mathbf{x}_i, f(\mathbf{x}_i)) \in \text{Epi}(f)$. If $\text{Epi}(f)$ is convex, then

$\lambda(\mathbf{x}_1, f(\mathbf{x}_1)) + (1 - \lambda)(\mathbf{x}_2, f(\mathbf{x}_2)) = (\lambda\mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2, \lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2)) \in \text{Epi}(f)$, and by definition this implies $\lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2) \geq f[\lambda\mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2]$. Let us now prove sufficiency. Let $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \in \text{Epi}(f)$. Then $f(\mathbf{x}_i) \leq y_i$. The following chain of inequalities holds

$$f[\lambda\mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2] \leq \lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2) \leq \lambda y_1 + (1 - \lambda)y_2, \quad (3.14)$$

and then the point $(\lambda\mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2, \lambda y_1 + (1 - \lambda)y_2)$ belongs to $\text{Epi}(f)$, which is then convex. \square

Jensen's inequality

As a consequence of the properties of convex functions, Jensen's inequality holds. This crucial result is stated in the following theorem.

Theorem 3.15 (Jensen's inequality). Let $f : D \rightarrow \mathbb{R}$ be a convex function. For every set of points $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subseteq D$ and any probability distribution (p_1, p_2, \dots, p_n) (namely $p_j \geq 0$ for every $1 \leq j \leq n$ and $\sum_{i=1}^n p_i = 1$) the following inequality holds

$$f\left(\sum_{i=1}^n \{p_i \mathbf{x}_i\}\right) \leq \sum_{i=1}^n \{p_i f(\mathbf{x}_i)\}. \quad (3.15)$$

Proof. The proof is by induction on n . We know that the thesis holds for $n = 2$. Let the thesis hold for $n = n_0$, and consider the case $n = n_0 + 1$. Let us define $q := p_{n_0+1}$, so that the l.h.s. in equation (3.15) can be written as

$$f\left(\sum_{i=1}^n \{p_i \mathbf{x}_i\}\right) = f\left[(1-q)\sum_{i=1}^{n_0} \left\{\frac{p_i}{1-q} \mathbf{x}_i\right\} + q \mathbf{x}_{n_0+1}\right].$$

Notice that $p_i/(1-q) \geq 0$ and $\sum_{i=1}^{n_0} \{p_i/(1-q)\} = 1$, and since the domain D is convex, the point $\sum_{i=1}^{n_0} \left\{\frac{p_i}{1-q} \mathbf{x}_i\right\}$ belongs to D . Thus by proposition 3.14 we have

$$f\left(\sum_{i=1}^n \{p_i \mathbf{x}_i\}\right) \leq (1-q)f\left(\sum_{i=1}^{n_0} \left\{\frac{p_i}{1-q} \mathbf{x}_i\right\}\right) + qf(\mathbf{x}_{n_0+1}).$$

Now, using the induction hypothesis, we have

$$f\left(\sum_{i=1}^{n_0} \left\{\frac{p_i}{1-q} \mathbf{x}_i\right\}\right) \leq \sum_{i=1}^{n_0} \left\{\frac{p_i}{1-q} f(\mathbf{x}_i)\right\},$$

and finally

$$f\left(\sum_{i=1}^n \{p_i \mathbf{x}_i\}\right) \leq (1-q)\sum_{i=1}^{n_0} \left\{\frac{p_i}{1-q} f(\mathbf{x}_i)\right\} + qf(\mathbf{x}_{n_0+1}) = \sum_{i=1}^{n_0+1} \{p_i f(\mathbf{x}_i)\},$$

which concludes the proof. \square

Jensen's inequality can be immediately re-phrased in the language of discrete random variables as follows

Corollary 3.16 (Jensen inequality for random variables). *Let X be a random variable with range $\text{Rng}(X) = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$. Let $f : D \rightarrow \mathbb{R}$ be a convex function, where D contains the convex hull of $\text{Rng}(X)$. Then one has*

$$f[\mathbb{E}(X)] \leq \mathbb{E}[f(X)]. \quad (3.16)$$

Log-sum inequality

We now apply Jensen's inequality to prove the *log-sum* inequality.

Lemma 3.17 (Log-sum inequality). *Let $\{a_i\}_{i=1}^n \subseteq \mathbb{R}$ satisfy $a_i \geq 0$ for all $1 \leq i \leq n$, and $\{b_i\}_{i=1}^n \subseteq \mathbb{R}$ satisfy $b_i > 0$ for all $1 \leq i \leq n$. Let us also define $a := \sum_{i=1}^n a_i$ and $b := \sum_{i=1}^n b_i$. Then the following inequality holds*

$$\sum_{i=1}^n \left\{a_i \log \frac{a_i}{b_i}\right\} \geq a \log \frac{a}{b}, \quad (3.17)$$

where the logarithm is taken in any basis $q > 1$.

Proof. First of all, we define $0 \log 0 := 0$ by continuity. Now, let us define the function $f : [0, 1] \rightarrow \mathbb{R}$ as $f(x) := x \log x$. It is easy to verify that f is convex, since it is twice differentiable in $(0, 1)$ and the second derivative is $f''(x) = 1/x \geq 0$. We can now write

$$\sum_{i=1}^n \left\{ a_i \log \frac{a_i}{b_i} \right\} = b \sum_{i=1}^n \left\{ \frac{b_i}{b} \frac{a_i}{b_i} \log \frac{a_i}{b_i} \right\} = b \sum_{i=1}^n \left\{ q_i f \left(\frac{a_i}{b_i} \right) \right\}, \quad q_i := \frac{b_i}{b}.$$

Applying Jensen's inequality to f we then obtain

$$\sum_{i=1}^n \left\{ a_i \log \frac{a_i}{b_i} \right\} \geq b f \left(\sum_{i=1}^n \left\{ q_i \frac{a_i}{b_i} \right\} \right) = b f \left(\frac{a}{b} \right) = a \log \frac{a}{b}.$$

This completes the proof. \square

Chapter 4

Lecture 4: Shannon entropy and typical sequences

4.1 Shannon entropy

In this section we introduce some preliminary notions that will serve the purpose of quantifying the information content of a source. We give some intuitive understanding of the main measures of information content, and derive their most relevant properties. However, we need to keep in mind that the precise theoretical and practical interpretation of such quantifications has to be found in theorems—primarily the two Shannon theorems and the data-processing inequality—which give precise operational meaning to the Shannon entropy and the mutual information.

In the following, we will often consider the probabilities of all the characters in the alphabet Σ as strictly positive. Reasoning in terms of the random variable X that maps the alphabet into $\text{Rng}(X)$, this means that $p_x > 0$ for all $x \in \text{Rng}(X)$, and treat the special case in which $p_x = 0$ for some x by the alphabet Σ' obtained by removing the symbol $X^{-1}(x)$ from the alphabet Σ , i.e. $\Sigma' := \Sigma \setminus X^{-1}(x)$. Since the range of X will be $\{1, 2, \dots, |\Sigma|\}$, in the above case we define a new random variable X' with $\text{Rng}(X') = \{1, 2, \dots, |\Sigma| - 1\}$. Finally, keep in mind that in the following we will often use the shorthand $\mathbb{P}_X(x)$ instead of $\mathbb{P}_X[\{x\}]$ or $\mathbb{P}_X[X = x]$.

4.1.1 Shannon information content and Shannon entropy

Definition 4.1 (Shannon information content). The *Shannon information content* of the discrete random variable X is the function $h_X : \text{Rng}(X) \rightarrow \mathbb{R}$ defined as

$$h_X(x) := \log_2 \frac{1}{\mathbb{P}_X(x)} . \quad (4.1)$$

When there will be no ambiguity in the specification of the random variable, we will omit it, and write $h(x)$ instead of $h_X(x)$. The Shannon information content was originally called by Shannon *self-information*, and it is often referred to as *surprisal*. Indeed, the smaller $\mathbb{P}_X(x)$, the larger is the surprisal upon occurrence of the event x .

Notice that since $0 < \mathbb{P}_X(x) \leq 1$, one has $h_X(x) \geq 0$ for every $x \in \text{Rng}(X)$.

Definition 4.2 (Shannon entropy). The Shannon entropy of a discrete random variable X is defined as

$$H(X) := \sum_{x \in \text{Rng}(X)} \left\{ \mathbb{P}_X(x) \log_2 \frac{1}{\mathbb{P}_X(x)} \right\} = - \sum_{x \in \text{Rng}(X)} \{ \mathbb{P}_X(x) \log_2 P_X(x) \}. \quad (4.2)$$

Notice that $H(X) = \sum_{x \in \text{Rng}(X)} \{ \mathbb{P}_X(X = x) h_X(x) \} = \mathbb{E}[h_X(X)]$, namely the Shannon entropy is equal to the Shannon information content averaged over all possible outcomes. We will sometimes write $H(\mathbf{p})$ instead of $H(X)$, when we need to emphasise the dependence of the Shannon entropy on the probability distribution. Indeed, notice that the definition never involves the range of X , only depending on the probability measure \mathbb{P}_X . This is a property that we expect from a quantity that measures the information content of a source: it must not depend on the symbols used for the encoding.

Being the expectation of the nonnegative function $h_X(X)$, the Shannon entropy is non-negative $H(X) \geq 0$. In the following we adopt the convention $0 \log_2 0 = 0$, which extends the function $x \log_2 x$ by continuity in 0, and thus $H(X) = 0$ for probability distributions \mathbf{p} that are concentrated at one value $x_i \in \text{Rng}(X)$, i.e. $p_j = \delta_{i,j}$.

Intuitive meaning of the information content

The actual meaning of the Shannon entropy will be given by the first Shannon theorem, in terms of compression ratio of a source. However, we can already provide a preliminary analysis of its most intuitive features.

Remark 1 (Additivity). What has $h_X(x)$ to do with the information content? First of all, notice that if one has two independent random variables X and Y , namely with probability distributions that are factorised $\mathbb{P}_{XY}(x, y) = \mathbb{P}_X(x)\mathbb{P}_Y(y)$, the function $h_{X,Y}(x, y)$ of the corresponding states is, by definition, additive, and this is what we expect from a quantity that measures the information content: when we read two independent messages the amount of information that we gain is the sum of the two amounts of information gained by reading each of the two messages separately. The additivity property is clearly inherited by the Shannon entropy. Therefore, for two *independent* variables X and Y , we have

$$h_{X,Y}(x, y) = h_X(x) + h_Y(y), \quad H(X, Y) = H(X) + H(Y).$$

Example 4.3 (Battleship). Consider a simplified battleship game, where a player must locate a submarine on a grid of 8×8 squares. Notice that all the cells of the checkerboard can be addressed by strings of 6 bits. At each round, the player can query whether the target is located in one square. They receive two possible answers: h (hit) and m (missed). At the first round the probability of finding the submarine at any square is $p_i = \frac{1}{64}$, and the Shannon information content of the event “hit” is

$$h_{(1)}(h) = \log_2 64 = 6 \text{ bits.}$$

Indeed, we have actually gained 6 bits of information, namely the 6 bits needed to specify the exact location of the submarine in the grid. On the other hand, the Shannon

information content of the event “missed” is

$$h_{(1)}(m) = \log_2 \frac{64}{63} = 0.0227 \text{ bits.}$$

Let us now play the second round. Now we have $p = \frac{1}{63}$ for all the remaining squares and $p = 0$ for the queried one. Then,

$$h_{(2)}(m) = \log_2 \frac{63}{62} = 0.0230 \text{ bits.}$$

Now, suppose that the player misses the submarine for 32 rounds. Then they gained

$$h_{(1)}(m) + h_{(2)}(m) + \cdots + h_{(32)}(m) = \log_2 \frac{64}{63} + \log_2 \frac{63}{62} + \cdots + \log_2 \frac{33}{32} = \log_2 2 = 1 \text{ bits.}$$

Indeed, the board is now divided into two parts: the queried squares and the remaining ones, that can be addressed by one bit. We know that the submarine is located in the second half, i. e. we have one bit of information. Now, it is easy to calculate that whenever the player hits the submarine, say at round k , the total information gained is

$$\begin{aligned} & h_{(1)}(m) + h_{(2)}(m) + \cdots + h_{(k-1)}(m) + h_{(k)}(h) \\ &= \log_2 \frac{64}{63} + \log_2 \frac{63}{62} + \cdots + \log_2 \frac{64-k+1}{64-k} + \log_2(64-k) \\ &= \log_2 64 = 6 \text{ bits,} \end{aligned}$$

independently of k . Notice that 6 is the number of binary questions (questions with yes/no answer) needed to identify a square on the grid, e.g. by subsequent bipartitions.

Example 4.4 (Binary Shannon entropy). For a binary random variable X with range $\text{Rng}(X) = \{0, 1\}$ and probability distribution $p(0) = p$ and $p(1) = 1 - p$, the Shannon entropy is a function of the parameter p . More precisely

$$H_2(p) := H(\{p, 1-p\}) = p \log_2 \frac{1}{p} + (1-p) \log_2 \frac{1}{1-p}. \quad (4.3)$$

From the plot of $H_2(p)$ in Fig. 4.1 we can see that the binary Shannon entropy is a concave function of p , with maximum at $p = \frac{1}{2}$, where $H_2(p) = 1$ bits. Indeed, the information that we gain from tossing a fair coin is 1 bit. On the other hand, for a biased coin the information gain from one toss is less than one bit, as one can understand analysing the extreme case in which the outcome “head” is certain: In this case we gain no information at all from one toss. More generally one can prove that $H(\mathbf{p})$ is a concave function of \mathbf{p} also for $\mathbf{p} \in \mathbb{R}^n$:

$$H(q\mathbf{p}_0 + (1-q)\mathbf{p}_1) \geq qH(\mathbf{p}_0) + (1-q)H(\mathbf{p}_1).$$

This means that when mixing two random variables with probabilities q and $1 - q$, one obtains a new random variable with an information content which at least equals the average information of each variable. In other words, we have more ignorance about the

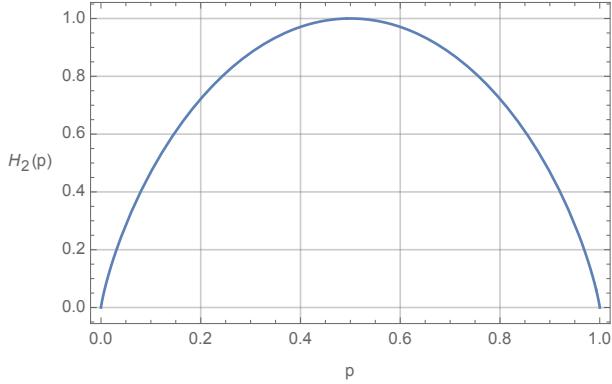


Figure 4.1 The binary entropy function $H_2(p)$ in Eq. (4.3).

new variable than the average ignorance about the original ones. For example, one can think of the new variable as obtained by the following procedure. First, a referee samples a binary variable C with $\text{Rng}(C) = \{0, 1\}$ and probabilities $(q, 1 - q)$. Then, if $c = 0$, they sample the variable B_0 with $\text{Rng}(B_0) = \{0, 1\}$ and distribution \mathbf{p}_0 , while if $c = 1$ they sample B_1 with $\text{Rng}(B_1) = \{0, 1\}$ and distribution \mathbf{p}_1 . Finally, after erasing the outcome c , they produce the outcome of the second sampled variable B_c . It is then clear that in this scenario, once we read the final outcome b_c , distributed according to $q\mathbf{p}_0 + (1 - q)\mathbf{p}_1$, we ignore the value of c , i.e. which of the two variables B_0 or B_1 was actually sampled in the given instance.

From what we have seen, we can say that the Shannon entropy $H(X)$ measures our prior uncertainty about X , or else $H(X)$ measures the amount of information that one gains from knowledge of X .

4.1.2 Relative entropy

We now introduce a quantity that has no immediate operational meaning, but turns out to be a powerful tool for proofs. This is the *relative entropy*, also known as *Kullback-Leibler divergence*.

Definition 4.5 (Kullback-Leibler divergence or relative entropy). Let \mathbf{p} and \mathbf{q} be two probability distributions for random variables whose range has cardinality n . Their *Kullback-Leibler divergence*, denoted as $H(\mathbf{p} \parallel \mathbf{q})$, is defined as

$$H(\mathbf{p} \parallel \mathbf{q}) := \begin{cases} \sum_{i=1}^n \left\{ p_i \log_2 \frac{p_i}{q_i} \right\} & q_i > 0 \text{ for all } i \text{ such that } p_i > 0, \\ +\infty & \text{otherwise,} \end{cases} \quad (4.4)$$

where we adopt the convention $0 \log_2 0 = 0 \log_2 \frac{0}{q} = 0 \log_2 \frac{0}{0} = 0$.

The relative entropy $H(X \parallel Y)$ of two random variables X and Y having probability distributions \mathbf{p} and \mathbf{q} , respectively, is defined as the Kullback-Leibler divergence $H(\mathbf{p} \parallel \mathbf{q})$. Intuitively speaking, the Kullback-Leibler divergence $H(\mathbf{p} \parallel \mathbf{q})$ measures how close the two

probability distributions \mathbf{p} and \mathbf{q} are. Nevertheless, it does not satisfy the requirements for a distance.

The main property of the Kullback-Leibler divergence is given by the following lemma

Lemma 4.6 (Gibbs' inequality). *The Kullback-Leibler divergence $H(\mathbf{p}\|\mathbf{q})$ is nonnegative and vanishes only if $\mathbf{p} = \mathbf{q}$.*

Proof. If $q_i = 0$ for some i with $p_i > 0$, the inequality is trivially satisfied. We then focus on the case where $q_i > 0$ for all i such that $p_i > 0$. Let us use the fundamental inequality of information theory in the following form

$$\ln \frac{1}{x} \geq 1 - x,$$

and remind that $\ln y = \log_2 y \ln 2$, which implies

$$\log_2 \frac{1}{x} \geq \frac{1 - x}{\ln 2}.$$

Now, we can write

$$H(\mathbf{p}\|\mathbf{q}) = \sum_{i=1}^n p_i \log_2 \frac{p_i}{q_i} \geq \frac{1}{\ln 2} \sum_{i=1}^n \left\{ p_i \left(1 - \frac{q_i}{p_i} \right) \right\} = \frac{1}{\ln 2} \sum_{i=1}^n \{p_i - q_i\} = 0.$$

Since the only inequality that we used is the fundamental inequality of information theory, which is tight only for $x = 1$, we conclude that also $H(\mathbf{p}\|\mathbf{q}) = 0$ only if $\frac{p_i}{q_i} = 1$ for all i , namely $\mathbf{p} = \mathbf{q}$. \square

Exercise 4.1

Let $M \in M_{m \times n}(\mathbb{R})$ be a stochastic (or Markov matrix), namely $M_{ij} \geq 0$ and $\sum_{i=1}^m M_{ij} = 1$ for all j . Prove that

$$H(M\mathbf{p}\|M\mathbf{q}) \leq H(\mathbf{p}\|\mathbf{q}) \tag{4.5}$$

Answer of exercise 4.1

Fix a value of i and apply the log-sum inequality to the expression

$$\sum_l M_{il} p_l \log \frac{\sum_k M_{ik} p_k}{\sum_s M_{is} q_s} \leq \sum_l M_{il} p_l \log \frac{M_{il} p_l}{M_{il} q_l} = \sum_l M_{il} p_l \log \frac{p_l}{q_l}.$$

Taking now the sum over i on both sides, we obtain $H(M\mathbf{p}\|M\mathbf{q}) \leq H(\mathbf{p}\|\mathbf{q})$.

The second property that we prove provides an interpretation of the Kullback-Leibler divergence as a sort of distance between probability distributions, despite it is not symmetric and therefore it cannot define a metric.

4.2 Bound on Shannon entropy

The relative entropy is not useful by itself, but only because other entropic quantities can be expressed in terms of it. As we will see, it is mostly a proving tool. The bounds that we prove here are important examples.

Lemma 4.7. *The Shannon entropy of a random variable X satisfies the following bound*

$$H(X) \leq \log_2 |\text{Rng}(X)|. \quad (4.6)$$

Moreover, $H(X) = \log_2 |\text{Rng}(X)|$ only for the uniform distribution $\mathbb{P}_X(x) = \frac{1}{|\text{Rng}(X)|}$ for all $x \in \text{Rng}(X)$.

Proof. Let us set $r := |\text{Rng}(X)|$, and define the random variable Y with $\text{Rng}(Y) = \text{Rng}(X)$ and $\mathbb{P}_Y(x) = \frac{1}{r}$ for all $x \in \text{Rng}(X)$. If we evaluate the relative entropy $H(X\|Y)$ we find that

$$\begin{aligned} H(X\|Y) &= \sum_{x \in \text{Rng}(X)} \left\{ \mathbb{P}_X(x) \log_2 \frac{\mathbb{P}_X(x)}{\mathbb{P}_Y(x)} \right\} = -H(X) + \sum_{x \in \text{Rng}(X)} \{ \mathbb{P}_X(x) \log_2 r \} \\ &= -H(X) + \log_2 r. \end{aligned}$$

Since, by Gibbs' inequality, $H(X\|Y) \geq 0$, we have the first statement of the thesis. Moreover, equality holds if and only if the two probability distributions are the same, namely $\mathbb{P}_X(x) = \frac{1}{|\text{Rng}(X)|}$. \square

4.2.1 Raw bit content

Definition 4.8 (Raw bit content). The *raw bit content* of a random variable X is defined as $H_0(X) := \log_2 |\text{Rng}(X)|$.

If we neglect the “surprisal” due to our prior expectations about the outcome of a random variable, a plain count of the number of bits needed for a perfect lossless encoding of the alphabet of X is precisely its raw bit content. Notice that the raw bit content is additive: For the variable Z = XY one has

$$H_0(X, Y) = \log_2 (|\text{Rng}(X)| |\text{Rng}(Y)|) = H_0(X) + H_0(Y).$$

The meaning of the raw bit content is then very intuitive: We need one string for every character in our alphabet, which means $|\text{Rng}(X)|$ different strings. Since N bits correspond to an alphabet with 2^N symbols, $|\text{Rng}(X)|$ strings require $H_0(X)$ bits for an exact encoding. The question is then: how is it possible to compress a string without loosing information?

4.3 Lossy and lossless compressors

Formally, a compression scheme $(\mathcal{E}, \mathcal{D})$ is defined in terms of a couple of channels: the encoding $\mathcal{E} : \text{Rng}(X) \rightarrow \bigcup_{N \in \mathbb{N}} \{0, 1\}^N$ and the decoding $\mathcal{D} : \bigcup_{N \in \mathbb{N}} \{0, 1\}^N \rightarrow \text{Rng}(X)$. The length $l[\mathcal{E}(x)]$ of the encoded character $x \in \text{Rng}(X)$ is denoted as $l_{\mathcal{E}}(x)$.

Clearly, by the simple counting argument that we used to derive the raw bit content, it is easy to see that it is impossible to make a reversible compression that reduces the size of every possible string. There are then only two kinds of compressors:

1. **Lossy compressor.** A lossy compressor compresses some strings to the same fixed length, but there are strings which after compression would be confusable. In some applications (e.g. image compression) the occurrence of confusable strings can be afforded. We will denote by δ the probability that the string belongs to the confusable set—the so-called probability of failure. The lossy compressor is practically useful if the probability of failure δ can be made small.
2. **Lossless compressor.** A lossless compressor maps all strings to a different encoding. Hence, while shortening some strings, it necessarily makes some other strings longer. The compressor is still useful if the probability δ that a string is lengthened is very small, and the probability $1 - \delta$ that it is shortened is high, as long as the *average* length of encoded strings is small.

In the present chapter we will focus on *lossy* compressors. In the next chapter we will analyse *lossless* compressors.

Example 4.9. Consider an alphabet of eight characters $A = \{a, b, c, d, e, f, g, h\}$, not equally distributed. A single-character string is described by the random variable X with

$$\text{Rng}(X) = \{1, 2, 3, 4, 5, 6, 7, 8\} \quad (4.7)$$

$$\mathbf{p} = \left\{ \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{3}{16}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64} \right\}. \quad (4.8)$$

The raw bit content is $H_0(X) = 3$ bits, but if we allow for a probability $\delta = \frac{1}{16}$ of loss, we can compress the subset $\{a, b, c, d\} \subset A$, loosing $\{e, f, g, h\}$. Then it is possible to compress with a small error probability to only 2 bits.

4.3.1 Smallest δ -sufficient set and essential bit content

The idea in the above example is captured by the notion of *essential bit content*, defined in the following. First of all, if we want a small probability of error—say $p_{\text{err}} \leq \delta$ —we must identify all possible subsets of our alphabet with probability larger than $1 - \delta$. A choice of such set will correspond to the set of strings that will be shortened without loss.

Definition 4.10 (δ -sufficient sets). Given a random variable X with probability distribution \mathbb{P}_X , for $0 \leq \delta \leq 1$ the family of δ -sufficient sets of X is the set

$$\mathcal{S}_\delta(X) := \{S \subseteq \text{Rng}(X) \mid \mathbb{P}_X[X \in S] \geq 1 - \delta\}. \quad (4.9)$$

Now, every choice of a δ -sufficient set will lead to a small error probability (in other words, the probability of confusable strings will be smaller than or equal to δ). However, in order to have an efficient compressor, among all δ -sufficient sets we would like to find the smallest one. Indeed, the smaller the set, the smaller is the number of bits that we will use to encode it.

Definition 4.11 (Smallest δ -sufficient set). Given a random variable X with probability distribution \mathbb{P}_X , for $0 \leq \delta \leq 1$ the *smallest δ -sufficient set* of X is the set

$$S_\delta(X) := \arg \min_{S \in \mathcal{S}_\delta(X)} |S| \quad (4.10)$$

In practice the subset $S_\delta(X)$ can be constructed by i) sorting the outcomes in order of decreasing probability, and ii) including elements in $S_\delta(X)$ starting from the most likely, until the total probability becomes larger than $1 - \delta$. Finding the smallest δ -sufficient set is not computationally easy, because it requires sorting the elements of $\text{Rng}(X)$ according to their probability. However, once the sorting has been completed, the smallest δ -sufficient set can be nicely visualized, as in the following Fig. 4.2 where $\text{Rng}(X) = \{1, 2, \dots, 7\}$ and $\mathbb{P}_X[x] = N/2^x$, with $N = 128/127$.

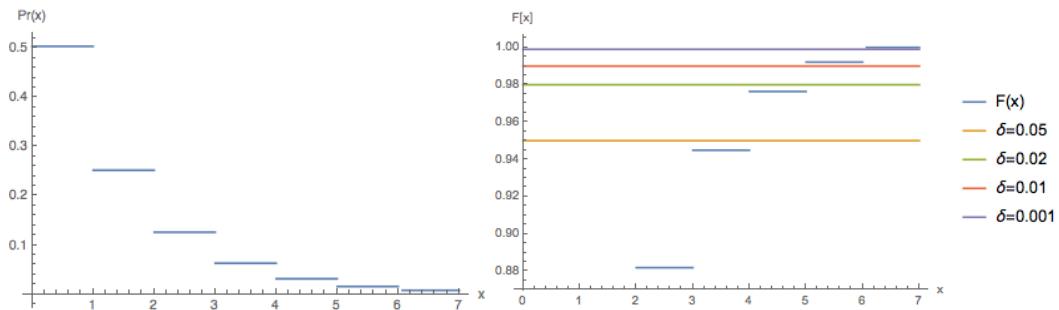


Figure 4.2 Plot on the left: Probability distribution $\mathbb{P}_X(x)$ over the sorted range of X . Right: Cumulative distribution F_X over the sorted range of X , evidencing the smallest δ -sufficient set: The first level above the line corresponding to a value of δ corresponds to the maximum value of x in $S_\delta(X)$. For example, for $\delta = 0.05$ one has $S_\delta = \{1, 2, 3, 4, 5\}$

We are finally in position to define the quantity that measures the number of bits needed for compression of the smallest δ -sufficient set.

Definition 4.12 (Essential bit content). Given a random variable X with probability distribution \mathbb{P}_X , for $0 \leq \delta \leq 1$ the *essential bit content* of X is defined as

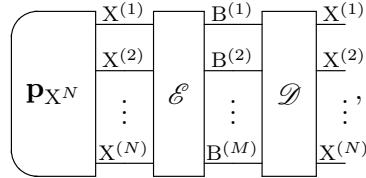
$$H_\delta(X) := \log_2 |S_\delta(X)|. \quad (4.11)$$

Remark 2. Notice that the notation for the raw bit content $H_0(X)$ is consistent with that for the essential bit content: indeed, the definition of raw bit content coincides with that of essential bit content in the special case $\delta = 0$.

From the above definitions we see that for affordable risk δ we can systematically achieve a compression length $H_\delta(X)$. In Example 4.9 we considered single-character strings. In the following we will more practically consider strings $x_i := x_{i_1} x_{i_2} \dots x_{i_N}$ of N characters, described by the random variable X^N made of N i.i.d. random variables equal to X . Then, for N very large, if X is not uniformly distributed, the least probable strings will be exponentially less probable than the least probable symbols, and this will make compression more and more reliable.

4.4 Source coding

A *lossy compression scheme* for strings of N i.i.d. random variables X^N is a pair of channels $(\mathcal{E}, \mathcal{D})$ with $\mathcal{E} : \text{Rng}(X)^N \rightarrow \text{Rng}(Y)$, and $\mathcal{D} : \text{Rng}(Y) \rightarrow \text{Rng}(X^N)$, with $\text{Rng}(Y) \subseteq \{0, 1\}^M$ for some integer M . Diagrammatically:



where $B^{(i)}$ is binary for every i . While X^N are i.i.d., the same is not true for B^M , thus we better think of the random variable Y .

Given a compression scheme $(\mathcal{E}, \mathcal{D})$, we define its *ratio* as follows.

Definition 4.13 (Data compression ratio). Let X^N be N i.i.d. random variables, and let $(\mathcal{E}, \mathcal{D})$ be a compression scheme for X^N . The *compression ratio* of the scheme $(\mathcal{E}, \mathcal{D})$ is

$$R_{\mathcal{E}}(X) := \frac{\log_2 |\mathcal{E}[\text{Rng}(X^N)]|}{N}. \quad (4.12)$$

The compression ratio quantifies the number of bits needed per character in a string. It is clear that this quantity accounts for the compression performances of the scheme $(\mathcal{E}, \mathcal{D})$. On the other hand, the performances of the compression scheme are also evaluated by the error probability, namely the probability that a string is compressed to a confusable string. The first Shannon theorem provides the smallest value for the compression ratio of a compression scheme for which the error probability can be made arbitrarily small.

4.4.1 Typical sets and asymptotic equipartition

The proof of Shannon's source coding theorem requires the notion of typical set, along with its properties.

Definition 4.14 (Typical set). Given a discrete random variable X with $\text{Rng}(X) = \{x_1, x_2, \dots, x_n\}$ distributed with probability $\{p_1, p_2, \dots, p_n\}$, we denote by $T_{N,\varepsilon}(X) \subseteq \text{Rng}(X^N)$ the *typical set to tolerance ε* defined as the following set of strings of length N

$$T_{N,\varepsilon}(X) := \left\{ x_i \in \text{Rng}(X^N) \mid \left| \frac{1}{N} \log_2 \frac{1}{P_{X^N}(x_i)} - H(X) \right| \leq \varepsilon \right\} \quad (4.13)$$

The elements of $T_{N,\varepsilon}(X)$ are called *typical strings*.

Using the weak law of large numbers we now see that the typical sequences have probabilities distributed according to the so-called asymptotic equipartition principle, which essentially states that the event $X^N = x_i = x_{i_1}x_{i_2} \dots x_{i_N}$ almost certainly belongs to a subset having only $2^{NH(X)}$ elements, each one having probability close to $2^{-NH(X)}$. More precisely the following theorem holds.

Theorem 4.15 (Asymptotic equipartition principle). *Let X be a random variable. The following three facts hold*

1. *The probability of strings $x_i \in T_{N,\varepsilon}$ is bounded as*

$$2^{-N[H(X)+\varepsilon]} \leq \mathbb{P}_{X^N}(x_i) \leq 2^{-N[H(X)-\varepsilon]}. \quad (4.14)$$

2. *The probability of the event $x_i \in T_{N,\varepsilon}(X)$ is lower bounded as*

$$\mathbb{P}_{X^N}[x_i \in T_{N,\varepsilon}] \geq 1 - \frac{\sigma_Y^2}{\varepsilon^2 N}, \quad (4.15)$$

where σ_Y^2 is the variance of the random variable $Y := h_X(X)$.

3. *For every $\delta > 0$, there exists N_0 such that for $N > N_0$ the cardinality of the typical set is bounded as*

$$(1 - \delta)2^{N[H(X)-\varepsilon]} \leq |T_{N,\varepsilon}(X)| \leq 2^{N[H(X)+\varepsilon]}. \quad (4.16)$$

Proof. Item 1: This is just a way of rewriting the defining condition of the typical set, indeed

$$\left| \frac{1}{N} \log_2 \frac{1}{\mathbb{P}_{X^N}(x_i)} - H(X) \right| \leq \varepsilon \Leftrightarrow -N[H(X) + \varepsilon] \leq \log_2 \mathbb{P}_{X^N}(x_i) \leq -N[H(X) - \varepsilon].$$

Item 2: Let us define the random variables $Y(X)$ with range $(h_X(x))_{x \in \text{Rng}(X)}$ and probability distribution $\mathbb{P}_Y[h_X(x)] = \mathbb{P}_X(x)$, and $Z(X^N) := \frac{1}{N} \sum_{i=1}^N Y(X^{(i)})$. Then we have

$$\begin{aligned} Z(x_i) &= \frac{1}{N} \log_2 \frac{1}{\mathbb{P}_Z(x_i)} = \frac{1}{N} \sum_{j=1}^N \log_2 \frac{1}{\mathbb{P}_X(x_{i_j})}, \\ \mathbb{E}[Z] &= \frac{1}{N} \sum_{j=1}^N \mathbb{E} \left[\log_2 \frac{1}{\mathbb{P}_X(x_{i_j})} \right] = H(X). \end{aligned}$$

The variable Z can then be used to rewrite the definition of typical set as follows

$$T_{N,\varepsilon} = \{x_i \mid [Z(x_i) - \mathbb{E}(Z)]^2 \leq \varepsilon^2\}.$$

Using the weak law of large numbers (lemma 3.9), we have

$$\mathbb{P}_{X^N}[x_i \in T_{N,\varepsilon}(X)] = \mathbb{P}_Z[(Z(x_i) - \mathbb{E}(Z))^2 \leq \varepsilon^2] \geq 1 - \frac{\sigma_Y^2}{N\varepsilon^2}.$$

Item 3: Given a value of δ , take N_0 such that $N_0 \geq \frac{\sigma_Y^2}{\delta\varepsilon^2}$. Then, by item 2, for $N > N_0$ the probability of the typical set $T_{N,\varepsilon}$ is bounded as

$$(1 - \delta) \leq \mathbb{P}_{X^N}[x_i \in T_{N,\varepsilon}(X)] \leq 1. \quad (4.17)$$

Thus, using the definition of typical set and the upper bound in equation (4.17), one can prove the following chain of inequalities

$$1 \geq \sum_{x_i \in T_{N,\varepsilon}(X)} \mathbb{P}_{X^N}(x_i) \geq \sum_{x_i \in T_{N,\varepsilon}(X)} 2^{-N[H(X+\varepsilon)]} = |T_{N,\varepsilon}(X)|2^{-N[H(X+\varepsilon)]}.$$

Thus $|T_{N,\varepsilon}(X)| \leq 2^{N[H(X+\varepsilon)]}$. Conversely, from the lower bound in equation (4.17) one can prove the following chain of inequalities

$$(1 - \delta) \leq \sum_{x_i \in T_{N,\varepsilon}(X)} \mathbb{P}_{X^N}(x_i) \leq \sum_{x_i \in T_{N,\varepsilon}(X)} 2^{-N[H(X)-\varepsilon]} = |T_{N,\varepsilon}(X)|2^{-N[H(X)-\varepsilon]}.$$

Thus $(1 - \delta)2^{N[H(X)-\varepsilon]} \leq |T_{N,\varepsilon}(X)|$. \square

4.4.2 Shannon source coding theorem

Using the asymptotic equipartition theorem, we can now prove the first Shannon theorem, known as *source coding theorem*. The theorem states that the essential bit content of N i.i.d. random variables X^N can be made arbitrarily close to $NH(X)$ for any value δ of the error probability, by taking long enough sequences. Since the essential bit content is the number of bits per character—or better, in the case of interest, per length N -string—in a reliable compression scheme with probability δ of error, the theorem provides the operational interpretation of the Shannon entropy, that is the smallest reliable asymptotic compression ratio for an information source represented by the random variable X . In other words, given a classical source whose state corresponds to a random variable X with probability distribution \mathbb{P}_X , one can find a family $(\mathcal{E}^N, \mathcal{D}^N)$ of encodings for strings of length N such that the error probability for these schemes is arbitrarily small for large N , while their compression ratio gets arbitrarily close to $H(X)$.

Let us now see the precise statement and the proof of the theorem.

Theorem 4.16 (Shannon's source coding theorem). *Let X a random variable with Shannon entropy $H(X)$. Given $\varepsilon > 0$ and $0 < \delta < 1$, there exists a positive integer N_0 such that for $N > N_0$ one has*

$$\left| \frac{1}{N} H_\delta(X^N) - H(X) \right| \leq \varepsilon \quad (4.18)$$

Proof. The statement is equivalent to the following chain of inequalities

$$H(X) - \varepsilon \leq \frac{1}{N} H_\delta(X^N) \leq H(X) + \varepsilon$$

The proof is split in two parts, in which we prove the two inequalities separately.

- *Proof of $H(X) - \varepsilon \leq \frac{1}{N} H_\delta(X^N)$.* Let $S \subseteq \text{Rng}(X^N)$ be any set with cardinality $|S| \leq 2^{N[H(X)-2\varepsilon]}$. We can bound the probability measure of such set as follows

$$\mathbb{P}_{X^N}[x_i \in S] = \mathbb{P}_{X^N}[x_i \in S \cap T_{N,\alpha}(X)] + \mathbb{P}_{X^N}[x_i \in S \cap \bar{T}_{N,\alpha}(X)],$$

where $\bar{\mathsf{T}}_{N,\alpha}(X) := \mathsf{Rng}(X^N) \setminus \mathsf{T}_{N,\alpha}(X)$ is the complement of the typical set $\mathsf{T}_{N,\alpha}(X)$. Since the probability of sequences in the typical set $\mathsf{T}_{N,\alpha}(X)$ is bounded by definition by $2^{-N[H(X)-\alpha]}$, one has

$$\begin{aligned}\mathbb{P}_{X^N}[x_i \in S \cap \mathsf{T}_{N,\alpha}(X)] &\leq \sum_{x_i \in S \cap \mathsf{T}_{N,\alpha}(X)} 2^{-N[H(X)-\alpha]} \leq \sum_{x_i \in S} 2^{-N[H(X)-\alpha]} \\ &\leq 2^{N[H(X)-2\alpha]} 2^{-N[H(X)-\alpha]} = 2^{-N\alpha}.\end{aligned}$$

On the other hand, by item 2 of the asymptotic equipartition theorem 4.15, one has

$$\begin{aligned}\mathbb{P}_{X^N}[x_i \in S \cap \bar{\mathsf{T}}_{N,\alpha}(X)] &\leq \mathbb{P}_{X^N}[x_i \in \bar{\mathsf{T}}_{N,\alpha}(X)] = 1 - \mathbb{P}_{X^N}[x_i \in \mathsf{T}_{N,\alpha}(X)] \\ &\leq 1 - \left(1 - \frac{\sigma^2}{N\alpha^2}\right) = \frac{\sigma^2}{N\alpha^2},\end{aligned}$$

where σ^2 is the variance of $\log_2 \frac{1}{\mathbb{P}_X(x)}$. Collecting the bounds on the two terms, we then have

$$\mathbb{P}_{X^N}[x_i \in S] \leq 2^{-N\alpha} + \frac{\sigma^2}{N\alpha^2}.$$

Setting $2\alpha = \varepsilon$ and $N_0 \geq \max\left\{\frac{8\sigma^2}{\varepsilon^2\gamma}, \frac{1}{2\varepsilon} \log_2 \frac{2}{\gamma}\right\}$, we conclude that for $N \geq N_0$, any set with cardinality $|S| \leq 2^{N[H(X)-\varepsilon]}$ has probability smaller than γ , and for suitably small γ it cannot be δ -sufficient. Then the smallest δ -sufficient set must have cardinality $|\mathsf{S}_\delta(X^N)|$ larger than $2^{N[H(X)-\varepsilon]}$. This implies that $\frac{1}{N}H_\delta(X^N) = \frac{1}{N} \log_2 |\mathsf{S}_\delta(X^N)| \geq H(X) - \varepsilon$

- *Proof of $\frac{1}{N}H_\delta(X^N) \leq H(X) + \varepsilon$.* From item 2 of the asymptotic equipartition theorem 4.15, setting $N > N_0$ in equation 4.15 and $N_0 \geq \frac{\sigma^2}{\varepsilon^2\delta}$, we observe that the typical set is δ -sufficient. The smallest δ -sufficient set has then cardinality not larger than the typical set, namely $|\mathsf{S}_\delta(X^N)| \leq |\mathsf{T}_{N,\varepsilon}(X)| \leq 2^{N[H(X)+\varepsilon]}$. Taking the logarithm we obtain $\frac{1}{N}H_\delta(X^N) \leq H(X) + \varepsilon$. \square

Intuitive consequences of the Shannon's source coding theorem are the following.

- For sufficiently large N the function $H_\delta(X^N)$ is essentially constant versus δ , for $0 < \delta < 1$.
- N i.i.d. random variables with Shannon entropy $H(X)$ can be compressed asymptotically for $N \rightarrow \infty$ into more than $NH(X)$ bits with negligible risk of information loss, whereas, if we try to compress them into fewer than $NH(X)$ bits, it is certain that some information will be lost.

Indeed, the first part of the proof tells us that even for very small probability of error δ , the number of bits per symbol $\frac{1}{N}H_\delta(X^N)$ required to specify a string x_i of N symbols does not need to exceed $H(X) + \varepsilon$ bits. In other words, for an arbitrarily small tolerance for error, the number of bits required drops from $H_0(X)$ to $H(X) + \varepsilon$. On the other hand, even if we tolerate an error δ very close to 1, the average number of bits per

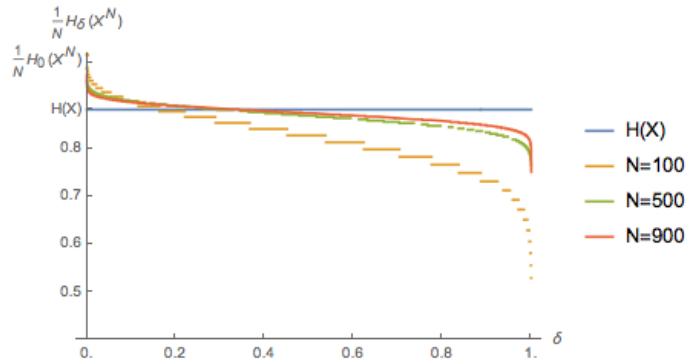


Figure 4.3 The essential bit content $\frac{1}{N}H_\delta(X^N)$ of N i.i.d. binary random variables X with probability distribution $\{p, 1-p\}$, as a function of the tolerated error probability δ . The plot represents the cases $N = 100, 500, 900$ for $p = 0.7$

symbol needed to specify x_i must still be at least $H(X) - \varepsilon$: The information becomes immediately unreliable if we take compression ratios smaller than $H(X)$. See figure 4.4.2.

We have thus far analysed the problem of source coding from a very general point of view, and Shannon's theorem summarises all the results that we can derive in the asymptotic case, exploiting the properties of typical sets. Shannon's theorem is proved by showing that, for large N , the smallest δ -sufficient set gets increasingly closer to the typical set. However, the ratio at which this phenomenon occurs is not analysed in detail, and thus a compression algorithm based on coding the typical set and messing up the remaining sequences may have performances that converge very slowly to the asymptotic regime. In the next lecture we will study a special class of coding algorithms, called *symbol codes*, that in the case of finite N already get close to the optimal compression ratio, represented by the Shannon entropy.

Chapter 5

Lecture 5: Symbol codes and stream codes

5.1 Symbol codes

Let $\Sigma = \{x_1, x_2, \dots, x_n\}$ be an alphabet. The set Σ^N denotes the set of length- N strings of characters from Σ , namely $\Sigma^N := \{x_{i_1}x_{i_2}\dots x_{i_N} \mid x_{i_j} \in \Sigma\}$. We then define

$$\Sigma^+ := \bigcup_{N \geq 1} \Sigma^N, \quad \Sigma^* := \{\lambda\} \cup \Sigma^N,$$

where λ denotes the empty string.

Example 5.1. For $\Sigma = \{0, 1\}$, we have

$$\begin{aligned} \Sigma^3 &= \{000, 001, 010, 011, 100, 101, 110, 111\}, \\ \Sigma^+ &= \{0, 1, 00, 01, 10, 11, 000, 001, 010, 011, 100, \dots\} \end{aligned}$$

Let us now consider a classical state of a system of type n , namely a random variable X with $|\text{Rng}(X)| = n$. In this case we have

$$\text{Rng}(X)^+ = \{x_1, x_2, \dots, x_n, x_1x_1, x_1x_2, \dots, x_nx_n, x_1x_1x_1, x_1x_1x_2, \dots\}.$$

A symbol code is a code that maps strings in $\text{Rng}(X)^+$ to strings of bits with variable length, by juxtaposing the codewords of the corresponding symbols.

Definition 5.2 (Symbol code). Let X be a random variable. A *symbol code* for X is a map $\mathcal{E} : \text{Rng}(X) \rightarrow \{0, 1\}^+$. The image of a *symbol* $x \in \text{Rng}(X)$, denoted as $\mathcal{E}(x) \in \{0, 1\}^+$, is a *codeword*, and the number of bits in $\mathcal{E}(x)$ is its length, denoted by $l_{\mathcal{E}}(x) = l(x)$. The codebook is the set of codewords $\mathcal{E}[\text{Rng}(X)]$.

Given a symbol code for X , one can use it to encode strings in $\text{Rng}(X)^+$, and this procedure defines the extended code.

Definition 5.3 (Extended code). Given a symbol code \mathcal{E} for X , we define its *extended code* as the map $\mathcal{E}^+ : \text{Rng}(X)^+ \rightarrow \{0, 1\}^+$, defined by concatenating codewords:

$$\mathcal{E}^+(x_{i_1}x_{i_2}\dots x_{i_k}) = \mathcal{E}(x_{i_1})\mathcal{E}(x_{i_2})\dots\mathcal{E}(x_{i_k}). \tag{5.1}$$

x	$\mathcal{E}_0(x)$	$l(x)$	$\mathbb{P}_X(x)$
1	00	2	$\frac{1}{2}$
2	01	2	$\frac{1}{4}$
3	10	2	$\frac{1}{4}$

Example 5.4. Let $\Omega = \{a, b, c\}$, $\text{Rng}(X) = \{1, 2, 3\}$ with $\mathbf{p} = \{\frac{1}{2}, \frac{1}{4}, \frac{1}{4}\}$. We can define the following symbol code

The following two examples are based on $\Omega = \{a, b, c, d\}$, with $\text{Rng}(X) = \{1, 2, 3, 4\}$ with $\mathbf{p} = \{\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}\}$.

x	$\mathcal{E}_1(x)$	$l(x)$	$\mathbb{P}_X(x)$	x	$\mathcal{E}_2(x)$	$l(x)$	$\mathbb{P}_X(x)$
1	0	1	$\frac{1}{2}$	1	0	1	$\frac{1}{2}$
2	10	2	$\frac{1}{4}$	2	10	2	$\frac{1}{4}$
3	110	3	$\frac{1}{8}$	3	100	3	$\frac{1}{8}$
4	111	3	$\frac{1}{8}$	4	111	3	$\frac{1}{8}$

Given a symbol code, we can define its expected length, that assesses its compression rate.

Definition 5.5 (Expected length). Given a symbol code \mathcal{E} for a random variable X with $\mathbf{p} = \{p_1, p_2, \dots, p_n\}$, the *expected length* $L(\mathcal{E}, X)$ of the code is defined as

$$L(\mathcal{E}, X) := \sum_{i=1}^n \{\mathbb{P}_X(x_i)l(x_i)\}. \quad (5.2)$$

As an example, consider the codes \mathcal{E}_1 and \mathcal{E}_2 introduced above, for which $L(\mathcal{E}_1, X) = L(\mathcal{E}_2, X) = 7/4$ bit/character. The length is not the only feature of symbol codes we are interested in. In particular, we need the code to be uniquely decodable. The following definitions will allow us to find the features that grant fulfilment of this requirement.

Definition 5.6 (Non-singular code). A symbol code \mathcal{E} for X is *non-singular* if for every two symbols $x \neq y$ in $\text{Rng}(X)$ one has $\mathcal{E}(x) \neq \mathcal{E}(y)$

The above condition is not sufficient yet to grant decodability. We need the following notion.

Definition 5.7 (Unique decodability). A symbol code \mathcal{E} for X is *uniquely decodable* if for every two strings $x \neq y$ in $\text{Rng}(X)^+$ one has $\mathcal{E}^+(x) \neq \mathcal{E}^+(y)$

Example 5.8. The codes \mathcal{E}_0 , \mathcal{E}_1 and \mathcal{E}_2 are non-singular. However, only \mathcal{E}_0^+ and \mathcal{E}_1^+ are uniquely decodable. Indeed, consider the following strings for the code \mathcal{E}_2 : $x = acbad$ and $y = ababad$. We have

$$\begin{aligned} \mathcal{E}_2^+[\text{X}^{(5)}(acbad)] &= 0100100111, \\ \mathcal{E}_2^+[\text{X}^{(6)}(ababad)] &= 0100100111. \end{aligned}$$

How can we ensure unique decodability? We will now define a property that provides a sufficient condition.

Definition 5.9 (Prefix). Let $p, w \in \{0, 1\}^+$. We say that p is a *prefix* for w if there is $t \in \{0, 1\}^*$ such that $w = pt$.

Definition 5.10 (Prefix code). We say that a code \mathcal{E} for X has the *prefix property*, or simply is *prefix*, if no codeword is a prefix for another codeword.

As an example, the codes \mathcal{E}_0 and \mathcal{E}_1 are prefix. Moreover, the code $\{0, 101\}$ for $\Omega = \{a, b\}$ is prefix, while $\{1, 101\}$ is not prefix, but it is uniquely decodable. The prefix property is sufficient for unique decodability, but it may seem too restrictive, as it is not necessary. However, the Kraft-McMillan theorem will show that we have no loss in performances if we restrict to prefix codes; on the other hand, prefix codes have the property of *self-punctuation*, namely one can decode the message on the fly (the so-called instantaneous decoding), without waiting for the end of transmission. Moreover, the prefix property is simply checkable, as it can be checked on $\text{Rng}(X)$.

Lemma 5.11. *If the symbol code \mathcal{E} is prefix the extended code \mathcal{E}^+ is uniquely decodable.*

Proof. If \mathcal{E} is prefix then $\mathcal{E}(x_{i_1}) = \mathcal{E}(x_{i_2})t$ iff $x_{i_1} = x_{i_2}$ and $t = \lambda$. Suppose now that \mathcal{E}^+ is not uniquely decodable, and let $x \neq y$ be two strings with $\mathcal{E}^+(x) = \mathcal{E}^+(y)$. Let $j := \min\{k | x_{i_k} \neq y_{i_k}\}$. Then we have

$$x = x_{i_1}x_{i_2} \dots x_{i_j}x_{i_{j+1}} \dots \quad y = x_{i_1}x_{i_2} \dots y_{i_j}y_{i_{j+1}} \dots \\ \mathcal{E}(x_{i_1})\mathcal{E}(x_{i_2}) \dots \mathcal{E}(x_{i_j})\mathcal{E}(x_{i_{j+1}}) \dots = \mathcal{E}(x_{i_1})\mathcal{E}(x_{i_2}) \dots \mathcal{E}(y_{i_j})\mathcal{E}(y_{i_{j+1}}) \dots$$

This implies that

$$\mathcal{E}(x_{i_j})\mathcal{E}(x_{i_{j+1}}) \dots = \mathcal{E}(y_{i_j})\mathcal{E}(y_{i_{j+1}}) \dots$$

We now have three possibilities:

1. $l_{\mathcal{E}}(x_{i_j}) = l_{\mathcal{E}}(y_{i_j})$: in this case, it must be $\mathcal{E}(x_{i_j}) = \mathcal{E}(y_{i_j})$.
2. $l_{\mathcal{E}}(x_{i_j}) > l_{\mathcal{E}}(y_{i_j})$: in this case, it must be $\mathcal{E}(x_{i_j}) = \mathcal{E}(y_{i_j})s$.
3. $l_{\mathcal{E}}(x_{i_j}) < l_{\mathcal{E}}(y_{i_j})$: in this case, it must be $\mathcal{E}(x_{i_j})s' = \mathcal{E}(y_{i_j})$.

In all three cases, by the prefix property of \mathcal{E} it must be $x_{i_j} = y_{i_j}$, contrarily to the hypothesis. \square

The prefix property is checkable on $\text{Rng}(X)$, contrarily to unique decodability: The code \mathcal{E}_2 in our example 5.4 indeed is uniquely decodable on $\text{Rng}(X)$, however $\mathcal{E}_2^+(ba) = \mathcal{E}_2^+(c)$.

5.2 Kraft-McMillan theorem

In this section we prove the most important result about symbol codes, that provides a necessary and sufficient condition for unique decodability, and allows one also to prove that for every uniquely decodable symbol code, there is a prefix one with the same codeword lengths.

Theorem 5.12 (Kraft-McMillan theorem). *Let X be a random variable.*

1. *A uniquely decodeable binary code \mathcal{E} for X must satisfy McMillan's bound*

$$\sum_{x_i \in \text{Rng}(X)} 2^{-l(x_i)} \leq 1. \quad (5.3)$$

2. *Conversely, if the bound 5.3 is satisfied by a set of lengths $\{l_1, l_2, \dots, l_{|\text{Rng}(X)|}\}$, then there exists a prefix code \mathcal{E} for X with $l_{\mathcal{E}}(x_i) = l_i$.*

Proof. Let us start from item 1. Define $S := \sum_{x_i \in \text{Rng}(X)} 2^{-l(x_i)} = \sum_{i=1}^R 2^{-l_i}$, where we set $|\text{Rng}(X)| = R$ and $l_i := l(x_i)$. Then we have

$$S^N = \sum_{i_1=1}^R \sum_{i_2=1}^R \dots \sum_{i_N=1}^R 2^{-(l_{i_1} + l_{i_2} + \dots + l_{i_N})}. \quad (5.4)$$

Since $l_{i_1} + l_{i_2} + \dots + l_{i_N} = l(x_{i_1} x_{i_2} \dots x_{i_N})$ for some string $x_i = x_{i_1} x_{i_2} \dots x_{i_N} \in \text{Rng}(X^N)$, we can re-write equation (5.4) as follows

$$S^N = \sum_{l=Nl_m}^{Nl_M} \{2^{-l} N(l)\}, \quad (5.5)$$

where $l_m := \min_{x_i \in \text{Rng}(X)} l(x_i)$, $l_M := \max_{x_i \in \text{Rng}(X)} l(x_i)$, and $N(l) := |\{x_i \in \text{Rng}(X^N) | l(x_i) = l\}|$ is the number of strings $x_i \in \text{Rng}(X^N)$ with $l(x_i) = l$. Since the number of bit strings of length l is 2^l , unique decodability enforces $N(l) \leq 2^l$. Thus,

$$S^N \leq \sum_{Nl_m}^{Nl_M} \{2^{-l} 2^l\} = N(l_M - l_m) + 1 \leq N(l_M - l_m + 1) = Nk, \quad k > 0.$$

Since the bound must hold for every N , it must be $S \leq 1$. Indeed, $S > 1$ would imply that, for sufficiently large N , one would have $S^N > Nk$. Let us now consider item 2. Let M denote the number of *different* values of the considered lengths—e.g., if $R = 5$ and $l_1 = 1, l_2 = 3, l_3 = 3, l_4 = 3, l_5 = 6$, then $M = 3$. Consider now the set of different lengths $\lambda_1 < \lambda_2 \dots < \lambda_M$, and let N_k denote the number of lengths l_i with the same value λ_k —e.g., in the same example as before, $\lambda_1 = 1, \lambda_2 = 3, \lambda_3 = 6$, and $N_1 = 1, N_2 = 3$, and $N_3 = 1$. We can express the hypothesis of item 2 as

$$S := \sum_{k=1}^M N_k 2^{-\lambda_k} \leq 1.$$

We now prove that one can construct a prefix code with N_k codewords of length λ_k for every $1 \leq k \leq M$, as follows. Consider any integer j with $1 \leq j \leq M$. One has

$$N_j 2^{-\lambda_j} + \sum_{k=m}^{j-1} \{N_k 2^{-\lambda_k}\} \leq S \leq 1,$$

and then

$$N_j \leq 2^{\lambda_j} - \sum_{k=1}^{j-1} \{N_k 2^{\lambda_j - \lambda_k}\}. \quad (5.6)$$

The proof now proceeds by induction on $j < M$. For $j = 1$ one has $N_1 \leq 2^{\lambda_1}$, namely the number N_1 is smaller than the number of strings of length l_1 , and one can thus encode N_1 different symbols in such strings. Clearly, no string of length λ_1 can be a prefix for another string of the same length. Suppose now that one can choose N_k codewords of length λ_k for every $k < j$, without using any assigned codeword as a prefix for another one. Since $2^{\lambda_j - \lambda_k}$ is the number of strings of length λ_j having a fixed prefix of length λ_k , $\sum_{k=1}^{j-1} N_k 2^{\lambda_j - \lambda_k}$ is precisely the number of strings of length λ_j with prefixes that have already been assigned as shorter codewords. If we subtract this number from the total number 2^{λ_j} of strings of length λ_j , inequality 5.6 ensures that we are left with a number of available strings of length λ_j that is not smaller than the number of strings we need to encode. Since this holds for every j , McMillan's bound grants the possibility to construct a prefix code. \square

Remark 3. Given a uniquely decodable code \mathcal{E} , its codeword lengths must satisfy McMillan's inequality. However, by Kraft-McMillan's theorem, this means that we can construct a prefix code with exactly the same lengths. This is the reason why it is not restrictive to consider prefix codes instead of uniquely decodable ones.

An important class of symbol codes is that of codes saturating McMillan's bound.

Definition 5.13 (Complete code). A uniquely decodable symbol code \mathcal{E} for X is called *complete* if its lengths satisfy

$$\sum_{x_i \in \text{Rng}(X)} 2^{-l(x_i)} = 1. \quad (5.7)$$

5.3 Source coding theorem for symbol codes

The task of compression consists in minimising the expected length of a code. For symbol codes, the figure of merit to be minimised is then

$$L(\mathcal{E}, X) := \sum_{x_i \in \text{Rng}(X)} \{\mathbb{P}_X(x_i) l(x_i)\}.$$

In this section, we will prove that the expected length is bounded from below by the Shannon entropy of X .

Theorem 5.14 (Lower bound for the expected code length). *The expected length $L(\mathcal{E}, \mathbf{X})$ of a uniquely decodable symbol code \mathcal{E} is bounded from below by $H(\mathbf{X})$.*

Proof. Let us define the two vectors

$$\begin{aligned}\mathbf{p} &:= (p_1, p_2, \dots, p_R), \quad p_i := \mathbb{P}_{\mathbf{X}}(x_i) \\ \mathbf{q} &:= (q_1, q_2, \dots, q_R), \quad q_i := \frac{2^{-l(x_i)}}{Z},\end{aligned}$$

where $l(x_i)$ are the lengths for \mathcal{E} and the *partition function* Z is $Z := \sum_{i=1}^R 2^{-l(x_i)} \leq 1$. The probabilities \mathbf{q} are called *implicit probabilities* of \mathcal{E} . Let us now apply Gibbs' inequality $0 \leq H(\mathbf{p} \parallel \mathbf{q})$ as follows

$$\sum_{i=1}^R \left\{ p_i \log_2 \frac{1}{q_i} \right\} \geq \sum_{i=1}^R \left\{ p_i \log_2 \frac{1}{p_i} \right\},$$

obtaining

$$\sum_{i=1}^R \{p_i l(x_i)\} = \sum_{i=1}^R \left\{ p_i \log_2 \frac{1}{q_i} \right\} - \log_2 Z \geq \sum_{i=1}^R \left\{ p_i \log_2 \frac{1}{p_i} \right\} = H(\mathbf{X}). \quad (5.8)$$

Notice that equality is achieved under two conditions: $Z = 1$ (the code is complete) and $q_i = p_i$, namely $l(x_i) = \log_2 \frac{1}{p_i}$. \square

Corollary 5.15 (Optimal source code-lengths). *The expected length of the code is minimum only if the code-length of a symbol equals its Shannon information content*

$$l(x_i) = \log_2 \frac{1}{\mathbb{P}_{\mathbf{X}}(x_i)} \quad (\text{Shannon's lengths}). \quad (5.9)$$

Conversely, any choice of code-lengths $l(x_i)$ defines a probability distribution $\mathbb{P}_L(x_i) = 2^{-l(x_i)}/Z$ for which the code-lengths would be optimal. If in addition the code is complete, then the expected length will equal the Shannon entropy $H(\mathbf{X})$, and $\mathbb{P}_{\mathbf{X}}(x_i) = 2^{-l(x_i)}$.

Notice however that the Shannon length $l(x_i)$ is generally not integer. Then, the question is: How close can we get to the lower bound for the expected length? This is provided by the following theorem

Theorem 5.16 (Source coding theorem for symbol codes). *For any random variable \mathbf{X} there exists a prefix code \mathcal{E} with expected length satisfying the source coding theorem for symbol codes*

$$H(\mathbf{X}) \leq L(\mathcal{E}, \mathbf{X}) < H(\mathbf{X}) + 1. \quad (5.10)$$

Proof. In general, the Shannon length $l(x_i) = -\log_2 \mathbb{P}_{\mathbf{X}}(x_i)$ cannot be exactly achieved because it is not an integer. Therefore, we use lengths slightly larger than the optimum, namely

$$l(x_i) = \left\lceil \log_2 \frac{1}{\mathbb{P}_{\mathbf{X}}(x_i)} \right\rceil, \quad (5.11)$$

where $\lceil a \rceil$ denotes the smallest integer greater than or equal to a . This choice of lengths satisfies McMillan's bound 5.3, since

$$\sum_{i=1}^R 2^{-l(x_i)} = \sum_{i=1}^R 2^{-\lceil \log_2 \frac{1}{p_i} \rceil} \leq \sum_{i=1}^R 2^{-\log_2 \frac{1}{p_i}} = \sum_{i=1}^R p_i = 1. \quad (5.12)$$

Then, according to theorem 5.12, there exists a prefix code with lengths $l(x_i)$. Moreover, we have

$$L(\mathcal{E}, X) = \sum_{i=1}^R \left\{ p_i \left\lceil \log_2 \frac{1}{p_i} \right\rceil \right\} < \sum_{i=1}^R \left\{ p_i \left(\log_2 \frac{1}{p_i} + 1 \right) \right\} = H(X) + 1. \quad (5.13)$$

The lower bound was already proved in theorem 5.14. \square

Remark 4. For complete codes, $Z = 1$, and by equation 5.8 one immediately sees that in that case the expected code length exceeds the Shannon entropy $H(X)$ exactly by the Kullback-Leibler divergence $H(\mathbf{p} \parallel \mathbf{q})$ between the actual distribution and that of implicit probabilities \mathbf{q} .

Chapter 6

Lecture 6: Huffman coding and stream codes

6.1 Optimal code: Huffman's coding

Given a probability distribution \mathbf{p} for the random variable X , can we construct the optimal symbol code? The answer is: Yes. We now show Huffman's code—constructed by a simple protocol—which is a prefix symbol code \mathcal{E} minimising the expected length $L(\mathcal{E}, X)$. The code is obtained by the following algorithm.

Huffman code $\mathcal{E}_H(X)$

1. Start by defining the random variable $X_0 := X$, with symbols sorted from the one with highest probability ($x_1^{(0)}$) to the one with smallest probability ($x_R^{(0)}$). Assign the empty string $\lambda \in \{0, 1\}^*$ to all symbols $x_i^{(0)} \in \text{Rng}(X_0)$.
2. If there is a unique symbol $x_1^{(n)}$, stop the algorithm. Otherwise take the two least probable symbols $x_{R-n-1}^{(n)}, x_{R-n}^{(n)}$ in the alphabet X_n , and append two different bits, 0 and 1, to the left of all strings associated to them.
3. Combine the two symbols $\{x_{R-n-1}^{(n)}, x_{R-n}^{(n)}\}$ into a single one $x'_{R-n-1}^{(n)} := \{x_{R-n-1}^{(n)}\} \cup \{x_{R-n}^{(n)}\}$, and define the new variable X_{n+1} with $x_j^{(n+1)} = x_j^{(n)}$ for $j < R-n-1$, and $x_{R-n-1}^{(n)} = x'_{R-n-1}^{(n)}$; and $p_j^{(n+1)} = p_j^{(n)}$ for $j < R-n-1$, and $p_{R-n-1}^{(n+1)} = p_{R-n-1}^{(n)} + p_{R-n}^{(n)}$. Sort X_{n+1} from the symbol with highest probability to the one with the lowest. Iterate the procedure from step 2.

Notice that by construction the Huffman algorithm associates the n -th longest codeword to the n -th least probable symbol. This feature is crucial in the proof of optimality of Huffman codes.

Example 6.1. Consider the alphabet:

$$X = \left\{ \begin{array}{l} \Omega = \{a, b, c, d, e\} \\ \mathbf{p} = \{0.25, 0.25, 0.2, 0.15, 0.15\} \end{array} \right.$$

The Huffman coding is obtained as follows

1.

$x_i^{(0)}$	$\mathbf{p}^{(0)}$	$\mathcal{E}_H(x)$	$x_i^{(1)}$	$\mathbf{p}^{(1)}$
a	0.25	λ	{d, e}	0.30
b	0.25	λ	a	0.25
c	0.20	λ	b	0.25
d	0.15	0	c	0.20
e	0.15	1		

2.

$x_i^{(1)}$	$\mathbf{p}^{(1)}$	$\mathcal{E}_H(x)$	$x_i^{(2)}$	$\mathbf{p}^{(2)}$
{d, e}	0.30	{0, 1}	{b, c}	0.45
a	0.25	λ	{d, e}	0.30
b	0.25	0	a	0.25
c	0.20	1		

3.

$x_i^{(2)}$	$\mathbf{p}^{(2)}$	$\mathcal{E}_H(x)$	$x_i^{(3)}$	$\mathbf{p}^{(3)}$
{b, c}	0.45	{0, 1}	{d, e, a}	0.55
{d, e}	0.30	{00, 01}	{b, c}	0.45
a	0.25	1		

4.

$x_i^{(3)}$	$\mathbf{p}^{(3)}$	$\mathcal{E}_H(x)$
{d, e, a}	0.55	{000, 001, 01}
{b, c}	0.45	{10, 11}

Thus we end up with the following code

x_i	p_i	$\mathcal{E}_H(x_i)$	$h(x_i)$
a	0.25	01	2.0
b	0.25	10	2.0
c	0.20	11	2.3
d	0.15	000	2.7
e	0.15	001	2.7

The expected length is $L = 2.40$ bits to be compared with the Shannon entropy $H(X) = 2.2855$ bits.

Example 6.2. Consider the alphabet:

$$X = \left\{ \begin{array}{l} \Omega = \{a, b, c, d, e, f, g\} \\ \mathbf{p} = \{0.01, 0.24, 0.05, 0.20, 0.47, 0.01, 0.02\} \end{array} \right.$$

The Huffman coding is obtained as follows

1.

$x_i^{(0)}$	$p_i^{(0)}$	$\mathcal{E}_H(x)$	$x_i^{(1)}$	$p_i^{(1)}$
e	0.47	λ	e	0.47
b	0.24	λ	b	0.24
d	0.20	λ	d	0.20
c	0.05	λ	c	0.05
g	0.02	λ	g	0.02
a	0.01	1	{a, f}	0.02
f	0.01	0		

2.

$x_i^{(1)}$	$p_i^{(1)}$	$\mathcal{E}_H(x)$	$x_i^{(2)}$	$p_i^{(2)}$
e	0.47	λ	e	0.47
b	0.24	λ	b	0.24
d	0.20	λ	d	0.20
c	0.05	λ	c	0.05
g	0.02	0	{g, a, f}	0.04
{a, f}	0.02	{10, 11}		

3.

$x_i^{(2)}$	$p_i^{(2)}$	$\mathcal{E}_H(x)$	$x_i^{(3)}$	$p_i^{(3)}$
e	0.47	λ	e	0.47
b	0.24	λ	b	0.24
d	0.20	λ	d	0.20
c	0.05	0	{c, g, a, f}	0.09
{g, a, f}	0.04	{10, 110, 111}		

4.

$x_i^{(3)}$	$p_i^{(3)}$	$\mathcal{E}_H(x)$	$x_i^{(4)}$	$p_i^{(4)}$
e	0.47	λ	e	0.47
b	0.24	λ	{d, c, g, a, f}	0.29
d	0.20	0	b	0.24
{c, g, a, f}	0.09	{10, 110, 1110, 1111}		

5.

$x_i^{(4)}$	$p_i^{(4)}$	$\mathcal{E}_H(x)$	$x_i^{(5)}$	$p_i^{(5)}$
e	0.47	λ	{d, c, g, a, f, b}	0.53
{d, c, g, a, f}	0.29	{00, 010, 0110, 01110, 01111}	e	0.47
b	0.24	1		

6.

$x_i^{(5)}$	$p_i^{(5)}$	$\mathcal{E}_H(x)$
{d, c, g, a, f, b}	0.53	{000, 0010, 00110, 001110, 001111, 01}
e	0.47	1

Thus we end up with the following code

x_i	p_i	$\mathcal{E}_H(x_i)$	$h(x_i)$
a	0.01	001110	6.64
b	0.24	01	2.06
c	0.05	0010	4.32
d	0.20	000	2.32
e	0.47	1	1.09
f	0.01	001111	6.64
g	0.02	00110	5.64

The expected length is $L = 1.97$ bits to be compared with the Shannon entropy $H(X) = 1.93$ bits.

From the examples we can see that the lengths of the code are close to the Shannon information content—sometimes smaller, sometimes larger. The Huffman code achieves an expected length that satisfies $H(X) \leq L(\mathcal{E}, X) < H(X) + 1$, as proven in theorem 6.5. This overhead between 0 and 1 bits per symbol can be in practice a big one, when $H(X)$ is itself 1 bit. Indeed, in practice one can achieve better compression not by *symbol codes*, but using *stream codes*.

We can prove optimality of the Huffman code by the following argument. First, we need a definition and a lemma.

Definition 6.3 (Sibling strings). Two binary strings $x, y \in \{0, 1\}^+$ are *sibling* if there exists $p \in \{0, 1\}^+$ such that $x = p0$ and $y = p1$.

Lemma 6.4. *The optimal uniquely decodable code \mathcal{E}_O must have the longest codewords occurring in pairs of sibling strings, with one pair encoding the two least probable symbols x_a, x_b .*

Proof. Since the optimal code can be chosen to be prefix without loss of generality, the longest codewords must occur in sibling pairs whose length we denote by l_M . Indeed, if one of the longest codewords had no sibling, since the code \mathcal{E}_B is prefix, one could shorten it without spoiling the prefix property, contradicting the optimality hypothesis. Suppose now that none of the longest sibling strings encodes one of the two least probable symbols—say x_a . Let then $l(x_a) < l(x_c) = l_M$, with $p(x_a) < p(x_c)$. If we now build a different code \mathcal{E}_C , exchanging the codewords for x_a, x_c , the difference of the expected lengths is

$$\Delta L := L(\mathcal{E}_C, X) - L(\mathcal{E}_B, X) = (p_c l_a + p_a l_c) - (p_c l_c + p_a l_a) = (p_c - p_a)(l_a - l_c) < 0.$$

This implies that the new code is better, contrarily to the optimality hypothesis. Thus, the best code must associate two sibling longest codewords to the two least probable symbols. \square

Optimality of Huffman codes can now be proved by induction on the number of symbols in the alphabet $\text{Rng}(X)$ as follows.

Theorem 6.5 (Optimality of the Huffman code). *The code constructed by Huffman's algorithm is optimal.*

Proof. The proof proceeds by induction. Indeed, for an alphabet with a unique symbol the encoding results in a code with $L(\mathcal{E}_H, X) = 0$, which is clearly optimal. Suppose now that the Huffman algorithm provides the optimal code $\mathcal{E}_{H,n}$ for alphabets with n symbols, and consider an alphabet $\text{Rng}(X)$ with $n+1$ symbols. By the above lemma the optimal code \mathcal{E}_O for X must associate two sibling strings of maximal length to the least probable symbols $x_n^{(0)}, x_{n+1}^{(0)}$. Let us then form the random variable X_1 with $x_j^{(1)} = x_j^{(0)}$ for $1 \leq j \leq n-1$ and $x_n^{(1)} := \{x_n^{(0)}, x_{n+1}^{(0)}\}$. The probabilities of X_1 are $p(x_i^{(1)}) = p(x_i^{(0)})$ for $1 \leq i \leq n-1$ and $p(x_n^{(1)}) := p(x_n^{(0)}) + p(x_{n+1}^{(0)})$. Consider now the code $\mathcal{E}_{O,n}$ obtained from \mathcal{E}_O by merging the symbols $x_n^{(0)}, x_{n+1}^{(0)}$, and associating to them the string $\mathcal{E}_{O,n}(x_n^{(1)}) := [\mathcal{E}_O(x_n^{(0)})]' = [\mathcal{E}_O(x_{n+1}^{(0)})]',$ where w' means the string w without the last bit. The code $\mathcal{E}_{O,n}$ is a code for the n -symbol random variable X_1 . Now, one has

$$\begin{aligned} L(\mathcal{E}_O, X) &= \sum_{i=1}^{n-1} p(x_i^{(0)}) l_{\mathcal{E}_O}(x_i^{(0)}) + p(x_n^{(0)}) l_{\mathcal{E}_O}(x_n^{(0)}) + p(x_{n+1}^{(0)}) l_{\mathcal{E}_O}(x_{n+1}^{(0)}) \\ &= \sum_{i=1}^{n-1} p(x_i^{(1)}) l_{\mathcal{E}_{O,n}}(x_i^{(1)}) + p(x_n^{(1)}) [l_{\mathcal{E}_{O,n}}(x_n^{(1)}) + 1] \\ &= L(\mathcal{E}_{O,n}, X_1) + p(x_n^{(0)}) + p(x_{n+1}^{(0)}). \end{aligned}$$

If we now consider the Huffman code $\mathcal{E}_{H,n+1}$, by an identical calculation we obtain

$$L(\mathcal{E}_{H,n+1}, X) = L(\mathcal{E}_{H,n}, X_1) + p(x_n^{(0)}) + p(x_{n+1}^{(0)}),$$

where $L(\mathcal{E}_{H,k}, Y)$ denotes the expected length of the Huffman code for the k -symbol random variable Y . Finally, by the induction hypothesis we have that $L(\mathcal{E}_{O,n}, X_1) \geq L(\mathcal{E}_{H,n}, X_1)$, which implies $L(\mathcal{E}_O, X) \geq L(\mathcal{E}_{H,n+1}, X)$. Since \mathcal{E}_O is optimal by hypothesis, it must be $L(\mathcal{E}_O, X) = L(\mathcal{E}_{H,n+1}, X)$, thus proving the optimality of the Huffman code. \square

6.2 Stream codes and redundancy

We have seen that symbol codes obey a source coding theorem that envisages a possible overhead of 1 bit/symbol. This may seem a little cost, however one should consider that e.g. a text written in the English language has an essential bit content of about 1 bit/symbol. This means that an overhead of 1 bit/symbol is huge. One possible way around the problem consists in treating long strings as symbols, namely encoding

$\text{Rng}(X^N)$ instead of $\text{Rng}(X)$. However, even this strategy can be very inefficient. Indeed, the Huffman algorithm constructs a full encoding table, which takes an exponential time to complete. In particular, sorting strings according to their decreasing probability is the part of the algorithm that takes exponential time versus the number of symbols. This makes the encoding necessarily rigid, since it is not convenient to update the table once the probabilities are updated. For this reason in practical applications where the i.i.d. hypothesis for symbols within long strings is not at all justified, it is convenient to go beyond symbol codes.

A better strategy with respect to the construction of Huffman's code tables for every update of probabilities is one which can be easily adapted to longer block lengths without an exponential computational complexity. This is the case of *stream codes*. We will analyse two main cases of stream codes: the arithmetic code and the Lempel-Ziv code. Arithmetic coding compresses data on the basis of a probabilistic modelling of the source. Since 1999 this is the best compression method for text files, and it is used also by some advanced image compressors. The Lempel-Ziv coding is a *universal* method, which is designed to compress any source. Even though it is not so effective as the arithmetic codes, it has a reasonable efficiency on most files, and it is proved to reach the Shannon limit asymptotically. The Lempel-Ziv algorithm is used in the popular `compress` and `gzip` routines.

The redundancy of the English language is a well known fact: in cross-puzzles, or in the game *hangman*, one can guess quite effectively what is the next letter in an incomplete word. The redundancy is not just due to the fact that characters are non equally distributed, but in particular it is due to the dependence of probabilities of symbols on the preceding character. More generally, entire strings of symbols have a likelihood that depends on the context. It is this kind of redundancy that can be exploited in compressing a text file. Indeed, the hang man example is particularly suited to our purpose, since we will now see how a similar guessing game can be used for compression.

Compression as a guessing game

Consider the following simple game in which each consecutive character of a sentence should be guessed by a player, and the number of attempts is recorded by a referee for each character. We restrict to the 27 characters of the English alphabet including only upper-case letters A,B,C,...,Z and the space '-'. For example, we report a sentence and the corresponding list of attempts below

```
I LOVE YOU
11 8211 311
```

Notice how easy it is to guess the right letter in just a single attempt, e.g. after Y. How can the above game be used for compression? Suppose that there is an exact twin of the guessing player, who makes exactly the same guesses as the first player when they are in the same context. Suppose that now there is a second referee, who claims to play the same game with the twin, as if they knew the sentence. However, all the second referee knows is the sequence of numbers of attempts from the first twin. The referee would

ask the second twin to guess the letter, and would say ‘yes, that’s right!’ whenever the number of guesses equals the number in the sequence. Since the second player is a perfect twin of the first one, exactly the same string would be obtained. It is easy to recognise that the first player has the role of the *encoder*, whereas his twin is the *decoder*. Indeed, we have not reduced the size of the alphabet, however, since the number of attempts is a small number most of the times, the new file has much less redundancy than the original one, and will be much easier to compress.

Do we have perfectly identical twins in practice? Yes: they would be simulated by a computer program with a rule for generating attempts for any possible preceding string (up to some maximum window length). This would be just a particular probabilistic language model—in our case this is the L -th order Markov model with maximum window length L . However, any probabilistic model can be used, more or less effectively, depending on the adequateness of the probabilistic model.

6.3 Arithmetic code

Let us consider a random variable X with $\text{Rng}(X) = \{x_1, x_2, \dots, x_n\}$, including the ‘end of file’ symbol $x_n := \odot$. Let us now consider the variable $X^{(N)}$, with $\text{Rng}(X^{(N)}) = \text{Rng}(X)^N$, but *not* i.i.d. Since the joint probabilities $P[x_i^{(N)}] := \mathbb{P}_{X^{(N)}}[x_{i_1}, x_{i_2}, \dots, x_{i_N}]$ of strings are not factorised, we need an efficient algorithm to compute them, in the first place. This is achieved through the computation of *conditional probabilities* $P[x_{i_k} | x_i^{(k-1)}]$, where $x_i^{(k)}$ denotes the string made of the first k symbols of $x_i^{(N)}$, namely $x_i^{(k)} = x_i^{(k-1)}x_{i_k}$. The i -th character in the string will be denoted by the random variable $X^{(i)}$. The same algorithm will be used to compute conditional probabilities, both at the encoding and at the decoding stage. This algorithm models the twin guessers.

Encoding and decoding

Every binary transmission can be regarded as the transmission of rational number belonging to the interval $[0, 1)$: the general string $x_i^{(k)}$ made of k binary digits is just the binary form of the fractional part of the number. For example the string 0101 corresponds to $\frac{1}{4} + \frac{1}{16} = 0.3125$. Let us denote the above map as

$$\begin{aligned}\text{Bin} : \{0, 1\}^+ &\rightarrow [0, 1), \\ \text{Bin} : x_i^{(k)} &\mapsto \sum_{j=1}^k x_{i_j} 2^{-j}.\end{aligned}$$

One can easily verify that the above map is bijective.

We now consider the following variation of the map Bin , that associates binary sequences $b \in \{0, 1\}^+$ to *binary subintervals* of $[0, 1)$ in the following way:

$$\begin{aligned}0 &\rightarrow [0, .5), 1 \rightarrow [.5, 1), \\ 00 &\rightarrow [0, .25), 01 \rightarrow [.25, .5), 10 \rightarrow [.5, .75), 11 \rightarrow [.75, 1),\end{aligned}$$

and rewriting the intervals in binary form

$$0 \rightarrow [0, .1), 1 \rightarrow [.1, 1),$$

$$00 \rightarrow [0, .01), 01 \rightarrow [.01, .1), 10 \rightarrow [.1, .11), 11 \rightarrow [.11, 1),$$

and so on. For example one has

$$01011 \rightarrow [.01011, .01011 + .00001) = [.01011, .01100).$$

Let us define the set BI of binary subintervals of $[0, 1]$, namely

$$\text{BI} := \{[0.x_{i_1}x_{i_2}\dots x_{i_k}, 0.x_{i_1}x_{i_2}\dots x_{i_k} + 0.\underbrace{00\dots 0}_{k-1 \text{ zeros}} 1) \mid x_{\mathbf{i}}^{(k)} \in \{0, 1\}^+\}$$

We will denote the above bijective map as

$$\text{Bintv} : \{0, 1\}^+ \leftrightarrow \text{BI},$$

where BI denotes the set of *binary intervals* in $[0, 1]$ defined above.

How can we encode information using a mapping on intervals? The idea is the following. We can divide the interval $[0, 1)$ also into n subintervals of lengths equal to the probabilities $\mathbb{P}_{X^{(N)}}[X^{(1)} = x_{i_1}], i = 1, 2, \dots, n$. We will denote these intervals by the symbols $I_{x_{i_1}}$. Then we can subdivide each interval x_{i_1} into subintervals denoted as $I_{x_{i_1}x_{i_2}}$, with lengths $|I_{x_{i_1}x_{i_2}}|$ proportional to $\mathbb{P}_{X^{(N)}}[X^{(2)} = x_{i_2}|X^{(1)} = x_{i_1}]$, or $P[x_{i_2}|x_{i_1}]$ in short notation. Indeed the length $|I_{x_{i_1}x_{i_2}}|$ of the interval $I_{x_{i_1}x_{i_2}}$ will be precisely the joint probability

$$|I_{x_{i_1}x_{i_2}}| = \mathbb{P}_{X^{(N)}}[X^{(1)} = x_{i_1}, X^{(2)} = x_{i_2}] = \mathbb{P}_{X^{(N)}}[X^{(1)} = x_{i_1}] \mathbb{P}_{X^{(N)}}[X^{(2)} = x_{i_2}|X^{(1)} = x_{i_1}].$$

Upon iterating the process, we will have a division into tiny intervals with lengths given by

$$|I_{x_{\mathbf{i}}^{(k)}}| = \mathbb{P}_{X^{(N)}}[x_{\mathbf{i}}^{(k)}] = \mathbb{P}_{X^{(N)}}[x_{\mathbf{i}}^{(k-1)}] \mathbb{P}_{X^{(N)}}[x_{i_k}|x_{\mathbf{i}}^{(k-1)}]. \quad (6.1)$$

The above map is one-to-one by construction, and we will denote it by

$$\text{Intv}_X : \text{Rng}(X)^+ \leftrightarrow I(X),$$

where $I(X)$ denotes the set of intervals $I_{x_{\mathbf{i}}^{(k)}}$ in $[0, 1)$ defined above. Notice that the set $\text{BI} \equiv \text{I}(B)$ can be regarded as a special case of $I(X)$ corresponding to the random variable $X = B$ with $\text{Rng}(B) = \{0, 1\}$ and uniform probability distribution. Similarly,

$$\text{Bintv} \equiv \text{Intv}_B : \text{Rng}(B)^+ \leftrightarrow I(B).$$

Now, consider the map

$$\text{Inc}_X : I(X) \rightarrow I(B), \quad \text{Inc}_X(I_{\mathbf{i}}) := \arg \max_{\substack{I_{b_j} \in I(B) \\ I_{b_j} \subseteq I_{x_{\mathbf{i}}}}} |I_{b_j}|$$

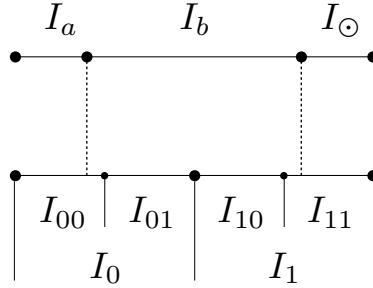


Figure 6.1 An example of interval $I_{x_i} = I_b$ with two largest binary subintervals $I_{b_{j_1}} = I_{01}$ and $I_{b_{j_2}} = I_{10}$.

which maps the interval $I_{x_i} \in \mathcal{I}(X)$ to the largest binary interval I_{b_j} contained in I_{x_i} . However, the map Inc_X is generally not well defined, because given a random variables $X^{(N)}$, it may happen that for some string x_i the interval I_{x_i} contains *two* largest intervals $I_{b_{j_1}}$ and $I_{b_{j_2}}$ of the same size. See e.g. Fig. 6.1. Notice that if the largest subinterval I_{b_j} contained in I_{x_i} is not unique, there can only be two largest included intervals. Indeed, if there are three included intervals $I_{b_{j_1}}, I_{b_{j_2}}$ and $I_{b_{j_3}}$ of size 2^{-k} , then either $I_{b_{j_1}} \cup I_{b_{j_2}}$ or $I_{b_{j_2}} \cup I_{b_{j_3}}$ would be an interval in $\mathcal{I}(B)$, included in I_{x_i} , of size 2^{-k+1} , and thus the three intervals $I_{b_{j_1}}, I_{b_{j_2}}$ and $I_{b_{j_3}}$ could be the largest intervals of $\mathcal{I}(B)$ included in I_{x_i} . Since the inclusion map Inc_X is precisely the one that we will use for arithmetic coding, we complete the definition by the rule that whenever there are two largest binary intervals $I_{b_{j_1}}$ and $I_{b_{j_2}}$ in $\mathcal{I}(B)$ included in I_{x_i} , then

$$\text{Inc}_X(I_{x_i}) := I_{\min\{b_{j_1}, b_{j_2}\}}.$$

Even though this case might seem rather exceptional, it will play a crucial role in the evaluation of the compression rate of the arithmetic coding. Having defined the map Inc_X , one can easily prove that it is invertible on its range, with

$$\text{Inc}_X^{-1} : \mathcal{I}(B) \rightarrow \mathcal{I}(X), \quad \text{Inc}_X^{-1}(I_{b_j}) := \arg \min_{\substack{I_{x_i} \in \mathcal{I}(X) \\ I_{x_i} \supseteq I_{b_j}}} |I_{x_i}|.$$

Now, to encode the string x_i we just locate the corresponding interval, and send the binary string of the largest interval (unique by construction) that lies in I_{x_i} . This corresponds to the one-to-one map

$$\mathcal{E}(x_i) = \text{Intv}_B^{-1} \circ \text{Inc}_X \circ \text{Intv}_X(x_i). \quad (6.2)$$

The fact that the encoding is a one-to-one map guarantees us that we can decode by just using the inverse map $\mathcal{D} := \mathcal{E}^{-1}$. In practice the decoding can be performed on the fly very efficiently by the following routines.

Encoding

The probabilities $P[x_i^{(1)}]$ for all symbols, including the separation and end-of-file characters, are included in the model. At step k , the encoder just evaluates the conditional probabilities $P[x_{i_k}|x_i^{(k-1)}]$ for the particular string $x_i^{(k-1)}$ that has been input up to step $k-1$. Then, they calculate the boundaries of the interval $I_{x_i^{(k)}}$. Considering now binary intervals, they find the largest one included in $I_{x_i^{(k)}}$, say $I_{b_j^{(k')}}$. Normally, $I_{b_j^{(k')}}$ is a unique binary subinterval of $I_{b_j^{(k-1)}}$, thus providing the next encoded bit. It may happen, however, that $I_{b_j^{(k)}}$ is not univocally identified. In this case, they proceed to the step $k+1$. If the last symbol x_{i_N} is such that the largest binary interval included in $I_{x_i^{(N)}}$ is not unique, they adopt the convention described above, taking the leftmost interval.

Decoding

The decoder uses the same algorithm as the one used by the encoder and calculates the probabilities $P[x_i^{(1)}]$ for all symbols. The decoder examines the first bit, and if the corresponding binary interval fits in the interval corresponding to any of the $x_i^{(1)}$'s, it starts decoding the first symbol. Otherwise they consider the binary interval corresponding to the first two bits, and so on increasing the number of bits, until an enclosing interval $I_{x_{i_1}}$ is found. Then they compute the probabilities $P[x_{i_2}|x_{i_1}^{(1)}]$ for all $x_{i_2} \in \text{Rng}(X)$, and deduce the boundaries of the intervals $x_i^{(2)}$. Next, they examine the subsequent encoded bits up to identification of the second Interval $I_{x_{i_2}^{(2)}}$, from which they can decode the second symbol x_{i_2} . They then compute the probabilities $P[x_{i_3}|x_{i_2}^{(2)}]$ for all $x_{i_3} \in \text{Rng}(X)$, and deduce the boundaries of the intervals $I_{x_{i_3}}$. The algorithm is repeated until the ‘end-of-file’ character \odot is decoded.

Compression rate

Arithmetic coding is very nearly optimal. Indeed, by construction the encoding map gives $\mathcal{E}(x_i^{(N)}) = b_j^{(k)}$ with $|I_{b_j^{(k)}}| = \frac{1}{2^k} \leq |I_{x_i^{(N)}}| = \mathbb{P}[x_i^{(N)}]$. Thus, the number k of bits used to encode the string $x_i^{(N)}$ is bounded by

$$k = \log_2 \frac{1}{|I_{b_j^{(k)}}|} \geq \log_2 \frac{1}{\mathbb{P}[x_i^{(N)}]},$$

which is precisely the Shannon information content of the encoded string. In order to find an upper bound to $l(x_i^{(N)})$, one has to find a lower bound to $|I_{b_j^{(k)}}|$, which requires the analysis of the situation in Fig. 6.1: in the case when there are two largest binary subintervals of $I_{x_i^{(N)}}$, indeed one has $|I_{b_j^{(k)}}| \geq |I_{x_i^{(N)}}|/4$, which thus implies

$$k = \log_2 \frac{1}{|I_{b_j^{(k)}}|} \leq \log_2 \frac{4}{|I_{x_i^{(N)}}|} = \log_2 \frac{1}{\mathbb{P}[x_i^{(N)}]} + 2,$$

thus leading to a possible overhead of two bits. Notice however that the overhead is now computed over the full string, namely we have $2/N$ bits overhead per symbol, to be compared with 1 bit/symbol of the Huffman code.

Complexity

Also the complexity of probability tables construction for the arithmetic coding is by far smaller than that of Huffman's code. Indeed, in order to communicate a string of N letters both the encoder and the decoder needed to compute only $N|\text{Rng}(X)|$ conditional probabilities, namely the probabilities of all possible letters, only in those (N) contexts that are actually encountered. This must be compared to the complexity of using Huffman's code with large blocks (N -symbols), where all sequences that may occur in principle must be considered, and all their probabilities must be evaluated. Their number amounts to $|\text{Rng}(X)|^N$, which grows exponentially as a function of the length N , instead of the linear growth required for the arithmetic code. Also, notice the flexibility of the algorithm: it can be used with any alphabet and any probability distribution, depending on the case.

The probabilistic model

For the purpose of evaluating the probabilities it is common to use a Bayesian model, in which the conditional probabilities are given by the Dirichlet rule

$$P[x_{i_k} | x_i^{(k-1)}] = \frac{F_{x_{i_k}}(x_i^{(k-1)}) + \alpha}{\sum_l [F_{x_{i_l}}(x_i^{(k-1)}) + \alpha]}, \quad (6.3)$$

$0 < \alpha \leq 1$, where $F_{x_{i_l}}(x_i^{(k-1)})$ denotes the number of times that x_{i_l} occurs in the string $x_i^{(k-1)}$. In other words, the probability is updated based on the number of occurrences in the previous steps. The parameter α sets an updating responsivity. The most responsive update rule is obtained by $\alpha = 0$. The special case $\alpha = 1$ is the most conservative model, and is known as Laplace rule.

6.4 Lempel-Ziv algorithm

Differently from the arithmetic coding, the Lempel-Ziv algorithm does not rely on probabilistic modelling. Roughly speaking, the compression method works on binary strings, and basically replaces a substring with a pointer to an earlier occurrence of the same substring, that has been suitably recorded in a *dictionary*. The incoming binary string $x_i^{(N)}$ is parsed into substrings, in such a way that every substring is composed by a prefix that has already been included in the dictionary plus an extra bit. For example, if the k -th substring is $= x^{(j)}x_{i_k}$, then the encoder outputs $(s(x^{(j)}), x_{i_k})$, namely the pointer $s(x^{(j)})$ to the previously recorded substring $x^{(j)}$ plus the tail bit x_{i_k} . In the output string the parentheses and comma separators are omitted. The algorithm starts setting by convention $x_0 = \lambda$ and $x^{(0)} = \lambda$, namely for the first string the output is (λ, x_{i_1}) .

The algorithm can be expressed in detail as follows.

1. Append the empty string λ in front of the input string $x_i^{(N)}$, as follows $x_i^{(N)} \mapsto \lambda x_i^{(N)}$.
2. Parse $\lambda x_i^{(N)}$ into substrings, from left to right, in such a way that the j -th substring did not appear in the first $j - 1$ instances.
3. Assign a binary pointer to every string, using ordered binary strings.
4. Re-write every substring—whose general length is denoted by l —as a couple, where the first element is the $l - 1$ -bits prefix string, and the second element is the l -th bit. Substitute the prefix with its pointer to the dictionary. When processing the k -th string, express the pointer using $\lceil \log_2 k \rceil$ bits.
5. Remove all parentheses and commas, thus obtaining the encoded string $\mathcal{E}(x_i^{(N)})$

The decoding algorithm works exactly in the same way, apart from the fact that the dictionary must be decoded.

1. Append the empty string λ in front of the input string $x_i^{(N)}$, as follows $\mathcal{E}(x_i^{(N)}) \mapsto \lambda \mathcal{E}(x_i^{(N)})$.
2. Parse $\lambda \mathcal{E}(x_i^{(N)})$ into substrings, from left to right, in such a way that the j -th substring has a prefix of $\lceil \log_2 j \rceil$ bits and a tail of 1 bit.
3. Assign a string to the binary pointer k by decoding the k -th couple, reading the string in the dictionary slot addressed by the prefix, and adding to it the tail bit.
4. Juxtapose all the dictionary entries from left to right thus obtaining the decoded string $x_i^{(N)}$

Let us see how the Lempel-Ziv algorithm works in some specific examples.

Example 6.6. Consider the input string

1011010100010...

1. We initialise the string appending the empty string λ to the left, thus obtaining $\lambda 1011010100010\dots$
2. We start parsing the string into an ordered *dictionary* of substrings, where the j -th string has not appeared in the first $j - 1$ instances, as follows:

$\lambda, 1, 0, 11, 01, 010, 00, 10, \dots$

3. We include the empty substring λ as the first substring in the dictionary, and assign ordered binary pointers to the substrings in the dictionary, following the order in which they appear.

source string: 1011010100010...									
$x^{(j)}$	λ	1	0	11	01	010	00	10	
$s(x^{(j)})$	λ	1	10	11	100	101	110	111	

4. We re-write every substring as a couple, expressing the pointer by $\lceil \log_2 k \rceil$ bits, where k is the number of strings previously included in the dictionary.

source string: 1011010100010....								
$x^{(j)}$	λ	1	0	11	01	010	00	10
$s(x^{(j)})$	λ	1	10	11	100	101	110	111
(pointer,bit)		($\lambda, 1$)	(0,0)	(01,1)	(10,1)	(100,0)	(010,0)	(001,0)

5. After removing parentheses and commas, we end up with the following encoded output: 1000111011000010000010....

In this example the output is actually longer than the input, because there is not much redundancy in the source string.

Example 6.7. Consider the input string 0000000000001000000000000.

1. Initialised string: $\lambda 0000000000001000000000000$

2. Parsed string

$$\lambda 0, 00, 000, 0000, 001, 00000, 000000$$

3. Dictionary

source string: 000000000000001000000000000....								
$x^{(j)}$	λ	0	00	000	0000	001	00000	000000
$s(x^{(j)})$	λ	1	10	11	100	101	110	111
(pointer,bit)		($\lambda, 0$)	(1,0)	(10,0)	(11,0)	(010,1)	(100,0)	(110,0)

4. Re-writing

source string: 000000000000001000000000000....								
$x^{(j)}$	λ	0	00	000	0000	001	00000	000000
$s(x^{(j)})$	λ	1	10	11	100	101	110	111
(pointer,bit)		($\lambda, 0$)	(1,0)	(10,0)	(11,0)	(010,1)	(100,0)	(110,0)

5. Encoded output: 010100110010110001100....

Exercise 6.1

Decode the string 00101011101100100100011010101000011.

Answer of exercise 6.1

1. Initialised string $\lambda 0010101110110010010001101010101000011$

2. Parsed string

$$\lambda, 0, 01, 010, 111, 0110, 0100, 1000, 1101, 01010, 00011$$

3. Dictionary

source string: 00101011101100100100011010101000011....									
(λ ,0)	(0,1)	(01,0)	(11,1)	(011,0)	(010,0)	(100,0)	(110,1)	(0101,0)	(0001,1)
λ	1	10	11	100	101	110	111	1000	1001
λ	0	1	00	001	000	10	0010	101	0000

4. Encoded output: 0100001000100010101000001

Even though the Lempel-Ziv coding doesn't rely on any probabilistic model, it is possible to show that for an ergodic source (i.e. memoryless on a sufficiently long scale) asymptotically it compresses up to the entropy of the source. The effectiveness of the Lempel-Ziv relies on the fact that in practice many files contain repetition of particularly short sequences of characters, a kind of redundancy which the algorithm is suited to.

Summary on codes

We have seen that symbol codes employ a variable-length code for each symbol of the source alphabet. The optimal code lengths are determined by the probabilities of the different symbols. The Huffman's code constructs an optimal symbol code for a given alphabet and probability distribution. The encoding is lossless, i.e. it is uniquely decodeable, and if the source has the assumed probability distribution, the expected length will stay in the interval $[H, H + 1]$, clearly with possible statistical fluctuations, that can give lengths longer or shorter than the expected length. If the source is not matched to the assumed probability distribution, then the expected length is increased by the relative entropy between the source distribution and the code's implicit distribution. For source with small entropy the code uses at least one bit per source symbol, and essentially no compression is achieved. One can, however, achieve compression by using bigger alphabets made of blocks of characters.

Stream codes, compared with symbol codes, are not constrained to emit at least one bit per input symbol, hence one can encode large strings into a smaller number of bits. Among stream codes, the arithmetic ones combine a probabilistic model with an encoding algorithm that identifies each string with a subinterval of $[0, 1)$ of size equal to the probability of that string under the model. Such coding is almost optimal in the sense that asymptotically the compressed length of a input string equals its Shannon information content, namely the expected length averaged over many strings achieves the Shannon limit. On the opposite side, Lempel-Ziv codes are independent of any probability model, and essentially compress based on the principle of building up a database of strings and encoding the strings into a list of pointers to the database (the length of the pointer will be logarithmically shorter than the string itself). Also the Lempel-Ziv coding achieves the Shannon compression limit asymptotically.

Both arithmetic and Lempel-Ziv encoding are very fragile to any alteration of the encoded string: a single bit error will completely destroy the remaining part of the compressed file. Therefore, if we need to communicate or store over noisy media, it will be essential to use error-correcting codes.

Chapter 7

Lecture 7: Mutual information and Markov chains

In the previous section we have considered encoding in order to achieve compression of information. However, we have implicitly assumed that the channel between the encoding and the decoding was noise-free. In real-life instances, channels are noisy, and even the tiniest error can ruin dramatically the recovery of the encoded file—as we have seen, especially when using the most efficient stream codes.

The aim of the next lectures is then to make the noisy channel effectively behave like a noiseless one. This will be achieved by error-correction coding. Therefore, the communication scheme will be the following. The input stream is first compressed, then encoded with redundancy for later error-correction. The output will then enter the noisy physical channel. At the output of the channel a first decoder will recover errors, and finally a second decoder will decompress the message and restore the original string. This communication scheme is summarised in Fig. 7.1.

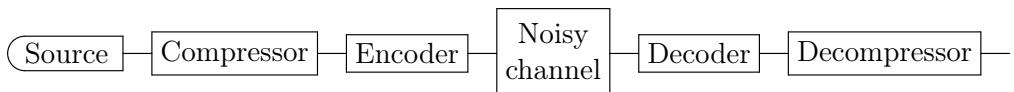


Figure 7.1 Schematic representation of a reliable and efficient communication channel

The problem with errors is that we do not know where they occurred in a given received string. For example one may think that a probability of bit-flip $f = 0.1$ in a symmetric binary channel would affect the rate of transmission of information by 10%, but this is not true: the situation is much worse. Indeed, for $f = 1/2$ there will be no transmitted information at all. A quantification of the transmitted information will be given by the mutual information between input and output. In the next lectures we will see that the channel capacity—i.e. the mutual information maximised over all prior probabilities of the input alphabet—provides the maximum rate at which we can transmit reliably through a noisy channel by a suitable encoding: this is the main content of the famous Shannon's noisy channel theorem.

7.1 Entropies for bi- and multi-variate random variables

7.1.1 Joint entropy and conditional entropy

In the remainder of this lecture we will consider the joint random variable $Z = (X, Y)$ whose range contains ordered couples (x_i, y_j) . As already discussed in section 3.1, one can define the two marginal random variables X and Y with marginal probabilities $p(x_i) = \sum_{y_j \in \text{Rng}(Y)} p(x_i, y_j)$ and $p(y_j) = \sum_{x_i \in \text{Rng}(X)} p(x_i, y_j)$. Notice that generally $p(x_i, y_j) \neq p(x_i)p(y_j)$, i.e. the two marginal random variables need not be independent. More generally, we will consider multivariate random variables $Z = (X_1, X_2, \dots, X_N)$. The marginal distribution for the M -tuple $(X_{i_1}, X_{i_2}, \dots, X_{i_M})$, obtained by excluding the random variables $X_{l_{M+1}}, X_{l_{M+2}}, \dots, X_{l_N}$ is defined by

$$p(x_{j_{i_1}}, x_{j_{i_2}}, \dots, x_{j_{i_M}}) := \sum_{\substack{x_{j_{l_{M+1}}} \in \text{Rng}(X_{l_{M+1}}) \\ x_{j_{l_{M+2}}} \in \text{Rng}(X_{l_{M+2}}) \\ \vdots \\ x_{j_{l_N}} \in \text{Rng}(X_{l_N})}} p(x_{j_1}, x_{j_2}, \dots, x_{j_N}).$$

Similarly, conditional probabilities are defined through Bayes' rule

$$p(x_{j_{i_1}}, x_{j_{i_2}}, \dots, x_{j_{i_M}} | x_{j_{l_{M+1}}}, x_{j_{l_{M+2}}}, \dots, x_{j_{l_N}}) := \frac{p(x_{j_1}, x_{j_2}, \dots, x_{j_N})}{p(x_{j_{l_{M+1}}}, x_{j_{l_{M+2}}}, \dots, x_{j_{l_N}})}.$$

Definition 7.1 (Joint entropy). The joint entropy of two random variables X and Y , denoted as $H(X, Y)$, is defined as

$$H(X, Y) := \sum_{\substack{x_i \in \text{Rng}(X) \\ y_j \in \text{Rng}(Y)}} \mathbb{P}_{X,Y}[X = x_i, Y = y_j] \log_2 \frac{1}{\mathbb{P}_{X,Y}[X = x_i, Y = y_j]} \quad (7.1)$$

$$= \sum_{\substack{x_i \in \text{Rng}(X) \\ y_j \in \text{Rng}(Y)}} p(x_i, y_j) \log_2 \frac{1}{p(x_i, y_j)} \quad (7.2)$$

The joint entropy is defined to be symmetric: $H(Y, X) := H(X, Y)$. Indeed, $\mathbb{P}_{X,Y}[Y = y_i, X = x_j] = \mathbb{P}_{X,Y}[X = x_j, Y = y_i] = p(x_i, y_j)$. Now that we have defined the joint entropy of the random variable X, Y , we will refer to the entropy $H(X)$ as the *marginal entropy* of X . The above definitions are straightforwardly extended to the case of multivariate random variables. We will then have

$$H(X_1, X_2, \dots, X_N) := \sum_{\substack{x_{1_{i_1}} \in \text{Rng}(X_1) \\ x_{2_{i_2}} \in \text{Rng}(X_2) \\ \vdots \\ x_{N_{i_N}} \in \text{Rng}(X_N)}} p(x_{1_{i_1}}, x_{2_{i_2}}, \dots, x_{N_{i_N}}) \log_2 \frac{1}{p(x_{1_{i_1}}, x_{2_{i_2}}, \dots, x_{N_{i_N}})}.$$

Theorem 7.2 (Subadditivity of the Shannon entropy). *For every couple of random variables X and Y the following inequality holds*

$$H(X, Y) \leq H(X) + H(Y), \quad (\text{with equality holding iff } X \text{ and } Y \text{ are independent}). \quad (7.3)$$

Proof. One has

$$\begin{aligned} H(p(x_i, y_j) \| p(x_i)p(y_j)) &= -H(p(x_i, y_j)) - \sum_{\substack{x_i \in \text{Rng}(X) \\ y_j \in \text{Rng}(Y)}} p(x_i, y_j) \log_2(p(x_i)p(y_j)) \\ &= -H(X, Y) + H(X) + H(Y). \end{aligned}$$

The statement is then just an application of Gibbs' inequality. Notice that the equal sign holds iff $p(x_i, y_j) = p(x_i)p(y_j) \forall x_i, y_j$, namely iff X and Y are independent. \square

The theorem then just asserts that the information content of the joint probability $p(x_i, y_j)$ is smaller than or equal to the sum of the two information contents of the two marginals. This can be intuitively understood, considering that the marginal distributions forget about correlations. For example, if we know that the face of a rolled fair die is even we have 1 bit of information. If we know that the face is \square , we have $\log_2 6 = 1 + \log_2 3 = 2.58$ bits. However, if we know both facts together, we do not have 3.58 bits. A more rigorous account for this observation is captured by the result that *conditioning reduces entropy*.

Definition 7.3 (Conditional entropy). The conditional entropy $H(X|Y = y_j)$ of X given $Y = y_j$ is the entropy of the conditional probability distribution

$$p(x_i|y_j) = \mathbb{P}_{X|y_j}[X = x_i|Y = y_j] := \frac{\mathbb{P}_{X,Y}[X = x_i, Y = y_j]}{\mathbb{P}_Y[Y = y_j]}.$$

Explicitly, using the notation $\mathbf{p}_{|y_j} = (p(x_1|y_j), p(x_2|y_j), \dots, p(x_n|y_j))$ one has

$$H(X|Y = y_j) := H(\mathbf{p}_{|y_j}) = \sum_{x_i \in \text{Rng}(X)} p(x_i|y_j) \log_2 \frac{1}{p(x_i|y_j)}.$$

The conditional entropy $H(X|Y)$ of X given Y is the expectation of $H(X|Y = y_j)$ over $\text{Rng}(Y)$, namely

$$\begin{aligned} H(X|Y) &:= \mathbb{E}[H(X|Y = y_j)] = \sum_{y_j \in \text{Rng}(Y)} p(y_j) H(\mathbf{p}_{|y_j}) = \sum_{y_j \in \text{Rng}(Y)} p(y_j) H(X|Y = y_j) \\ &= \sum_{\substack{x_i \in \text{Rng}(X) \\ y_j \in \text{Rng}(Y)}} p(x_i, y_j) \log_2 \frac{1}{p(x_i|y_j)}, \end{aligned}$$

where we used Bayes' rule $p(x_i, y_j) = p(y_j)p(x_i|y_j)$.

We stress that the above definition can be applied also to multivariate random variables $Z = (X_1, X_2, \dots, X_N)$, provided that one subdivides the X_i 's in two multivariate random variables $Y_1 = (X_{i_1}, X_{i_2}, \dots, X_{i_k})$ and $Y_2 = (X_{j_1}, X_{j_2}, \dots, X_{j_l})$, with $\{1, 2, \dots, N\} = \{i_1, i_2, \dots, i_k\} \cup \{j_1, j_2, \dots, j_l\}$, and then

$$H(X_{i_1}, X_{i_2}, \dots, X_{i_k} | X_{j_1}, X_{j_2}, \dots, X_{j_l}) := H(Y_1 | Y_2).$$

Notice that since $H(X|Y = y_j) \geq 0$ for all y_j , also $H(X|Y) \geq 0$.

Lemma 7.4 (Chain rules for entropies). *For the information content and for the conditional entropy one has the chain rules*

$$\log_2 \frac{1}{p(x_i, y_j)} = \log_2 \frac{1}{p(x_i)} + \log_2 \frac{1}{p(y_j|x_i)} = \log_2 \frac{1}{p(y_j)} + \log_2 \frac{1}{p(x_i|y_j)}, \quad (7.4)$$

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y). \quad (7.5)$$

Proof. The identity for the information content is a trivial consequence of its definition. The second identity is obtained from the first one by just taking the expectation value with $\mathbb{P}_{X,Y}$ on both sides. \square

Corollary 7.5. *The following bound holds*

$$H(X, Y) \geq H(X). \quad (7.6)$$

Proof. From the chain rule $H(X, Y) = H(X) + H(Y|X)$ and non negativity of $H(Y|X)$ we can easily establish the bound $H(X, Y) \geq H(X)$. \square

Theorem 7.6 (Conditioning reduces entropy). *The conditional entropy is bounded as follows*

$$0 \leq H(X|Y) \leq H(X), \quad (7.7)$$

and is exactly zero if and only if $X = f(Y)$, whereas it achieves the upper bound if and only if X and Y are independent.

Proof. Non negativity of $H(X|Y)$ follows from its definition. Now, if $X = f(Y)$, then $p(x_i = f(y_j)|y_j) = 1$ and $p(x_i \neq f(y_j)|y_j) = 0$, hence $H(X|Y = y_j) = H(\mathbf{p}_{|y_j}) = 0$ for every $y_j \in \text{Rng}(Y)$, and $H(X|Y) = 0$. Conversely, if $0 < p(x_i|y_j) < 1$ strictly for some x_i, y_j , then $\log_2 \frac{1}{p(x_i|y_j)} > 0$, and consequently $H(X|Y) > 0$ strictly. Thus, $H(X|Y) = 0$ implies that $p(x_i|y_j) = \delta_{x_i, f(y_j)}$, namely $x_i = f(y_j)$. Finally, the upper bound along with its achievement follow from the subadditivity of Shannon entropy in Theorem 7.2, and the chain rules of Lemma 7.4. Indeed, one has $H(X|Y) = H(X, Y) - H(Y)$, and then

$$H(X|Y) - H(X) = H(X, Y) - H(X) - H(Y) \leq 0,$$

with equality holding iff X and Y are independent. \square

The above theorem has the immediate interpretation that on average *conditioning reduces the entropy*, it does not reduce it if the conditional variable is independent of the conditioning one, and it reduces it to zero when the conditional variable is functionally dependent on the conditioning one.

The conditional entropy $H(X|Y)$ has a nice interpretation in quantifying the quality of inference of X from Y , in the sense that, intuitively speaking the smaller the conditional entropy, the better the inference: indeed, the uncertainty about X is decreased by the knowledge of Y . The intuition is strengthened by considering the limiting case $H(X|Y) = 0$ where one has $X = f(Y)$. This interpretation is made more rigorous by Fano's inequality, that we will prove in the next lecture.

Notice that the inequality $H(Y|X) \leq H(Y)$ holds only in average. Indeed, occurrence of a special outcome $Y = y_j$ may increase our uncertainty about X , namely there may exist y_j such that

$$H(X|Y = y_j) > H(X).$$

This is the case, for example, for the following distribution for two bits X, Y .

	$p(0, y)$	$p(1, y)$
$p(x, 0)$	$\frac{1}{8}$	$\frac{1}{8}$
$p(x, 1)$	0	$\frac{3}{4}$

In this case the marginal distribution for Y is obtained by summing the entries in the two rows: $\mathbb{P}_Y[Y = 0] = \frac{1}{4}$ and $\mathbb{P}_Y[Y = 1] = \frac{3}{4}$. Thus, $\mathbb{P}_{X|y=0}[X = 0|Y = 0] = \frac{1}{2}$, $\mathbb{P}_{X|y=0}[X = 1|Y = 0] = \frac{1}{2}$, while $\mathbb{P}_{X|y=1}[X = 0|Y = 1] = 0$, $\mathbb{P}_{X|y=1}[X = 1|Y = 1] = 1$. The marginal distribution for X , instead, is obtained by summing the entries in the two columns. In this way we obtain $\mathbb{P}_X[X = 0] = \frac{1}{8}$, $\mathbb{P}_X[X = 1] = \frac{7}{8}$. One can then easily compute

$$\begin{aligned} H(X) &= H_2\left(\frac{1}{8}\right) = 0.54 \\ H(X|Y = 0) &= H_2\left(\frac{1}{2}\right) = 1, \\ H(X|Y = 1) &= H_2(0) = 0, \\ H(X|Y) &= 0.25. \end{aligned}$$

While, clearly, $H(X|Y) = 0.25 < H(X) = 0.54$, for the case $Y = 0$ we have $H(X|Y = 0) = 1 > H(X)$.

7.1.2 Mutual information

We are now in position to define the mutual information.

Definition 7.7 (Mutual information). The *mutual information between X and Y , denoted by $I(X : Y)$, is defined as*

$$\begin{aligned} I(X : Y) &= H(X) - H(X|Y) = H(Y) - H(Y|X) \\ &= H(X) + H(Y) - H(X, Y). \end{aligned} \tag{7.8}$$

Theorem 7.8 (Properties of the mutual information). *The mutual information satisfies the following properties:*

1. $I(X : Y) \geq 0$;
2. $I(X : Y) = I(Y : X)$;
3. $I(X : Y) \leq H(Y)$ (*equality holds iff $Y = f(X)$*).

Proof. Non negativity of mutual information is a direct consequence of subadditivity of the Shannon entropy. The symmetry is a direct consequence of the definition and of symmetry of $H(X, Y)$. Finally the last item is just a restatement of Theorem 7.6. Indeed,

$$H(X) + H(Y) - H(X, Y) = H(Y) - H(Y|X) \leq H(Y),$$

with equality holding iff $H(Y|X) = 0$. \square

Inferring the output from the input

The mutual information $I(X : Y)$ represents the amount of information that one can infer about X by knowing Y . This can be seen with the following steps

1. Before reading Y our knowledge about X is given by the marginal probability distribution $p(x_i)$, and uncertainty about X is quantified by the Shannon entropy $H(X)$.
2. We now read Y , getting outcome $Y = y_j$. We know the conditional probability $p(x_i|y_j)$. For example, this is the case if X and Y represent the input and the output of a channel, since the probability $p(y_j|x_i)$ of reading y_j at the output for input x_i characterises the channel: it is the so-called *transfer matrix*. Hence, the probability that the input was $X = x_i$ given that we read $Y = y_j$ is given by the Bayes rule

$$p(x_i|y_j) = \frac{p(x_i, y_j)}{p(y_j)} = \frac{p(x_i)p(y_j|x_i)}{\sum_{x_k \in \text{Rng}(X)} p(x_k)p(y_j|x_k)}. \quad (7.9)$$

In the process we have updated the *prior* probability to the *posterior* one:

$$\mathbf{P} \longrightarrow \mathbf{P}_{|y_j}.$$

3. How much information have we gained about X from reading $Y = y_j$? This is quantified by the "decrease" (possibly negative) of our uncertainty from $H(X) \equiv H(\mathbf{P})$ to $H(X|Y = y_j) \equiv H(\mathbf{P}_{|y_j})$, namely the difference $H(X) - H(X|Y = y_j)$.
4. Upon averaging over all possible outcomes $y_j \in \text{Rng}(Y)$ we get

$$\sum_{y_j \in \text{Rng}(Y)} p(y_j)[H(X) - H(X|Y = y_j)] = H(X) - H(X|Y) = I(X : Y),$$

namely, the amount of information that we gained in average on X from Y is given by the mutual information $I(X : Y)$.

From the above considerations, we see that the mutual information quantifies the amount of information that the output of a channel conveys about the input. Our aim is now to find ways of using the channel in a reliable way, with negligible error probability in recovering the transmitted bits. Shannon's noisy channel theorem shows that this is possible, and the mutual information exactly quantifies the amount of information that can be reliably transmitted.

Regarding the practical problem of evaluating numerically the mutual information for a channel, we notice that using $I(X : Y) = H(Y) - H(Y|X)$ instead of $I(X : Y) = H(X) - H(X|Y)$ is more convenient, as it avoids evaluating the posterior probabilities in Eq. (7.9).

7.2 Markov chains

In the analysis of channel coding, an important role is played by Markov chains. This notion is conceptually very relevant, as it will allow us to prove the *data processing theorem*, and it turns out to be technically useful to characterise the extreme cases where the inequalities of strong subadditivity hold tightly.

Definition 7.9 (Markov chain). We say that the multivariate random variable

$$Z = (X_1, X_2, X_3, \dots, X_n, X_{n+1}, \dots)$$

is a Markov chain, denoted as

$$X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow \dots \rightarrow X_n \rightarrow X_{n+1} \rightarrow \dots,$$

if for all $n \geq 1$

$$\begin{aligned} & \mathbb{P}(X_{n+1} = x_{i_{n+1}} | X_n = x_{i_n}, X_{n-1} = x_{i_{n-1}}, \dots, X_1 = x_{i_1}) = \\ & \mathbb{P}(X_{n+1} = x_{i_{n+1}} | X_n = x_{i_n}). \end{aligned} \tag{7.10}$$

Lemma 7.10. *Given three random variables X, Y and Z, one has that $X \rightarrow Y \rightarrow Z$ iff $Z \rightarrow Y \rightarrow X$.*

Proof. Indeed, one has

$$p(x_i | y_j, z_k) = \frac{p(x_i, y_j, z_k)}{p(y_j, z_k)} = \frac{p(z_k | x_i, y_j)p(x_i, y_j)}{p(z_k | y_j)p(y_j)}.$$

Now, using the definition of Markov chain, $p(z_k | x_i, y_j) = p(z_k | y_j)$, and then

$$p(x_i | y_j, z_k) = \frac{p(z_k | y_j)p(x_i, y_j)}{p(z_k | y_j)p(y_j)} = p(x_i | y_j). \quad \square$$

For this reason Markov chains are also denoted as

$$X_1 \leftrightarrow X_2 \leftrightarrow X_3 \leftrightarrow \dots \leftrightarrow X_n \leftrightarrow X_{n+1} \leftrightarrow \dots,$$

The “Markovianity” property is the same as conditional independence. More precisely, we have

Lemma 7.11 (Conditional independence). $X \leftrightarrow Y \leftrightarrow Z$ if and only if

$$p(x_i, z_k | y_j) = p(x_i | y_j)p(z_k | y_j). \quad (7.11)$$

Proof. If $X \leftrightarrow Y \leftrightarrow Z$ one has

$$p(x_i, z_k | y_j) = \frac{p(x_i, y_j, z_k)}{p(y_j)} = \frac{p(z_k | x_i, y_j)p(x_i, y_j)}{p(y_j)} = \frac{p(z_k | y_j)p(x_i, y_j)}{p(y_j)} = p(x_i | y_j)p(z_k | y_j).$$

On the other hand, if $p(x_i, z_k | y_j) = p(x_i | y_j)p(z_k | y_j)$ then one has

$$p(x_i | y_j, z_k) = \frac{p(x_i, y_j, z_k)}{p(y_j, z_k)} = \frac{p(x_i, z_k | y_j)p(y_j)}{p(y_j, z_k)} = \frac{p(x_i | y_j)p(z_k | y_j)p(y_j)}{p(y_j, z_k)} = p(x_i | y_j). \quad \square$$

From the above lemma the following theorem follows.

Theorem 7.12 (Strong subadditivity of the joint entropy). *The following inequality holds*

$$H(X, Y, Z) + H(Y) \leq H(X, Y) + H(Y, Z), \quad (\text{equality holds iff } X \leftrightarrow Y \leftrightarrow Z). \quad (7.12)$$

Proof. If we treat $p(x_i, z_k | y_j)$ as a usual joint distribution, where y_j is just a parameter, we can apply subadditivity of Shannon entropy to obtain

$$H(X|Y = y_j) + H(Z|Y = y_j) - H(X, Z|Y = y_j) \geq 0.$$

Averaging over y_j one has

$$\begin{aligned} & H(X|Y) + H(Z|Y) - H(X, Z|Y) \\ &= \sum_{y_j \in \text{Rng}(Y)} p(y_j)[H(X|Y = y_j) + H(Z|Y = y_j) - H(X, Z|Y = y_j)] \geq 0. \end{aligned} \quad (7.13)$$

Then one has

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

namely

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

Considering equation (7.13), one easily realises that equality holds iff $H(X|Y = y_j) + H(Z|Y = y_j) - H(X, Z|Y = y_j) = 0$ for every y_j , and by the subadditivity of Shannon entropy in theorem 7.2 this happens iff

$$p(x_i, z_k | y_j) = p(x_i | y_j)p(z_k | y_j),$$

namely iff conditional independence holds, which in turn is equivalent of the condition that X, Y, Z is a Markov chain (see Lemma 7.11). \square

From the above theorem we can generalize the assertion from theorem 7.6 that conditioning reduces the entropy

Corollary 7.13 (Conditioning reduces the entropy). *The following bound holds*

$$H(X|Y, Z) \leq H(X|Y), \quad (7.14)$$

with the equality iff $X \leftrightarrow Y \leftrightarrow Z$.

Proof. In fact, using the definition of conditional entropy, the statement can be rewritten as follows

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

which is just strong subadditivity of theorem 7.12. The case of equality is just the definition of Markov chain.

Alternatively, one can just observe that the statement is the thesis of theorem 7.6 for $p(x_i, z_k | y_j)$, averaged over y_j . \square

7.3 Data Processing theorem

We now prove the data processing theorem for Markov chains, and other preliminary results that will be used in the proof of Shannon's noisy channel coding theorem. In particular, we will prove Fano's inequality, which provides a quantitative evidence for the intuitive notion that "the better is the guess of Y given the knowledge of X , the smaller is $H(Y|X)$ ".

Theorem 7.14 (Data-processing theorem). *If $X \leftrightarrow Y \leftrightarrow Z$ is a Markov chain the following chain of bounds holds*

$$I(X : Z) \leq I(X : Y) \leq H(X), \quad (7.15)$$

with the last bound saturated iff $X = f(Y)$. Equivalently,

$$H(X|Z) \geq H(X|Y) \geq 0, \quad (7.16)$$

with the last bound saturated iff $X = f(Y)$

Proof. Equivalence of equations (7.15) and (7.16) is obtained by considering that

$$\begin{aligned} I(X : Z) \leq I(X : Y) &\Leftrightarrow H(X) - I(X : Z) \geq H(X) - I(X : Y) \\ I(X : Y) \leq H(X) &\Leftrightarrow H(X) - I(X : Y) \geq 0, \end{aligned}$$

and reminding the definition of $I(X : Y) = H(X) - H(X|Y)$. The second bound, along with conditions for its saturation is thus already proved in theorem 7.6, as $H(X|Y) \geq 0$ and $H(X|Y) = 0$ iff $X = f(Y)$. Let us now consider the first bound. Since $X \leftrightarrow Y \leftrightarrow Z$ is a Markov chain, by the strong subadditivity theorem 7.12, and in particular by corollary 7.13, we have

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

This allows us to rewrite the first bound, in the form of equation 7.16, as

$$H(X|Z) \geq H(X|Y, Z), \quad (7.17)$$

which is equivalent to strong subadditivity [see equation (7.14)]. \square

Notice that, despite $X \leftrightarrow Y \leftrightarrow Z$ is a Markov chain, the bound in Eq. (7.17) cannot be saturated, because this would require the Markov condition to hold with the ordering $X \leftrightarrow Z \leftrightarrow Y$.

The data-processing theorem states that *forgetful* data-processing can only destroy information. This is a relevant and non trivial assertion, even though we have to be careful with both words “data processing” and “information”. Since the theorem holds for a Markov chain of random variables, we can regard the “data-processing” also as the input-output correlation in an analog device, such as an amplifier, or a telephone line, etc.

Example 7.15 (Old vinyl music). As an example, consider the case of amplification of vibrations by the stylus of a gramophon. At first sight one can think that the amplifier increases the amount of information encoded on the vinyl record. However, by the data processing theorem it can only degrade it. Why does the amplifier work, then? The stylus vibrations are too faint to be heard, and the amplifier enhances them above our audibility threshold. However, there is more information on the tracks of the vinyl record than that coming to our ears.

Example 7.16 (Recovering information from noise by digitalisation). Can one digitalise the sound from a gramophon, and then enhance it by a smart software on a computer? Of course, this is possible, however one cannot recover the information lost in the noisy amplifier: The post-processing software is just a forgetful data-processing device, because it is completely independent of the original track. This way of processing corresponds to a step in a Markov chain. According to the data-processing theorem, it will be useless in re-gaining the information lost. However, one may hear a better sound (e.g. after applying a Dolby filter), just because the processing optimises the sound for the characteristics of the human ear.

Example 7.17 (Saturable amplifiers). There is no linear analog signal amplifier: for high input power, the output will saturate. The high-power signal cut due to saturation cuts high-energy fluctuations, thus effectively improving the signal-to-noise ratio (i.e. $\mathbb{E}(Y^2)/\sigma_Y^2$). Does this mean that saturation is a positive feature? Of course not. Indeed, improving the signal to noise ratio does not increase the mutual information with the source if the output signal has little to do with the input. The data-processing theorem immediately helps clarifying this point.

Chapter 8

Lecture 8: Channel capacity and joint typicality

8.1 Fano's inequality

We now prove a key result for noisy channel coding: Fano's inequality. This result will provide a more quantitative account of the intuition that conditional entropy measures the guessing uncertainty of the random variable X upon knowledge of Y , and will be used in the proof of the Shannon noisy channel theorem. First, let us prove the following chaining rules for conditional entropy.

Theorem 8.1 (Chaining rules). *The following rules hold for every $N \geq 1$*

$$H(X_1, X_2, \dots, X_N | Y) = \sum_{i=1}^N H(X_i | Y, X_1, X_2, \dots, X_{i-1}), \quad (8.1)$$

where $Y, X_1, X_2, \dots, X_{i-1}$ is by definition Y for $i = 1$.

Proof. The proof is by induction. For $N = 1$ the statement is trivial, since it reduces to $H(X_1, | Y) = H(X_1, | Y)$. Let us now suppose it is true for $N = n$. Then

$$\begin{aligned} H(X_1, X_2, \dots, X_{n+1} | Y) &= H(X_1, X_2, \dots, X_{n+1}, Y) - H(Y) \\ &= H(X_1, X_2, \dots, X_{n+1}, Y) - H(X_1, X_2, \dots, X_n, Y) \\ &\quad + H(X_1, X_2, \dots, X_n, Y) - H(Y) \\ &= H(X_{n+1} | Y, X_1, X_2, \dots, X_n) + H(X_1, X_2, \dots, X_n | Y). \end{aligned}$$

If we now use the induction hypothesis, we have

$$\begin{aligned} H(X_1, X_2, \dots, X_{n+1} | Y) &= H(X_{n+1} | Y, X_1, X_2, \dots, X_n) \\ &\quad + \sum_{j=1}^n H(X_j | Y, X_1, X_2, \dots, X_{j-1}) \\ &= \sum_{j=1}^{n+1} H(X_j | Y, X_1, X_2, \dots, X_{j-1}) \quad \square \end{aligned}$$

We can now prove Fano's theorem. The situation to which the theorem refers is the following. Suppose that there are two random variables, X and Y , X representing Alice's source. Bob has only access to Y , e.g. because it is the output of a noisy channel fed with input X , and he tries to guess the value of X . Whatever strategy Bob uses, his guess \hat{X} is a random variable that represents the third element of a Markov chain $X \leftrightarrow Y \leftrightarrow \hat{X}$, because Bob does not have direct access to X , and he can only produce his guess out of a—possibly random—algorithm that processes Y . Thus, $p(\hat{x}|y, x) = p(\hat{x}|y)$.

Theorem 8.2 (Fano's inequality). *Let X, Y be a pair of random variables. Let \hat{X} be a third random variable, representing a guess of the value of X provided that one knows Y . Clearly, $X \leftrightarrow Y \leftrightarrow \hat{X}$. The following bound holds*

$$H(X|Y) \leq H(X|\hat{X}) \leq H_2(p_e) + p_e \log_2(|\text{Rng}(X)| - 1), \quad (8.2)$$

where p_e denotes the error probability in the guess, namely $p_e := 1 - \mathbb{P}_{X,\hat{X}}[X = \hat{X}]$.

Proof. Let us define the binary random variable E with $\text{Rng}(E) = \{0, 1\}$ as a function of X and \hat{X} as follows

$$E := \begin{cases} 0 & \hat{x} = x, \\ 1 & \hat{x} \neq x. \end{cases}$$

One then has

$$\begin{aligned} H(E) &= \mathbb{P}_E(E = 0) \log_2 \frac{1}{\mathbb{P}_E(E = 0)} + \mathbb{P}_E(E = 1) \log_2 \frac{1}{\mathbb{P}_E(E = 1)} \\ &= (1 - p_e) \log_2 \frac{1}{1 - p_e} + p_e \log_2 \frac{1}{p_e} = H_2(p_e). \end{aligned}$$

Applying the chaining rule to $H(E, X|\hat{X})$ one obtains

$$H(E, X|\hat{X}) = H(E|X, \hat{X}) + H(X|\hat{X}).$$

Since $E = f(X, \hat{X})$, by theorem 7.6 the conditional entropy $H(E|X, \hat{X})$ is null, thus

$$H(E, X|\hat{X}) = H(X|\hat{X}).$$

On the other hand, applying a different chaining rule, we have

$$H(E, X|\hat{X}) = H(X|E, \hat{X}) + H(E|\hat{X}).$$

By theorem 7.6, conditioning reduces entropy, then $H(E|\hat{X}) \leq H(E) = H_2(p_e)$. Thus, we have

$$H(E, X|\hat{X}) \leq H(X|E, \hat{X}) + H_2(p_e).$$

Moreover, we can bound the first term on r.h.s. by the observing that

$$H(X|E, \hat{X}) = p_e H(X|\hat{X}, E = 1) + (1 - p_e) H(X|\hat{X}, E = 0),$$

and since $\mathbb{P}(X = x | \hat{X} = \hat{x}, E = 0) = \delta_{x,\hat{x}}$, we have $H(X|\hat{X}, E = 0) = 0$. Then we can write

$$\begin{aligned} H(X|E, \hat{X}) &= p_e H(X|\hat{X}, E = 1) \\ &= p_e \sum_{\hat{x} \in \text{Rng}(X)} p(\hat{x}|1) \sum_{x \in \text{Rng}(X) \setminus \{\hat{x}\}} p(x|\hat{x}, 1) \log_2 \frac{1}{p(x|\hat{x}, 1)} \\ &\leq p_e \sum_{\hat{x} \in \text{Rng}(X)} p(\hat{x}|1) \log_2 [|\text{Rng}(X)| - 1] \\ &= p_e \log_2 [|\text{Rng}(X)| - 1]. \end{aligned}$$

We then proved that

$$H(X|\hat{X}) \leq H_2(p_e) + p_e \log_2 [|\text{Rng}(X)| - 1].$$

As to the first inequality in the thesis, we use the fact that $X \leftrightarrow Y \leftrightarrow \hat{X}$ is a Markov chain, and by the data processing theorem

$$I(X : \hat{X}) \leq I(X : Y) \Leftrightarrow H(X|Y) \leq H(X|\hat{X}). \quad \square$$

8.2 Discrete memoryless channels

In section 2.6 we introduced classical systems, their types, their states and their transformations. Transformations describe transmission media, or more generally any information processing algorithm, which is a sequence of physical or logical operations, producing an *output* state when fed with an *input* state. Also transformations have a type: If the input is the state of a system of type n and the output is a system of type m the transformation has type $n \rightarrow m$.

Transformations allow us to determine the state of the output system provided that we know the state of the input one. As we discussed in lecture 3, states of a system X of type n represent the (probabilistic) preparation procedure for the system—a register of cells with n levels. States of a system of type n are thus probability distributions for a random variable X with $|\text{Rng}(X)| = n$. A transformation of type $n \rightarrow m$ turns states of system X into state of system Y of type m . The diagrammatic representation of a transformation \mathcal{C} of type $n \rightarrow m$ is the following

$$\xrightarrow[X]{\mathcal{C}} \xrightarrow[Y]{\mathcal{C}}, \quad \xrightarrow[\mathbf{q}]{\mathcal{C}} \xrightarrow[Y]{\mathcal{C}} = \xrightarrow[\mathbf{p}]{\mathcal{C}} \xrightarrow[X]{\mathcal{C}} \xrightarrow[Y]{\mathcal{C}}.$$

The probability distribution \mathbf{p} , corresponding to the state of system X describes a random preparation algorithm that produces an excitation of the j -th level of the cell X with probability p_j . A transformation represents a second algorithm that reads the content of the cell X and randomly prepares system Y in a new state \mathbf{q} .

Clearly, the transformation algorithm cannot behave differently depending on the probability distribution \mathbf{p} . For example, suppose that X is a bit B , prepared by reading the outcome of a coin toss. Let Y be a second bit B' . A transformation \mathcal{C} of type $2 \rightarrow 2$ represents a transmission of the state of B to the state of B' , which cannot be sensitive to the bias p of the coin used to prepare the input system B . The transformation will just

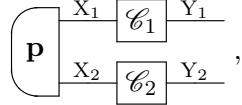
read the content of the input cell, and upon reading 0 performs a protocol, while upon reading 1 performs another one.

The above requirement on the transformation \mathcal{C} , of behaving in a fixed way, independently of the distribution \mathbf{p} , is mathematically expressed through the property of *linearity*. The state \mathbf{q} of Y_m is a linear function of the state \mathbf{p} of X_n . Then, the transformation \mathcal{C} corresponds to a $m \times n$ matrix $Q = [Q_{i,j}]$, such that $q_i = \sum_{j=1}^n Q_{i,j} p_j$. In order to respect positivity and (sub-)normalisation, the linear transformation must have positive elements $Q_{ij} \geq 0$, and every column must be (sub-)normalised: $\sum_{i=1}^m Q_{i,j} \leq 1$ for every j .

A *deterministic* transformation is a normalised one, i.e. $\sum_{i=1}^m Q_{i,j} = 1$ for every j . Normalised transformations are usually referred to as *channels*, and the corresponding matrices are *stochastic* matrices, whose entries $Q_{i,j}$ can be interpreted as a conditional probability distribution $Q_{i,j} = p(y_i|x_j)$, and the output state $\mathbf{q} = (q(y_1), q(y_2), \dots, q(y_m))$ is obtained as the marginal distribution

$$q(y_i) = \sum_{j=1}^n q(y_i|x_j)p(x_j).$$

The parallel action of a channel \mathcal{C}_1 from X_1 to Y_1 and a channel \mathcal{C}_2 from X_2 to Y_2 , as in the following diagram



is represented by the tensor product of matrices $Q^{(1)}$ and $Q^{(2)}$, namely

$$p(y_{i_1}, y_{i_2}) = \sum_{j_1, j_2} Q_{i_1, j_1}^{(1)} Q_{i_2, j_2}^{(2)} p(x_{j_1}, x_{j_2}).$$

The memoryless channel is then defined as follows.

Definition 8.3 (Discrete memoryless channel). A discrete memoryless channel \mathcal{C} of type $n \rightarrow m$ is characterised by an input and an output system X of type n and Y of type m , respectively, and by a complete set of conditional probabilities $Q_{i,j} = p(y_i|x_j)$ for all $x_j \in \text{Rng}(X)$ and $y_i \in \text{Rng}(Y)$, such that $\mathbf{p}_Y = Q\mathbf{p}_X$.

The definition can be generalised to the case of N uses of the same channel \mathcal{C} as follows.

Definition 8.4 (Extended channel). For a discrete memoryless channel \mathcal{C} we consider the channel \mathcal{C}^N corresponding to N uses of the channel \mathcal{C} , having input and output systems of type n^N and m^N with ranges $\text{Rng}(X)^N$ and $\text{Rng}(Y)^N$, respectively, and with conditional probability given by

$$p(y_{\mathbf{j}}^{(N)}|x_{\mathbf{i}}^{(N)}) = \prod_{j=1}^N p(y_{i_j}|x_{i_j}). \quad (8.3)$$

Alternatively, the transfer matrix of \mathcal{C}^N is given by the N -th tensor power $Q^{\otimes N}$ of the transfer matrix Q , where

$$(Q^{\otimes N})_{\mathbf{i}, \mathbf{j}} = Q_{i_1, j_1} Q_{i_2, j_2} \dots Q_{i_N, j_N}, \quad (8.4)$$

$\mathbf{i} = (i_1, i_2, \dots, i_N)$, and $\mathbf{j} = (j_1, j_2, \dots, j_N)$ denoting poly-indices.

Remark 5. This model of channel actually hides an important underlying assumption: All the uses of the channel are considered to be independent and identical, pretty much in the same way as N uses of a source are assumed to prepare N i.i.d. random variables. This assumption is based on the hypothesis that the channel has no internal memory, allowing it to count the uses. Otherwise, every use could end up in a different behaviour, possibly influenced by the previous input states. These models of channels have been studied, also in the quantum case, under the name of *channels with memory* or *memory channels*, for obvious reasons. However, in the present course we will focus on the case of *memoryless channels*.

Let us now have a look at a few examples of discrete memoryless channels.

Example 8.5 (Binary symmetric channel). The input and output system types are 2, thus $\text{Rng}(X) = \text{Rng}(Y) = \{0, 1\}$. The transfer matrix is given by

$$Q = \begin{pmatrix} 1-f & f \\ f & 1-f \end{pmatrix},$$

where $0 < f < 1/2$ is the error probability for every input.

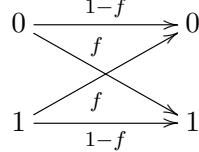


Figure 8.1 Graphical representation of the binary symmetric channel.

Example 8.6 (Binary erasure channel). The input system type is 2, thus the input random variable X is a bit with $\text{Rng}(X) = \{0, 1\}$, while the output system type is 3, namely Y has a range with cardinality 3: we choose $\text{Rng}(Y) = \{0, ?, 1\}$, which makes the physical interpretation of the channel clear. The transfer matrix of the binary erasure channel is given by

$$Q = \begin{pmatrix} 1-f & 0 \\ f & f \\ 0 & 1-f \end{pmatrix},$$

where $0 < f < 1/2$ describes the probability of losing the input, and is independent of the particular input. In other words, there is a probability of loss, however errors are

always detected because they produce the outcome “?”. This makes the error probability exactly null, while it decreases the rate, because in general more than one transmission is required in order to deliver the input bit.

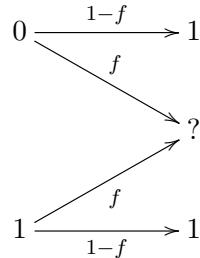


Figure 8.2 Graphical representation of the binary erasure channel.

Example 8.6 is particularly interesting, as it represents the simplest case of nonideal channel with no error.

Example 8.7 (Binary Z channel). In this case $\text{Rng}(X) = \text{Rng}(Y) = \{0, 1\}$ and the transfer matrix is

$$Q = \begin{pmatrix} 1 & f \\ 0 & 1-f \end{pmatrix}$$

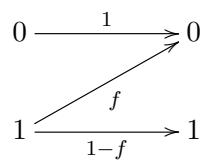


Figure 8.3 Graphical representation of the binary Z channel.

Example 8.8 (Noisy typewriter). In this case $\text{Rng}(X) = \text{Rng}(Y) = \{A, B, \dots, Z, -\}$ One can think of the letters as if they were arranged on a circle, and when the typist presses one button $X = x_i$ then one of the three letters consisting in the input x_i itself and its two neighbours x_{i+1}, x_{i-1} is produced, where \oplus and \ominus denote sum and difference modulo 27. So if the button B is pressed, the output can be A, B or C, each one with probability

$\frac{1}{3}$. The transfer matrix is given by

$$Q = \begin{pmatrix} \frac{1}{3} & \frac{1}{3} & 0 & \dots & 0 & 0 & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & \dots & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & \dots & 0 & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} & \dots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & 0 & 0 & \dots & 0 & \frac{1}{3} & \frac{1}{3} \end{pmatrix} \quad (8.5)$$

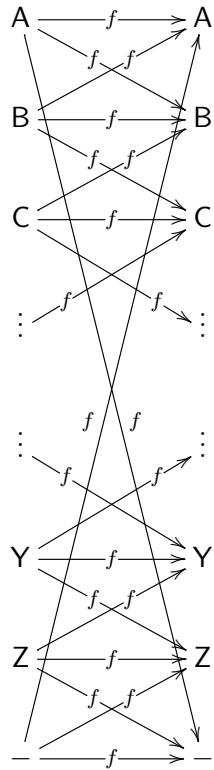


Figure 8.4 Graphical representation of the noisy typewriter channel, where $f = 1/3$.

8.3 Capacity of a channel

The mutual information between input and output of a channel \mathcal{C} is not just a function of the channel Q —from now on we assume that the channel is completely specified by its transfer matrix—but depends also on the *prior* probability distribution \mathbb{P}_X of the random variable representing the state of the input system X . Indeed, by Eq. (8.3), we have

$$\mathbb{P}_{X,Y}(x_i, y_j) = Q_{j,i} \mathbb{P}_X(x_i). \quad (8.6)$$

In order to define a quantity that only depends on Q , we then maximise the mutual information over such prior distribution: This quantity is the so-called *channel capacity*.

Definition 8.9 (Channel capacity). For a discrete memoryless channel Q , the channel capacity is defined as

$$C(Q) = \max_{\mathbb{P}_X} I(X : Y). \quad (8.7)$$

Exercise 8.1

Evaluate the channel capacity of the noisy typewriter channel.

Answer of exercise 8.1

We first evaluate the mutual information for an arbitrary prior. The conditional probabilities are given by the only non-null elements

$$\mathbb{P}_{Y|X=x}(Y = x|X = x) = \mathbb{P}_{Y|X=x}(Y = x \oplus 1|X = x) = \mathbb{P}_{Y|X=x}(Y = x \ominus 1|X = x) = \frac{1}{3},$$

which gives

$$\begin{aligned} H(Y|X = x) &\equiv H(\mathbf{p}|_x) = \sum_{i=1}^{27} \mathbb{P}_{Y|X=x}(Y = y_i|X = x) \log_2 \frac{1}{\mathbb{P}_{Y|X=x}(Y = y_i|X = x)} \\ &= 3 \times \frac{1}{3} \log_2 3 = \log_2 3. \end{aligned}$$

independently of x , hence $H(Y|X) = \log_2 3$. The mutual information is then bounded as

$$I(X : Y) = H(Y) - \log_2 3 \leq \log_2 27 - \log_2 3 = \log_2 9.$$

The question is now whether we can saturate the bound for a suitable choice of \mathbb{P}_X , namely whether we can make the distribution \mathbb{P}_Y uniform. The marginal probability \mathbb{P}_Y is given by

$$\mathbb{P}_Y(y) = \sum_i \mathbb{P}_{Y|X=x_i}(Y = y|X = x_i) \mathbb{P}_X(x_i) = \frac{1}{3} (\mathbb{P}_X(y) + \mathbb{P}_X(y \oplus 1) + \mathbb{P}_X(y \ominus 1)),$$

which can be made uniform, e.g. for uniform \mathbb{P}_X . In such case, then, $H(Y) = \log_2 27$, and we have

$$C(Q) = \max_{\mathbb{P}_X} \{H(Y) - H(Y|X)\} = \log_2 9.$$

For the noisy typewriter channel we can actually achieve the capacity by coding only on a *non-confusable set* of characters, picking out one symbol out of three, e.g. B, E, H... as in figure 8.5. In this way, when the receiver gets A, B or C he/she knows for sure that the letter B was transmitted, and similarly for all the remaining groups of three subsequent letters. The communication by non-confusable sets is then equivalent to a new ideal, noiseless channel for a system of type 9, namely with a range of nine letters. Thus, the new channel has capacity $\log_2 9$.

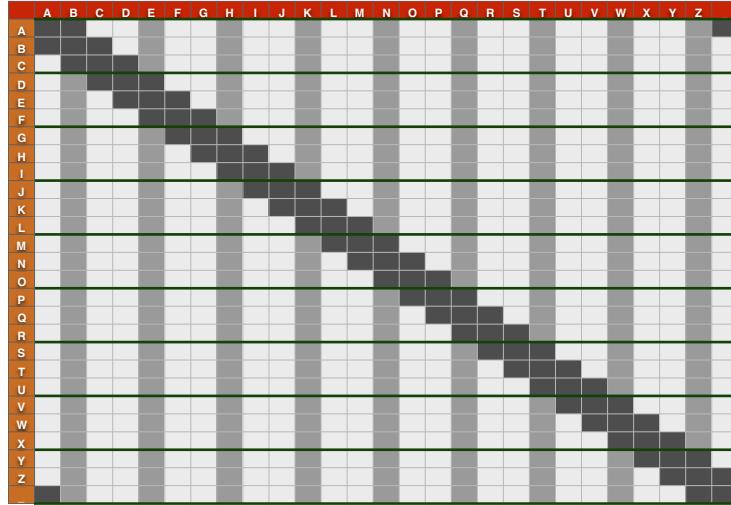


Figure 8.5 A non-confusable subset of inputs for the noisy typewriter.

Exercise 8.2

Evaluate the channel capacity of the binary symmetric channel.

Answer of exercise 8.2

We evaluate the mutual information through the following expression

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

The output marginal probability is given by

$$\begin{aligned} \mathbb{P}_Y(Y = 0) &= \mathbb{P}_{Y|X=0}(Y = 0|X = 0)\mathbb{P}_X(0) + \mathbb{P}_{Y|X=1}(Y = 0|X = 1)\mathbb{P}_X(1) \\ &= (1 - f)p_0 + f(1 - p_0), \end{aligned}$$

where $p_0 := \mathbb{P}_X(0)$. Therefore, we have

$$H(Y) = H_2((1 - f)p_0 + (1 - p_0)f).$$

On the other hand we have

$$H(Y|X) = \sum_x p_x H(Y|X = x) = \sum_x p_x H_2(f) = H_2(f),$$

hence, overall we have

$$I(X : Y) = H_2((1 - f)p_0 + (1 - p_0)f) - H_2(f).$$

Since $H_2(x) \leq 1$ is maximal for $x = 1/2$, the channel capacity is achieved for $p_0 = \frac{1}{2}$, and is given by

$$C(f) = 1 - H_2(f).$$

According to the result of Exercise 8.2, for $f = 0.1$ one has $C = 0.53$. According to the Shannon noisy-channel theorem, which will be presented in the next chapter, it is possible to have reliable communication at the rate equal to the channel capacity with vanishingly small probability of error for sufficiently long strings. Therefore, as already mentioned when discussing the majority vote error-correction, if you have a disk drive whose error probability is that of a binary symmetric channel with $f = 0.1$, in order to achieve a bit error probability 10^{-15} using a majority-voting code you would need *sixty* disks, whereas according to the second Shannon theorem *only two* disks are needed.

Exercise 8.3

Evaluate the channel capacity of the binary erasure channel.

Answer of exercise 8.3

The mutual information $I(X : Y)$ is evaluated as

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

The conditional probability $\mathbb{P}_{Y|X=x}$ is given by

$$\begin{array}{ll} \mathbb{P}_{Y|X=0}(Y = 0|X = 0) = 1 - f & \mathbb{P}_{Y|X=1}(Y = 0|X = 1) = 0 \\ \mathbb{P}_{Y|X=0}(Y = ?|X = 0) = f & \mathbb{P}_{Y|X=1}(Y = ?|X = 1) = f \\ \mathbb{P}_{Y|X=0}(Y = 1|X = 0) = 0 & \mathbb{P}_{Y|X=1}(Y = 1|X = 1) = 1 - f. \end{array}$$

Then, the conditional entropy is

$$H(Y|X = 0) = H(Y|X = 1) = H_2(f), \Rightarrow H(Y|X) = H_2(f).$$

Now, the joint probability $\mathbb{P}_{X,Y}$ is given by

$$\begin{array}{ll} \mathbb{P}_{Y,X}(Y = 0, X = 0) = p_0(1 - f) & \mathbb{P}_{Y,X}(Y = 0, X = 1) = 0 \\ \mathbb{P}_{Y,X}(Y = ?, X = 0) = p_0f & \mathbb{P}_{Y,X}(Y = ?, X = 1) = p_1f \\ \mathbb{P}_{Y,X}(Y = 1, X = 0) = 0 & \mathbb{P}_{Y,X}(Y = 1, X = 1) = p_1(1 - f), \end{array}$$

And the marginal probability \mathbb{P}_Y is then

$$\mathbb{P}_Y(0) = p_0(1 - f), \quad \mathbb{P}_Y(?) = f, \quad \mathbb{P}_Y(1) = p_1(1 - f).$$

Finally, this allows us to calculate $H(Y)$ that amounts to

$$\begin{aligned} H(Y) &= -p_0(1 - f) \log_2[p_0(1 - f)] - f \log_2 f - p_1(1 - f) \log_2[p_1(1 - f)] \\ &= (1 - f)H_2(p_0) + H_2(f) = (1 - f)H(X) + H_2(f). \end{aligned}$$

Finally, we have

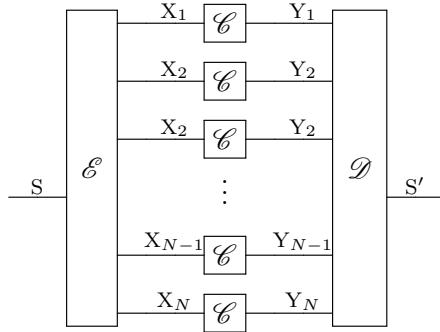
$$I(X : Y) = (1 - f)H(X).$$

The optimisation is now straightforward: The optimal choice is the uniform distribution for X , which yields

$$C(Q) = (1 - f).$$

Notice that in the case of a binary erasure channel, the fact that we know when the errors occur makes the channel capacity, namely the rate as we will show in Shannon's noisy channel coding theorem, exactly equal to the complement of the error probability.

When we analysed the noisy typewriter channel, we realised that it allows for a set of non-confusable characters. The basic idea of channel coding is that all extended channels—namely channels consisting in multiple uses of the same discrete memoryless channel—get increasingly close to a sort of noisy typewriter channel, as the number of uses increases. This makes it possible to find sets of strings that are approximately non-confusable, reaching channel capacity with vanishingly small error probability. The scheme representing extended channels is the following



In the remainder we introduce the tools used for the analysis of extended channels.

8.4 Block codes

Definition 8.10 (Block code). A (N, K) block code for a channel \mathcal{C} with input system X and output Y is a map $\mathcal{E} : S \rightarrow \text{Rng}(X)^N$, where S is the set $S = \{1, 2, \dots, 2^K\}$ of cardinality $|S| = 2^K$. The images of the elements $s \in S$ under \mathcal{E} are 2^K codewords denoted by

$$\{x_{\mathbf{i}(1)}, x_{\mathbf{i}(2)}, \dots, x_{\mathbf{i}(2^K)}\}, \quad x_{\mathbf{i}(s)} \in \text{Rng}(X)^N, \quad (8.8)$$

each codeword of length N .

The rate of a code is defined as

Definition 8.11 (Rate of a code). Let \mathcal{E} be a (N, K) block code for a channel \mathcal{C} . The rate of the code is $R := K/N$ bits per channel use.

Using a code \mathcal{E} one encodes the “signal” $s \in S$ into the codeword $\mathcal{E}(s) = x_{\mathbf{i}(s)}$. Notice that the number of codewords 2^K is an integer, whereas K (the number of bits necessary to specify s) generally is not. Notice also that “per use” refers to N , which is not a base 2 logarithm. The reader should keep in mind such asymmetry in the definition of *rate* as number of bits per channel use.

Definition 8.12 (Decoder). A decoder for a (N, K) block code \mathcal{E} is a map $\mathcal{D} : \text{Rng}(Y)^N \rightarrow S'$ from the set of length- N strings of channel outputs Y^N , to the set $S' = \{0, 1, 2, \dots, 2^K\}$ of codeword labels $\hat{s} \in S'$. The extra label $\hat{s} = 0$ is useful for indicating a “failure”.

In Fig. 8.6 the optimal decoder for a binary erasure channel is illustrated: in this particular case the decoding is error-free, and there is only a small probability of failure. The rate is $R = K/N = 1/2$.

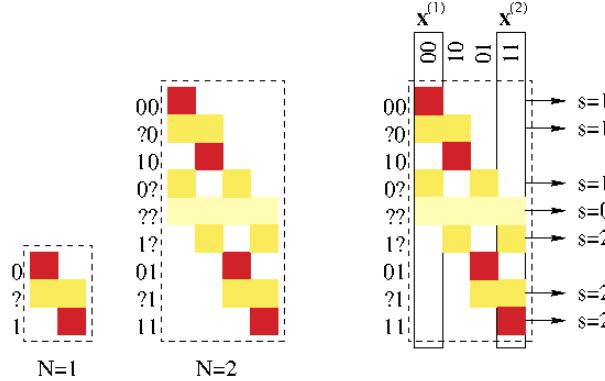


Figure 8.6 (a) Extended channel of the binary erasure channel with $f = 0.15$: each column corresponds to a possible input and each row to a possible output. (b) A block code consisting of the two codewords 00 and 11. (c) The optimal decoder for this code: notice that the decoding is error-less, and there is only a small probability of failure.

8.5 Error probabilities

In the following we will need the precise definitions of different kinds of error probability. We will treat the input signal s and the output s' as the value of two random variables S and S' , with ranges $\text{Rng}(S) = S$ and $\text{Rng}(S') = S' = S \cup \{0\}$.

Definition 8.13 (Probability of error for a codeword). For given channel \mathcal{C} , code \mathcal{E} , and decoder \mathcal{D} , the *probability of error for the codeword* $s \in S$ is given by

$$\begin{aligned} p_s(\mathcal{E}) &= \mathbb{P}_{S'|S=s}(S' \neq s | S = s, \mathcal{E}, \mathcal{D}) \\ &= \sum_{s' \neq s} \sum_{y_j \in \text{Rng}(Y)^N} \underbrace{p(s'|y_j)}_{\mathcal{D}} \underbrace{p(y_j | \tilde{x}_{i^{(s)}})}_{Q}. \end{aligned} \quad (8.9)$$

Definition 8.14 (Probability of block error). For a given channel \mathcal{C} , and for given code \mathcal{E} , decoder \mathcal{D} , and prior probability $\mathbb{P}_S(s)$ of the signal, the *probability of block error* is given by

$$p_B(\mathcal{E}) = \sum_{s \in S} \mathbb{P}_S(s) \mathbb{P}_{S'|S=s}(S' \neq s | S = s, \mathcal{E}, \mathcal{D}). \quad (8.10)$$

Definition 8.15 (Maximum probability of block error). For a given channel \mathcal{C} , and for given code \mathcal{E} , decoder \mathcal{D} , and prior probability $\mathbb{P}_S(s)$ of the signal, the *maximum probability of block error* is

$$p_{\text{BM}}(\mathcal{E}) = \max_s \mathbb{P}_{S'|S=s}(S' \neq s | S = s, \mathcal{E}, \mathcal{D}), \quad (8.11)$$

8.6 Optimal decoder

Given a channel \mathcal{C} , a prior distribution $\mathbb{P}_S(s)$, and an encoding \mathcal{E} , the optimal decoder is uniquely identified.

Definition 8.16 (Optimal decoder). The optimal decoder for a channel code \mathcal{E} is the one that minimises the probability of block error. It decodes an output y_j as the symbol s_0 that has maximum posterior probability $p(s_0|y_j)$, where the prior $p(s_0)$ is $p(s_0) := \mathbb{P}_S(S = s_0)$

$$p(s_0|y_j) = \frac{p(y_j|s_0)p(s_0)}{\sum_t p(y_j|t)p(t)}, \quad s_{\text{opt}}(y_j) := \arg \max_s p(s|y_j). \quad (8.12)$$

When the maximum is not unique, then $s_{\text{opt}}(y_j) := 0$. In other words, \mathcal{D} has transfer matrix $p(s'|y_j) = \mathbb{P}_{S'|Y^N=y_j}[S' = s' | Y^N = y_j] := \delta_{s', s_{\text{opt}}(y_j)}$

For uniform prior probability distribution $\mathbb{P}_S(s)$ the optimal decoder is also the *maximum likelihood decoder*, namely it maps an output y_j to the input s that has maximum likelihood $P(y_j|s)$.

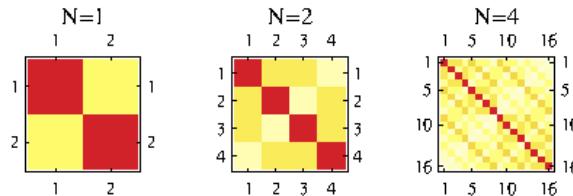


Figure 8.7 Extended channel of the binary symmetric channel with $f = 0.15$. Each column corresponds to a possible input and each row to a possible output.

Example 8.17 (Block decoders for the binary symmetric channel). In Fig. 8.7 the Q matrix (columns correspond to possible inputs and rows to possible outputs) of the extended channel of the binary symmetric channel with $f = 0.15$ is reported for $N = 1, 2$ and 4 . Notice how the extended channel looks more and more similar to a noisy typewriter channel. The same phenomenon occurs essentially for any noisy channel for increasing N . For $N = 2$ we can devise an optimal decoding for $s \in \{1, 2\}$ with codewords $x^{(1)} = 00$ and $x^{(2)} = 11$ corresponding to $\hat{s}(00) = 1, \hat{s}(11) = 1$ and $\hat{s}(01) = \hat{s}(10) = 0$. We can see that in this way we achieve a rate $R = .5$ with an average block error probability $2f - f^2 = .28$.

For $N = 4$ instead we can devise an optimal decoding for $s \in \{1, 2\}$ with codewords $x^{(1)} = 0000, x^{(2)} = 1111$, and decoding 1 if the majority of bits is 0, and 2 if the

majority is 1 (and all other codewords corresponding to the failure value $\hat{s} = 0$). In this way we can achieve a rate $R = 1/4 = .25$ bits with an average error probability $1 - (1 + 3f)(1 - f)^3 = .11$. We can see that by increasing N we have improved the error probability, however we had to decrease the rate. As we will see, according to the second Shannon theorem, for sufficiently large N it is possible to achieve a rate equal to the channel capacity with a vanishingly small probability of error. In the present case the channel capacity is $C = 1 - H_2(f) = .61$ (see Example 8.2).

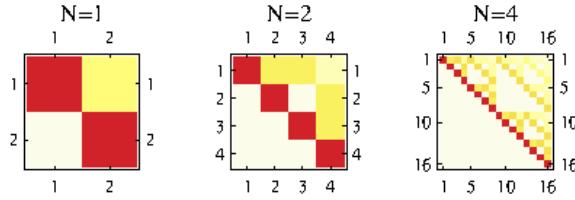


Figure 8.8 Extended channel of the binary Z channel with $f = 0.15$. Each column corresponds to a possible input and each row to a possible output.

8.7 Joint typicality

As in the case of Shannon's first theorem, also in the case of noisy channel coding the proof exploits the property of typical sets. Differently from the case of source coding, however, here two random variables are involved: The input and the output. We then need an extension of the notion of typical sets, given by the following definition.

Definition 8.18 (Jointly typical sequences). A pair of sequences (x_i, y_j) that are values of N i.i.d. pairs of random variables $(X, Y)^N$ are called jointly typical to tolerance ε if $x_i \in T_{N,\varepsilon}(X)$, $y_j \in T_{N,\varepsilon}(Y)$, and $(x_i, y_j) \in T_{N,\varepsilon}(X, Y)$, the third condition meaning that

$$\left| \frac{1}{N} \log_2 \frac{1}{P_{(X,Y)^N}(x_i, y_j)} - H(X, Y) \right| = \left| \frac{1}{N} \sum_{l=1}^N \log_2 \frac{1}{P_{(X,Y)}(x_{i_l}, y_{j_l})} - H(X, Y) \right| \leq \varepsilon. \quad (8.13)$$

Definition 8.19 (Jointly typical set). We denote the set of all jointly typical sequences for X and Y by $J_{N,\varepsilon}(X, Y)$. This set is the *Jointly typical set*.

Having introduced the notion of joint typicality, we can now prove the joint typicality theorem.

Theorem 8.20 (Joint typicality). Consider the values x_i, y_j of $(X, Y)^N = (X^N, Y^N)$, where the couples are i.i.d., namely they are distributed according to the probability distribution $P_{(X,Y)^N}(x_i, y_j) = \prod_{k=1}^N P_{(X,Y)}(x_{i_k}, y_{j_k})$. One has

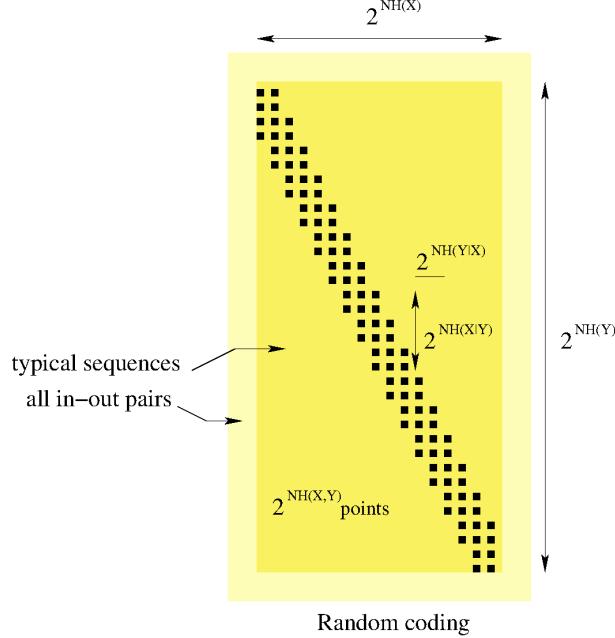


Figure 8.9 The jointly typical set. The horizontal direction represents $\text{Rng}(X)^N$, the set of all input strings of length N , and the vertical one represents $\text{Rng}(Y)^N$, the set of all output strings of length N . The outer box contains all input-output pairs. The inside box contains all pairs made of typical sequences. Dots represent jointly typical pairs of sequences (x_i, y_j) . The total number of jointly-typical sequences is approximately $2^{NH(X,Y)}$.

1. For every $\delta > 0$ there exists N_0 such that for every $N \geq N_0$ the probability $\mathbb{P}_{(X,Y)^N}[(x_i, y_j) \in J_{N,\varepsilon}(X, Y)]$ of the jointly typical set is bounded as

$$\mathbb{P}_{(X,Y)^N}[(x_i, y_j) \in J_{N,\varepsilon}(X, Y)] \geq 1 - \delta. \quad (8.14)$$

2. The number of jointly typical sequences $|J_{N,\varepsilon}(X, Y)|$ is close to $2^{NH(X,Y)}$, in the sense that

$$(1 - \delta)2^{N(H(X,Y) - \varepsilon)} \leq |J_{N,\varepsilon}(X, Y)| \leq 2^{N(H(X,Y) + \varepsilon)} \quad (8.15)$$

3. Let X' and Y' be independent random variables with $\text{Rng}(X') = \text{Rng}(X)$, $\text{Rng}(Y') = \text{Rng}(Y)$, and distributed as the marginals X and Y of (X, Y) , namely

$$\mathbb{P}_{X',Y'}(x_i, y_j) = \mathbb{P}_X(x_i)\mathbb{P}_Y(y_j). \quad (8.16)$$

Then the probability $\mathbb{P}_{(X',Y')^N}[(x_i, y_j) \in J_{N,\varepsilon}(X, Y)]$ that (X', Y') occurs in the jointly typical set of (X, Y) is bounded as

$$\mathbb{P}_{(X',Y')^N}[(x_i, y_j) \in J_{N,\varepsilon}(X, Y)] \leq 2^{-N(I(X:Y) - 3\varepsilon)}. \quad (8.17)$$

Proof. The proof of items 1 and 2 are exactly the same as the analogous ones in the asymptotic equipartition theorem 4.15. Let

$$\begin{aligned} A &:= \left\{ (x_i, y_j) \mid \left| \frac{1}{N} \log_2 \frac{1}{P_{X^N}(x_i)} - H(X) \right| > \varepsilon \right\}, \\ B &:= \left\{ (x_i, y_j) \mid \left| \frac{1}{N} \log_2 \frac{1}{P_{Y^N}(y_j)} - H(Y) \right| > \varepsilon \right\}, \\ C &:= \left\{ (x_i, y_j) \mid \left| \frac{1}{N} \log_2 \frac{1}{P_{(X,Y)^N}(x_i, y_j)} - H(X, Y) \right| > \varepsilon \right\}. \end{aligned}$$

By definition, we have that

$$J_{N,\varepsilon}(X, Y) = \text{Rng}(X, Y)^N \setminus [A \cup B \cup C],$$

and thus

$$\begin{aligned} P_{(X,Y)^N}[(x_i, y_j) \in J_{N,\varepsilon}(X, Y)] &= 1 - P_{(X,Y)^N}[(x_i, y_j) \in (A \cup B \cup C)] \\ &\geq 1 - \{P_{(X,Y)^N}[(x_i, y_j) \in A] + P_{(X,Y)^N}[(x_i, y_j) \in B] + P_{(X,Y)^N}[(x_i, y_j) \in C]\}. \end{aligned}$$

From the asymptotic equipartition theorem 4.15 we have that, for sufficiently large N_0 , for every $N \geq N_0$

$$\begin{aligned} P_{(X,Y)^N}[(x_i, y_j) \in A] &= 1 - P_{X^N}[x_i \in T_{N,\varepsilon}(X)] \leq \frac{\sigma_X^2}{\varepsilon^2 N_0}, \\ P_{(X,Y)^N}[(x_i, y_j) \in B] &= 1 - P_{Y^N}[y_j \in T_{N,\varepsilon}(Y)] \leq \frac{\sigma_Y^2}{\varepsilon^2 N_0}, \\ P_{(X,Y)^N}[(x_i, y_j) \in C] &\leq \frac{\sigma_{(X,Y)}^2}{\varepsilon^2 N_0} = \frac{\sigma_X^2}{\varepsilon^2 N_0} + \frac{\sigma_Y^2}{\varepsilon^2 N_0}, \end{aligned}$$

where σ_W^2 is the variance of $\log_2 \frac{1}{P_W(W=w)}$. Thus, for $N \geq N_0$, and upon suitable choice of N_0 ,

$$P_{(X,Y)^N}[(x_i, y_j) \in A] + P_{(X,Y)^N}[(x_i, y_j) \in B] + P_{(X,Y)^N}[(x_i, y_j) \in C] \leq \delta.$$

Finally, this implies that

$$P_{(X,Y)^N}[(x_i, y_j) \in J_{N,\varepsilon}(X, Y)] \geq 1 - \delta.$$

As to item 2, by definition of jointly typical sequence we have that for $(x_i, y_j) \in J_{N,\varepsilon}(X, Y)$

$$2^{-N[H(X,Y)+\varepsilon]} \leq P_{(X,Y)^N}(x_i, y_j) \leq 2^{-N[H(X,Y)-\varepsilon]},$$

and then

$$\begin{aligned} 1 &= \sum_{(x_i, y_j) \in \text{Rng}(X, Y)^N} P_{(X,Y)^N}(x_i, y_j) \\ &\geq \sum_{(x_i, y_j) \in J_{N,\varepsilon}(X, Y)} P_{(X,Y)^N}(x_i, y_j) \geq |J_{N,\varepsilon}(X, Y)| 2^{-N[H(X,Y)+\varepsilon]}, \end{aligned}$$

which implies the upper bound in item 2. Analogously, since

$$\begin{aligned} 1 - \delta &\leq \mathbb{P}_{(X,Y)^N}[(x_i, y_j) \in J_{N,\varepsilon}(X, Y)] \\ &= \sum_{(x_i, y_j) \in J_{N,\varepsilon}(X, Y)} \mathbb{P}_{(X,Y)^N}(x_i, y_j) \leq |J_{N,\varepsilon}(X, Y)| 2^{-N[H(X, Y) - \varepsilon]}, \end{aligned}$$

we also have the lower bound. Item 3 is proved through the following steps

$$\begin{aligned} \mathbb{P}_{(X',Y')}[(x_i, y_j) \in J_{N,\varepsilon}(X, Y)] &= \sum_{(x_i, y_j) \in J_{N,\varepsilon}(X, Y)} \mathbb{P}_{X'^N}(x_i) \mathbb{P}_{Y'^N}(y_j) \\ &\leq |J_{N,\varepsilon}(X, Y)| 2^{-N[H(X) - \varepsilon]} 2^{-N[H(Y) - \varepsilon]} \\ &\leq 2^{N[H(X, Y) + \varepsilon] - N[H(X) + H(Y) - 2\varepsilon]} = 2^{-N[I(X;Y) - 3\varepsilon]}. \quad \square \end{aligned}$$

Chapter 9

Lecture 9: Shannon noisy channel coding theorem

9.1 Random coding and typical-set decoding

In this section we describe a coding scheme, named *random coding*, along with a decoder that is sub-optimal, and defined as *typical set decoding*. Surprisingly, these very lame procedures for identification of $(\mathcal{E}, \mathcal{D})$ are sufficient to provide a proof of the noisy channel coding theorem.

For a given prior probability distribution for the alphabet \mathbb{P}_X , a random (N, K) code \mathcal{E} is generated as follows

1. Generate 2^K codewords of a (N, K) code at random, sampling from the probability distribution

$$\mathbb{P}(x_{\mathbf{i}}) = \prod_{n=1}^N \mathbb{P}_X(x_{i_n}). \quad (9.1)$$

2. The code is shared by sender and receiver.
3. The signal s is chosen with uniform probability distribution $\mathbb{P}_S(s) = \frac{1}{2^K}$. Upon transmission of $x_{\mathbf{i}^{(s)}}$, the output probability distribution is given by

$$\mathbb{P}_{Y^N|X^N=x_{\mathbf{i}^{(s)}}}(y_{\mathbf{j}}|x_{\mathbf{i}^{(s)}}) = \prod_{n=1}^N p(y_{j_n}|x_{i_n^{(s)}}). \quad (9.2)$$

4. Decode by the *typical-set decoding*, i.e. decode $\mathcal{D}(y_{\mathbf{j}}) = \hat{s}$ if $(x_{\mathbf{i}^{(\hat{s})}}, y_{\mathbf{j}})$ is jointly typical *and* there is no other \hat{s}' such that $(x_{\mathbf{i}^{(\hat{s}')}}, y_{\mathbf{j}})$ is jointly typical—otherwise declare failure ($\hat{s} = 0$).

A random coding is illustrated in Fig. 9.1.

For random coding and typical set decoding, the following result holds.

Lemma 9.1. *The average block error $\langle p_B \rangle$ over all codes, defined as*

$$\langle p_B \rangle := \sum_{\mathcal{E}} \mathbb{P}(\mathcal{E}) p_B(\mathcal{E}) = \sum_{\mathcal{E}} \mathbb{P}(\mathcal{E}) \sum_{s \in S} \mathbb{P}_S(s) P(s' \neq s | s, \mathcal{E}), \quad (9.3)$$

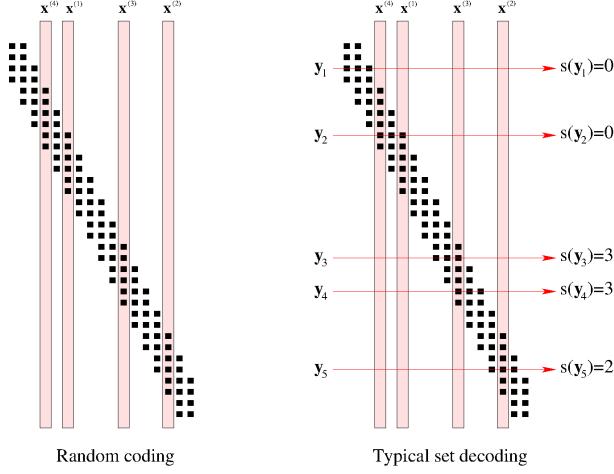


Figure 9.1 (a) A random code. (b) Typical-set decoding. A sequence that is not jointly-typical with any codeword such as y_1 is decoded to $\hat{s} = 0$. A sequence that is jointly-typical with a single codeword $x^{(s)}$ (e.g. y_3 and y_4) is decoded to $\hat{s} = s$. A sequence that is jointly-typical with more than one codeword such as y_2 is decoded to $\hat{s} = 0$.

where $P_S(s) = 1/2^K$, is equal to

$$\langle p_B \rangle = \sum_{\mathcal{E}} \mathbb{P}(\mathcal{E}) P(s' \neq 1 | s = 1, \mathcal{E}). \quad (9.4)$$

Proof. We can write a code \mathcal{E} as a $2^K \times N$ matrix B , the *codebook matrix* whose lines are the codewords $x_{i^{(s)}}$, namely $B_{s,j} = x_{i_j^{(s)}}$. Now, since codewords are randomly extracted with i.i.d. characters distributed according to \mathbb{P}_X , for every code \mathcal{E} , the probability

$$\mathbb{P}(\mathcal{E}) = \prod_{s=1}^{2^K} \prod_{j=1}^N \mathbb{P}_X(x_{i_j^{(s)}}), \quad (9.5)$$

is independent of permutations of the rows. We can then substitute a code \mathcal{E} with the couple (B_0, τ) , where B_0 is a fixed *reference codebook* with the same codewords as \mathcal{E} , and τ is the permutation of S by which the rows of the codebook matrix B_0 are mapped to those of the codebook matrix B for \mathcal{E} . Notice that in general one can have a codebook B_0 with repeated instances. In this case, the non-trivial permutations $\tilde{\tau}$ are not all the possible permutations of 2^K elements. In particular, if the codebook B_0 contains the m words $x_{i^{(1)}}, x_{i^{(2)}}, \dots, x_{i^{(m)}}$ repeated n_1, n_2, \dots, n_m times, respectively, then for every relevant permutation $\tilde{\tau}$ there are $n_1!n_2!\dots n_m!$ permutations τ of the 2^K elements that are equivalent to $\tilde{\tau}$, as can be verified by simple combinatorial arguments. This implies that the number of relevant permutations $\tilde{\tau}$ is $2^K!/(n_1!n_2!\dots n_m!)$. Now, re-expressing \mathcal{E} as a couple $(B_0, \tilde{\tau})$, according to equation 9.5, one has that the probability of every code with the same reference codebook B_0 is the same. Thus, by the counting of relevant

permutations $\tilde{\tau}$ for a given B_0 , we have

$$\mathbb{P}(\mathcal{E}) = \mathbb{P}(B_0, \tilde{\tau}) = \mathbb{P}(B_0)\mathbb{P}(\tilde{\tau}|B_0) = \frac{n_1!n_2!\dots n_m!}{2^K!}\mathbb{P}(B_0). \quad (9.6)$$

Moreover, for a general code \mathcal{E} the decoding of the symbol s depends on the set of strings y_j that are jointly typical with $x_{j(s)}$. Notice that the set of strings y_j that are jointly typical with x_i is independent of the particular label s associated to the string x_i . Let us then consider two codes \mathcal{E} and \mathcal{E}' with the same set of codewords, and differing only by the (non-trivial) permutation $\tilde{\sigma}$ of labels, namely $x_i = \mathcal{E}'(s) = \mathcal{E}[\tilde{\sigma}(s)]$. One has

$$\mathbb{P}_{S'|S=s}(S' \neq s|S = s, \mathcal{E}') = \mathbb{P}_{S'|S=\tilde{\sigma}(s)}(S' \neq \tilde{\sigma}(s)|S = \tilde{\sigma}(s), \mathcal{E}). \quad (9.7)$$

Now, considering the code \mathcal{E}_0 corresponding to the reference codebook B_0 , $\mathcal{E}_0 \leftrightarrow (B_0, \text{id})$, if we associate $\mathcal{E} \leftrightarrow (B_0, \tilde{\tau})$ we have $\mathcal{E}(s) = \mathcal{E}_0[\tilde{\tau}(s)]$, and thus $\mathcal{E}'(t) = \mathcal{E}[\tilde{\sigma}(t)] = \mathcal{E}_0[\tilde{\tau}\tilde{\sigma}(t)]$. This implies that $\mathcal{E}' \leftrightarrow (B_0, \tilde{\tau}\tilde{\sigma})$, and one can rewrite the identity in equation (9.7) as

$$\mathbb{P}_{S'|S=s}(S' \neq s|S = s, B_0, \tilde{\tau}\tilde{\sigma}) = \mathbb{P}_{S'|S=\tilde{\sigma}(s)}(S' \neq \tilde{\sigma}(s)|S = \tilde{\sigma}(s), B_0, \tilde{\tau}).$$

In the following, in order to make the notation lighter, we omit the parameters of the probability distribution. Since the number of permutations sending 1 to s is $(2^K - 1)!$, one clearly has

$$\begin{aligned} \frac{1}{2^K} \sum_{s \in S} \mathbb{P}(S' \neq s|S = s, B_0, \tilde{\tau}) &= \frac{1}{2^K} \sum_{s \in S} \frac{1}{(2^K - 1)!} \sum_{\sigma(1)=s} \mathbb{P}(S' \neq \sigma(1)|S = \sigma(1), B_0, \tilde{\tau}) \\ &= \frac{1}{2^K!} \sum_{\sigma} \mathbb{P}(S' \neq \sigma(1)|S = \sigma(1), B_0, \tilde{\tau}) \\ &= \frac{n_1!n_2!\dots n_m!}{2^K!} \sum_{\tilde{\sigma}} \mathbb{P}(S' \neq \tilde{\sigma}(1)|S = \tilde{\sigma}(1), B_0, \tilde{\tau}) \\ &= \frac{n_1!n_2!\dots n_m!}{2^K!} \sum_{\tilde{\sigma}} \mathbb{P}(S' \neq 1|S = 1, B_0, \tilde{\tau}\tilde{\sigma}), \end{aligned}$$

where in the second last equality we used the fact that for every non-trivial permutation $\tilde{\sigma}$ there are $n_1!n_2!\dots n_m!$ equivalent permutations σ . We can then rewrite equation 9.3 as follows

$$\begin{aligned} \langle p_B \rangle &= \sum_{B_0, \tilde{\tau}} \mathbb{P}(B_0, \tilde{\tau}) \frac{n_1!n_2!\dots n_m!}{2^K!} \sum_{\tilde{\sigma}} \mathbb{P}(S' \neq 1|S = 1, B_0, \tilde{\tau}\tilde{\sigma}) \\ &= \sum_{B_0, \tilde{\tau}} \mathbb{P}(B_0, \tilde{\tau}) \frac{n_1!n_2!\dots n_m!}{2^K!} \sum_{\tilde{\sigma}'} \mathbb{P}(S' \neq 1|S = 1, B_0, \tilde{\sigma}') \\ &= \sum_{B_0} \mathbb{P}(B_0) \frac{n_1!n_2!\dots n_m!}{2^K!} \sum_{\tilde{\sigma}'} \mathbb{P}(S' \neq 1|S = 1, B_0, \tilde{\sigma}') \\ &= \sum_{B_0, \tilde{\sigma}'} \mathbb{P}(B_0, \tilde{\sigma}') \mathbb{P}(S' \neq 1|S = 1, B_0, \tilde{\sigma}') \\ &= \sum_{\mathcal{E}} \mathbb{P}(\mathcal{E}) \mathbb{P}(S' \neq 1|S = 1, \mathcal{E}), \end{aligned}$$

where in the second identity we used the fact that $\sum_{\tilde{\sigma}} f(\tilde{\sigma}) = \sum_{\tilde{\sigma}'} f(\tilde{\sigma}')$, and in the fourth we used the fact that $\mathbb{P}(B_0, \tilde{\sigma}) = \mathbb{P}(B_0) \frac{n_1! n_2! \dots n_m!}{2^K!}$ as in equation (9.6). \square

9.2 The second theorem of Shannon's

In the proof of the noisy-channel coding theorem we will use the following lemma.

Lemma 9.2 (Expurgation). *For a uniformly distributed random variable X , and a non-negative (measurable) function $f : \text{Rng}(X) \rightarrow \mathbb{R}_+$, the subset $A \subseteq \text{Rng}(X)$ of values $x_i \in \text{Rng}(X)$ with $f(x_i) < 2\mathbb{E}[f(X)]$, namely*

$$A := \{x_i \in \text{Rng}(X) | f(x_i) < 2\mathbb{E}[f(X)]\},$$

has probability $\mathbb{P}_X[x \in A] \geq \frac{1}{2}$ and cardinality $|A| \geq \frac{1}{2}|\text{Rng}(X)|$.

Proof. Using the Markov inequality (3.8)

$$\mathbb{P}_X[f(X) \geq a] \leq \frac{\mathbb{E}[f(X)]}{a},$$

with $a = 2\mathbb{E}[f(X)]$, one has

$$\mathbb{P}_X[x_i \in A] = \mathbb{P}_X[f(X) < 2\mathbb{E}[f(X)]] = 1 - \mathbb{P}_X[f(X) \geq 2\mathbb{E}[f(X)]] \geq \frac{1}{2}.$$

Moreover, since X is uniformly distributed, one has

$$\frac{1}{2} \leq \mathbb{P}_X[x_i \in A] = \sum_{x_i \in A} \mathbb{P}_X[x_i] = \sum_{x_i \in A} \frac{1}{|\text{Rng}(X)|} = \frac{|A|}{|\text{Rng}(X)|},$$

from which we conclude $|A| \geq \frac{1}{2}|\text{Rng}(X)|$. \square

We will now prove the celebrated noisy-channel coding theorem of Shannon. The proof will use all the following previously analyzed ingredients:

1. joint typicality: item (3) of Theorem 8.20;
2. random coding and typical-set decoding: lemma 9.1;
3. the Fano inequality in Theorem 8.2;
4. the data-processing Theorem 7.14;
5. the expurgation technique (see Lemma 9.2);
6. the definition of channel capacity for a single use;
7. the probability of errors defined in 8.13 and 8.14.

Theorem 9.3 (Noisy-channel coding theorem). *Let \mathcal{C} be a discrete memoryless channel with transfer matrix Q . Then we have the following.*

1. *For every $R_0 < C(Q)$, for every $\varepsilon > 0$ there exists N_0 such that for $N \geq N_0$ there exists a block code \mathcal{E} of length N and rate R with $R_0 < R < C(Q)$, and a decoding \mathcal{D} such that the maximal probability of block error is $p_{\text{BM}}(\mathcal{E}) < \varepsilon$.*
2. *Conversely, for any block code \mathcal{E} with rate $R > C(Q)$, the probability $p_{\text{BM}}(\mathcal{E})$ is bounded from below by $k > 0$, independently of N .*

Sketch of the proof

Before giving the proof of the theorem, we want to focus attention on the sketch of the logical sequence that constitutes the proof of item 1. The point is to prove existence of a code that has maximum probability of block error bounded by ε . Here are the steps for this proof.

1. Since it is very difficult to evaluate the probability of block error $p_B(\mathcal{E}) := \sum_{s \in \text{Rng}(S)} \mathbb{P}(S' \neq s | S = s, \mathcal{E}) \mathbb{P}_S(s)$ for a particular code, while it is easy to evaluate its average over all randomly chosen codes $\langle p_B \rangle = \sum_{\mathcal{E}} \mathbb{P}(\mathcal{E}) p_B(\mathcal{E})$, Shannon first bounds the above average.
2. Clearly, there must exist a code \mathcal{E} for which $p_B(\mathcal{E}) \leq \langle p_B \rangle$.
3. However, we need to bound the maximal block error probability, given by the expression $p_{BM}(\mathcal{E}) := \max_s \mathbb{P}(S' \neq s | S = s, \mathcal{E})$, and the fact that we have found a code \mathcal{E} for which $p_B(\mathcal{E})$ is bounded does not mean that $p_{BM}(\mathcal{E})$ is bounded. This code can be modified by throwing away the worst 50% of its codewords using the Expurgation of Lemma 9.2, yielding bounded maximal error without affecting the rate by any significant amount.

The proof

Now let us see the real proof of Shannon's noisy channel coding theorem.

Proof. We use random coding and typical set decoding. We remind that this implies extracting the code at random, with i.i.d. characters distributed according to the distribution $\mathbb{P}_X(x_{i_k^{(s)}})$, independent of s, k .

Let us then evaluate the probability of block error averaged over all codes, namely

$$\langle p_B \rangle = \sum_{\mathcal{E}} \mathbb{P}(\mathcal{E}) p_B(\mathcal{E})$$

The typical set decoding is subject to two different kinds of errors: *type 1*: the output y_j is not jointly typical with the input $x_{i^{(s)}}$. This case includes both the cases when y_j is jointly typical with no codeword, and when $(x_{i^{(s')}}, y_j)$ is jointly typical with $s' \neq s$; *type 2*: $(x_{i^{(s)}}, y_j)$ may or may not be jointly typical, but there is another codeword s' in \mathcal{E} that is actually jointly typical with y_j . Notice that the two types of errors are not independent, as it may happen that both $(y_j, x_{i^{(s)}}) \notin J_{N,\varepsilon}(X, Y)$ and $(y_j, x_{i^{(s')}}) \in J_{N,\varepsilon}(X, Y)$ for $s' \neq s$. According to lemma 9.1, we can use the expression of eq. 9.4 for the average block error probability, and restrict to the case $s = 1$ without loss of generality. The situation is thus the following: one generates the code by random coding, and then the output y_j is

generated from input $x_{\mathbf{i}^{(1)}}$. The joint probability for this sequence of events is

$$\begin{aligned} p(\mathcal{E}, y_{\mathbf{j}} | S = 1) &= \mathbb{P}_{X^{2^K N}, Y^N}(x_{\mathbf{i}^{(1)}}, x_{\mathbf{i}^{(2)}}, \dots, x_{\mathbf{i}^{(2^K)}}, y_{\mathbf{j}}) \\ &= P_{X^N}(x_{\mathbf{i}^{(1)}}) \mathbb{P}_{Y^N | X^N = x_{\mathbf{i}^{(1)}}}(y_{\mathbf{j}} | x_{\mathbf{i}^{(1)}}) \prod_{j=2}^{2^K} \mathbb{P}_{X^N}(x_{\mathbf{i}^{(j)}}), \\ \mathbb{P}_{Y^N | X^N = x^{(1)}}(y_{\mathbf{j}} | x^{(1)}) &= \prod_{l=1}^N p(y_{j_l} | x_{i_l^{(1)}}). \end{aligned} \quad (9.8)$$

We then define the following sets related to the two error types

$$\begin{aligned} A &:= \{(x_{\mathbf{i}^{(1)}}, \dots, x_{\mathbf{i}^{(2^K)}}, y_{\mathbf{j}}) \mid (x_{\mathbf{i}^{(1)}}, y_{\mathbf{j}}) \notin J_{N,\beta}(X, Y)\}, \\ B &:= \{(x_{\mathbf{i}^{(1)}}, \dots, x_{\mathbf{i}^{(2^K)}}, y_{\mathbf{j}}) \mid (x_{\mathbf{i}^{(s)}}, y_{\mathbf{j}}) \in J_{N,\beta}(X, Y), s \neq 1\}. \end{aligned}$$

Since the two types of error are not independent, it may happen that $A \cap B \neq \emptyset$. Anyway, we can bound the average block error probability as

$$\begin{aligned} \langle p_B \rangle &= \mathbb{P}_{X^{2^K N}, Y^N}[A \cup B] \\ &\leq \mathbb{P}_{X^{2^K N}, Y^N}[A] + \mathbb{P}_{X^{2^K N}, Y^N}[B] \\ &= \langle p_B^1 \rangle + \langle p_B^2 \rangle, \end{aligned}$$

where we used the *union bound* $\mathbb{P}(C \cup D) \leq \mathbb{P}(C) + \mathbb{P}(D)$, holding for general probability distributions and general sets C, D , and $\langle p_B^n \rangle$ denotes the contribution from type n error.

Let us first evaluate $\langle p_B^1 \rangle$. According to item 1 of theorem 8.20, for a given code \mathcal{E} the probability that the input $x_{\mathbf{i}^{(1)}}$ is jointly typical with the output $y_{\mathbf{j}}$ with $x_{\mathbf{i}^{(1)}}$ and $y_{\mathbf{j}}$ distributed according to the distribution $\mathbb{P}_{X^N, Y^N}(x_{\mathbf{i}^{(1)}}, y_{\mathbf{j}}) = \mathbb{P}_X(x_{\mathbf{i}^{(1)}}) \mathbb{P}_{Y^N | X^N = x_{\mathbf{i}^{(1)}}}(y_{\mathbf{j}} | x_{\mathbf{i}^{(1)}})$ —which is obtained as the marginal of the joint distribution in equation (9.8)—converges asymptotically to 1 for $N \rightarrow \infty$. More precisely, for any desired $\delta > 0$ there exists a block length N_δ such that, for $N \geq N_\delta$, it holds that $\mathbb{P}_{X^N, Y^N}[(x_{\mathbf{i}^{(1)}}, y_{\mathbf{j}}) \notin J_{N,\beta}(X, Y)] < \delta$. Thus one has $\langle p_B^1 \rangle = \mathbb{P}_{X^N, Y^N}[A] = \mathbb{P}_{X^N, Y^N}[(x_{\mathbf{i}^{(1)}}, y_{\mathbf{j}}) \notin J_{N,\beta}(X, Y)] < \delta$ for $N \geq N_\delta$.

Let us then evaluate the second term $\langle p_B^2 \rangle$. One has

$$\begin{aligned} \langle p_B^2 \rangle &= \mathbb{P}_{X^{2^K N}, Y^N}[\bigcup_{s \neq 1} \{(x_{\mathbf{i}^{(1)}}, y_{\mathbf{j}}) \mid (x_{\mathbf{i}^{(s)}}, y_{\mathbf{j}}) \in J_{N,\beta}(X, Y)\}] \\ &\leq \sum_{s \neq 1} \mathbb{P}_{X^{2^K N}, Y^N}[\{(x_{\mathbf{i}^{(1)}}, y_{\mathbf{j}}) \mid (x_{\mathbf{i}^{(s)}}, y_{\mathbf{j}}) \in J_{N,\beta}(X, Y)\}]. \end{aligned}$$

Since the average is over all random codes, we can rewrite the last expression as

$$\begin{aligned}
\langle p_B^2 \rangle &\leq \sum_{s \neq 1} \sum_{x_{\mathbf{i}(1)} \in \text{Rng}(X^N)} \overbrace{\sum_{x_{\mathbf{i}(2)} \in \text{Rng}(X^N)} \cdots \sum_{x_{\mathbf{i}(2^K)} \in \text{Rng}(X^N)}^{\text{2}^K - 2 \text{ sums}} \mathbb{P}_{X^N}(x_{\mathbf{i}(1)}) \mathbb{P}_{X^N}(x_{\mathbf{i}(2)}) \dots \mathbb{P}_{X^N}(x_{\mathbf{i}(2^K)})} \\
&\times \left\{ \sum_{(x_{\mathbf{i}(s)}, y_{\mathbf{j}}) \in J_{N,\beta}(X, Y)} \mathbb{P}_{X^N}(x_{\mathbf{i}(s)}) \mathbb{P}_{Y^N|X^N=x_{\mathbf{i}(1)}}(y_{\mathbf{j}} | x_{\mathbf{i}(1)}) \right\} \\
&= \sum_{s \neq 1} \sum_{x_{\mathbf{i}(1)} \in \text{Rng}(X^N)} \sum_{(x_{\mathbf{i}(s)}, y_{\mathbf{j}}) \in J_{N,\beta}(X^N, Y^N)} \mathbb{P}_{X^N}(x_{\mathbf{i}(s)}) \mathbb{P}_{Y^N|X^N=x_{\mathbf{i}(1)}}(y_{\mathbf{j}} | x_{\mathbf{i}(1)}) \mathbb{P}_{X^N}(x_{\mathbf{i}(1)}) \\
&= \sum_{s \neq 1} \sum_{x_{\mathbf{i}(1)} \in \text{Rng}(X^N)} \sum_{(x_{\mathbf{i}(s)}, y_{\mathbf{j}}) \in J_{N,\beta}(X^N, Y^N)} \mathbb{P}_{X^N}(x_{\mathbf{i}(s)}) \mathbb{P}_{Y^N, X^N}(y_{\mathbf{j}}, x_{\mathbf{i}(1)}) \\
&= \sum_{s \neq 1} \sum_{(x_{\mathbf{i}(s)}, y_{\mathbf{j}}) \in J_{N,\beta}(X^N, Y^N)} \mathbb{P}_{X^N}(x_{\mathbf{i}(s)}) \mathbb{P}_{Y^N}(y_{\mathbf{j}}).
\end{aligned}$$

According to item 3 of Theorem 8.20 we can bound every term of the sum over s as

$$\sum_{(x_{\mathbf{i}(s)}, y_{\mathbf{j}}) \in J_{N,\beta}(X^N, Y^N)} \mathbb{P}_{X^N}(x_{\mathbf{i}(s)}) \mathbb{P}_{Y^N}(y_{\mathbf{j}}) \leq 2^{-N(I(X:Y)-3\beta)},$$

and we thus have

$$\langle p_B^2 \rangle \leq \sum_{s=2}^{2^K} 2^{-N(I(X:Y)-3\beta)} = (2^K - 1) 2^{-N(I(X:Y)-3\beta)} < 2^{-N(I(X:Y)-R-3\beta)},$$

since there are $2^K - 1 = 2^{NR} - 1$ terms in the sum. Therefore, the total average block error probability is bounded as

$$\langle p_B \rangle < \delta + 2^{-N(I(X:Y)-R-3\beta)}.$$

Now, if the following condition holds

$$R < I(X : Y) - 3\beta, \tag{9.9}$$

for sufficiently large N we can make

$$\langle p_B \rangle < 2\delta, \tag{9.10}$$

Now, if we choose \mathbb{P}_X to be that particular probability distribution maximising the mutual information, condition (9.9) becomes

$$R < C(Q) - 3\beta. \tag{9.11}$$

Since the error probability averaged over all codes is smaller than 2δ , there exists a code \mathcal{E}_0 with mean probability of block error $p_B(\mathcal{E}_0) < 2\delta$. However, such code may have

maximum probability of error $p_{BM}(\mathcal{E}) > 2\delta$. In order to show that a code exists having $p_{BM}(\mathcal{E})$ arbitrarily small, we use the *expurgation* technique to build up a new \mathcal{E}' from \mathcal{E}_0 , by removing the half of the codewords that are most likely to produce errors, i.e. those with the highest probability of error. Since the codewords s are uniformly distributed by the hypothesis of random coding, if we apply Lemma 9.2 to the random variable $p_s(\mathcal{E}_0)$, the remaining codewords must have probability of error $p_s(\mathcal{E}') = p_s(\mathcal{E}_0) \leq 2p_B(\mathcal{E}_0) < 4\delta$ for every s , and thus $p_{BM}(\mathcal{E}') < 4\delta$. The expurgated code \mathcal{E}' has $2^{NR'} = 2^{NR-1}$ codewords, namely the rate has been reduced from R to $R' = R - \frac{1}{N}$, which is a negligible reduction for large N . In conclusion, for $R < C(Q)$ we have proved the existence of a code with rate $R' = R - \frac{1}{N} < C(Q) - 3\beta$ and with maximum probability of error $p_{BM} < 4\delta$. The statement of the theorem is finally obtained by setting $R := (R_0 + C(Q))/2 - 1/N \geq R_0$, $\delta = \varepsilon/4$, $\beta < (C(Q) - R)/3$, and N sufficiently large to satisfy the remaining conditions. We have proved in this way the first part of the theorem.

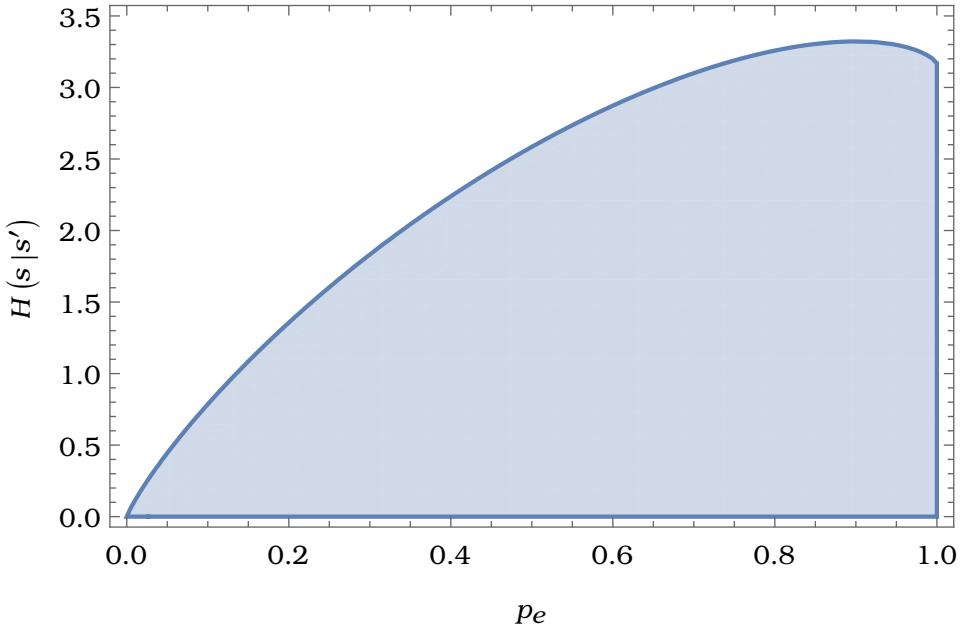


Figure 9.2 Plot of the Fano bound in Eq. (9.13) for $S = 10$.

To prove the second part of the theorem, consider any coding-decoding, and denote by p_e the corresponding error probability. The Data Processing Theorem 7.14 for the Markov chain $s \leftrightarrow x_{i(s)} \leftrightarrow y_j \leftrightarrow s'$ guarantees that

$$I(s : s') \leq I(s : y_j) \leq I(x_{i(s)} : y_j) \leq NC(Q), \quad (9.12)$$

where the last bound follows by the definition of channel capacity. According to the Fano inequality (8.2)

$$H_2(p_e) + p_e \log_2(|S| - 1) \geq H(s|s'), \quad (9.13)$$

For error probability $p_e = 0$ one has $H(s|s') = 0$, and $I(s : s') = H(s)$. Thus, for uniform probability distribution $p(s) = 1/2^K$ one has $I(s : s') = \log_2 |S| = K = NR$ (see the

definition 8.10 of block code), and inequality (9.12) must be satisfied in this special case, namely

$$NR \leq NC(Q) \Leftrightarrow R \leq C(Q). \quad (9.14)$$

On the other hand, if $R - C(Q) > 0$ one has

$$H(s) - H(s|s') \leq NC(Q) \Rightarrow H(s|s') \geq H(s) - NC(Q) = N(R - C(Q)) > 0.$$

Thus, by Fano's theorem

$$\begin{aligned} N(R - C(Q)) &\leq H(s|s') \\ &\leq H_2(p_e) + p_e \log_2(|\mathcal{S}| - 1) \\ &\leq 1 + p_e \log_2 |\mathcal{S}| \\ &= 1 + p_e K \\ &= 1 + p_e NR, \end{aligned}$$

and finally one has

$$p_e \geq \frac{N(R - C(Q)) - 1}{NR} = 1 - \frac{C(Q)}{R} - \frac{1}{NR}, \quad (9.15)$$

which for large enough N implies $p_e > k$ with $k > 0$ independent of N . \square

Notice the precise form of the statement of the theorem. In the first part—the *direct* statement—it is not claimed that the maximum rate is the channel capacity, but that for every $R < C(Q)$ and for any $\epsilon > 0$ there exists a rate $R' \geq R$ such that the maximum probability of block error is $p_{BM} < \epsilon$. Indeed, in principle, rates larger than the channel capacity are allowed, as long as a non-vanishing error probability is accepted. The second part—the *converse* statement—specifies this situation in terms of the allowed bit error probability. The extension of information theory at non-zero error probability is the objective of the so-called *rate distortion theory*.

9.3 Communication above capacity [extra]

What can we say about error probability in a communication with rate above capacity? In general we have:

$$R \leq \frac{C(Q)}{1 - H_2(p_b)},$$

where p_b is called *bit error probability*. What is p_b precisely? It is a measure of communication error that we now define. In a (N, K) block code one encodes 2^K strings, namely $K = NR$ bits. Let us then represent the input and output symbols $s, s' \in \mathcal{S}$ by the binary representation

$$b : \mathcal{S} \rightarrow \{0, 1\}^K :: s \mapsto \mathbf{b}_s = b_{s1} b_{s2} \dots b_{sK}.$$

Whenever a communication error occurs, one has $\mathbf{b}_s \neq \mathbf{b}_{s'}$, namely there exists $1 \leq l \leq K$ such that $b_{sl} \neq b_{s'l}$. We can then define p_l as the probability of occurrence of such error. We then define the bit error probability as follows.

Definition 9.4 (Bit error probability). The bit error probability for a (N, K) block code is defined as

$$p_b := \frac{1}{K} \sum_{l=1}^K p_l. \quad (9.16)$$

We now prove the following lemma that will be used to bound the attainable rates above capacity.

Lemma 9.5. *Let the input symbols of \mathbf{b}_s be i.i.d., and let $\mathbf{b}_{s'}$ denote the output string for a (N, K) block code for channel \mathcal{C} . Then the following bound holds:*

$$I(\mathbf{b}_s : \mathbf{b}_{s'}) \geq \sum_{l=1}^K I(b_{sl} : b_{s'l}). \quad (9.17)$$

Proof. Let us recall that

$$\begin{aligned} I(\mathbf{b}_s : \mathbf{b}_{s'}) &= H(\mathbf{b}_s) - H(\mathbf{b}_s | \mathbf{b}_{s'}) = \sum_{\mathbf{b}_s, \mathbf{b}_{s'}} p(\mathbf{b}_s, \mathbf{b}_{s'}) \log_2 \frac{p(\mathbf{b}_s | \mathbf{b}_{s'})}{p(\mathbf{b}_s)} \\ &= \sum_{\mathbf{b}_s, \mathbf{b}_{s'}} p(\mathbf{b}_s, \mathbf{b}_{s'}) \log_2 \frac{p(\mathbf{b}_s | \mathbf{b}_{s'})}{p(b_{s1})p(b_{s2}) \dots p(b_{sK})}. \end{aligned}$$

On the other hand, one has

$$\begin{aligned} \sum_{l=1}^K I(b_{sl} : b_{s'l}) &= \sum_{l=1}^K \sum_{b_{sl}, b_{s'l}} p(b_{sl}, b_{s'l}) \log_2 \frac{p(b_{sl} | b_{s'l})}{p(b_{sl})} \\ &= \sum_{l=1}^K \sum_{\mathbf{b}_s, \mathbf{b}_{s'}} p(\mathbf{b}_s, \mathbf{b}_{s'}) \log_2 \frac{p(b_{sl} | b_{s'l})}{p(b_{sl})} \\ &= \sum_{\mathbf{b}_s, \mathbf{b}_{s'}} \sum_{l=1}^K p(\mathbf{b}_s, \mathbf{b}_{s'}) \log_2 \frac{p(b_{sl} | b_{s'l})}{p(b_{sl})} \\ &= \sum_{\mathbf{b}_s, \mathbf{b}_{s'}} p(\mathbf{b}_s, \mathbf{b}_{s'}) \log_2 \frac{p(b_{s1} | b_{s'1})p(b_{s2} | b_{s'2}) \dots p(b_{sK} | b_{s'K})}{p(b_{s1})p(b_{s2}) \dots p(b_{sK})}. \end{aligned}$$

Now, one can write

$$\begin{aligned}
I(\mathbf{b}_s : \mathbf{b}_{s'}) - \sum_{l=1}^K I(b_{sl} : b_{s'l}) &= \sum_{\mathbf{b}_s, \mathbf{b}_{s'}} p(\mathbf{b}_s, \mathbf{b}_{s'}) \log_2 \frac{p(\mathbf{b}_s | \mathbf{b}_{s'})}{p(b_{s1} | b_{s'1}) p(b_{s2} | b_{s'2}) \dots p(b_{sK} | b_{s'K})} \\
&= - \sum_{\mathbf{b}_s, \mathbf{b}_{s'}} p(\mathbf{b}_s, \mathbf{b}_{s'}) \log_2 \frac{p(b_{s1} | b_{s'1}) p(b_{s2} | b_{s'2}) \dots p(b_{sK} | b_{s'K})}{p(\mathbf{b}_s | \mathbf{b}_{s'})} \\
&= - \mathbb{E}_{\mathbf{b}_s, \mathbf{b}_{s'}} \left(\log_2 \frac{p(b_{s1} | b_{s'1}) p(b_{s2} | b_{s'2}) \dots p(b_{sK} | b_{s'K})}{p(\mathbf{b}_s | \mathbf{b}_{s'})} \right) \\
&\geq - \log_2 \mathbb{E}_{\mathbf{b}_s, \mathbf{b}_{s'}} \left(\frac{p(b_{s1} | b_{s'1}) p(b_{s2} | b_{s'2}) \dots p(b_{sK} | b_{s'K})}{p(\mathbf{b}_s | \mathbf{b}_{s'})} p(\mathbf{b}_{s'}) \right) \\
&= - \log_2 \left[\sum_{\mathbf{b}_s, \mathbf{b}_{s'}} p(b_{s1} | b_{s'1}) p(b_{s2} | b_{s'2}) \dots p(b_{sK} | b_{s'K}) p(\mathbf{b}_{s'}) \right] = 0. \quad \square
\end{aligned}$$

We are now in position to prove the main result in this section, that provides a lower bound for the rate above the channel capacity as a function of the bit error probability.

Theorem 9.6. *For a discrete memory channel \mathcal{C} with capacity $C(Q)$, every block code with bit error probability p_b must have a rate R that satisfies the bound*

$$R \leq \frac{C(Q)}{1 - H_2(p_b)}. \quad (9.18)$$

Proof. In the first place, by the second Shannon theorem one has the following inequalities

$$I(\mathbf{b}_s : \mathbf{b}_{s'}) = I(s : s') \leq I(x_{\mathbf{i}^{(s)}} : y_{\mathbf{j}}) \leq NC(Q).$$

By lemma 9.5 we then have

$$I(\mathbf{b}_s : \mathbf{b}_{s'}) \geq \sum_{l=1}^K I(b_{sl} : b_{s'l}) = \sum_{l=1}^K \{H(b_{sl}) - H(b_{sl} | b_{s'l})\}.$$

Now, since the symbols s are uniformly distributed, one has $H(b_{sl}) = H_2(1/2) = 1$. Finally, by Fano's inequality one has

$$H(b_{sl} | b_{s'l}) \leq H_2(p_l) + p_l \log_2 1 = H_2(p_l).$$

This provides the bound

$$I(\mathbf{b}_s : \mathbf{b}_{s'}) \geq NR \left[1 - \frac{1}{NR} \sum_{l=1}^K H_2(p_l) \right].$$

Now, by convexity of the function $-H_2(x)$, Jensen inequality gives

$$-\frac{1}{NR} \sum_{l=1}^K H_2(p_l) \geq -H_2(p_b).$$

Recollecting all the inequalities, we obtain

$$NC(Q) \geq NR[1 - H_2(p_b)]. \quad \square$$

The inequality in Eq. (9.18) can be recast inverting the function H_2 (on the range $0 \leq p_b \leq 1/2$), obtaining

$$p_b \geq H_2^{-1} \left(1 - \frac{C(Q)}{R} \right). \quad (9.19)$$

The plot of the above function in the plane (R, p_e) is given in Fig. 9.3. All the points

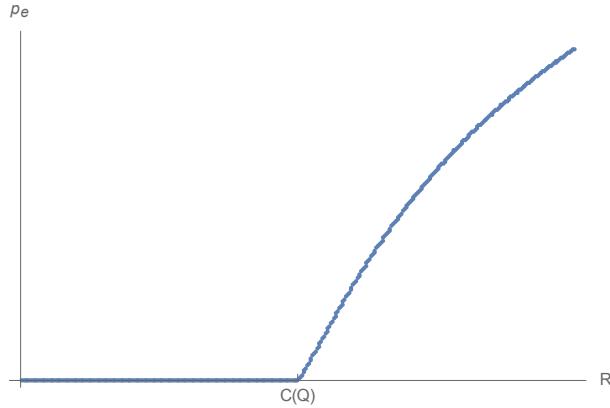


Figure 9.3 The plot of the boundary of the attainable region of block codes for a discrete memoryless channel with capacity $C(Q)$. Points above the blue line represent pairs (R, p_e) of values of rate and maximal block error probability that can be achieved by a suitable (N, K) block code.

above the bound are actually attainable.

Part II

Quantum information theory

Chapter 10

Lecture 10: Elements of linear algebra

We will use in the following the theory of finite dimensional Hilbert spaces. We review here the basic facts about linear algebra on complex Hilbert spaces, avoiding the most advanced topics and sticking to the structures that are required in the development of the quantum theory of information.

10.1 Hilbert spaces

A complex Hilbert space is a vector space \mathcal{H} on \mathbb{C} , equipped with a sesquilinear form $(,) : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{C}$ with the following properties

$$\begin{aligned} (\varphi, \psi) &= (\psi, \varphi)^* \\ (\varphi, a\psi + b\eta) &= a(\varphi, \psi) + b(\varphi, \eta) \\ (\varphi, \varphi) &\geq 0, \quad (\varphi, \varphi) = 0 \Leftrightarrow \varphi = 0, \end{aligned}$$

and closed in the metric $d(\varphi, \psi) := \|\varphi - \psi\|$ induced by the sesquilinear form by $\|\eta\| := (\eta, \eta)^{\frac{1}{2}}$. The dual space \mathcal{H}^* of continuous linear functionals on \mathcal{H} is isomorphic to \mathcal{H} by the Riesz-Fréchet representation theorem

$$D : \mathcal{H}^* \rightarrow \mathcal{H} \quad \tilde{\eta} \mapsto D\tilde{\eta} : \tilde{\eta}(\psi) = (D\tilde{\eta}, \psi).$$

We will adopt the Dirac notation, using *kets* $|\psi\rangle$ to denote vectors $\psi \in \mathcal{H}$ and *bras* $\langle\varphi|$ to denote vectors in \mathcal{H}^* , thus

$$(\varphi, \psi) = \langle\varphi|\psi\rangle. \tag{10.1}$$

It is often useful to represent vectors by expanding them on an *orthonormal basis*, i.e. a collection of orthonormal vectors $\{e_i\}_{i=1}^d$ such that

$$\langle e_i | e_j \rangle = \delta_{i,j}.$$

Indeed, one can prove that such a collection is a basis, namely it is *complete*

$$\psi \in \mathcal{H} \Rightarrow \exists \{\psi_i\}_{i=1}^d \in \mathbb{C}^d \quad |\psi\rangle = \sum_{i=1}^d \psi_i |e_i\rangle,$$

and linearly independent

$$\sum_{i=1}^d \eta_i |e_i\rangle = 0 \Leftrightarrow \eta_i = 0 \forall i.$$

Thus, every complex Hilbert space \mathcal{H} of dimension d is isomorphic to \mathbb{C}^d . Notice that for $d = \infty$ the sums are replaced by series, and convergence is meant in the norm induced by the sesquilinear product. In the following we will restrict attention to the *finite dimensional* case $d < \infty$.

10.1.1 Subspaces

Let now $\mathcal{K} \subseteq \mathcal{H}$ be a subset closed under linear combinations, namely

$$\varphi, \psi \in \mathcal{K} \rightarrow \forall a, b \in \mathbb{C} a|\varphi\rangle + b|\psi\rangle \in \mathcal{K}.$$

We call such \mathcal{K} *subspace* of \mathcal{H} ¹. The dimension if \mathcal{K} is the maximum cardinality of a set of linearly independent vectors in \mathcal{K} . Clearly, 0 belongs to every subspace \mathcal{K} . Special cases are $\mathcal{K} = \{0\}$ and $\mathcal{K} = \mathcal{H}$. For every subspace $\mathcal{K} \subseteq \mathcal{H}$ one can define the set

$$\mathcal{K}^\perp := \{\psi \in \mathcal{H} | \langle \varphi | \psi \rangle = 0 \forall \varphi \in \mathcal{K}\}.$$

One can easily prove that \mathcal{K}^\perp is itself a subspace, and $(\mathcal{K}^\perp)^\perp = \mathcal{K}$. Actually, \mathcal{K}^\perp is a subspace even if \mathcal{K} is not. Thus, in general $(\mathcal{K}^\perp)^\perp$ is the *linear closure* of the subset \mathcal{K} , namely the smallest subspace that contains \mathcal{K} .

Given two subspaces $\mathcal{K}_1, \mathcal{K}_2 \subseteq \mathcal{H}$ such that for every $\psi \in \mathcal{H}$ there exists a unique pair $(\varphi_1, \varphi_2) \in \mathcal{K}_1 \times \mathcal{K}_2$ such that $\psi = \varphi_1 + \varphi_2$, we write $\mathcal{H} = \mathcal{K}_1 \oplus \mathcal{K}_2$. It is easy to prove, in particular, that for every subspace $\mathcal{K} \subseteq \mathcal{H}$ one has

$$\mathcal{H} = \mathcal{K} \oplus \mathcal{K}^\perp.$$

10.2 Linear operators on complex Hilbert spaces

Let $\mathcal{H}_A, \mathcal{H}_B$ be two Hilbert spaces with dimension d_A, d_B , respectively. The set of *linear operators* $\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_B)$ from \mathcal{H}_A to \mathcal{H}_B is the set of linear maps $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$, namely

$$\begin{aligned} T|\psi\rangle &\in \mathcal{H}_B \quad \forall \psi \in \mathcal{H}_A, \\ T(a|\psi\rangle + b|\varphi\rangle) &= aT|\psi\rangle + bT|\varphi\rangle. \end{aligned}$$

Two linear operators $S, T : \mathcal{H}_A \rightarrow \mathcal{H}_B$ are equal if and only if $T|\psi\rangle = S|\psi\rangle$ for all $\psi \in \mathcal{H}_A$. Clearly, given any orthonormal basis $\{e_i\}_{i=1}^{d_A}$ we have

$$T = S \Leftrightarrow T|e_i\rangle = S|e_i\rangle, \quad \forall 1 \leq i \leq d_A.$$

¹In the infinite dimensional case this definition corresponds to a linear manifold. A subspace in this case needs to be topologically closed, in addition to being closed under linear combinations. This means that Cauchy sequences in \mathcal{K} must have their limit within \mathcal{K} . The results for subspaces in the finite dimensional case hold exactly also in the infinite dimensional case, but they might fail if manifolds are considered instead of subspaces.

A relevant class of linear operators in $\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_B)$ is that of operators of the form

$$T_{\eta,\varphi}|\psi\rangle := |\eta\rangle\langle\varphi|\psi\rangle,$$

for $\eta \in \mathcal{H}_B$ and $\varphi \in \mathcal{H}_A$. We can consistently denote $T_{\eta,\varphi}$ as $|\eta\rangle\langle\varphi|$, so that $(|\eta\rangle\langle\varphi|)|\psi\rangle = |\eta\rangle\langle\varphi|\psi\rangle$. Linear operators $\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_B)$ form a complex vector space, provided that we define

$$(aT + bS)|\psi\rangle := aT|\psi\rangle + bS|\psi\rangle.$$

A special linear operator in $\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_A)$ is the *identical operator*, or *identity* I_A :

$$I_A|\psi\rangle = |\psi\rangle \quad \forall \psi \in \mathcal{H}_A.$$

Notice that given the orthonormal basis $\{e_i\}_{i=1}^{d_A}$ one has

$$\sum_{i=1}^{d_A} |e_i\rangle\langle e_i| = I_A. \quad (10.2)$$

Two linear operators $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$ and $S : \mathcal{H}_B \rightarrow \mathcal{H}_C$ can be composed, obtaining the liner operator $R := ST : \mathcal{H}_A \rightarrow \mathcal{H}_C$. The product satisfies the following properties

$$\begin{aligned} (AB)C &= A(BC), \\ A(B+C) &= AB+AC, \\ (A+B)C &= AC+BC. \end{aligned}$$

The above properties imply that $\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_A)$ is an *associative algebra*. Notice that, for any $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$ one has $I_B T = T I_A = T$, thus

$$T = I_B T I_A = \sum_{i=1}^{d_B} \sum_{j=1}^{d_A} |e_i\rangle\langle e_i| T |f_j\rangle\langle f_j| = \sum_{i=1}^{d_B} \sum_{j=1}^{d_A} T_{i,j} |e_i\rangle\langle f_j|, \quad (10.3)$$

where $\{e_i\}_{i=1}^{d_A}$ and $\{f_j\}_{j=1}^{d_B}$ are orthonormal bases in \mathcal{H}_A and \mathcal{H}_B , respectively. Thus, chosen a basis in every space, operators $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$ are in one-to-one correspondence with matrices $\mathbb{M}_{d_B \times d_A}(\mathbb{C})$ as $T \leftrightarrow (T_{i,j})$.

The matrix representation of operators satisfies the following properties

$$\begin{aligned} (aA + bB)_{i,j} &= aA_{i,j} + bB_{i,j} \\ (AB)_{i,j} &= \sum_{k=1}^{d_B} A_{i,k} B_{k,j} \end{aligned}$$

For every operator $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$ we define the *adjoint* $T^\dagger : \mathcal{H}_B \rightarrow \mathcal{H}_A$, defined by the following equation

$$|\eta\rangle = T^\dagger|\psi\rangle \quad \Leftrightarrow \quad \langle\eta|\varphi\rangle = \langle\psi|T|\varphi\rangle, \quad \forall \varphi \in \mathcal{H}_A. \quad (10.4)$$

By definition one then has

$$\begin{aligned}\langle \psi | T | \varphi \rangle^* &= \langle \eta | \varphi \rangle^* \\ &= \langle \varphi | \eta \rangle \\ &= \langle \varphi | T^\dagger | \psi \rangle.\end{aligned}$$

The matrix corresponding to the adjoint of T has then matrix elements

$$T_{i,j}^\dagger = T_{j,i}^*$$

This implies the following identities

$$\begin{aligned}(aT_1 + bT_2)^\dagger &= a^*T_1^\dagger + b^*T_2^\dagger, \quad \forall ab, \in \mathbb{C}, \quad T_1, T_2 : \mathcal{H}_A \rightarrow \mathcal{H}_B \\ (AB)^\dagger &= B^\dagger A^\dagger, \quad \forall B : \mathcal{H}_A \rightarrow \mathcal{H}_B, \quad A : \mathcal{H}_B \rightarrow \mathcal{H}_C, \\ (A^\dagger)^\dagger &= A, \quad \forall A : \mathcal{H}_A \rightarrow \mathcal{H}_B.\end{aligned}$$

Finally, notice that applying the definition in equation (10.4) one can easily prove that $T_{\eta,\varphi}^\dagger = T_{\varphi,\eta}$, i.e.

$$(|\eta\rangle\langle\varphi|)^\dagger = |\varphi\rangle\langle\eta|. \quad (10.5)$$

Let us now consider $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$, and define the *kernel* of T as the set

$$\text{Ker}(T) := \{\psi \in \mathcal{H}_A \mid T|\psi\rangle = 0\}.$$

It is straightforward to prove that $\text{Ker}(T)$ is a subspace of \mathcal{H}_A ². We then define the *support* as

$$\text{Supp}(T) := \text{Ker}(T)^\perp,$$

and then $\mathcal{H}_A = \text{Ker}(T) \oplus \text{Supp}(T)$. Similarly, we define the *range* of T as

$$\text{Rng}(T) := \{\varphi \in \mathcal{H}_B \mid \exists \psi \in \mathcal{H}_A, |\varphi\rangle = T|\psi\rangle\}.$$

Also $\text{Rng}(T)$ is a subspace³, and we can define the *co-kernel* as $\text{Co-Ker}(T) := \text{Rng}(T)^\perp$. It is remarkable that

$$\dim(\text{Rng}(T)) = \dim(\text{Supp}(T)).$$

The above dimension is called *rank* of T , denoted as $\text{rank}(T) := \dim(\text{Supp}(T))$. Notice that $\varphi \in \text{Ker}(T^\dagger)$ if and only if

$$T^\dagger|\varphi\rangle = 0 \Leftrightarrow \langle\varphi|T|\psi\rangle = 0 \quad \forall \psi \in \mathcal{H}_A,$$

²In the infinite dimensional case, this is true of bounded operators T (whose domain can always be extended to all \mathcal{H}_A). Indeed, being a linear operator bounded if and only if it is continuous, then $\text{Ker}(T) = T^{-1}(0) = \mathcal{H}_A \setminus T^{-1}(\mathcal{H}_B \setminus \{0\})$, which is (topologically) closed, being the complement of the open set $T^{-1}(\mathcal{H}_B \setminus \{0\})$.

³In the infinite dimensional case, again this is true if T is bounded.

which implies $\text{Ker}(T^\dagger) = \text{Rng}(T)^\perp$, namely $\text{Supp}(T^\dagger) = \text{Rng}(T)$. An operator $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$ is *invertible* iff it is injective, namely $T|\psi\rangle = 0$ iff $|\psi\rangle = 0$. In other words, T is invertible iff $\text{Ker}(T) = \{0\}$, equivalently $\text{Supp}(T) = \mathcal{H}_A$. In this case, we can define $T^{-1} : \mathcal{H}_B \rightarrow \mathcal{H}_A$ such that $T^{-1}T = I_A$. Indeed, let $\varphi \in \text{Rng}(T)$. Suppose that $|\varphi\rangle = T|\psi\rangle = T|\eta\rangle$, which implies $T(|\psi\rangle - |\eta\rangle) = 0$. Then, it must be $|\psi\rangle = |\eta\rangle$. Thus, there is a unique ψ such that $|\varphi\rangle = T|\psi\rangle$. One can then define $T^{-1}|\varphi\rangle = |\psi\rangle$, and e.g. $\text{Ker}(T^{-1}) = \text{Rng}(T)^\perp$.

Polarisation identity

Let $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$. Then, for every $\psi, \varphi \in \mathcal{H}_A$ one has

$$\langle \psi | T | \varphi \rangle = \frac{1}{4} \sum_{l=\pm 1, \pm i} l \langle \varphi + l\psi | T | \varphi + l\psi \rangle, \quad (10.6)$$

where $|\varphi + l\psi\rangle = |\varphi\rangle + l|\psi\rangle$, and thus $\langle \varphi + l\psi | = \langle \varphi | + l^* \langle \psi |$. Equation (10.6) implies that

$$S = T \Leftrightarrow \langle \psi | T | \psi \rangle = \langle \psi | S | \psi \rangle \quad \forall \psi \in \mathcal{H}_A$$

10.2.1 Unitary operators

A special class of invertible operators in $\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_A)$ is that of *unitary operators*. An operator U is unitary if

$$U^\dagger U = UU^\dagger = I_A.$$

By the polarisation identity, it is immediate to observe that the following conditions are equivalent to unitarity

$$\forall \psi \in \mathcal{H}_A, \quad \|U\psi\| = \|\psi\|, \quad \|U^\dagger\psi\| = \|\psi\|.$$

Since by definition a unitary operator preserves the sesquilinear products, being

$$\langle \varphi | U^\dagger U | \psi \rangle = \langle \varphi | \psi \rangle,$$

a unitary operator maps orthonormal bases to orthonormal bases. Viceversa, an operator that maps an orthonormal basis to another is unitary. Indeed, we can write

$$U = \sum_{i=1}^{d_A} |f_i\rangle \langle e_i|,$$

and since $|\varphi\rangle \langle \psi|^\dagger = |\psi\rangle \langle \varphi|$, we have

$$U^\dagger U = \sum_{i,i'=1}^{d_A} |e_{i'}\rangle \langle f_{i'}| f_i \rangle \langle e_i| = \sum_{i=1}^{d_A} |e_i\rangle \langle e_i| = I_A.$$

Thus, every change of orthonormal basis is represented by a unitary operator, and viceversa every unitary operator sends an orthonormal basis to another one.

Trace

The quantities that are invariant under change of basis, like the determinant, are intrinsic properties of the operators. One of these quantities that plays a distinguished role in quantum theory is the trace $\text{Tr}[T]$ of $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$, defined as the unique linear functional on the space $\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_A)$ such that

$$\text{Tr}[\langle \eta | \varphi \rangle] = \langle \varphi | \eta \rangle. \quad (10.7)$$

One can easily prove that

$$\text{Tr}[T] := \sum_{i=1}^{d_A} \langle e_i | T | e_i \rangle, \quad (10.8)$$

where $\{e_i\}_{i=1}^{d_A}$ is an orthonormal basis in \mathcal{H}_A . Let us now consider N operators $T_1 : \mathcal{H}_{A_1} \rightarrow \mathcal{H}_{A_2}$, $T_2 : \mathcal{H}_{A_2} \rightarrow \mathcal{H}_{A_3}, \dots, T_N : \mathcal{H}_{A_N} \rightarrow \mathcal{H}_{A_1}$. Then the following property of the trace, called *invariance under cyclic permutations* holds

$$\text{Tr}[T_N T_{N-1} \dots T_2 T_1] = \text{Tr}[T_{N-1} \dots T_2 T_1 T_N].$$

Indeed, one has

$$\begin{aligned} \text{Tr}[T_N T_{N-1} \dots T_2 T_1] &= \sum_{e_i=1}^{d_{A_1}} \langle e_i | T_N T_{N-1} \dots T_2 T_1 | e_i \rangle \\ &= \sum_{i=1}^{d_{A_1}} \sum_{j=1}^{d_{A_N}} \langle e_i | T_N | f_j \rangle \langle f_j | T_{N-1} \dots T_2 T_1 | e_i \rangle \\ &= \sum_{e_i=1}^{d_{A_1}} \sum_{j=1}^{d_{A_N}} \langle f_j | T_{N-1} \dots T_2 T_1 | e_i \rangle \langle e_i | T_N | f_j \rangle \\ &= \sum_{j=1}^{d_{A_N}} \langle f_j | T_{N-1} \dots T_2 T_1 T_N | f_j \rangle \\ &= \text{Tr}[T_{N-1} \dots T_2 T_1 T_N]. \end{aligned}$$

Now, since every change of basis is represented by a unitary operator $U : \mathcal{H}_A \rightarrow \mathcal{H}_A$, we can prove that the trace does not depend on the basis in which equation (10.8) is expressed. Indeed, let us write the trace of $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$ in a different basis $\{f_i\}$ with $|f_i\rangle = U|e_i\rangle$

$$\begin{aligned} \sum_{i=1}^{d_A} \langle f_i | T | f_i \rangle &= \sum_{i=1}^{d_A} \langle e_i | U^\dagger T U | e_i \rangle \\ &= \text{Tr}[U^\dagger T U] \\ &= \text{Tr}[T U U^\dagger] \\ &= \text{Tr}[T]. \end{aligned}$$

10.2.2 Selfadjoint operators

A special class of operators $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$ that is relevant to quantum theory is that of *selfadjoint operators*, namely those T such that

$$T^\dagger = T$$

Clearly, a necessary and sufficient condition for T to be selfadjoint is that

$$\langle \varphi | T | \psi \rangle = \langle \psi | T | \varphi \rangle^*, \quad \forall \psi, \varphi \in \mathcal{H}_A.$$

As to the matrix elements, T is selfadjoint if and only if

$$T_{i,j} = T_{j,i}^*$$

in any basis. The identity I_A is thus selfadjoint. Notice that the adjoint is basis independent, however it can be obtained as the composition of two involutive maps that, on the contrary, are both basis dependent: the *transpose*

$$\tau : \mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_B) \rightarrow \mathcal{L}(\mathcal{H}_B \rightarrow \mathcal{H}_A), \quad \tau : T \mapsto T^T, \quad (T^T)_{i,j} := T_{j,i},$$

and the *complex conjugate*

$$* : \mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_B) \rightarrow \mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_B), \quad * : T \mapsto T^*, \quad (T^*)_{i,j} := (T_{i,j})^*.$$

Notice that, since $\text{Supp}(T^\dagger) = \text{Rng}(T)$, for a selfadjoint T it holds that

$$\begin{aligned} \text{Rng}(T) &= \text{Supp}(T), \\ \text{Ker}(T) &= \text{Co} - \text{Ker}(T). \end{aligned}$$

Remark 6. The trace of T is equal to the trace of T^T . Indeed,

$$\text{Tr}[T] = \sum_{i=1}^{d_A} T_{i,i}, \tag{10.9}$$

and $(T^T)_{i,i} = T_{i,i}$. Similarly, it is straightforward to verify that $\text{Tr}[T^*] = \text{Tr}[T]^*$, and thus $\text{Tr}[T^\dagger] = \text{Tr}[T^*] = \text{Tr}[T]^*$.

Every operator $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$ can be decomposed in a *real* and *imaginary* part, both selfadjoint—the names come from the analogy with complex numbers. The decomposition is the following

$$T = A + iB, \quad A = A^\dagger, \quad B = B^\dagger, \tag{10.10}$$

$$A := \frac{1}{2}(T + T^\dagger), \quad B := \frac{1}{2i}(T - T^\dagger). \tag{10.11}$$

Positive operators

A very relevant subclass of selfadjoint operators for quantum theory is that of *positive semidefinite* or *non-negative definite* operators. These are defined as follows

Definition 10.1 (Positive semidefinite operator). An operator $P : \mathcal{H}_A \rightarrow \mathcal{H}_A$ is *positive semidefinite* if

$$\langle \psi | P | \psi \rangle \geq 0 \quad \forall \psi \in \mathcal{H}_A. \quad (10.12)$$

We denote positive semidefiniteness of P as $P \geq 0$.

It is immediate to verify that the set $\mathcal{P}(A) := \{P : \mathcal{H}_A \rightarrow \mathcal{H}_A | P \geq 0\}$ is a *cone*. Indeed, for any $\lambda \geq 0$ and $P \in \mathcal{P}(A)$ one has $\lambda P \in \mathcal{P}(A)$.

Substituting strict inequality in the definition, one obtains the definition of *positive definite* operators $P > 0$. Notice that selfadjointness is not required in the definition. Indeed, it comes as a consequence of the definition itself, through the polarisation identity

Lemma 10.2. *If an operator $P : \mathcal{H}_A \rightarrow \mathcal{H}_A$ is positive semidefinite, then it is selfadjoint.*

Proof. Let us evaluate through the polarisation identity (10.6) the general quantity

$$\langle \psi | P | \varphi \rangle = \frac{1}{4} \sum_{l=\pm 1, \pm i} l \langle \varphi + l\psi | P | \varphi + l\psi \rangle.$$

If $P \geq 0$ then we have

$$\langle \psi | P | \varphi \rangle^* = \frac{1}{4} \sum_{l=\pm 1, \pm i} l^* \langle \varphi + l\psi | P | \varphi + l\psi \rangle,$$

and since

$$\begin{aligned} |\varphi + l\psi\rangle &= |\varphi\rangle + l|\psi\rangle \\ &= l(|\psi\rangle + l^*|\varphi\rangle) \\ &= l|\psi + l^*\varphi\rangle, \end{aligned}$$

and $\langle \varphi + l\psi | = l^* \langle \psi + l^*\varphi |$, we have

$$\begin{aligned} \langle \psi | P | \varphi \rangle^* &= \frac{1}{4} \sum_{l=\pm 1, \pm i} l^* \langle \psi + l^*\varphi | P | \psi + l\varphi \rangle \\ &= \frac{1}{4} \sum_{l'=\pm 1, \pm i} l' \langle \psi + l'\varphi | P | \psi + l'\varphi \rangle \\ &= \langle \varphi | P | \psi \rangle. \end{aligned} \quad \square$$

In the following, “positive semidefinite” will be often replaced by “non-negative”, or with an even worse abuse of terminology, “positive”.

Projections

A further restriction of the class of positive semidefinite operators, that also plays a major role in quantum theory, is that of *projections*.

Definition 10.3 (Projection). An operator $P : \mathcal{H}_A \rightarrow \mathcal{H}_A$ is a *projection* if $P = P^\dagger$ and $P^2 = P$.

First of all, we prove that a projection is indeed positive semidefinite.

Lemma 10.4. *Let $P : \mathcal{H}_A \rightarrow \mathcal{H}_A$ be a projection. Then $P \geq 0$.*

Proof. Let us consider the expectation

$$\begin{aligned}\langle \psi | P | \psi \rangle &= \langle \psi | P^2 | \psi \rangle \\ &= \langle \psi | P^\dagger P | \psi \rangle \\ &= \|P\psi\|^2 \geq 0, \quad \forall \psi \in \mathcal{H}_A.\end{aligned}$$
 \square

Moreover, if P is a projection, also $I_A - P$ is, since $(I_A - P)^\dagger = (I_A - P)$, and $(I_A - P)^2 = (I_A - P)$. A projection P projects any vector in the subspace $\text{Rng}(P) = \text{Supp}(P)$. Indeed, suppose that $\psi \in \text{Supp}(P) = \text{Rng}(P)$. Then there exists $\varphi \in \mathcal{H}_A$ such that

$$|\psi\rangle = P|\varphi\rangle,$$

and then

$$\begin{aligned}P|\psi\rangle &= P^2|\varphi\rangle \\ &= P|\varphi\rangle \\ &= |\psi\rangle\end{aligned}$$

Thus, since $\mathcal{H}_A = \text{Supp}(P) \oplus \text{Ker}(P)$, if one decomposes a general vector $\psi \in \mathcal{H}_A$ accordingly, it is

$$\psi = \psi_S + \psi_K,$$

and by definition of kernel one has

$$\begin{aligned}P|\psi\rangle &= P|\psi_S\rangle + P|\psi_K\rangle \\ &= P|\psi_S\rangle \\ &= |\psi_S\rangle.\end{aligned}$$

By the above observation, it is then clear $\text{Supp}(P) = \text{Ker}(I_A - P)$, and thus also $\text{Ker}(P) = \text{Supp}(I_A - P)$. This implies that for any projection P there exists $I_A - P$ such that P projects on $\text{Supp}(P)$ and $I_A - P$ projects on $\text{Supp}(P)^\perp = \text{Ker}(P)$. This statement can be reversed as follows.

Lemma 10.5. *Let $\mathcal{K} \subseteq \mathcal{H}_A$ be a subspace. Then there exists a projection P with $\text{Supp}(P) = \mathcal{K}$.*

Proof. Let us construct the operator P as follows. First, let $\psi \in \mathcal{H}$, and consider the decomposition

$$|\psi\rangle = |\psi_{\mathcal{K}}\rangle + |\psi_{\mathcal{K}^\perp}\rangle.$$

Then we define P by

$$P|\psi\rangle := |\psi_{\mathcal{K}}\rangle, \quad \forall \psi \in \mathcal{H}.$$

The operator P is well defined, since for every $\psi \in \mathcal{H}_A$ the decomposition $\psi = \psi_{\mathcal{K}} + \psi_{\mathcal{K}^\perp}$ is unique, and moreover P is linear, since

$$\begin{aligned} a\psi + b\varphi &= a\psi_{\mathcal{K}} + a\psi_{\mathcal{K}^\perp} + b\varphi_{\mathcal{K}} + b\varphi_{\mathcal{K}^\perp} \\ &= (a\psi_{\mathcal{K}} + b\varphi_{\mathcal{K}}) + (a\psi_{\mathcal{K}^\perp} + b\varphi_{\mathcal{K}^\perp}), \end{aligned}$$

which implies

$$\begin{aligned} P(a|\psi\rangle + b|\varphi\rangle) &= a|\psi_{\mathcal{K}}\rangle + b|\varphi_{\mathcal{K}}\rangle \\ &= aP|\psi\rangle + bP|\varphi\rangle. \end{aligned}$$

Finally, P is a projection. Indeed, P is positive semidefinite, since

$$\begin{aligned} \langle \psi | P | \psi \rangle &= (\langle \psi_{\mathcal{K}} | + \langle \psi_{\mathcal{K}^\perp} |) |\psi_{\mathcal{K}}\rangle \\ &= \langle \psi_{\mathcal{K}} | \psi_{\mathcal{K}} \rangle \geq 0, \end{aligned}$$

and thus $P^\dagger = P$. Moreover,

$$\begin{aligned} P^2|\psi\rangle &= P|\psi_{\mathcal{K}}\rangle \\ &= |\psi_{\mathcal{K}}\rangle \\ &= P|\psi\rangle. \end{aligned} \quad \square$$

Let now $\eta \in \text{Ker}(P)$. Then $\eta_{\mathcal{K}} = 0$, i.e. $\eta \in \mathcal{K}^\perp$. Viceversa, if $\eta \in \mathcal{K}^\perp$ then $P|\eta\rangle = 0$. This shows that $\text{Ker}(P) = \mathcal{K}^\perp$, i.e. $\text{Supp}(P) = \mathcal{K}$.

In conclusion, this result shows that there is a one-to-one correspondence between subspaces of \mathcal{H}_A and projections in $\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_A)$.

10.2.3 Normal operators and the spectral theorem

All the classes of operators $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$ that we considered so far belong to a single general class, the *normal operators*. What makes normal operators special is that they are the most general class of operators that satisfy the *spectral theorem*.

Definition 10.6 (Normal operator). An operator $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$ is *normal* if it commutes with its adjoint, namely

$$[T, T^\dagger] = 0. \quad (10.13)$$

Notice that the *commutator* of $T, S : \mathcal{H}_A \rightarrow \mathcal{H}_A$ is defined as

$$[S, T] := ST - TS = -[T, S].$$

We say that two operators *commute* if their commutator is null.

Lemma 10.7. *An operator $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$ is normal iff its real and imaginary parts $A = \frac{1}{2}(T + T^\dagger)$ and $B = \frac{1}{2i}(T - T^\dagger)$ commute.*

Proof. Let $T = A + iB$ with $A = A^\dagger$ and $B = B^\dagger$. Then

$$\begin{aligned} [T, T^\dagger] &= [A + iB, A - iB] \\ &= [A, A] + [B, B] + i[B, A] - i[A, B] \\ &= 2i[B, A]. \end{aligned}$$

Thus, $[T, T^\dagger] = 0$ iff $[A, B] = 0$. \square

Definition 10.8 (Diagonalisable operator). We say that $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$ is diagonalisable if there are projections P_k satisfying $P_k P_{k'} = \delta_{kk'} P_k$ and $\{\lambda_k\} \subseteq \mathbb{C}$ with $\lambda_{k_1} \neq \lambda_{k_2}$ such that

$$T = \sum_k \lambda_k P_k, \quad I_A = \sum_k P_k. \tag{10.14}$$

The support of P_k in equation (10.14) is called *eigenspace* of T .

Lemma 10.9. *Let $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$ be diagonalisable. The decomposition of equation (10.14) for T is unique.*

Proof. If $T = 0$, then $T = 0I_A$, and any other decomposition would be forbidden by the condition $\lambda_{k_1} \neq \lambda_{k_2}$. Let then $T \neq 0$, and

$$T = \sum_k \lambda_k P_k = \sum_j \mu_j Q_j.$$

Then

$$P_k T Q_j = \lambda_k P_k Q_j = \mu_j P_k Q_j.$$

Thus, either $P_k Q_j = 0$ or $\mu_j = \lambda_k$. Now, $P_k Q_j$ cannot be null for every j, k , otherwise $\lambda_k P_k = P_k T = \sum_j \mu_j P_k Q_j = 0$ for every k , which means $T = 0$. Then there must exist j such that $P_k Q_j \neq 0$. Such a value of j is unique, otherwise $\mu_{j_1} = \lambda_k = \mu_{j_2}$, contrarily to the hypotheses. Let us then define the function $\varphi(k)$ such that $\lambda_k = \mu_{\varphi(k)}$. The function is injective, otherwise $\varphi(k_1) = \varphi(k_2)$ would imply $\lambda_{k_1} = \mu_{\varphi(k_1)} = \mu_{\varphi(k_2)} = \lambda_{k_2}$. Thus,

$$\begin{aligned} \lambda_k P_k &= P_k T \\ &= \mu_{\varphi(k)} P_k Q_{\varphi(k)} \\ &= \lambda_k P_k Q_{\varphi(k)} \\ &= T Q_{\varphi(k)} \\ &= \mu_{\varphi(k)} Q_{\varphi(k)}. \end{aligned}$$

Finally, this implies $Q_{\varphi(k)} = P_k$ and $\mu_{\varphi(k)} = \lambda_k$, for $\lambda_k \neq 0$. Finally, if there is k_0 such that $\lambda_{k_0} = \mu_{\varphi(k_0)} = 0$, one has

$$\begin{aligned} P_{k_0} &= I - \sum_{k \neq k_0} P_k \\ &= I - \sum_{k \neq k_0} Q_{\varphi(k)} \\ &= Q_{\varphi(k_0)}. \end{aligned} \quad \square$$

Definition 10.10. Let $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$. We say that $0 \neq \psi \in \mathcal{H}_A$ is an *eigenvector* of T with *eigenvalue* λ if

$$T|\psi\rangle = \lambda|\psi\rangle. \quad (10.15)$$

The set of eigenvalues of T is called *spectrum* of T , and denoted as $\text{Spec}(T)$.

Let T be diagonalisable. By the Gram-Schmidt construction, one can take a basis $\{\varphi_i^{(k)}\}_{i=1}^{\dim \mathcal{K}_k}$ for each eigenspace $\mathcal{K}_k := \text{Supp}(P_k)$, and thus $\{\psi_j\}_{j=1}^{d_A} = \bigcup_k \{\varphi_i^{(k)}\}_{i=1}^{\dim \mathcal{K}_k}$ constitutes an orthonormal basis for \mathcal{H}_A . Moreover, one can easily verify that for every k

$$P_k = \sum_{i=1}^{\dim \mathcal{K}_k} |\varphi_i^{(k)}\rangle \langle \varphi_i^{(k)}|.$$

Thus, an operator T is diagonalisable iff

$$T = \sum_k \sum_{i=1}^{\dim \mathcal{K}_k} \lambda_k |\varphi_i^{(k)}\rangle \langle \varphi_i^{(k)}|,$$

which implies $T|\varphi_i^{(k)}\rangle = \lambda_k |\varphi_i^{(k)}\rangle$. Namely, if T is diagonalisable there is an orthonormal basis of eigenvectors of T , and viceversa if there is an orthonormal basis $\{\psi_i\}_{i=1}^{d_A}$ of eigenvectors of T one has

$$T|\psi_i\rangle = \lambda_i |\psi_i\rangle,$$

and one can easily verify that this implies

$$T = \sum_{i=1}^{d_A} \lambda_i |\psi_i\rangle \langle \psi_i|.$$

Finally, an eigenspace \mathcal{K}_k contains all eigenvectors corresponding to the same eigenvalue λ_k . The dimension $\dim \mathcal{K}_k = \text{Tr}[P_k]$ is called *degeneracy* of λ_k .

We now state the spectral theorem without proving it, as follows.

Theorem 10.11 (Spectral theorem). *The operator $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$ is diagonalisable iff T is normal.*

Corollary 10.12. *Unitary and selfadjoint operators are diagonalisable.*

Lemma 10.13. *Let T be diagonalisable. The spectrum $\text{Spec}(T^\dagger)$ of the adjoint operator of T is the complex conjugate of $\text{Spec}(T)$, and the eigenvectors are the same.*

Proof. Let

$$T = \sum_k \lambda_k P_k.$$

Then

$$T^\dagger = \sum_k \lambda_k^* P_k,$$

and clearly $\text{Spec}(T^\dagger) = \text{Spec}(T)^*$. Moreover, the projections for the decomposition of T and T^\dagger are the same, thus they have the same eigenvectors. \square

Corollary 10.14. 1. A normal operator U is unitary iff $\text{Spec}(U)$ is contained in the unit circle \mathbb{C} in the complex plane \mathbb{C} .

2. A normal operator T is selfadjoint iff $\text{Spec}(T)$ is contained in the real axis.

3. A normal operator P is positive semidefinite or positive definite iff $\text{Spec}(P)$ is non-negative or positive, respectively.

4. A normal operator P is a projection iff $\text{Spec}(P)$ is contained in $\{0, 1\}$.

Proof. 1. Let U be normal, and $\{\psi_i\}_{i=1}^{d_A}$ be an orthonormal basis of eigenvectors $U|\psi_i\rangle = \lambda_i|\psi_i\rangle$. Then by lemma 10.13 U is unitary iff

$$\begin{aligned} |\psi_i\rangle &= U^\dagger U |\psi_i\rangle \\ &= \lambda_i U^\dagger |\psi_i\rangle \\ &= \lambda_i \lambda_i^* |\psi_i\rangle \\ &= |\lambda_i|^2 |\psi_i\rangle. \end{aligned}$$

Thus, U is unitary iff $\text{Spec}(U) \subseteq \mathbb{C}$.

2. Let T be normal, and $\{\psi_i\}_{i=1}^{d_A}$ be an orthonormal basis of eigenvectors $T|\psi_i\rangle = \lambda_i|\psi_i\rangle$. Then T is selfadjoint iff

$$\begin{aligned} \lambda_i |\psi_i\rangle &= T |\psi_i\rangle \\ &= T^\dagger |\psi_i\rangle \\ &= \lambda_i^* |\psi_i\rangle, \end{aligned}$$

Thus, T is selfadjoint iff $\text{Spec}(T) \subseteq \mathbb{R}$.

3. Let P be normal, and $\{\psi_i\}_{i=1}^{d_A}$ be an orthonormal basis of eigenvectors $P|\psi_i\rangle = \lambda_i|\psi_i\rangle$. Then P is positive semidefinite iff

$$\begin{aligned} 0 &\leq \langle \psi_i | P | \psi_i \rangle \\ &= \lambda_i \langle \psi_i | \psi_i \rangle \end{aligned}$$

Thus, P is positive semidefinite iff $\text{Spec}(P) \geq 0$. For positive definite P , it is sufficient to replace the first inequality by its strict version.

4. Let P be normal and $\{\psi_i\}_{i=1}^{d_A}$ be an orthonormal basis of eigenvectors $P|\psi_i\rangle = \lambda_i|\psi_i\rangle$. Then P is a projection iff

$$\begin{aligned} \lambda_i |\psi_i\rangle &= P|\psi_i\rangle \\ &= P^2|\psi_i\rangle \\ &= \lambda_i^2 |\psi_i\rangle. \end{aligned}$$

Thus, P is a projection iff $\text{Spec}(P) \subseteq \{0, 1\}$. \square

The above results allow us to prove two relevant properties of positive operators. These results are related to the observation that, for an arbitrary $A : \mathcal{H}_A \rightarrow \mathcal{H}_B$, one has

$$A^\dagger A \geq 0.$$

Indeed, we have

$$\langle \psi | A^\dagger A | \psi \rangle = \|A\psi\|^2 \geq 0, \quad \forall \psi \in \mathcal{H}_A. \quad (10.16)$$

The above result can also be reversed: P is positive iff $P = A^\dagger A$. Let us now see how.

Theorem 10.15 (Existence and uniqueness of the square root). *The operator $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$ is positive iff there exists a positive operator $S : \mathcal{H}_A \rightarrow \mathcal{H}_A$ such that $S^2 = T$. Such an operator S is unique.*

Proof. Let $S \geq 0$ and $T = S^2$. Then $T = S^\dagger S \geq 0$. On the other hand, let $T \geq 0$. Then

$$T = \sum_k \lambda_k P_k, \quad \lambda_k \geq 0 \ \forall k.$$

Let us define

$$S := \sum_k \sqrt{\lambda_k} P_k.$$

Then clearly $S \geq 0$ and $S^2 = T$. Let us now prove uniqueness of S . Indeed, suppose that $S_1^2 = S_2^2 = T \geq 0$ and $S_i \geq 0$ for both $i = 1$ and $i = 2$. Then

$$T = \sum_k (\lambda_k^{(1)})^2 P_k^{(1)} = \sum_k (\lambda_k^{(2)})^2 P_k^{(2)}.$$

However, by the uniqueness of the diagonalisation of T proved in lemma 10.9, it must be $P_k^{(2)} = P_{\varphi(k)}^{(1)}$ and since $S_i \geq 0$ for $i = 1, 2$, the eigenvalues of S_i must be positive for $i = 1, 2$, then $\lambda_k^{(2)} = \lambda_{\varphi(k)}^{(1)}$. \square

The unique positive operator $S \geq 0$ such that $S^2 = T$ for $T \geq 0$ will be denoted by \sqrt{T} or $T^{\frac{1}{2}}$.

Corollary 10.16. *The operator $P : \mathcal{H}_A \rightarrow \mathcal{H}_A$ is positive iff there exists \mathcal{H}_B and $A : \mathcal{H}_A \rightarrow \mathcal{H}_B$ such that $P = A^\dagger A$.*

Corollary 10.17. *Let $P : \mathcal{H}_A \rightarrow \mathcal{H}_A$ be positive semidefinite. Then*

$$\langle \psi | P | \psi \rangle = 0 \quad (10.17)$$

iff $\psi \in \text{Ker}(P)$.

Proof. Since $P = A^\dagger A$, we have

$$\langle \psi | P | \psi \rangle = 0 \Leftrightarrow 0 = \langle \psi | A^\dagger A | \psi \rangle = \|A\psi\|^2 = 0$$

which is equivalent to $A|\psi\rangle = 0$. Finally, this implies $A^\dagger A|\psi\rangle = P|\psi\rangle = 0$. The converse is trivial. \square

Corollary 10.18. *Let $P : \mathcal{H}_A \rightarrow \mathcal{H}_A$ be positive. Then*

$$\text{Tr}[P] = 0 \quad (10.18)$$

iff $P = 0$

Proof. By definition of trace, given any orthonormal basis $\{\varphi_i\}_{i=1}^{d_A}$ one has

$$\text{Tr}[P] = \sum_{i=1}^{d_A} \langle \varphi_i | P | \varphi_i \rangle, \quad (10.19)$$

and since $P \geq 0$ implies that $\langle \varphi_i | P | \varphi_i \rangle \geq 0$ for every i , the trace is null iff $\langle \varphi_i | P | \varphi_i \rangle = 0$ for every i . By corollary 10.17, this is equivalent to say that for every basis $\{\varphi_i\}_{i=1}^{d_A}$, every φ_i belongs to $\text{Ker}(P)$. Thus, $\text{Ker}(P) = \mathcal{H}_A$ and $P = 0$. \square

Notice that the set $\mathcal{P}(A) := \{P : \mathcal{H}_A \rightarrow \mathcal{H}_A \mid P \geq 0\}$ is complete in $\mathcal{L}(\mathcal{H}_A)$. Indeed, the following lemmas hold.

Lemma 10.19. *Let $A : \mathcal{H}_A \rightarrow \mathcal{H}_A$ be selfadjoint. Then $A = A_+ - A_-$ where $A_\pm \geq 0$ and $A_+ A_- = A_- A_+ = 0$.*

Proof. If $A = A^\dagger$ then A is diagonalisable, and has real eigenvalues $\{\lambda_k\}$. If we divide the spectrum as $\{\lambda_k\} = \mathbb{S}_+ \cup \mathbb{S}_-$ such that $\lambda_k \in \mathbb{S}_+$ if $\lambda_k \geq 0$ and $\lambda_k \in \mathbb{S}_-$ if $\lambda_k < 0$, we have

$$\begin{aligned} A &= A_+ - A_- \\ A_+ &:= \sum_{\lambda_k \in \mathbb{S}_+} \lambda_k P_k \geq 0 \\ A_- &:= \sum_{\lambda_k \in \mathbb{S}_-} |\lambda_k| P_k \geq 0. \end{aligned}$$

Moreover, the projections on eigenspaces satisfy $P_k P_{k'} = 0$ for $\lambda_k \in \mathbb{S}_+$ and $\lambda_{k'} \in \mathbb{S}_-$. Thus, $A_+ A_- = A_- A_+ = 0$. \square

Moreover, every operator $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$ can be expressed as in equation (10.11) as $T = A + iB$, thus we have the following corollary.

Corollary 10.20. *Let $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$. Then T can be expanded as a linear combination of positive operators.*

Proof. It is sufficient to write

$$\begin{aligned} T &= A + iB \\ &= A_+ - A_- + iB_+ - iB_-. \end{aligned} \quad \square$$

Finally, since for every positive operator $A \geq 0$ one has $\text{Tr}[A] = 0$ iff $A = 0$, we have $A = \rho_A \text{Tr}[A]$, with $\rho_A \geq 0$ and $\text{Tr}[\rho_A] = 1$. Thus, positive operators with unit trace are complete in $\mathcal{L}(\mathcal{H}_A)$, as one can write

$$\begin{aligned} T &= A_+ - A_- + iB_+ - iB_- \\ &= \text{Tr}[A_+] \rho_{A_+} - \text{Tr}[A_-] \rho_{A_-} + i \text{Tr}[B_+] \rho_{B_+} - i \text{Tr}[B_-] \rho_{B_-}. \end{aligned}$$

Definition 10.21 (Modulus of an operator). Let $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$. Then we define the *modulus* of T as the operator

$$|T| := (T^\dagger T)^{\frac{1}{2}}. \quad (10.20)$$

Remark 7. The above definition of \sqrt{T} is just a special case of functional calculus for normal operators T , which can be consistently defined. Indeed, let $T = \sum_k \lambda_k P_k$. Now, for a general function $f : \mathbb{C} \rightarrow \mathbb{C}$ we can define

$$f(T) := \sum_k f(\lambda_k) P_k.$$

One can easily verify that $(g \circ f)(T) = g[f(T)]$. Moreover, functional calculus can be defined for selfadjoint operators also for functions $f : \mathbb{R} \rightarrow \mathbb{R}$. Finally, for positive semidefinite operators P one can define $f(P)$ even for functions $f : [0, \infty) \rightarrow \mathbb{R}$. This is actually the case of interest here, as we will extensively use the function $f(x) := x \log_2 x$ which can be defined on non-negative reals.

10.3 Isometric operators

The last class of operators that we introduce is that of *isometric operators* or *isometries*.

Definition 10.22 (Isometric operator). Let $V : \mathcal{H}_A \rightarrow \mathcal{H}_B$ with $d_B \geq d_A$. We say that V is an *isometric operator* or *isometry* if

$$V^\dagger V = I_A. \quad (10.21)$$

Remark 8. By the polarisation identity, it is immediate to observe that the following condition is equivalent to isometry

$$\forall \psi \in \mathcal{H}_A, \quad \|V\psi\| = \|\psi\|.$$

Since by definition an isometric operator preserves the sesquilinear products, being

$$\langle \varphi | V^\dagger V | \psi \rangle = \langle \varphi | \psi \rangle,$$

an isometric operator maps orthonormal bases to orthonormal sets, generally not complete.

Remark 9. Notice that, differently from the case of unitary operators we do not require $VV^\dagger = I_B$. In the general case, however, VV^\dagger is a projection. Indeed, $(VV^\dagger)^\dagger = VV^\dagger$, and $(VV^\dagger)(VV^\dagger) = V(V^\dagger V)V^\dagger = VV^\dagger$. Since for an orthonormal basis $\{\varphi_i\}_{i=1}^{d_A}$ one has that $\{V\varphi_i\}_{i=1}^{d_A}$ is an orthonormal set, this set will be a basis for $\text{Rng}(V)$. Moreover, let $\psi \in \text{Rng}(V)$. Then there exists $\varphi \in \mathcal{H}_A$ such that $\psi = V\varphi$. Thus

$$\begin{aligned} VV^\dagger |\psi\rangle &= VV^\dagger V|\varphi\rangle \\ &= V|\varphi\rangle \\ &= |\psi\rangle. \end{aligned}$$

Since VV^\dagger is a projection, it follows that $\text{Rng}(V) \subseteq \text{Supp}(VV^\dagger)$. Moreover, let $\eta \in \text{Rng}(V)^\perp$. Then for all $\psi \in \mathcal{H}_A$

$$\begin{aligned} 0 &= \langle \eta | V | \psi \rangle \\ &= \langle \psi | V^\dagger | \eta \rangle^*, \end{aligned}$$

namely $\eta \in \text{Ker}(V^\dagger)$, and thus $\eta \in \text{Ker}(VV^\dagger)$. Then $\text{Rng}(V)^\perp \subseteq \text{Ker}(VV^\dagger)$, and finally $\text{Supp}(VV^\dagger) \subseteq \text{Rng}(V)$. This implies that $\text{Supp}(VV^\dagger) = \text{Rng}(V)$, and VV^\dagger is the projection on $\text{Rng}(V)$. Thus, $\dim \text{Supp}(VV^\dagger) = \dim \text{Rng}(V) = d_A$, and if $d_B = d_A$ one has $VV^\dagger = I_B$, namely V is unitary.

Remark 10. The adjoint V^\dagger of an isometry V is called co-isometry, and is a special case of a *partial isometry*. In general, an operator $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$ is called a *partial isometry* if $V^\dagger V = P_{\text{Supp}(V)}$.

10.4 Polar decomposition

Theorem 10.23 (Polar decomposition). *Let $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$, with $d_A \leq d_B$. Then there exists an isometry $S : \mathcal{H}_A \rightarrow \mathcal{H}_B$ such that*

$$T = S|T|. \tag{10.22}$$

On the other hand, if $d_B < d_A$ there exists a co-isometry $S : \mathcal{H}_A \rightarrow \mathcal{H}_B$ such that equation (10.22) holds, and the projection $S^\dagger S$ satisfies $S^\dagger S \geq P_{\text{Supp}(|T|)}$, namely S is isometric on $\text{Supp}(|T|)$.

Proof. For $T = 0$ the thesis is trivial. Let us then consider $T \neq 0$. We observe that for every $\psi \in \mathcal{H}_A$

$$\begin{aligned}\|T\psi\|^2 &= \langle\psi|T^\dagger T|\psi\rangle \\ &= \langle\psi|(|T|^2)|\psi\rangle \\ &= \|(|T|\psi)\|^2.\end{aligned}\tag{10.23}$$

Let us now define $S_0 : \mathcal{H}_A \rightarrow \mathcal{H}_B$ with $\text{Ker}(S_0) = \text{Ker}(|T|)$, and for $\psi \in \text{Rng}(|T|) = \text{Supp}(|T|)$, namely $\psi = |T|\varphi$

$$S_0|\psi\rangle := T|\varphi\rangle.$$

We need to prove that the operator S_0 is well defined. For this purpose, let

$$|\psi\rangle = |T||\varphi_0\rangle = |T||\varphi_1\rangle.$$

Then by equation (10.23) we have

$$\begin{aligned}0 &= \|\{|T|(\varphi_0 - \varphi_1)\}\|^2 \\ &= \|T(\varphi_0 - \varphi_1)\|^2,\end{aligned}$$

and finally $S_0|\psi\rangle = T|\varphi_0\rangle = T|\varphi_1\rangle$. Thus, S_0 is well defined. We now show that S_0 is linear. Indeed, let $\psi = a\psi_1 + b\psi_2$. Then

$$\begin{aligned}P_{\text{Supp}(|T|)}|\psi\rangle &= aP_{\text{Supp}(|T|)}|\psi_1\rangle + bP_{\text{Supp}(|T|)}|\psi_2\rangle \\ &= a|T||\varphi_1\rangle + b|T||\varphi_2\rangle \\ &= |T|(a|\varphi_1\rangle + b|\varphi_2\rangle),\end{aligned}$$

and thus

$$\begin{aligned}S_0(a|\psi_1\rangle + b|\psi_2\rangle) &= S_0|\psi\rangle \\ &= S_0P_{\text{Supp}(|T|)}|\psi\rangle \\ &= T(a|\varphi_1\rangle + b|\varphi_2\rangle) \\ &= aS_0|\psi_1\rangle + bS_0|\psi_2\rangle.\end{aligned}$$

Since $|T|$ is selfadjoint, $\text{Supp}(|T|) = \text{Rng}(|T|)$, therefore for any $\psi \in \mathcal{H}_A$ there exists φ such that $P_{\text{Supp}(|T|)}\psi = |T|\varphi$. By equation (10.23) we then have for every $\psi \in \mathcal{H}_A$

$$\begin{aligned}\langle\psi|(S_0^\dagger S_0 - P_{\text{Supp}(|T|)})|\psi\rangle &= \langle\psi|P_{\text{Supp}(|T|)}(S_0^\dagger S_0 - P_{\text{Supp}(|T|)})P_{\text{Supp}(|T|)}|\psi\rangle \\ &= \langle\varphi||T|(S_0^\dagger S_0 - P_{\text{Supp}(|T|)})|T||\varphi\rangle \\ &= 0.\end{aligned}$$

Finally, by the polarisation identity this implies that $S_0^\dagger S_0 = P_{\text{Supp}(|T|)}$. Since $\text{Rng}(|T|) = \text{Supp}(|T|)$, we have $k := \dim \text{Ker}(|T|) = d_A - \dim \text{Supp}(|T|) \leq d_B - \dim \text{Supp}(|T|) = d_B - \dim \text{Rng}(S_0) = \dim \text{Co} - \text{Ker}(S_0)$. Thus, for an orthonormal basis $\{f_j\}_{j=1}^k$ of $\text{Ker}(|T|)$

one can find k orthonormal vectors $\{f'_j\}_{j=1}^k$ in $\text{Co} - \text{Ker}(S_0)$. We then construct the operator

$$S := S_0 + S_1, \quad S_1 := \sum_{i=1}^k |f'_i\rangle\langle f_i|.$$

Now, it is immediate to see that $S_1^\dagger S_0 = 0$, since for every $\psi \in \mathcal{H}$ one has $\langle f'_i | S_0 | \psi \rangle = 0$ for all $1 \leq i \leq k$, and thus

$$S_1^\dagger S_0 = \sum_{i=1}^k |f_i\rangle\langle f'_i| S_0 = 0.$$

Moreover, $S_1^\dagger S_1 = I_A - P_{\text{Supp}(|T|)}$, since

$$S_1^\dagger S_1 = \sum_{i=1}^k |f_i\rangle\langle f_i| = P_{\text{Ker}(|T|)}.$$

Then one can easily verify that

$$\begin{aligned} S^\dagger S &= (S_0 + S_1)^\dagger (S_0 + S_1) \\ &= S_0^\dagger S_0 + S_0^\dagger S_1 + S_1^\dagger S_0 + S_1^\dagger S_1 \\ &= S_0^\dagger S_0 + S_1^\dagger S_1 = I_A. \end{aligned}$$

In conclusion, S is isometric and for every $\psi \in \mathcal{H}_A$ one has

$$S|T||\psi\rangle = T|\psi\rangle.$$

Now, if $d_B < d_A$ we can use the first part of the theorem to prove that $T^\dagger = S|T^\dagger|$. Since $|T^\dagger|^2 = TT^\dagger$, we have

$$\begin{aligned} |T|^2 &= T^\dagger T \\ &= S|T^\dagger|^2 S^\dagger \\ &= S|T^\dagger|S^\dagger S|T^\dagger|S^\dagger, \end{aligned}$$

and then, since $\langle \psi | S|T^\dagger|S^\dagger | \psi \rangle \geq 0$ for every $\psi \in \mathcal{H}_A$, one has

$$S|T^\dagger|S^\dagger = |T|.$$

Thus, since $T = |T^\dagger|S^\dagger$, we conclude that

$$T = S^\dagger S|T^\dagger|S^\dagger = S^\dagger|T|.$$

Notice that since $SS^\dagger|T| = S|T^\dagger|S^\dagger = |T|$, one has $SS^\dagger \geq P_{\text{Supp}(|T|)}$. \square

Corollary 10.24. *Let $A : \mathcal{H}_A \rightarrow \mathcal{H}_B$ and $B : \mathcal{H}_A \rightarrow \mathcal{H}_C$. Then*

$$A^\dagger A = B^\dagger B \tag{10.24}$$

iff $A = VQ$ and $B = WQ$, with $Q := |A| = |B|$ and $P_{\text{Supp}(Q)} \leq V^\dagger V \leq I_A$ and $P_{\text{Supp}(Q)} \leq W^\dagger W \leq I_A$.

10.4.1 Singular value decomposition

Theorem 10.25 (Singular value decomposition). *Let $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$, with $d_A \leq d_B$. Then there exists an isometry $V : \mathcal{H}_A \rightarrow \mathcal{H}_B$ and a unitary U such that*

$$T = V\Sigma U, \quad (10.25)$$

where $\Sigma_{i,j} = \sigma_i \delta_{i,j}$, and $\{\sigma_i\} = \text{Spec}(|T|)$. On the other hand, if $d_B < d_A$ there exists a co-isometry $V : \mathcal{H}_A \rightarrow \mathcal{H}_B$ such that equation (10.25) holds and $V^\dagger V \geq P_{\text{Supp}(\Sigma)}$.

Proof. The proof uses the polar decomposition

$$T = S|T|,$$

and since $|T| \geq 0$ one has

$$|T| = \sum_{i=1}^{d_A} \sigma_i |\psi_i\rangle\langle\psi_i|,$$

with $\sigma_i \geq 0$ for all values of i . Upon defining $U|\psi_i\rangle = |e_i\rangle$, where $\{e_i\}_{i=1}^{d_A}$ is the canonical basis, we have $U|T|U^\dagger = \Sigma$, and then

$$T = SU^\dagger\Sigma U.$$

Finally, the thesis is proved by setting $V := SU^\dagger$. □

Corollary 10.26. *For every $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$ one has*

$$T = \sum_{i=1}^{d_A} \sigma_i |\varphi_i\rangle\langle\psi_i|, \quad (10.26)$$

where $\{\sigma_i\} = \text{Spec}(|T|) \subseteq \mathbb{R}_+$ and $\langle\varphi_i|\varphi_j\rangle = \langle\psi_i|\psi_j\rangle = \delta_{i,j}$.

10.5 Tensor product

We will now construct the tensor product of two Hilbert spaces, and review a few of its properties. We also introduce the *double-ket* notation, a very useful tool for calculations, which makes it very practical to use tensor product spaces.

Let $\mathcal{H}_A \simeq \mathbb{C}^{d_A}$ and $\mathcal{H}_B \simeq \mathbb{C}^{d_B}$. We formally introduce free linear span \mathcal{F} of formal vectors $\psi \otimes \varphi$ with $(\psi, \varphi) \in \mathcal{H}_A \times \mathcal{H}_B$, namely

$$\mathcal{F} = \left\{ \sum_l c_l (\psi_l \otimes \varphi_l) \mid (\psi_l, \varphi_l) \in \mathcal{H}_A \times \mathcal{H}_B, \{c_l\} \subseteq \mathbb{C} \right\}.$$

We define the following equivalence relations in \mathcal{F}

$$\begin{aligned} (a\psi_1 + b\psi_2) \otimes \varphi &\sim a(\psi_1 \otimes \varphi) + b(\psi_2 \otimes \varphi), \\ \psi \otimes (a\varphi_1 + b\varphi_2) &\sim a(\psi \otimes \varphi_1) + b(\psi \otimes \varphi_2). \end{aligned}$$

We then define the tensor product $\mathcal{H}_A \otimes \mathcal{H}_B$ as the quotient \mathcal{F}/\sim . Notice that, taking $b = 0$ one can easily prove that $[(a\psi) \otimes \varphi] = [\psi \otimes (a\varphi)] = [a(\psi \otimes \varphi)]$. Thus, the null element in $\mathcal{H}_A \otimes \mathcal{H}_B$ is the class that contains $0 \otimes \varphi$ for any $\varphi \in \mathcal{H}_B$ and $\psi \otimes 0$ for any $\psi \in \mathcal{H}_A$. The sesquilinear product needs to be defined only on classes of *factorised vectors*, represented by $|\psi\rangle \otimes |\varphi\rangle := [\psi \otimes \varphi]$, and then just extended by sesquilinearity. We then set

$$(\langle \psi_1 | \otimes \langle \varphi_1 |)(|\psi_2\rangle \otimes |\varphi_2\rangle) := \langle \psi_1 | \psi_2 \rangle \langle \varphi_1 | \varphi_2 \rangle.$$

If we now consider two orthonormal bases $\{e_i\}_{i=1}^{d_A} \subseteq \mathcal{H}_A$ and $\{f_j\}_{j=1}^{d_B} \subseteq \mathcal{H}_B$, one can easily prove that the set $\{e_i \otimes f_j\}$ is orthonormal in $\mathcal{H}_A \otimes \mathcal{H}_B$. Moreover, it is complete by construction. Thus $\mathcal{H}_A \otimes \mathcal{H}_B \simeq \mathbb{C}^{d_A d_B}$. Moreover, by construction one has

$$|\psi\rangle \otimes |\varphi\rangle = \sum_{i=1}^{d_A} \sum_{j=1}^{d_B} \psi_i \varphi_j (|e_i\rangle \otimes |f_j\rangle) \quad (10.27)$$

10.5.1 Linear operators

The linear operators in $\mathcal{H}_A \otimes \mathcal{H}_B$ are simply the linear operators on $\mathbb{C}^{d_A d_B}$. However, due to the construction, there is a special class of linear operators in $\mathcal{H}_A \otimes \mathcal{H}_B$: *factorised operators* of the form $A \otimes B$, defined by

$$(A \otimes B)(\psi \otimes \varphi) := (A\psi) \otimes (B\varphi).$$

Notice that being $\mathcal{H}_A \otimes \mathcal{H}_B$ spanned by factorised vectors this definition is sufficient to specify the action of $A \otimes B$ on the whole $\mathcal{H}_A \otimes \mathcal{H}_B$. Given two factorised operators $A \otimes B$ and $C \otimes D$, one can then easily verify that

$$(A \otimes B)(C \otimes D) = AC \otimes BD, \quad (10.28)$$

and that the tensor product distributes over sums:

$$(A + B) \otimes C = A \otimes C + B \otimes C, \quad A \otimes (B + C) = A \otimes B + A \otimes C. \quad (10.29)$$

We can also observe that

$$|\psi_1\rangle \langle \psi_2| \otimes |\varphi_1\rangle \langle \varphi_2| = (|\psi_1\rangle \otimes |\varphi_1\rangle)(\langle \psi_2| \otimes \langle \varphi_2|).$$

Reminding the equality in (10.2), the identity operator I_{AB} on $\mathcal{H}_A \otimes \mathcal{H}_B$ is then equal to

$$\begin{aligned} I_{AB} &= \sum_{i=1}^{d_A} \sum_{j=1}^{d_B} (|e_i\rangle \otimes |f_j\rangle)(\langle e_i | \otimes \langle f_j |) \\ &= \sum_{i=1}^{d_A} \sum_{j=1}^{d_B} |e_i\rangle \langle e_i| \otimes |f_j\rangle \langle f_j| \\ &= I_A \otimes I_B. \end{aligned} \quad (10.30)$$

The definition of $A \otimes B$ can be straightforwardly generalised to the case where $A : \mathcal{H}_A \rightarrow \mathcal{H}_C$ and $B : \mathcal{H}_B \rightarrow \mathcal{H}_D$. It is immediate to verify that two operators of the form $A \otimes I$ and $I \otimes B$ with $A : \mathcal{H}_A \rightarrow \mathcal{H}_B$ and $B : \mathcal{H}_C \rightarrow \mathcal{H}_D$ commute, and satisfy the identity

$$(A \otimes I_D)(I_A \otimes B) = (I_B \otimes B)(A \otimes I_C) = A \otimes B \quad (10.31)$$

By the same procedure that we adopted in equation (10.3), we can then see that the matrix elements of $T : \mathcal{H}_A \otimes \mathcal{H}_B \rightarrow \mathcal{H}_C \otimes \mathcal{H}_D$ are defined by

$$\begin{aligned} T &= I_{CD}TI_{AB} \\ &= \sum_{i,j} \sum_{i',j'} (|e_i\rangle\langle e_i| \otimes |f_j\rangle\langle f_j|)T(|e_{i'}\rangle\langle e_{i'}| \otimes |f_{j'}\rangle\langle f_{j'}|) \\ &= \sum_{i,j} \sum_{i',j'} T_{(i,j),(i',j')} |e_i\rangle\langle e_{i'}| \otimes |f_j\rangle\langle f_{j'}|. \end{aligned}$$

Thus, it holds that

$$(A \otimes B)_{(i,j),(i',j')} = A_{i,j}B_{i',j'}.$$

As the set of operators between any pair of Hilbert spaces, also $\mathcal{L}(\mathcal{H}_A \otimes \mathcal{H}_B \rightarrow \mathcal{H}_C \otimes \mathcal{H}_D)$ is a complex vector space. In this space, factorised operators are complete. Indeed, we have the following lemma.

Lemma 10.27. *Let $T : \mathcal{H}_A \otimes \mathcal{H}_B \rightarrow \mathcal{H}_C \otimes \mathcal{H}_D$. Then there exist two finite sets $\{A_i\}_{i \in X} \subseteq \mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_C)$ and $\{B_i\}_{i \in X} \subseteq \mathcal{L}(\mathcal{H}_B \rightarrow \mathcal{H}_D)$ such that*

$$T = \sum_{i \in X} A_i \otimes B_i. \quad (10.32)$$

Proof. Indeed, one has

$$\begin{aligned} T &= \sum_{i,j} \sum_{i',j'} T_{(i,j),(i',j')} |e_i\rangle\langle e_{i'}| \otimes |f_j\rangle\langle f_{j'}| \\ &= \sum_{i,i'} |e_i\rangle\langle e_{i'}| \otimes B_{i,i'}, \\ B_{i,i'} &:= \sum_{j,j'} T_{(i,j),(i',j')} |f_j\rangle\langle f_{j'}|. \end{aligned} \quad \square$$

Before concluding this section, we introduce a relevant unitary operator on $\mathcal{H}_A \otimes \mathcal{H}_B$, called *swap*.

Definition 10.28 (Swap operator). For every tensor product $\mathcal{H}_A \otimes \mathcal{H}_B$ with $\mathcal{H}_A \simeq \mathcal{H}_B$ we define the *swap* operator

$$E(|\psi\rangle \otimes |\varphi\rangle) = |\varphi\rangle \otimes |\psi\rangle. \quad (10.33)$$

Lemma 10.29. *The swap operator E is selfadjoint.*

Proof. Let us write the adjoint equation for a factorised basis $|\psi_i\rangle \otimes |\varphi_j\rangle$

$$\begin{aligned}\langle\theta_{i,j}|(|\psi\rangle \otimes |\varphi\rangle) &= (\langle\psi_i| \otimes \langle\varphi_j|)E(|\psi\rangle \otimes |\varphi\rangle) \\ &= (\langle\psi_i| \otimes \langle\varphi_j|)(|\varphi\rangle \otimes |\psi\rangle) \\ &= \langle\psi_i|\varphi\rangle \langle\varphi_j|\psi\rangle \\ &= (\langle\varphi_j| \otimes \langle\psi_i|)(|\psi\rangle \otimes |\varphi\rangle).\end{aligned}$$

Thus, the equation is satisfied by $|\theta_{i,j}\rangle = |\varphi_j\rangle \otimes |\psi_i\rangle$, namely $E^\dagger(|\psi_i\rangle \otimes |\varphi_j\rangle) = |\varphi_j\rangle \otimes |\psi_i\rangle$, and finally $E^\dagger = E$. \square

Lemma 10.30. *The square of E is the identity; $E^2 = I_{AB}$.*

Proof. It is sufficient to prove the statement on factorised states. However, in this case the proof is trivial. \square

Corollary 10.31. *The swap operator is unitary.*

Proof. It is sufficient to collect the two previous results, getting

$$\begin{aligned}I_{AB} &= E^2 \\ &= E^\dagger E \\ &= EE^\dagger.\end{aligned}\quad \square$$

10.6 The partial trace

Notice that $\langle\psi|$ can be thought of as a special kind of operator from \mathcal{H}_A to $\mathcal{H}_I \simeq \mathbb{C}$, namely $d_I = 1$. Clearly, $\mathcal{H}_A \otimes \mathcal{H}_I = \mathcal{H}_I \otimes \mathcal{H}_A = \mathcal{H}_A$. Similarly, $|\psi\rangle$ can be thought of as a special kind of operator from \mathcal{H}_I to \mathcal{H}_A . Thus, for $T : \mathcal{H}_A \rightarrow \mathcal{H}_B$ and $\psi \in \mathcal{H}_C$, we can define the following factorised operators

$$T \otimes \langle\psi| : \mathcal{H}_A \otimes \mathcal{H}_C \rightarrow \mathcal{H}_B, \quad T \otimes \langle\psi|(|\eta\rangle \otimes |\varphi\rangle) := (\langle\psi|\varphi\rangle)T|\eta\rangle, \quad (10.34)$$

$$T \otimes |\psi\rangle : \mathcal{H}_A \rightarrow \mathcal{H}_B \otimes \mathcal{H}_C, \quad (T \otimes |\psi\rangle)|\eta\rangle := (T|\eta\rangle) \otimes |\psi\rangle. \quad (10.35)$$

We can thus define the *partial trace* as follows.

Definition 10.32 (Partial trace). Let $T : \mathcal{H}_A \otimes \mathcal{H}_B \rightarrow \mathcal{H}_A \rightarrow \mathcal{H}_C$. Then we define its *partial trace* on \mathcal{H}_A , denoted as $\text{Tr}_A[T]$, as follows

$$\text{Tr}_A[T] := \sum_{i=1}^{d_A} (\langle e_i| \otimes I_C)T(|e_i\rangle \otimes I_B). \quad (10.36)$$

Lemma 10.33. *Let $T : \mathcal{H}_A \otimes \mathcal{H}_B \rightarrow \mathcal{H}_A \otimes \mathcal{H}_C$, and write $T = \sum_i A_i \otimes B_i$. The partial trace of T on \mathcal{H}_A is given by*

$$\text{Tr}_A[T] = \sum_i \text{Tr}[A_i]B_i. \quad (10.37)$$

Proof. Applying the definition one has

$$\begin{aligned}\mathrm{Tr}_A[T] &= \sum_i \sum_{j=1}^{d_A} (\langle e_j | \otimes I_C)(A_i \otimes B_i)(|e_j\rangle \otimes I_B) \\ &= \sum_i \sum_{j=1}^{d_A} (\langle e_j | A_i | e_j \rangle) I_C B_i I_B \\ &= \sum_i \mathrm{Tr}[A_i] B_i.\end{aligned}$$

□

By the above lemma and by the independence of the trace on the basis, we have the following corollary

Corollary 10.34. *The partial trace on \mathcal{H}_A is independent of the basis $\{e_i\}_{i=1}^{d_A}$ in the defining equation (10.36).*

Moreover, using equation (10.28) and the decomposition of T in equation (10.32), one can easily prove the following corollary.

Corollary 10.35. *The following identity holds for operators T, X, Y defined on suitable Hilbert spaces*

$$\mathrm{Tr}_A[(I_A \otimes X)T(I_A \otimes Y)] = X \mathrm{Tr}_A[T]Y. \quad (10.38)$$

Lemma 10.36. *The partial traces Tr_A and Tr_B commute on operators $T : \mathcal{H}_A \otimes \mathcal{H}_B \rightarrow \mathcal{H}_A \otimes \mathcal{H}_B$, and satisfy the following identity*

$$\mathrm{Tr}_A[\mathrm{Tr}_B[T]] = \mathrm{Tr}_B[\mathrm{Tr}_A[T]] = \mathrm{Tr}[T]. \quad (10.39)$$

Proof. The thesis follows straightforwardly from the definition. Indeed, reminding equation 10.31, along with the fact that $\mathcal{H}_X \otimes \mathcal{H}_I = \mathcal{H}_I \otimes \mathcal{H}_X = \mathcal{H}_X$ and $I_I = 1$, we have

$$\begin{aligned}(I_A \otimes |\psi\rangle)|\varphi\rangle &= (|\varphi\rangle \otimes I_B)|\psi\rangle = |\varphi\rangle \otimes |\psi\rangle \\ \langle\varphi|(I_A \otimes \langle\psi|) &= \langle\psi|(\langle\varphi| \otimes I_B) = \langle\varphi| \otimes \langle\psi|.\end{aligned}$$

Thus the thesis follows. □

10.7 The Vec isomorphism: double-ket notation

In the remainder we will often simplify the notation by posing

$$|\psi\rangle|\varphi\rangle := |\psi\rangle \otimes |\varphi\rangle \quad (10.40)$$

The isomorphism between \mathcal{H} and \mathcal{H}^* and the consequent isomorphism between $\mathcal{H} \rightarrow \mathcal{K}$ and $\mathcal{K} \otimes \mathcal{H}^*$ is well known in tensor analysis, and it is extensively used under the operation of *lowering* or *raising* indices. Here we introduce it with a special formalism that is extremely convenient for calculations.

The formalism is based on the following way of denoting vectors $\Psi \in \mathcal{H}_B \otimes \mathcal{H}_A$

$$|\Psi\rangle\langle\Psi| := \sum_{i=1}^{d_A} \sum_{j=1}^{d_B} \Psi_{j,i} |f_j\rangle\langle e_i|. \quad (10.41)$$

Definition 10.37. Let $\{e_i\}_{i=1}^{d_A}$ and $\{f_j\}_{j=1}^{d_B}$ be orthonormal bases in \mathcal{H}_A and \mathcal{H}_B , respectively. Let

$$A = \sum_{i=1}^{d_A} \sum_{j=1}^{d_B} A_{j,i} |f_j\rangle\langle e_i| \quad (10.42)$$

$$|\Psi\rangle\langle\Psi| = \sum_{i=1}^{d_A} \sum_{j=1}^{d_B} \Psi_{j,i} |f_j\rangle\langle e_i| \quad (10.43)$$

Then we define the isomorphism $\text{Vec} : \mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_B) \rightarrow \mathcal{H}_B \otimes \mathcal{H}_A$ along with its inverse $\text{Vec}^{-1} : \mathcal{H}_B \otimes \mathcal{H}_A \rightarrow \mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_B)$ as follows

$$\text{Vec} : \quad A \mapsto |A\rangle\langle A| := \sum_{i=1}^{d_A} \sum_{j=1}^{d_B} A_{j,i} |f_j\rangle\langle e_i|, \quad (10.44)$$

$$\text{Vec}^{-1} : \quad |\Psi\rangle\langle\Psi| \mapsto \Psi := \sum_{i=1}^{d_A} \sum_{j=1}^{d_B} \Psi_{j,i} |f_j\rangle\langle e_i| \quad (10.45)$$

Theorem 10.38. *The following identity holds for $C : \mathcal{H}_A \rightarrow \mathcal{H}_B$, $A : \mathcal{H}_B \rightarrow \mathcal{H}_C$ and $B : \mathcal{H}_A \rightarrow \mathcal{H}_D$.*

$$A \otimes B |C\rangle\langle C| = |ACB^T\rangle\langle C| \quad (10.46)$$

Moreover, for every pair $A, B : \mathcal{H}_A \rightarrow \mathcal{H}_B$ one has

$$\text{Tr}_A[|A\rangle\langle B|] = AB^\dagger, \quad (10.47)$$

$$\text{Tr}_B[|A\rangle\langle B|] = A^T B^*, \quad (10.48)$$

$$\text{Tr}[|A\rangle\langle B|] = \langle B|A\rangle = \text{Tr}[B^\dagger A] = \text{Tr}[B^* A^T]. \quad (10.49)$$

Proof. First of all, we prove the identity of equation (10.46). Indeed, by definition we have

$$|C\rangle\rangle = \sum_{i=1}^{d_A} \sum_{j=1}^{d_B} C_{j,i} |f_j\rangle\langle e_i|.$$

Thus, if we apply $A \otimes B$ on both sides we obtain

$$A \otimes B |C\rangle\rangle = \sum_{i=1}^{d_A} \sum_{j=1}^{d_B} C_{j,i} A |f_j\rangle\otimes B |e_i\rangle.$$

Finally, applying $I_B \otimes I_D$ to the left on both sides, using the form of equation (10.30), we get

$$A \otimes B |C\rangle\rangle = \sum_{j'=1}^{d_C} \sum_{i'=1}^{d_D} \sum_{i=1}^{d_A} \sum_{j=1}^{d_B} C_{j,i} A_{j',j} B_{i',i} |f'_{j'}\rangle\langle e'_{i'}|,$$

where $\{e'_i\}_{i=1}^{d_D}$ and $\{f'_j\}_{j=1}^{d_C}$ are orthonormal bases in \mathcal{H}_D and \mathcal{H}_C , respectively, and

$$A_{j',j} := \langle f'_{j'} | A | f_j \rangle, \quad B_{i',i} := \langle e'_{i'} | B | e_i \rangle.$$

Finally, reminding that $(B^T)_{i,i'} = B_{i',i}$, we obtain

$$A \otimes B |C\rangle\rangle = \sum_{j'=1}^{d_C} \sum_{i'=1}^{d_D} (ACB^T)_{j',i'} |f'_{j'}\rangle\langle e'_{i'}|,$$

namely equation (10.46). Let us then consider the vector $|I_A\rangle\rangle \in \mathcal{H}_A \otimes \mathcal{H}_{A'}$, with $\mathcal{H}_A \simeq \mathcal{H}_{A'}$. Now, one simply needs to calculate $\text{Tr}_A [|I_A\rangle\rangle \langle I_A|]$, since using equation (10.36)

$$\begin{aligned} \text{Tr}_A [|A\rangle\rangle \langle B|] &= \text{Tr}_A [(A \otimes I_A) |I_A\rangle\rangle \langle I_A| (B^\dagger \otimes I_A)] \\ &= A \text{Tr}_A [|I_A\rangle\rangle \langle I_A|] B^\dagger, \end{aligned} \tag{10.50}$$

$$\begin{aligned} \text{Tr}_B [|A\rangle\rangle \langle B|] &= \text{Tr}_B [(I_B \otimes A^T) |I_B\rangle\rangle \langle I_B| (I_B \otimes B^*)] \\ &= A^T \text{Tr}_B [|I_B\rangle\rangle \langle I_B|] B^*. \end{aligned} \tag{10.51}$$

Now, consider that

$$\begin{aligned} |I_A\rangle\rangle \langle I_A| &= \sum_{i,i'=1}^{d_A} |e_i\rangle\langle e_i| |e_{i'}\rangle\langle e_{i'}| \\ &= \sum_{i,i'=1}^{d_A} |e_i\rangle\langle e_{i'}|_A \otimes |e_i\rangle\langle e_{i'}|_{A'}, \end{aligned}$$

and taking Tr_A we obtain

$$\begin{aligned}\text{Tr}_A[|I_A\rangle\langle I_A|] &= \sum_{i,i'=1}^{d_A} \text{Tr}[|e_i\rangle\langle e_{i'}|]|e_i\rangle\langle e_{i'}|_{A'} \\ &= \sum_{i=1}^{d_A} |e_i\rangle\langle e_i|_{A'} \\ &= I_{A'}.\end{aligned}$$

The same argument leads to $\text{Tr}_B[|I_B\rangle\langle I_B|] = I_{B'}$. Thus, equations (10.50) and (10.51) lead to

$$\text{Tr}_A[|A\rangle\langle B|] = AB^\dagger, \quad \text{Tr}_B[|A\rangle\langle B|] = A^T B^*.$$

Finally, reminding lemma 10.36 we have

$$\begin{aligned}\text{Tr}[|A\rangle\langle B|] &= \text{Tr}_B[\text{Tr}_A[|A\rangle\langle B|]] = \text{Tr}_A[\text{Tr}_B[|A\rangle\langle B|]] \\ &= \text{Tr}[AB^\dagger] = \text{Tr}[A^T B^*].\end{aligned}\quad \square$$

Using the double-ket notation it is easy to prove the so-called Schmidt decomposition of vectors in $\mathcal{H}_A \otimes \mathcal{H}_B$.

Theorem 10.39 (Schmidt decomposition). *Let $|\Psi\rangle\rangle \in \mathcal{H}_B \otimes \mathcal{H}_A$. Then one can find two orthonormal sets $\{\psi_i\}_{i=1}^{d_B}$ and $\{\varphi_j\}_{j=1}^{d_A}$ and a set of positive reals $\{\sigma_i\}_{i=1}^s$ with $s \leq \min(d_A, d_B)$ such that*

$$|\Psi\rangle\rangle = \sum_{i=1}^s \sigma_i |\psi_i\rangle |\varphi_i\rangle. \quad (10.52)$$

The integer number s is called Schmidt rank of $|\Psi\rangle\rangle$.

Proof. We use the inverse of the Vec isomorphism to study the operator $\Psi : \mathcal{H}_A \rightarrow \mathcal{H}_B$

$$\Psi = \text{Vec}^{-1}|\Psi\rangle\rangle.$$

Now, by the singular value decomposition 10.25 and its corollary 10.26 we know that we can write

$$\Psi = V\Sigma U,$$

with $V : \mathcal{H}_A \rightarrow \mathcal{H}_B$ isometric on $\text{Rng}(\Sigma)$ and $U : \mathcal{H}_A \rightarrow \mathcal{H}_A$ unitary. Thus we have

$$\begin{aligned}|\Psi\rangle\rangle &= V \otimes U^T |\Sigma\rangle\rangle \\ &= V \otimes U^T \sum_{i=1}^s \sigma_i |e_i\rangle |e_i\rangle \\ &= \sum_{i=1}^s \sigma_i |\psi_i\rangle |\varphi_i\rangle,\end{aligned}$$

where $|\psi_i\rangle = V|e_i\rangle$ and $|\varphi_j\rangle = U^T|e_j\rangle$. Since V is isometric on $\text{Span}(\{e_i\}_{i=1}^s)$ and U^T is unitary, the sets $\{\psi_i\}_{i=1}^s$ and $\{\varphi_j\}_{j=1}^s$ are orthonormal, and they can be completed to orthonormal bases in the corresponding spaces. \square

10.8 Linear maps on $\mathcal{L}(\mathcal{H})$

We remind that $\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_B)$ is a complex Hilbert space of dimension $d_A d_B$, as witnessed by the Vec isomorphism. The sesquilinear product is the Hilbert-Schmidt product $\langle\langle A|B \rangle\rangle = \text{Tr}[A^\dagger B]$. We can thus define linear maps $\mathcal{L}(\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_A) \rightarrow \mathcal{L}(\mathcal{H}_B \rightarrow \mathcal{H}_B))$. For simplicity, and for other reasons that will be clearer later, we will denote the set of such linear maps as $A \rightarrow B$. Moreover, since in $\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_A)$ one has an ordering relation $A \geq B$ defined by

$$A \geq B \Leftrightarrow A - B \geq 0. \quad (10.53)$$

The space $\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_A)$ equipped with the ordering \geq is an *ordered linear space*. Having a partial order, we can also define *order-preserving* linear maps.

Definition 10.40 (Linear map). A map $\mathcal{M} : A \rightarrow B$ is *linear* if for all $A, B : \mathcal{H}_A \rightarrow \mathcal{H}_A$ and $a, b \in \mathbb{C}$ one has

$$\mathcal{M}(aA + bB) = a\mathcal{M}(A) + b\mathcal{M}(B). \quad (10.54)$$

Being a set of linear maps on a complex vector space, also the set of maps $A \rightarrow B$ is a complex vector space.

Definition 10.41 (Identity map). The identity map $\mathcal{I}_A : A \rightarrow A$ is defined as

$$\mathcal{I}_A(A) = A, \quad \forall A : \mathcal{H}_A \rightarrow \mathcal{H}_A. \quad (10.55)$$

Since unit trace positive operators on \mathcal{H}_A are complete in $\mathcal{L}(\mathcal{H}_A \rightarrow \mathcal{H}_A)$, we have the following lemma.

Lemma 10.42. Let $\mathcal{M}, \mathcal{N} : A \rightarrow B$. Then $\mathcal{M} = \mathcal{N}$ iff

$$\mathcal{M}(\rho) = \mathcal{N}(\rho), \quad \forall \rho : \mathcal{H}_A \rightarrow \mathcal{H}_A, \rho \geq 0, \text{Tr}[\rho] = 1. \quad (10.56)$$

Definition 10.43 (Adjoint preserving). A map $\mathcal{M} : A \rightarrow B$ is *adjoint-preserving* if

$$\mathcal{M}(A^\dagger) = \mathcal{M}(A)^\dagger, \quad , \forall A : \mathcal{H}_A \rightarrow \mathcal{H}_A. \quad (10.57)$$

Lemma 10.44. The map $\mathcal{M} : A \rightarrow B$ is adjoint-preserving iff

$$\mathcal{M}(A) = \mathcal{M}(A)^\dagger, \quad \forall A : \mathcal{H}_A \rightarrow \mathcal{H}_A, A = A^\dagger. \quad (10.58)$$

Proof. Let us write $T = A + iB$. Then

$$\mathcal{M}(T) = \mathcal{M}(A) + i\mathcal{M}(B).$$

Now, if $\mathcal{M}(A) = \mathcal{M}(A)^\dagger$ for every selfadjoint A , we have

$$\begin{aligned} \mathcal{M}(T^\dagger) &= \mathcal{M}(A) - i\mathcal{M}(B) \\ &= \mathcal{M}(T)^\dagger, \end{aligned}$$

for every $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$. □

10.8.1 Positive maps

Order-preserving maps are called *positive*, because they map positive operators to positive operators. Here is the precise definition.

Definition 10.45 (Positive map). A map $\mathcal{M} : A \rightarrow B$ is *positive* if

$$\mathcal{M}(P) \geq 0, \quad \forall P \in \mathcal{P}(A). \quad (10.59)$$

10.9 Maps on tensor products

The tensor product structure, that is lifted on linear operators, is also lifted to the level of linear maps on operators. In the following we will denote by $\mathcal{M} : AB \rightarrow CD$ a map from $\mathcal{L}(\mathcal{H}_A \otimes \mathcal{H}_B)$ to $\mathcal{L}(\mathcal{H}_C \otimes \mathcal{H}_D)$. In the set of maps $\{\mathcal{M} : AB \rightarrow CD\}$ there is a special class that we define here.

Definition 10.46 (Tensor-product maps). Let $\mathcal{M} : A \rightarrow C$ and $\mathcal{N} : B \rightarrow D$. Then $\mathcal{M} \otimes \mathcal{N} : AB \rightarrow CD$ is the map defined as

$$\mathcal{M} \otimes \mathcal{N}(A \otimes B) = \mathcal{M}(A) \otimes \mathcal{N}(B), \quad \forall A : \mathcal{H}_A \rightarrow \mathcal{H}_A, \quad B : \mathcal{H}_B \rightarrow \mathcal{H}_B. \quad (10.60)$$

Remark 11. Notice that since every operator $T : \mathcal{H}_A \otimes \mathcal{H}_B \rightarrow \mathcal{H}_A \otimes \mathcal{H}_B$ can be expanded as in equation (10.32) $T = \sum_i A_i \otimes B_i$, tensor-product maps are well defined, indeed

$$\mathcal{M} \otimes \mathcal{N}(T) = \sum_i \mathcal{M}(A_i) \otimes \mathcal{N}(B_i). \quad (10.61)$$

The most important property of the tensor product of maps is its interplay with sequential composition.

Lemma 10.47. *Let $\mathcal{M} : A \rightarrow B$, $\mathcal{N} : D \rightarrow E$, $\mathcal{M}' : B \rightarrow C$, $\mathcal{N}' : E \rightarrow F$. Then the following identity holds*

$$(\mathcal{M}' \otimes \mathcal{N}')(\mathcal{M} \otimes \mathcal{N}) = (\mathcal{M}' \mathcal{M}) \otimes (\mathcal{N}' \mathcal{N}). \quad (10.62)$$

Proof. Since tensor product operators are complete in $\mathcal{L}(\mathcal{H}_A \otimes \mathcal{H}_B)$, we just need to prove that the identity holds when both members are applied to an operator of the form $A \otimes B$. In this case we have

$$\begin{aligned} [(\mathcal{M}' \otimes \mathcal{N}')(\mathcal{M} \otimes \mathcal{N})](A \otimes B) &= (\mathcal{M}' \otimes \mathcal{N}')[\mathcal{M}(A) \otimes \mathcal{N}(B)] \\ &= [(\mathcal{M}' \mathcal{M})(A)] \otimes [\mathcal{N}' \mathcal{N}(B)] \\ &= [(\mathcal{M}' \mathcal{M}) \otimes (\mathcal{N}' \mathcal{N})](A \otimes B), \end{aligned}$$

where in both equalities we used the defining equation (10.60). \square

Corollary 10.48. *Let $\mathcal{M} : A \rightarrow B$ and $\mathcal{N} : C \rightarrow D$. Then the following identity holds*

$$(\mathcal{M} \otimes \mathcal{I}_D)(\mathcal{I}_A \otimes \mathcal{N}) = (\mathcal{I}_B \otimes \mathcal{N})(\mathcal{M} \otimes \mathcal{I}_C) = \mathcal{M} \otimes \mathcal{N}. \quad (10.63)$$

Remark 12. Notice that the trace is a special linear map from $A \rightarrow I$, and thus corollary 10.48 in the special case where $\mathcal{M} = \mathcal{N} = \text{Tr}$ gives lemma 10.36 as a corollary.

Corollary 10.49. *Let $\mathcal{N} : B \rightarrow C$. Then for every $T : \mathcal{H}_A \otimes \mathcal{H}_B$ one has*

$$\text{Tr}_A[(\mathcal{I}_A \otimes \mathcal{N})(T)] = \mathcal{N}(\text{Tr}_A[T]). \quad (10.64)$$

Proof. The thesis is just the statement of corollary 10.48 with $\mathcal{M} = \text{Tr} : A \rightarrow I$. Indeed

$$(\text{Tr}_A \otimes \mathcal{I}_C)(\mathcal{I}_A \otimes \mathcal{N}) = \mathcal{N}(\text{Tr}_A \otimes \mathcal{I}_B). \quad \square$$

Finally, we prove that the tensor product distributes over linear combinations.

Lemma 10.50. *Let $\mathcal{M}, \mathcal{N} : A \rightarrow B$. Then for every $a, b \in \mathbb{C}$ and every \mathcal{H}_C , the following identities hold*

$$(a\mathcal{M} + b\mathcal{N}) \otimes \mathcal{I}_C = a\mathcal{M} \otimes \mathcal{I}_C + b\mathcal{N} \otimes \mathcal{I}_C, \quad (10.65)$$

$$\mathcal{I}_C \otimes (a\mathcal{M} + b\mathcal{N}) = a\mathcal{I}_C \otimes \mathcal{M} + b\mathcal{I}_C \otimes \mathcal{N}. \quad (10.66)$$

Proof. We prove the first identity, the proof for the second being straightforwardly analogous. It is sufficient to consider the action of the map on the l.h.s. on factorised operators $A \otimes B$. In this case

$$\begin{aligned} [(a\mathcal{M} + b\mathcal{N}) \otimes \mathcal{I}_C](A \otimes B) &= [(a\mathcal{M} + b\mathcal{N})(A)] \otimes B \\ &= [a\mathcal{M}(A) + b\mathcal{N}(A)] \otimes B \\ &= a(\mathcal{M} \otimes \mathcal{I}_C)(A \otimes B) + b(\mathcal{N} \otimes \mathcal{I}_C)(A \otimes B). \end{aligned} \quad \square$$

The desired result is obtained invoking corollary 10.48, in the form of the following corollary.

Corollary 10.51. *Let $\mathcal{M}, \mathcal{N} : A \rightarrow B$ and $\mathcal{L} : C \rightarrow D$. Then for every $a, b \in \mathbb{C}$ the following identity holds*

$$(a\mathcal{M} + b\mathcal{N}) \otimes \mathcal{L} = a\mathcal{M} \otimes \mathcal{L} + b\mathcal{N} \otimes \mathcal{L}.$$

10.10 Complete positivity

We now introduce a class of maps that have a crucial importance in quantum theory: completely positive maps. We need a preliminary definition.

Definition 10.52 (*n*-positive map). Let $\mathcal{M} : A \rightarrow B$. We say that \mathcal{M} is *n*-positive if for $\mathcal{H}_C \simeq \mathbb{C}^n$ the map $\mathcal{M} \otimes \mathcal{I}_C$ is positive.

Definition 10.53 (Completely positive map). Let $\mathcal{M} : A \rightarrow B$. We say that \mathcal{M} is *completely positive (CP)* if it is *n* positive for every $n \in \mathbb{N}$.

Remark 13. Notice that the notion of complete positivity is strictly stronger than simple positivity. There are indeed examples of maps that are positive but not completely positive. The typical example is the transpose \mathcal{T} . Indeed the *partial transpose* $\mathcal{I}_A \otimes \mathcal{T}_B$, defined for $T = \sum_i A_i \otimes B_i : \mathcal{H}_A \otimes \mathcal{H}_B \rightarrow \mathcal{H}_A \otimes \mathcal{H}_B$ as

$$\mathcal{I}_A \otimes \mathcal{T}_B(T) = \sum_i A_i \otimes B_i^T, \quad (10.67)$$

is not positive (see example 11.9).

Lemma 10.54. Let $\mathcal{M} : B \rightarrow C$ and $\mathcal{N} : A \rightarrow B$ be both CP. Then $\mathcal{M}\mathcal{N}$ is CP.

Proof. Since it is

$$\begin{aligned} [\mathcal{M}\mathcal{N} \otimes \mathcal{I}_C](P) &= (\mathcal{M} \otimes \mathcal{I}_C)(\mathcal{N} \otimes \mathcal{I}_C)(P) \\ &= (\mathcal{M} \otimes \mathcal{I}_C)[(\mathcal{N} \otimes \mathcal{I}_C)(P)], \end{aligned}$$

then clearly if $0 \leq P : \mathcal{H}_A \rightarrow \mathcal{H}_A$ is positive then

$$[\mathcal{M}\mathcal{N} \otimes \mathcal{I}_C](P) \geq 0. \quad \square$$

Lemma 10.55. Let $\mathcal{M}, \mathcal{N} : A \rightarrow B$ be CP, and $a, b \geq 0$. Then $a\mathcal{M} + b\mathcal{N}$ is CP.

Proof. Let $0 \leq P : \mathcal{H}_A \otimes \mathcal{H}_C$. By lemma 10.50 one has

$$[(a\mathcal{M} + b\mathcal{N}) \otimes \mathcal{I}_C](P) = a(\mathcal{M} \otimes \mathcal{I}_C)(P) + b(\mathcal{N} \otimes \mathcal{I}_C)(P).$$

Then, since \mathcal{M}, \mathcal{N} are CP, $\mathcal{M} \otimes \mathcal{I}_C(P) \geq 0$ and $\mathcal{N} \otimes \mathcal{I}_C(P) \geq 0$, and finally

$$[(a\mathcal{M} + b\mathcal{N}) \otimes \mathcal{I}_C](P) \geq 0. \quad \square$$

Definition 10.56. Let $X : \mathcal{H}_A \rightarrow \mathcal{H}_B$. We define $\mathcal{N}_X : A \rightarrow B$ as

$$\mathcal{N}_X(A) := XAX^\dagger. \quad (10.68)$$

Lemma 10.57. For every $\mathcal{N}_X : A \rightarrow B$ and for every \mathcal{H}_C one has $\mathcal{N}_X \otimes \mathcal{I}_C = \mathcal{N}_{X \otimes I_C}$.

Proof. It is sufficient to prove the statement for the action on factorised operators $A \otimes B : \mathcal{H}_A \otimes \mathcal{H}_C \rightarrow \mathcal{H}_A \otimes \mathcal{H}_C$. In this case one has

$$\begin{aligned}\mathcal{N}_X \otimes \mathcal{I}_C(A \otimes B) &= \mathcal{N}_X(A) \otimes B \\ &= XAX^\dagger \otimes B \\ &= (X \otimes I_C)(A \otimes B)(X^\dagger \otimes I_C) \\ &= \mathcal{N}_{X \otimes I_C}(A \otimes B).\end{aligned}\quad \square$$

Lemma 10.58. *Every map $\mathcal{N}_X : A \rightarrow B$ is positive.*

Proof. Let $P \geq 0$. By corollary 10.16 one has $P = A^\dagger A$ for some $A : \mathcal{H}_A \rightarrow \mathcal{H}_C$. Then

$$\begin{aligned}\mathcal{N}_X(P) &= XA^\dagger AX^\dagger \\ &= (AX^\dagger)^\dagger(AX^\dagger),\end{aligned}$$

which by corollary 10.16 is non-negative. \square

Corollary 10.59. *Every map $\mathcal{N}_X : A \rightarrow B$ is CP.*

Proof. Indeed, for every \mathcal{H}_C the map $\mathcal{N}_X \otimes \mathcal{I}_C$ is $\mathcal{N}_{X \otimes I_C}$ by lemma 10.57, and thus by lemma 10.58 $\mathcal{N}_X \otimes \mathcal{I}_C$ is positive. \square

Corollary 10.60. *For every collection $\{\mathcal{N}_{X_i}\}_{i \in Y}$ and for coefficients $a_i \geq 0$ the map $\sum_{i \in Y} a_i \mathcal{N}_{X_i}$ is CP.*

Proof. This is a consequence of corollary 10.59 and lemma 10.55. \square

10.11 Trace-preserving and unital maps

Definition 10.61 (Trace-preserving maps). Let $\mathcal{M} : A \rightarrow B$. We say that \mathcal{M} is *trace-preserving* if for all $T : \mathcal{H}_A \rightarrow \mathcal{H}_A$

$$\text{Tr}[\mathcal{M}(T)] = \text{Tr}[T]. \quad (10.69)$$

Definition 10.62 (Trace non-increasing maps). Let $\mathcal{M} : A \rightarrow B$. We say that \mathcal{M} is *trace non-increasing* if for all $0 \leq T : \mathcal{H}_A \rightarrow \mathcal{H}_A$

$$\text{Tr}[\mathcal{M}(T)] \leq \text{Tr}[T]. \quad (10.70)$$

Definition 10.63 (Dual map). Let $\mathcal{M} : A \rightarrow B$. We define the *dual map* \mathcal{M}^\dagger as the map satisfying the following identity for every $X : \mathcal{H}_A \rightarrow \mathcal{H}_A$ and $Y : \mathcal{H}_B \rightarrow \mathcal{H}_B$

$$\text{Tr}[Y \mathcal{M}(X)] = \text{Tr}[\mathcal{M}^\dagger(Y)X]. \quad (10.71)$$

Lemma 10.64. *The dual map \mathcal{M}^\dagger is well defined and linear.*

Proof. Let us consider the set of operators $|\varphi_i\rangle\langle\varphi_j|$, where $\{\varphi_i\}_{i=1}^{d_A}$ is an orthonormal basis in \mathcal{H}_A . Then for every $Y : \mathcal{H}_B \rightarrow \mathcal{H}_B$ we have

$$\begin{aligned}\mathcal{M}^\dagger(Y)_{j,i} &= \text{Tr}[\mathcal{M}^\dagger(Y)|\varphi_i\rangle\langle\varphi_j|] \\ &= \text{Tr}[Y\mathcal{M}(|\varphi_i\rangle\langle\varphi_j|)].\end{aligned}$$

Thus, if $Y = 0$ one has $\mathcal{M}^\dagger(Y) = 0$. Moreover the map \mathcal{M}^\dagger is linear. Indeed

$$\begin{aligned}\text{Tr}[\mathcal{M}^\dagger(aY_1 + bY_2)X] &= \text{Tr}[(aY_1 + bY_2)\mathcal{M}(X)] \\ &= a\text{Tr}[Y_1\mathcal{M}(X)] + b\text{Tr}[Y_2\mathcal{M}(X)] \\ &= a\text{Tr}[\mathcal{M}^\dagger(Y_1)X] + b\text{Tr}[\mathcal{M}^\dagger(Y_2)X] \\ &= \text{Tr}[(a\mathcal{M}^\dagger(Y_1) + b\mathcal{M}^\dagger(Y_2))X].\end{aligned}$$

Thus, $\mathcal{M}^\dagger(aY_1 + bY_2) = a\mathcal{M}^\dagger(Y_1) + b\mathcal{M}^\dagger(Y_2)$. \square

The dual map is relevant for conditions of trace-preservation. Indeed, we have the following lemmas.

Definition 10.65 (Unital map). We say that the map $\mathcal{M} : A \rightarrow B$ is *unital* if $\mathcal{M}(I_A) = I_B$.

Lemma 10.66. *The map \mathcal{M} is trace-preserving iff \mathcal{M}^\dagger is unital.*

Proof. It is clear that if $I_A = \mathcal{M}^\dagger(I_B)$ then for every X

$$\begin{aligned}\text{Tr}[X] &= \text{Tr}[I_AX] \\ &= \text{Tr}[\mathcal{M}^\dagger(I_B)X] \\ &= \text{Tr}[I_B\mathcal{M}(X)] \\ &= \text{Tr}[\mathcal{M}(X)],\end{aligned}$$

namely \mathcal{M} is trace-preserving. Viceversa, if $\text{Tr}[X] = \text{Tr}[\mathcal{M}(X)]$ for every X , then

$$\begin{aligned}\text{Tr}[\mathcal{M}^\dagger(I_B)X] &= \text{Tr}[\mathcal{M}(X)] \\ &= \text{Tr}[X] \\ &= \text{Tr}[I_AX],\end{aligned}$$

namely $\mathcal{M}^\dagger(I_B) = I_A$. \square

Remark 14. Notice that we used the property that for $A, B : \mathcal{H}_A \rightarrow \mathcal{H}_A$ one has $A = B$ iff $\text{Tr}[AX] = \text{Tr}[BX]$ for every $X : \mathcal{H}_A \rightarrow \mathcal{H}_A$. Actually, it is sufficient to consider the case of unit-trace positive X . Indeed, this case includes all operators of the form $X = |\varphi\rangle\langle\varphi|$ for every $\varphi \in \mathcal{H}_A$ with $\|\varphi\| = 1$. Thus, the condition $\text{Tr}[AX] = \text{Tr}[BX]$ for $X = |\varphi\rangle\langle\varphi|$ implies

$$\begin{aligned}\langle\psi|A|\psi\rangle &= \|\psi\|^2\langle\varphi|A|\varphi\rangle \\ &= \|\psi\|^2\text{Tr}[|\varphi\rangle\langle\varphi|A] \\ &= \|\psi\|^2\text{Tr}[|\varphi\rangle\langle\varphi|B] \\ &= \|\psi\|^2\langle\varphi|B|\varphi\rangle \\ &= \langle\psi|B|\psi\rangle \quad \forall \psi \in \mathcal{H}_A,\end{aligned}\tag{10.72}$$

where $|\psi\rangle\langle\psi| = \|\psi\|^2|\varphi\rangle\langle\varphi|$ for $\|\varphi\| = 1$. Then, by the polarisation identity, equation 10.72 this implies $A = B$. Analogously, one can easily prove that if $\text{Tr}[X\rho] \geq 0$ for every $\rho \geq 0$, then $X \geq 0$. Indeed, all the operators $\rho = |\varphi\rangle\langle\varphi|$ are positive semidefinite.

Finally, we have the following lemma

Lemma 10.67. *The map $\mathcal{M} : A \rightarrow B$ is trace non-increasing iff $\mathcal{M}^\dagger(I_B) \leq I_A$.*

Proof. The map is trace non-increasing iff for every $P \in \mathcal{P}(A)$

$$\begin{aligned}\text{Tr}[\mathcal{M}(P)] &= \text{Tr}[I_B \mathcal{M}(P)] \\ &= \text{Tr}[\mathcal{M}^\dagger(I_B)P] \\ &\leq \text{Tr}[P] \\ &= \text{Tr}[I_A P],\end{aligned}$$

which is equivalent to $\text{Tr}[(I_A - \mathcal{M}^\dagger(I_B))P] \geq 0$ for all $P \geq 0$, and then to

$$I_A - \mathcal{M}^\dagger(I_B) \geq 0.$$

□

Chapter 11

Lecture 11: Quantum Theory and the Choi isomorphism

11.1 Quantum theory

Analogously to classical information theory, quantum theory is a theory dealing with *quantum* information carrying systems, that can be thought of as quantum n -levels systems. A quantum system A is in correspondence with a complex Hilbert space \mathcal{H}_A , having d_A different levels, but also allowing for the storage of superpositions, and for the measurement of complementary quantities. The *type* of a quantum system A is the dimension d_A of the corresponding space. Also in the quantum case, for the purpose of information encoding, all systems of the same type are equivalent.

11.2 States

The encoding of information in a system or an array of systems corresponds to the preparation of a special *state* of the system. For a given system A, states are represented by positive operators ρ with $\text{Tr}[\rho] \leq 1$. The preparation may occur according to a probabilistic algorithm, and $\text{Tr}[\rho]$ represents the preparation probability. The set of states of A is denoted by $\text{St}(A)$. A state ρ is *deterministic* if and only if it is *normalised*, namely $\text{Tr}[\rho] = 1$, and the set of deterministic states is denoted by $\text{St}_1(A)$. In the general case, a preparation of system A corresponds to an *ensemble*, that is a collection of states $\{\rho_i\}_{i \in X}$ such that $\sum_{i \in X} \rho_i$ is deterministic. The ensemble represents a preparation algorithm where the occurrence of a particular classical outcome $i \in X$ —that can be a digit on a display, or any other classical signal from a device—heralds preparation of the specific state ρ_i in the ensemble. A state ρ of system A is represented by the following diagram

$$\boxed{\rho} \xrightarrow{\quad} A \tag{11.1}$$

Composing two systems A of type d_A and B of type d_B *in parallel* allows one to encode states of a new system C whose type is $d_A d_B$ and $\mathcal{H}_{AB} = \mathcal{H}_A \otimes \mathcal{H}_B$. The parallel run of preparation algorithms ρ and σ for A and B, respectively, is represented by the state $\rho \otimes \sigma$.

Notice that both the set of states $\text{St}(A)$ and the set of deterministic states $\text{St}_1(A)$ are convex. Let us consider in particular $\text{St}_1(A)$. It is clear from the spectral theorem that every element $\rho \in \text{St}_1(A)$ can be decomposed as a convex combination of elements of the form $\rho_i = |\psi_i\rangle\langle\psi_i|$, namely

$$\rho = \sum_{i=1}^{\text{rank}(\rho)} p_i |\psi_i\rangle\langle\psi_i|.$$

In mathematical terms, the elements $|\psi_i\rangle\langle\psi_i|$ correspond to *extremal elements* of the convex set $\text{St}_1(A)$. We call these states *pure states* of system A. The reason is that they correspond to preparations where we have maximal control of the system. Every other preparation can be perfectly simulated via a random preparation algorithm.

What is the relation of the above formalism with the traditional representation of quantum mechanical systems? The connection is through pure states. In the traditional terminology, the vector $|\psi\rangle$ is used to denote the state $\rho_\psi = |\psi\rangle\langle\psi|$. Differently from ρ_ψ , however, $|\psi\rangle$ carries irrelevant information, e.g. any overall phase factor. This makes the correspondence $\psi \mapsto \rho_\psi$ not invertible. On the other hand, the full information about preparation of system A can be summarised in the specification of a state $\rho \in \text{St}_1(A)$.

11.2.1 Quantum operations and channels

As in the classical case, also transformations have a type $A \rightarrow B$, that accounts for the type A of the input system and the type B of the output.

A transformation of type $A \rightarrow B$ maps states of system A into states of system B. The diagrammatic representation of a transformation \mathcal{C} of type $A \rightarrow B$ is the following

$$\xrightarrow{A} \boxed{\mathcal{C}} \xrightarrow{B}, \quad \xrightarrow{\rho'} \xrightarrow{B} = \xrightarrow{\rho} \xrightarrow{A} \boxed{\mathcal{C}} \xrightarrow{B} .$$

Also in the quantum case, the transformation algorithm cannot behave differently depending on the probability distribution $\mathbb{P}_X(i) := \text{Tr}[\rho_i]$ of an input ensemble $\{\rho_i\}_{i \in X}$. We already know from the classical case that the requirement on the transformation \mathcal{C} , of behaving in a fixed way independently of the probability distribution of the ensemble, is mathematically expressed through the property of *linearity*: The state ρ' of B is a linear function of the state ρ of A. Then, the transformation \mathcal{C} is a linear map $\mathcal{C} : A \rightarrow B$.

In order to respect positivity, also when applied in parallel with the identical channel \mathcal{I}_C on an arbitrary independent system C, as in the following diagram

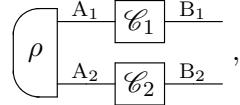
$$\xrightarrow{\rho} \xrightarrow{A} \boxed{\mathcal{C}} \xrightarrow{B} \xrightarrow{C}$$

the linear transformation must preserve the positivity of joint states $\text{St}(AC)$. Thus, a transformation $\mathcal{C} : A \rightarrow B$ must be a linear CP map.

Finally, in order to preserve (sub-)normalisation the map $\mathcal{C} : A \rightarrow B$ must be trace non-increasing. Linear CP trace non-increasing maps are called *quantum operations*, and their set will be denoted by $\text{QO}(A \rightarrow B)$. Deterministic quantum operations are called *channels*, and they are trace-preserving. The set of quantum channels is denoted by

$\text{QC}(A \rightarrow B)$. A test in quantum theory is a collection $\{\mathcal{A}_i\}_{i \in X}$ of quantum operations summing to a channel. In quantum theory a test is called a *quantum instrument*, and represents a collection of possible events in a measurement, with outcome $i \in X$ heralding the occurrence of the event corresponding to the quantum operation \mathcal{A}_i .

The parallel action of a channel (or operation) \mathcal{C}_1 from A_1 to B_1 and of a channel (or operation) \mathcal{C}_2 from A_2 to B_2 , as in the following diagram



is represented by the tensor product channel (or operation)

$$\mathcal{D} = \mathcal{C}_1 \otimes \mathcal{C}_2.$$

11.2.2 Effects

One can think of states $\text{St}(A)$ as special quantum operations $\text{QO}(I \rightarrow A)$, where the system I has $d_I = 1$. The system I is such that $AI = IA = A$. Thanks to the system I we can define another special type of operations, $\text{Eff}(A) := \text{QO}(A \rightarrow I)$. These operations are called *effects*, and map states into probabilities given by the *Born rule*

$$p(P, \rho) = \text{Tr}[\rho P]. \quad (11.2)$$

In circuit notation, we will write

$$p(P, \rho) = \text{Tr}[\rho P] = (\rho \square^A P). \quad (11.3)$$

There is clearly a unique deterministic effect P such that $\text{Tr}[P\rho] = 1$ for all $\rho \in \text{St}_1(A)$, namely I_A . A complete instrument in this case is thus represented by a *POVM*, i.e. a collection of positive operators $\{P_i\}$ such that $\sum_i P_i = I_A$.

11.3 The qubit

A system type of particular relevance is the qubit, i.e. the type of systems whose dimension is $d = 2$. In this case $\mathcal{H} = \mathbb{C}^2$. By the *Vec* isomorphism, the space $\mathcal{L}(\mathcal{H} \rightarrow \mathcal{H})$ of linear operators on \mathcal{H} is isomorphic to $\mathcal{H} \otimes \mathcal{H}$, and one can then expand such operators on any set that corresponds to an orthonormal basis in $\mathcal{H} \otimes \mathcal{H}$. A basis that is particularly suited for the expansion of operators on \mathcal{H} is represented by the identity I along with the three Pauli sigma operators

$$\sigma_1 = \sigma_x = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad \sigma_2 = \sigma_y = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}, \quad \sigma_3 = \sigma_z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}.$$

Indeed, defining $\sigma_0 := I$, since, for $i = 1, 2, 3$, $\sigma_i^\dagger = \sigma_i$ with $\text{Tr}[\sigma_i] = 0$ and $\sigma_j \sigma_k = \delta_{jk} I + i \varepsilon_{jkl} \sigma_l$, one has that $\{(1/\sqrt{2})|\sigma_i\rangle\}_{i=0}^3$ is an orthonormal basis in $\mathcal{H} \otimes \mathcal{H}$. Thus, for every operator T we can write

$$T = t_0 I + \mathbf{t} \cdot \boldsymbol{\sigma}. \quad (11.4)$$

Moreover, T is selfadjoint if and only if $t_0, t_1, t_2, t_3 \in \mathbb{R}$. One can easily prove that, for $t_k \in \mathbb{R}$, $k = 1, 2, 3$, the spectrum of $\mathbf{t} \cdot \boldsymbol{\sigma}$ is $\{-\|\mathbf{t}\|, \|\mathbf{t}\|\}$. As a consequence, the spectrum of a selfadjoint operator expressed as in Eq. (11.4) is $\{t_0 + \|\mathbf{t}\|, t_0 - \|\mathbf{t}\|\}$. It is now easy to conclude that a necessary and sufficient condition for T to be positive is $t_0 \geq \|\mathbf{t}\|$. Moreover, $\text{Tr}[T] = 2t_0$. Thus, the set of deterministic states of a qubit system is of the form

$$T = \frac{1}{2}(I + \mathbf{n} \cdot \boldsymbol{\sigma}), \quad \|\mathbf{n}\| \leq 1. \quad (11.5)$$

In other words, the set of deterministic states is a 3-dimensional sphere, with pure states corresponding to vectors on the boundary ($\|\mathbf{n}\| = 1$), and the center represented by the maximally mixed state $I/2$. Such a representation of the set of deterministic states is known as the *Bloch sphere* representation.

11.4 Purification

The most important feature of quantum theory is a property of its states called *purification*. Indeed, purification can be used as an axiom to derive the Hilbert space structure, along with other five axioms that are shared by classical information theory. In this respect, one can say that purification is not only one of the properties of quantum theory, but it really captures its essence.

Definition 11.1 (Purification of a state). Let $\rho \in \text{St}(A)$. We say that $R = |\Psi\rangle\langle\Psi| \in \text{St}(AB)$ is a purification of ρ if ρ is the marginal of R on system A. In formula

$$\rho = \text{Tr}_B[R] \quad (11.6)$$

The property of purification is expressed by the following theorem.

Theorem 11.2 (Purification). *Every state $\rho \in \text{St}(A)$ has a purification $R \in \text{St}(AB)$ with $d_B = \text{rank}(\rho)$. If $R' \in \text{St}(AC)$ is another purification of ρ , then there is an isometry $V : \mathcal{H}_B \rightarrow \mathcal{H}_C$ such that*

$$R' = (\mathcal{J}_A \otimes \mathcal{N}_V)(R). \quad (11.7)$$

Proof. First of all, let $|\Psi\rangle\langle\Psi| \in \text{St}(AB)$ be a purification of ρ . By theorem 10.38, one has

$$\rho = \text{Tr}_B[|\Psi\rangle\langle\Psi|] = \Psi\Psi^\dagger. \quad (11.8)$$

We observe that since $\text{rank}(AB) \leq \min\{\text{rank}(A), \text{rank}(B)\}$, it must be $d_B \geq \text{rank}(\Psi) \geq \text{rank}(\rho)$. Existence of a *minimal* purification—that is a purification $R \in \text{St}(AB)$ with $d_B = \text{rank}(\rho)$ —is proved by exhibiting R . For this purpose, let us diagonalise ρ as

$$\rho = \sum_{i=1}^{\text{rank}(\rho)} p_i |\psi_i\rangle\langle\psi_i|.$$

By definition $p_i > 0$ and $\sum_{i=1}^{\text{rank}(\rho)} p_i \leq 1$, with equality holding iff $\text{Tr}[\rho] = 1$. We then define $R := |\Psi\rangle\langle\Psi|$, with

$$|\Psi\rangle := \sum_{i=1}^{\text{rank}(\rho)} \sqrt{p_i} |\psi_i\rangle |\varphi_i\rangle, \quad (11.9)$$

where $\{\varphi_i\}_{i=1}^{\text{rank}(\rho)}$ is an orthonormal complete set in $\mathcal{H}_B \simeq \mathbb{C}^{\text{rank}(\rho)}$. It is easy to check that $\text{Tr}_B[R] = \rho$. Notice that by the Vec isomorphism (with respect to the bases $\{\psi_i\}_{i=1}^{d_A}$ and $\{\varphi_i\}_{i=1}^{d_B}$) it is

$$\Psi = \sum_{i=1}^{\text{rank}(\rho)} \sqrt{p_i} |\psi_i\rangle \langle \varphi_i|,$$

and thus $\text{rank}(\Psi) = \text{rank}(\rho) = d_B$. In order to prove the second statement, i.e. that there is an isometry $V : \mathcal{H}_B \rightarrow \mathcal{H}_C$ such that $R' = (\mathcal{J}_A \otimes \mathcal{N}_V)(R)$, we start observing that

$$|\Psi^\dagger| = \sum_{i=1}^{\text{rank}(\rho)} \sqrt{p_i} |\psi_i\rangle \langle \psi_i|, \quad |\Psi| = \sum_{i=1}^{d_B} \sqrt{p_i} |\varphi_i\rangle \langle \varphi_i|.$$

namely $|\Psi^\dagger|$ and $|\Psi|$ have the same spectrum (neglecting the eigenvalue 0), with the same degeneracy of non-null eigenvalues. Let us now suppose that $R' = |\Phi\rangle\langle\Phi| \in \text{St}(AC)$ is a purification of ρ . Then we have

$$|\Psi^\dagger|^2 = \Psi \Psi^\dagger = \rho = \text{Tr}_C[R'] = \Phi \Phi^\dagger = |\Phi^\dagger|^2.$$

Also $|\Phi^\dagger|$ has then the same spectrum as $|\Psi|$, and by the polar decomposition theorem 10.23 we can write

$$\Phi^\dagger = U |\Phi^\dagger| = \sum_{i=1}^{d_B} \sqrt{p_i} |\eta_i\rangle \langle \psi_i|,$$

with the (co-)isometry U such that $U|\psi_i\rangle = |\eta_i\rangle$. Now, we have

$$\Phi = |\Phi^\dagger| U^\dagger = \sum_{i=1}^{d_B} \sqrt{p_i} |\psi_i\rangle \langle \eta_i|. \quad (11.10)$$

If we define the isometry $V^* : \mathcal{H}_B \rightarrow \mathcal{H}_C$ such that $V^*|\varphi_i\rangle = |\eta_i\rangle$, we can write

$$|\Phi\rangle\langle\Phi| = |\Psi V^T\rangle\langle\Psi| = (I \otimes V)|\Psi\rangle\langle\Psi|.$$

Finally, we have

$$R' = (I \otimes V) R (I \otimes V^\dagger) = (\mathcal{J}_A \otimes \mathcal{N}_V)(R). \quad \square$$

Theorem 11.3 (Essential uniqueness of purification). *Let $R_0, R_1 \in \text{St}(AC)$ be two purifications of ρ . Then there is a unitary $U : \mathcal{H}_C \rightarrow \mathcal{H}_C$ such that*

$$R_1 = (\mathcal{J}_A \otimes \mathcal{N}_U)(R_0). \quad (11.11)$$

Proof. By the proof of theorem 11.3, we know that R_0 and R_1 must be of the form

$$R_i = |\Phi_i\rangle\langle\Phi_i|, \quad , i = 0, 1,$$

with

$$\Phi_i = \sum_{j=1}^{\text{rank}(\rho)} \sqrt{p_j} |\psi_j\rangle\langle\eta_j^{(i)}|.$$

One can define a unitary U^* such that $U^*|\eta_j^{(0)}\rangle = |\eta_j^{(1)}\rangle$, and then

$$|\Phi_1\rangle\langle\Phi_1| = I \otimes U|\Phi_0\rangle\langle\Phi_0|.$$

Finally, this implies that

$$R_1 = (\mathcal{I} \otimes \mathcal{N}_U)(R_0). \quad \square$$

11.5 The Choi correspondence

We now introduce the Choi correspondence, which allows us to treat the complete positivity constraint in a very convenient way. The correspondence is actually a general correspondence between linear maps of type $A \rightarrow B$ (we remind that this means $\mathcal{L}(\mathcal{H}_A) \rightarrow \mathcal{L}(\mathcal{H}_B)$) and linear operators on $\mathcal{H}_B \otimes \mathcal{H}_A$ (namely elements of $\mathcal{L}(\mathcal{H}_B \otimes \mathcal{H}_A)$).

Definition 11.4 (Choi correspondence). Let $\mathcal{M} : A \rightarrow B$ be a linear map. Its Choi representative (also called Choi state, matrix, operator) is

$$\mathbf{Ch}(\mathcal{M}) := \mathcal{M} \otimes \mathcal{I}_A(|I_A\rangle\langle I_A|) \in \mathcal{L}(\mathcal{H}_B \otimes \mathcal{H}_A). \quad (11.12)$$

Notice that according to the above definition one has

$$\mathbf{Ch}(\mathcal{N}_X) = |X\rangle\langle X|, \quad (11.13)$$

for every X , as one can straightforwardly verify using the properties of the \mathbf{Vec} map. The first property that we show is that the Choi correspondence is linear and invertible.

Theorem 11.5. *The correspondence $\mathbf{Ch} : \mathcal{M} \mapsto \mathbf{Ch}(\mathcal{M})$ is linear and invertible. The inverse map $\mathbf{Ch}^{-1} : \mathcal{L}(\mathcal{H}_B \otimes \mathcal{H}_A) \rightarrow \mathcal{L}[\mathcal{L}(\mathcal{H}_A) \rightarrow \mathcal{L}(\mathcal{H}_B)]$ is defined through the identity*

$$\mathbf{Ch}^{-1}[R](X) = \text{Tr}_A[(I_B \otimes X^T)R] \in \mathcal{L}(\mathcal{H}_B). \quad (11.14)$$

Proof. Let us first prove linearity. Consider the map $\alpha\mathcal{M} + \beta\mathcal{N}$ with $\alpha, \beta \in \mathbb{C}$ and $\mathcal{M}, \mathcal{N} : A \rightarrow B$. Then by definition

$$\begin{aligned} \mathbf{Ch}(\alpha\mathcal{M} + \beta\mathcal{N}) &= (\alpha\mathcal{M} + \beta\mathcal{N}) \otimes \mathcal{I}_A(|I_A\rangle\langle I_A|) \\ &= \alpha\mathcal{M} \otimes \mathcal{I}_A(|I_A\rangle\langle I_A|) + \beta\mathcal{N} \otimes \mathcal{I}_A(|I_A\rangle\langle I_A|) \\ &= \alpha\mathbf{Ch}(\mathcal{M}) + \beta\mathbf{Ch}(\mathcal{N}). \end{aligned}$$

Let us now consider the expression in equation (11.14). First of all, this defines a linear correspondence, since

$$\begin{aligned}\text{Ch}^{-1}[\alpha R + \beta S](X) &= \text{Tr}_A[(I_B \otimes X^T)(\alpha R + \beta S)] \\ &= \alpha \text{Tr}_A[(I_B \otimes X^T)R] + \beta \text{Tr}_A[(I_B \otimes X^T)S] \\ &= \alpha \text{Ch}^{-1}(R)(X) + \beta \text{Ch}^{-1}(S)(X).\end{aligned}$$

Let us now calculate $\text{Ch}^{-1}[\text{Ch}(\mathcal{M})](X)$ explicitly.

$$\begin{aligned}\text{Ch}^{-1}[\text{Ch}(\mathcal{M})](X) &= \text{Tr}_A[(I_B \otimes X^T)\text{Ch}(\mathcal{M})] \\ &= \text{Tr}_A[(\mathcal{I}_B \otimes \mathcal{L}_{X^T})\text{Ch}(\mathcal{M})] \\ &= \text{Tr}_A[(\mathcal{M} \otimes \mathcal{I}_A)(\mathcal{I}_A \otimes \mathcal{L}_{X^T})(|I_A\rangle\langle I_A|)] \\ &= \text{Tr}_A[(\mathcal{M} \otimes \mathcal{I}_A)(I \otimes X^T|I_A\rangle\langle I_A|)] \\ &= \mathcal{M}(\text{Tr}_A[|X\rangle\langle I|]) \\ &= \mathcal{M}(X),\end{aligned}$$

where for every $X : \mathcal{H}_A \rightarrow \mathcal{H}_A$ the linear map \mathcal{L}_X acts as $\mathcal{L}_X(Y) := XY$. Thus, we can conclude that $\text{Ch}^{-1}[\text{Ch}(\mathcal{M})](X) = \mathcal{M}(X)$ for every $X \in \mathcal{L}(\mathcal{H}_A)$, namely

$$\text{Ch}^{-1}[\text{Ch}(\mathcal{M})] = \mathcal{M}. \quad (11.15)$$

Moreover, by equation 11.14 if $\text{Ch}^{-1}(R) = 0$ one has $R = 0$, which implies invertibility of the map Ch^{-1} . Thus, since by equation 11.15 one has

$$\text{Ch}^{-1}\{\text{Ch}[\text{Ch}^{-1}(R)]\} = \text{Ch}^{-1}(R),$$

we must conclude that $\text{Ch}[\text{Ch}^{-1}(R)] = R$. \square

We can now prove the most important property of the Choi correspondence, which is given by the next theorem.

Theorem 11.6. *Let $\mathcal{M} : A \rightarrow B$ be a linear map. Then \mathcal{M} is CP iff*

$$\text{Ch}(\mathcal{M}) \geq 0. \quad (11.16)$$

Proof. We first prove that equation (11.16) is necessary for \mathcal{M} to be CP. Indeed, since the operator $|I_A\rangle\langle I_A|$ is non-negative definite, if \mathcal{M} is CP one has

$$\text{Ch}(\mathcal{M}) = \mathcal{M} \otimes \mathcal{I}_A(|I_A\rangle\langle I_A|) \geq 0.$$

On the other hand, the condition is sufficient. Indeed, let $\text{Ch}(\mathcal{M}) \geq 0$. Then for every $|\Psi\rangle \in \mathcal{H}_A \otimes \mathcal{H}_C$ we have

$$\begin{aligned}\mathcal{M} \otimes \mathcal{I}_C(|\Psi\rangle\langle\Psi|) &= (\mathcal{M} \otimes \mathcal{I}_C)(\mathcal{I}_A \otimes \mathcal{N}_{\Psi^T})(|I_A\rangle\langle I_A|) \\ &= (\mathcal{I}_B \otimes \mathcal{N}_{\Psi^T})(\mathcal{M} \otimes \mathcal{I}_A)(|I_A\rangle\langle I_A|) \\ &= (\mathcal{I}_B \otimes \mathcal{N}_{\Psi^T})\text{Ch}(\mathcal{M}).\end{aligned}$$

Finally, since \mathcal{N}_X is CP for every X and $\text{Ch}(\mathcal{M})$ is non-negative definite by hypothesis, we have

$$\mathcal{M} \otimes \mathcal{I}_C(|\Psi\rangle\langle\Psi|) \geq 0,$$

for any C and $|\Psi\rangle\langle\Psi| \in \mathcal{H}_A \otimes \mathcal{H}_C$. Finally, since every positive operator $\rho \in \mathcal{L}(\mathcal{H}_A \otimes \mathcal{H}_C)$ can be diagonalised and has positive eigenvalues, $\rho = \sum_{i=1}^{d_A d_C} p_i |\Psi_i\rangle\langle\Psi_i|$, the result can be straightforwardly extended from rank-one positive operators $|\Psi\rangle\langle\Psi|$ to general positive operators $\rho \geq 0$. Thus, $\mathcal{M} \otimes \mathcal{I}_C$ is positive for every C , namely \mathcal{M} is CP. \square

This result bears two important consequences, given by the following two corollaries.

Corollary 11.7. *Let $\mathcal{M} : A \rightarrow B$. Then \mathcal{M} is CP iff it is d_A -positive.*

The second one is a celebrated theorem known as Kraus theorem, leading to the so-called *operator-sum representation* of CP maps.

Corollary 11.8 (Kraus' theorem). *Let $\mathcal{M} : A \rightarrow B$. Then \mathcal{M} is CP iff there exists a set of operators $\{K_i\}_{i=1}^{d_A d_B}$ such that*

$$\mathcal{M}(X) = \sum_{i=1}^{d_A d_B} K_i X K_i^\dagger. \quad (11.17)$$

The operators K_i are called *Kraus operators* of \mathcal{M} , and equation 11.17 is called *Kraus form* or *operator-sum representation* of \mathcal{M} .

Proof. A map \mathcal{M} defined as in equation (11.17) is of the form

$$\mathcal{M} = \sum_{i=1}^{d_A d_B} \mathcal{N}_{K_i},$$

and since \mathcal{N}_{K_i} is CP for any K_i , and $\alpha \mathcal{N}_1 + \beta \mathcal{N}_2$ —with $\alpha, \beta \geq 0$ and \mathcal{N}_i CP—is itself CP, we conclude that such a map \mathcal{M} is CP. On the other hand, let $\mathcal{M} : A \rightarrow B$ be CP. By theorem 11.6 $\text{Ch}(\mathcal{M}) \geq 0$, and then $\text{Ch}(\mathcal{M}) \in \mathcal{L}(\mathcal{H}_B \otimes \mathcal{H}_A)$ is diagonalisable with non-negative eigenvalues. Then we have

$$\begin{aligned} \text{Ch}(\mathcal{M}) &= \sum_{i=1}^{d_A d_B} q_i |H_i\rangle\langle H_i| \\ &= \sum_{i=1}^{d_A d_B} |K_i\rangle\langle K_i|, \end{aligned}$$

where $\langle\langle H_i | H_j \rangle\rangle = \delta_{i,j}$, while $K_i := \sqrt{q_i} H_i$. Thus,

$$\text{Ch}(\mathcal{M}) = \sum_{i=1}^{d_A d_B} \text{Ch}(\mathcal{N}_{K_i}),$$

and by linearity of \mathbf{Ch}^{-1} we conclude

$$\mathcal{M} = \sum_{i=1}^{d_A d_B} \mathcal{N}_{K_i}. \quad \square$$

As a consequence of theorem 11.6, one can easily prove that the transpose map, which is positive, is not completely positive.

Example 11.9. Let the transpose map $\mathcal{T} : A \rightarrow A$ be defined by $\mathcal{T}(X) = X^T$ for every $X \in \mathcal{L}(\mathcal{H}_A)$. The positive map \mathcal{T} is not CP. Indeed, its Choi representative is

$$\begin{aligned} \mathbf{Ch}(\mathcal{T}) &= \mathcal{T} \otimes \mathcal{I}_A(|I_A\rangle\langle I_A|) \\ &= \sum_{m,n=1}^{d_A} |e_m\rangle\langle e_n| \otimes |e_n\rangle\langle e_m| = E, \end{aligned}$$

where E is the swap operator $E(|\varphi\rangle \otimes |\psi\rangle) = |\psi\rangle \otimes |\varphi\rangle$. Clearly, the swap operator E has negative eigenvalue -1 , as one can immediately see applying E to an antisymmetric vector $|\Psi\rangle\rangle = (|\varphi\rangle \otimes |\psi\rangle - |\psi\rangle \otimes |\varphi\rangle)$. One indeed obtains

$$E|\Psi\rangle\rangle = -|\Psi\rangle\rangle.$$

This proves that $\mathbf{Ch}(\mathcal{T}) \not\succeq 0$.

We can now move on to the study of the second relevant property of quantum channels, namely the property of being trace-preserving. For this purpose, it is useful to prove the following consequence of the Kraus theorem 11.8.

Corollary 11.10. *Let $\mathcal{M} : A \rightarrow B$ be CP, with Kraus form*

$$\mathcal{M}(X) = \sum_{i=1}^{d_A d_B} K_i X K_i^\dagger.$$

then the dual map \mathcal{M}^\dagger is CP and has Kraus form

$$\mathcal{M}^\dagger(Y) = \sum_{i=1}^{d_A d_B} K_i^\dagger Y K_i. \quad (11.18)$$

We can now prove the following lemma.

Lemma 11.11. *Let $\mathcal{M} : A \rightarrow B$ be CP. Then \mathcal{M} is trace-preserving iff*

$$\mathrm{Tr}_B[\mathbf{Ch}(\mathcal{M})] = I_A. \quad (11.19)$$

Proof. We proved that a necessary and sufficient condition for \mathcal{M} to be TP is $\mathcal{M}^\dagger(I_B) = I_A$. Now, for a map \mathcal{M} of the form of equation (11.17), by corollary 11.10 one has

$$\mathcal{M}^\dagger(Y) = \sum_i K_i^\dagger Y K_i.$$

This implies $\mathcal{M}^\dagger(I_B) = \sum_i K_i^\dagger K_i$. In other words, one has

$$\mathcal{M}^\dagger(I_B) = I_A \Leftrightarrow \sum_i K_i^\dagger K_i = I_A,$$

and finally, considering that $\sum_i K_i^\dagger K_i = \sum_i \text{Tr}_B[\mathbf{Ch}(\mathcal{N}_{K_i})]^T = \text{Tr}_B[\mathbf{Ch}(\mathcal{M})]^T$, we conclude that

$$\mathcal{M}^\dagger(I_B) = I_A \Leftrightarrow \text{Tr}_B[\mathbf{Ch}(\mathcal{M})] = I_A. \quad \square$$

This result has as an immediate consequence a characterisation of Choi representatives of quantum operations, given by the following corollary.

Corollary 11.12. *Let $\mathcal{M} : A \rightarrow B$ be CP. Then $\mathcal{M} \in \text{QO}(A \rightarrow B)$ iff*

$$\text{Tr}_B[\mathbf{Ch}(\mathcal{M})] \leq I_A. \quad (11.20)$$

Proof. As we showed in the proof of lemma 11.11, one has

$$\mathcal{M}^\dagger(I_B) = \sum_i K_i^\dagger K_i.$$

Now, a CP map \mathcal{M} is trace-non-increasing iff

$$\text{Tr}[\mathcal{M}(\rho)] = \text{Tr}[\mathcal{M}^\dagger(I_B)\rho] \leq \text{Tr}[\rho], \quad \forall \rho \geq 0.$$

This is equivalent to

$$\text{Tr}[(I_A - \mathcal{M}^\dagger(I_B))\rho] \geq 0 \quad \forall \rho \geq 0,$$

and this clearly implies

$$\begin{aligned} I_A &\geq \mathcal{M}^\dagger(I_B) \\ &= \sum_i K_i^\dagger K_i \\ &= \text{Tr}_B[\mathbf{Ch}(\mathcal{M})]^T. \end{aligned} \quad \square$$

Corollary 11.13. *Let $\mathcal{N}_V : A \rightarrow B$ be a quantum channel. Then $V : \mathcal{H}_A \rightarrow \mathcal{H}_B$ is an isometry.*

Proof. The map \mathcal{N}_V is CP. Now, it is TP iff

$$\mathcal{N}_V^\dagger(I_B) = \mathcal{N}_{V^\dagger}(I_B) = V^\dagger V = I_A. \quad \square$$

11.5.1 Effects

We now consider the particular case of quantum operations $\text{QO}(A \rightarrow I)$, the effects. Since $\mathcal{H}_I = \mathbb{C}$, we have that the Choi representative of an effect is a positive operator $\text{Ch}(a) = P \geq 0$, with the property that

$$\text{Tr}_I[P] = P \leq I_A.$$

Notice that application of the Choi map in the case of effects would require to define $\text{Ch}(a) = P^T$, since

$$\begin{aligned} \text{Tr}_A[\rho^T \text{Ch}(a)] &= \text{Tr}[\rho^T P] \\ &= \text{Tr}[\rho P^T]. \end{aligned}$$

However, we will stick to the notation of the Born rule (11.2), and omit the transpose from the effect.

11.6 Unitary dilation of channels

We now prove one of the most important theorems of quantum information theory, that is a consequence of purification.

Theorem 11.14. *Let $\mathcal{M} \in \text{QC}(A \rightarrow B)$. Then there exist*

- Two systems C and D with $\mathcal{H}_A \otimes \mathcal{H}_C \simeq \mathcal{H}_B \otimes \mathcal{H}_D$;
- A unitary operator $U \in \mathcal{L}(\mathcal{H}_A \otimes \mathcal{H}_C)$;
- A pure state $\eta \in \text{St}_1(C)$;

such that

$$\boxed{\mathcal{M}}_{AB} = \boxed{\eta}_{C} \boxed{N_U}_{AB} \boxed{I}_{CD} . \quad (11.21)$$

In formula

$$\mathcal{M}(\rho) = \text{Tr}_D[U(\rho_A \otimes \eta_C)U^\dagger]. \quad (11.22)$$

Proof. Let us denote by $\Omega \in \text{St}_1(AA')$ the pure state $\frac{1}{d_A}|I_A\rangle\langle I_A|$. Then we have $\frac{1}{d_A}\text{Ch}(\mathcal{M}) = (\mathcal{M} \otimes \mathcal{I}_{A'})\Omega$. Let now $\Sigma \in \text{St}_1(ABA'B')$ be a purification of $\frac{1}{d_A}\text{Ch}(\mathcal{M})$. Then

$$\frac{1}{d_A} \boxed{\text{Ch}(\mathcal{M})}_{AB} = \boxed{\Omega}_{AB} \boxed{\mathcal{M}}_{AB} = \boxed{\Sigma}_{ABA'B'} \boxed{I}_{B'A} .$$

Since the channel \mathcal{M} is normalised, we have $\text{Tr}_B[\text{Ch}(\mathcal{M})] = I_A = d_A \text{Tr}_A[\Omega]$, thus, introducing $C = BB'$ and a pure state $\eta \in \text{St}_1(C)$, we have

$$\begin{array}{c} (\eta \xrightarrow{\text{BB'}} I) \\ \Omega \xrightarrow{\text{A}} I \end{array} = \Sigma \xrightarrow{\text{B'A}} I , \quad \text{A'}$$

and by the essential uniqueness of purification 11.3 we obtain

$$\Sigma \xrightarrow{\text{B'A}} I = \begin{array}{c} (\eta \xrightarrow{\text{BB'}} \mathcal{N}_U \xrightarrow{\text{B'A}} I) \\ \Omega \xrightarrow{\text{A}} \mathcal{N}_U \xrightarrow{\text{B'A}} I \end{array} , \quad \text{B} \quad \text{A'}$$

for some unitary operator $U : \mathcal{H}_{B'AB} \rightarrow \mathcal{H}_{B'A}$. Now, applying the deterministic effect I_D with $D = B'A$, we have

$$\Omega \xrightarrow{\text{A}} \mathcal{M} \xrightarrow{\text{B}} I = \begin{array}{c} (\eta \xrightarrow{\text{BB'}} \mathcal{N}_U \xrightarrow{\text{B'A}} I) \\ \Omega \xrightarrow{\text{A}} \mathcal{N}_U \xrightarrow{\text{B'A}} I \end{array} .$$

Finally, considering that

$$d_A \Omega \xrightarrow{\text{A}} \mathcal{X} \xrightarrow{\text{B}} I = \text{Ch}(\mathcal{X}) \xrightarrow{\text{B}} I , \quad \text{A'}$$

by the invertibility of the map Ch we conclude that

$$\begin{array}{c} A \xrightarrow{\mathcal{M}} B \\ \Omega \end{array} = \begin{array}{c} A \xrightarrow{\mathcal{N}_U} B \\ \eta \xrightarrow{C} \mathcal{N}_U \xrightarrow{D} I \end{array} . \quad \square$$

Chapter 12

Lecture 13: von Neumann entropy

12.1 Quantum transmission of classical information

In a quantum transmission of classical information, such as in a low-signal communication along an optical fibre, the input character, say x_i , is transmitted by sending the quantum state ρ_i^{in} . At the output of the channel, described by \mathcal{C} , the state is degraded to $\rho_i = \mathcal{C}(\rho_i^{\text{in}})$. The receiver then performs a measurement for the purpose of discriminating the possible input states ρ_i . Such a measurement will be described in general by a POVM $\mathbf{P} = \{P_j\}$ where the element P_j corresponds to the effect giving the probability of detection of the output character y_j . Therefore, the joint probability of the channel will be given by the Born rule

$$\mathbb{P}_{XY}(x_i, y_j) = \text{Tr}[\rho_i P_j] = p_i \text{Tr}[\tilde{\rho}_i P_j], \quad (12.1)$$

where $p_i := \text{Tr}[\rho_i] = \text{Tr}[\rho_i^{\text{in}}]$ is the probability of preparing the state ρ_i^{in} , namely the probability of sending the symbol x_i . This corresponds to the following scheme

$$\left(\begin{array}{c} \{\rho_i^{\text{in}}\} \\ \end{array} \right) \xrightarrow{\text{A}} \boxed{\mathcal{C}} \xrightarrow{\text{B}} \left(\begin{array}{c} \{P_j\} \\ \end{array} \right) = \left(\begin{array}{c} \{\rho_i\} \\ \end{array} \right) \xrightarrow{\text{B}} \left(\begin{array}{c} \{P_j\} \\ \end{array} \right). \quad (12.2)$$

In summary, the ingredients for the quantum communication scheme are:

1. an ensemble of states $\mathsf{E} = \{\rho_i\} \equiv \{p_i \tilde{\rho}_i\}$
2. a POVM $\mathbf{P} = \{P_j\}$.

where we remind that the quantum ensemble E could be equivalently described by a set of normalised positive operators $\tilde{\rho}_j$ along with their probabilities p_j , so that $\rho_j = p_j \tilde{\rho}_j$, and

$$\mathbb{P}_{Y|X=x_i}[y_j|x_i] = \text{Tr}[P_j \tilde{\rho}_i] = \frac{\text{Tr}[P_j \rho_i]}{\text{Tr}[\rho_i]}. \quad (12.3)$$

We can now regard the mutual information of a quantum communication scheme as a functional of the ensemble E and the POVM \mathbf{P} , and we will then write equivalently

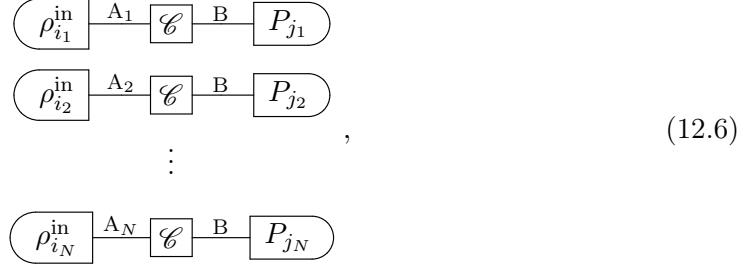
$$I(\mathsf{E}, \mathbf{P}) \equiv I(X : Y). \quad (12.4)$$

This scenario naturally leads to the following definition

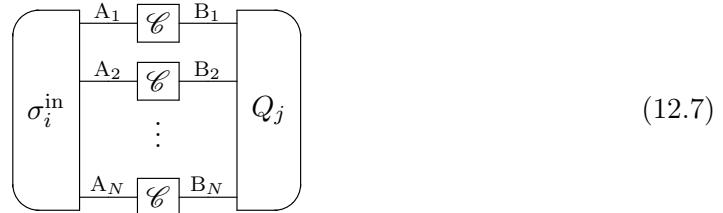
Definition 12.1 (Accessible information of an ensemble). Given an ensemble $\mathsf{E} = \{\rho_i\} \subseteq \text{St}(A)$, the *accessible information* of E is

$$\text{Acc}(\mathsf{E}) := \max_{\mathbf{P}} I(X : Y) = \max_{\mathbf{P}} I(\mathsf{E}, \mathbf{P}). \quad (12.5)$$

If we use the channel N times independently on N copies of the ensemble, as in the following scheme



the accessible information clearly provides a bound on the classical capacity of \mathcal{C} . However, differently from the classical case, the existence of quantum entangled states (and non-local measurements) offers the possibility of achieving better performances. Indeed, the N uses of \mathcal{C} can be exploited as follows



This possibility makes it harder to find a closed formula—what is usually called a *single letter formula*—for the classical capacity of a quantum channel. However, *in the case of separable input states* the Holevo-Schumacher-Westmoreland theorem characterises the classical capacity in terms of the so-called Holevo χ -quantity, which represents an upper bound on the accessible information of an ensemble.

The purpose of the next two chapters is precisely to prove the latter relevant bound, due to Holevo.

12.2 Von Neumann entropy

We will see that the role of the Shannon entropy here is played by the von Neumann entropy $S(\rho)$ of the quantum state ρ , defined as

$$S(\rho) := -\text{Tr}[\rho \log_2 \rho]. \quad (12.8)$$

Indeed, the von Neumann entropy plays a dual role. Not only it quantifies the maximum amount of *classical* information that one can gain about a classical symbol encoded in the input state by performing the best possible measurement, but it also

quantifies the maximum amount of *quantum* information content per state of the ensemble \mathbb{E} , i.e. the minimum number of *qubits* needed to reliably encode the quantum state.

Thus quantification of quantum information is largely concerned with von Neumann entropy, in the same way as classical information uses the Shannon entropy. And, indeed, much of the mathematical framework of quantum information theory resembles the classical one, even though the conceptual context is radically different. Clearly, the central issue in quantum information is the existence of entangled states, which bear features with no classical analogue.

Remark 15. Since every state $\rho \in \text{St}_1(\mathbf{A})$ is positive by definition, the spectral theorem holds, and one can write

$$\rho = \sum_k p_k |\varphi_k\rangle\langle\varphi_k|,$$

and for every function $f : \mathbb{R} \rightarrow \mathbb{R}$ one has

$$f(\rho) := \sum_k f(p_k) |\varphi_k\rangle\langle\varphi_k|.$$

In particular, the function f defining the von Neumann entropy is

$$f(x) := \begin{cases} -x \log_2 x & x > 0, \\ 0 & x = 0. \end{cases}$$

Thus, if $\{p_k\}_{k=1}^{d_{\mathbf{A}}}$ is the spectrum of ρ , one has $p_k \geq 0$ and $\sum_k p_k = \text{Tr}[\rho] = 1$, and then

$$S(\rho) = H(\mathbf{p}).$$

This allows us to observe that *the von Neumann entropy only depends on the spectrum of ρ , and every map \mathcal{M} that does not change the spectrum leaves the von Neumann entropy invariant, $S(\rho) = S(\mathcal{M}(\rho))$.* In particular, we have

$$\rho = \sum_i p_i |\psi_i\rangle\langle\psi_i|, \Rightarrow \rho^T = \sum_i p_i |\psi_i^*\rangle\langle\psi_i^*|,$$

and since $\langle\psi_i^*|\psi_j^*\rangle = \langle\psi_j|\psi_i\rangle = \delta_{i,j}$, the transpose leaves the spectrum of ρ invariant, then

$$S(\rho^T) = S(\rho). \quad (12.9)$$

Remark 16. As regards notation, we will often write $S(\mathbf{A})$ instead of $S(\rho)$ when system \mathbf{A} is in state $\rho \in \text{St}_1(\mathbf{A})$. In particular, when dealing with bipartite systems \mathbf{AB} we will write $\rho_{\mathbf{AB}} \in \text{St}_1(\mathbf{AB})$, and

$$\rho_{\mathbf{A}} := \text{Tr}_{\mathbf{B}}[\rho_{\mathbf{AB}}], \quad \rho_{\mathbf{B}} := \text{Tr}_{\mathbf{A}}[\rho_{\mathbf{AB}}],$$

thus

$$S(\mathbf{AB}) = S(\rho_{\mathbf{AB}}), \quad S(\mathbf{A}) = S(\rho_{\mathbf{A}}), \quad S(\mathbf{B}) = S(\rho_{\mathbf{B}}).$$

Remark 17. It is very natural to describe encoding of classical information on quantum states. In particular, given a classical random variable of type n , one has $\text{St}_1(\mathbf{X}) = \{(p_1, p_2, \dots, p_n) | p_i \geq 0, \sum_i p_i = 1\}$. Given a quantum system A of type n , we can fix an orthonormal basis $\{\varphi_i\}_{i=1}^n$, and given a state (p_1, p_2, \dots, p_n) one can encode it in the quantum state $\rho \in \text{St}_1(\mathbf{A})$ as follows

$$(p_1, p_2, \dots, p_n) \mapsto \rho = \sum_{i=1}^n p_i |\varphi_i\rangle\langle\varphi_i|.$$

In particular, for a bit X and a qubit A one has

$$\text{St}_1(\mathbf{X}) = \{(p, 1-p) | 0 \leq p \leq 1\}, \quad (p, 1-p) \mapsto p|0\rangle\langle 0| + (1-p)|1\rangle\langle 1| \in \text{St}_1(\mathbf{A}).$$

However, the set of states of A contains also projections on arbitrary superpositions of $|0\rangle$ and $|1\rangle$ —and mixtures thereof—which are not simply encodings of classical sources under the above encoding map. Conversely, one might in general encode classical information into a set of states that are not mutually orthogonal, and not even pure.

We now define the analogues of relative and conditional entropy, and mutual information.

Definition 12.2 (Quantum relative entropy). We define the *quantum relative entropy* of ρ and σ in $\text{St}(\mathbf{A})$ as follows

$$S(\rho\|\sigma) := \begin{cases} \text{Tr}[\rho \log_2 \rho] - \text{Tr}[\rho \log_2 \sigma] & \text{Ker}(\sigma) \subseteq \text{Ker}(\rho) \\ +\infty & \text{otherwise.} \end{cases} \quad (12.10)$$

Similarly to the Kullback-Leibler divergence, also the quantum relative entropy is a useful proving tool.

Definition 12.3 (Quantum conditional entropy). We define the *quantum conditional entropy* of system A conditional on system B, with AB in state $\rho \in \text{St}_1(\mathbf{AB})$, as follows

$$S(A|B) := S(AB) - S(B). \quad (12.11)$$

Definition 12.4 (Quantum mutual information). We define the *quantum mutual information* of systems A and B, with AB in state $\rho \in \text{St}_1(\mathbf{AB})$, as follows

$$\begin{aligned} S(A : B) &:= S(A) + S(B) - S(AB) \\ &= S(A) - S(A|B) \\ &= S(B) - S(B|A). \end{aligned} \quad (12.12)$$

Notice that, differently from the classical case, the joint von Neumann entropy $S(AB)$ is not necessarily bounded from below by the two marginals. Indeed, for a pure maximally entangled state $\rho_{AB} = \frac{1}{d_A} |I_A\rangle\langle I_A|$ one has $\rho_A = I_A/d_A$ and $\rho_B = I_B/d_B$. Then, $S(AB) = 0$ and $S(A) = \log d_A$, hence $S(B|A) < 0$. In particular, if ρ_{AB} is pure, then $S(B|A)$ is negative if and only if ρ_{AB} is entangled.

For the special case of diagonal factorised states we have

$$\rho_{AB} = \sum_{(x_i, y_j) \in \text{Rng}(X) \times \text{Rng}(Y)} p_{ij}^{AB} |i\rangle\langle i|_A \otimes |j\rangle\langle j|_B,$$

and

$$\begin{aligned} \rho_A &= \sum_{x_i \in \text{Rng}(X)} p_i^A |i\rangle\langle i|_A, & p_i^A &= \sum_{y_j \in \text{Rng}(Y)} p_{ij}^{AB}, \\ \rho_B &= \sum_{y_j \in \text{Rng}(Y)} p_j^B |j\rangle\langle j|_B, & p_j^B &= \sum_{x_i \in \text{Rng}(X)} p_{ij}^{AB}. \end{aligned}$$

In this case

$$S(A : B) = I(X : Y), \quad \mathbb{P}_{X,Y}[x_i, y_j] = p_{ij}$$

In consideration of the relevant role played by the von Neumann entropy, we now devote a subsection to analyse its properties.

12.2.1 Klein's inequality

We first prove the main property of quantum relative entropy.

Theorem 12.5 (Klein's inequality). *Let the states $\rho, \sigma \in \text{St}(A)$ have the same probability $\text{Tr}[\rho] = \text{Tr}[\sigma]$. The quantum relative entropy of ρ with respect to σ is non-negative*

$$S(\rho\|\sigma) \geq 0, \quad \text{with equality iff } \rho = \sigma. \quad (12.13)$$

Proof. The thesis is straightforward if $\text{Ker}(\sigma) \not\subseteq \text{Ker}(\rho)$. Let us focus then on the complementary case. One can diagonalise both ρ and σ , as

$$\rho = \sum_i p_i P_i, \quad \sigma = \sum_j q_j Q_j,$$

where P_i and Q_j are the projections on the i -th and j -th eigenspace of ρ and σ , respectively, and $p_i \neq p_{i'}$ for $i \neq i'$, as well as $q_j \neq q_{j'}$ for $j \neq j'$. Let $\text{Tr}[P_i] = n_i$ and $\text{Tr}[Q_j] = m_j$. We then write the relative entropy as follows

$$\begin{aligned} S(\rho\|\sigma) &= \sum_i n_i p_i \log_2 p_i - \sum_{i,j} p_i \log_2 q_j \text{Tr}[P_i Q_j] \\ &= \sum_i n_i p_i \left(\log_2 p_i - \sum_j S_{ji} \log_2 q_j \right), \end{aligned}$$

where the matrix with elements $S_{ji} := \text{Tr}[P_i Q_j]/n_i$ is stochastic, namely $S_{ji} \geq 0$ for every i, j , and $\sum_j S_{ji} = \text{Tr}[P_i]/n_i = 1$, for every i . Since the logarithm is a strictly concave function, by Jensen's inequality we then have

$$\sum_{j=1}^{d_A} S_{ji} \log_2 q_j \leq \log_2 \left(\sum_{j=1}^{d_A} S_{ji} q_j \right) = \log_2 r_i, \quad r_i := \sum_{j=1}^{d_A} S_{ji} q_j. \quad (12.14)$$

Notice that

$$\sum_i n_i r_i = \sum_{i,j} n_i S_{ji} q_j = \sum_{i,j} \text{Tr}[P_i Q_j] q_j = \text{Tr}[\sigma] = \text{Tr}[\rho],$$

thus the vector \mathbf{r} where r_i has the same multiplicity n_i as p_i , has the same sum as \mathbf{p} . For deterministic states ($\text{Tr}[\rho] = \text{Tr}[\sigma] = 1$), we can then write

$$S(\rho\|\sigma) \geq \sum_i n_i p_i \log_2 \frac{p_i}{r_i} \equiv H(\mathbf{p}\|\mathbf{r}),$$

while for sub-normalised states one can define $\mathbf{p}' = \mathbf{p}/\text{Tr}[\rho]$ and $\mathbf{r}' = \mathbf{r}/\text{Tr}[\sigma]$, and

$$S(\rho\|\sigma) \geq \frac{1}{\text{Tr}[\rho]} \sum_i n_i p_i \log_2 \frac{p_i}{r_i} \equiv \frac{1}{\text{Tr}[\rho]} H(\mathbf{p}'\|\mathbf{r}'),$$

From the Gibbs' inequality for $H(\mathbf{p}\|\mathbf{r})$ we thus get

$$S(\rho\|\sigma) \geq 0.$$

Equality holds upon two conditions. The first one is saturation of the bound in equation (12.14), which is Jensen's inequality for the \log_2 function, and by strict concavity of \log_2 one must have $S_{ji} = \delta_{j,\varphi(i)}$ for every i . By making this condition explicit, we have

$$\text{Tr}[P_i Q_{\varphi(i)}] = n_i, \quad \text{Tr}[P_i Q_j] = 0, \quad j \neq \varphi(i).$$

Then, since $\text{Tr}[P_i Q_j] = \text{Tr}[P_i Q_j P_i]$ is the trace of the non-negative definite operator $P_i Q_j P_i$, the trace is null iff $P_i Q_j P_i = 0$. Now, since $P_i Q_j P_i = (Q_j P_i)^\dagger (Q_j P_i)$, one has $Q_j P_i = 0$ for $j \neq \varphi(i)$. Also, $P_i Q_j = (Q_j P_i)^\dagger = 0$ for $j \neq \varphi(i)$. Then

$$\begin{aligned} P_i &= \sum_j Q_j P_i = Q_{\varphi(i)} P_i \\ &= P_i^\dagger = P_i Q_{\varphi(i)}. \end{aligned}$$

Thus, $[P_i, Q_j] = 0$ for every i, j and $P_i Q_{\varphi(i)} = Q_{\varphi(i)} P_i = P_i$. Moreover, since $S_{ji} = \delta_{j,\varphi(i)}$, one has $r_i = q_{\varphi(i)}$. Suppose now that for some j there is no i such that $j = \varphi(i)$. This implies that

$$Q_j = Q_j \sum_i P_i = \sum_i (Q_j P_i) = 0,$$

namely $Q_j = 0$. This implies that there cannot be any such j , and φ is thus surjective on the set of values of j , namely the spectrum of σ is the collection of the r_i 's.

The second condition is saturation of Gibbs' inequality for $H(\mathbf{p}\|\mathbf{r})$ [or $H(\mathbf{p}'\|\mathbf{r}')$], namely $\mathbf{p} = \mathbf{r}$. Then it must be $p_i = r_i = q_{\varphi(i)}$, and thus, ρ and σ must have the same spectrum. Moreover, for $i \neq i'$ one has $q_{\varphi(i)} = r_i = p_i \neq p_{i'} = r_{i'} = q_{\varphi(i')}$, which implies $\varphi(i) \neq \varphi(i')$, thus φ is an injective function. Then, one can write

$$\begin{aligned} Q_{\varphi(j)} &= \sum_i Q_{\varphi(j)} P_i = Q_{\varphi(j)} P_j \\ &= \sum_l Q_l P_j = P_j. \end{aligned}$$

In conclusion, ρ and σ have the same eigenspaces and the same spectrum, which implies

$$S(\rho\|\sigma) = 0 \Leftrightarrow \rho = \sigma. \quad \square$$

12.2.2 Preliminary lemmas

We now need two lemmas that we use in the following.

Lemma 12.6. *Let $0 \leq X, Y \in \mathcal{L}(\mathcal{H}_A)$ with $[X, Y] = 0$. Then*

$$\log_2 XY = \log_2 X + \log_2 Y. \quad (12.15)$$

Proof. The operators X and Y are simultaneously diagonalisable, with

$$X = \sum_{i=1}^{d_A} x_i |\psi_i\rangle\langle\psi_i|, \quad Y = \sum_{i=1}^{d_A} y_i |\psi_i\rangle\langle\psi_i|.$$

Then $XY = \sum_{i=1}^{d_A} x_i y_i |\psi_i\rangle\langle\psi_i|$, and

$$\begin{aligned} \log_2 XY &= \sum_{i=1}^{d_A} \log_2(x_i y_i) |\psi_i\rangle\langle\psi_i| \\ &= \sum_{i=1}^{d_A} (\log_2 x_i + \log_2 y_i) |\psi_i\rangle\langle\psi_i| \\ &= \log_2 X + \log_2 Y. \end{aligned} \quad \square$$

Lemma 12.7. *Let $0 \leq X \in \mathcal{L}(\mathcal{H}_A)$. Then one has*

$$\log_2(X \otimes I_B) = (\log_2 X) \otimes I_B \quad (12.16)$$

Proof. One has

$$X = \sum_{i=1}^{d_A} x_i |\psi_i\rangle\langle\psi_i| \Rightarrow X \otimes I_B = \sum_{i=1}^{d_A} \sum_{j=1}^{d_B} x_i |\psi_i\rangle\langle\psi_i| \otimes |\varphi_j\rangle\langle\varphi_j|.$$

Then

$$\begin{aligned} \log_2(X \otimes I_B) &= \log_2 \left(\sum_{i=1}^{d_A} \sum_{j=1}^{d_B} x_i |\psi_i\rangle\langle\psi_i| \otimes |\varphi_j\rangle\langle\varphi_j| \right) \\ &= \sum_{i=1}^{d_A} \log_2(x_i) |\psi_i\rangle\langle\psi_i| \otimes \sum_{j=1}^{d_B} |\varphi_j\rangle\langle\varphi_j| \\ &= (\log_2 X) \otimes I_B. \end{aligned} \quad \square$$

12.3 Mathematical properties of the von Neumann entropy

1. **Entropy of pure states.** The entropy of a pure state $\rho = |\psi\rangle\langle\psi|$ is $S(\rho) = 0$. We notice that in the definition of the von Neumann entropy we follow the same rule as for the Shannon entropy, namely we put $0 \log 0 := 0$. Indeed, the von Neumann entropy is equal to the Shannon entropy $S(\rho) = H(\mathbf{p})$ of the eigenvalues $\{p_i\}_{i=1}^{d_A}$ of the state $\rho = \sum_i p_i |\psi_i\rangle\langle\psi_i|$, and for a pure state the eigenvalues are all zero, apart from a single one which is unit.
2. **Maximum value.** The entropy is upper bounded as $S(\rho) \leq \log_2 d_A$, where $\rho \in \text{St}_1(A)$. This clearly follows from the fact that $S(\rho) = H(\mathbf{p})$, where $\{p_i\}_{i=1}^{d_A}$ is the spectrum of ρ , and the maximum value of the Shannon entropy is achieved when the probability distribution is uniform: $p_i = 1/d_A$, i.e. for $\rho = I_A/d_A$.
3. **Isometric invariance.** The entropy is left unchanged by an isometric transformation, i.e.

$$S(V\rho V^\dagger) = S(\rho), \quad (12.17)$$

with $V : \mathcal{H}_A \rightarrow \mathcal{H}_B$ is an isometry, namely $V^\dagger V = I_A$. This is just a consequence of the fact that $S(\rho)$ depends only on the spectrum of ρ , and the isometric transformation preserves the diagonal form of the state:

$$\rho = \sum_{i=1}^{d_A} p_i |\psi_i\rangle\langle\psi_i| \Rightarrow V\rho V^\dagger = \sum_{i=1}^{d_A} p_i V|\psi_i\rangle\langle\psi_i|V^\dagger, \quad (12.18)$$

and by the properties of isometries $\langle\psi_i|V^\dagger V|\psi_j\rangle = \langle\psi_i|\psi_j\rangle = \delta_{i,j}$.

4. **Marginal entropies.** For a composite system AB in a pure state $|\Psi\rangle\langle\Psi| \in \text{St}_1(AB)$, the marginal entropies are equal, i.e. $S(A) = S(B)$.

Proof. Let us fix $d_A \leq d_B$, which is clearly not restrictive, as the role of A and B can be interchanged. For a pure joint state $|\Psi\rangle\langle\Psi|$ with $\Psi = V|\Psi|$ we have

$$\begin{aligned} \rho_A &= \text{Tr}_B[|\Psi\rangle\langle\Psi|] = V|\Psi|^2V^\dagger, \\ \rho_B &= \text{Tr}_A[|\Psi\rangle\langle\Psi|] = |\Psi|^T V^T V^* |\Psi|^T = (|\Psi|^2)^T \end{aligned}$$

namely $\rho_A = V\rho_B^T V^\dagger$. Since the transpose as well as isometries do not change the spectrum, by equations (12.9) and (12.17), the two states have the same von Neumann entropy. \square

5. **Entropy of statistically independent systems.**

$$S(\rho \otimes \sigma) = S(\rho) + S(\sigma). \quad (12.19)$$

Proof. The statement is an immediate consequence of the identity

$$\rho \otimes \sigma = (\rho \otimes I_B)(I_A \otimes \sigma),$$

which by lemmas 12.6 and 12.7 implies

$$\log_2(\rho \otimes \sigma) = \log_2 \rho \otimes I_B + I_A \otimes \log_2 \sigma. \quad (12.20)$$

More explicitly, the latter follows from $\log_2(XY) = \log_2(X) + \log_2(Y)$ for $[X, Y] = 0$, and from $\log_2(X \otimes I_B) = \log_2(X) \otimes I_B$ [and $\log_2(I_A \otimes Y) = I_A \otimes \log_2(Y)$]. \square

6. **Subadditivity.** For any state $\rho_{AB} \in \text{St}_1(AB)$ one has

$$S(AB) \leq S(A) + S(B), \quad \text{with equality iff A and B are factorised.} \quad (12.21)$$

Proof. Using the Klein's inequality for $S(\rho_{AB} \| \rho_A \otimes \rho_B)$ and identity (12.20), one has

$$\begin{aligned} S(\rho_{AB}) &\leq -\text{Tr}[\rho_{AB} \log(\rho_A \otimes \rho_B)] \quad (\text{with equality iff } \rho_{AB} = \rho_A \otimes \rho_B) \\ &= -\text{Tr}[\rho_{AB}(\log \rho_A \otimes I_B)] - \text{Tr}[\rho_{AB}(I_A \otimes \log \rho_B)] \\ &= S(A) + S(B). \end{aligned} \quad \square$$

7. **States with orthogonal supports.** For a set of states $\{\rho_i\}$ having orthogonal supports, i.e. $\rho_i \rho_j = 0$ for $i \neq j$

$$S\left(\sum_i p_i \tilde{\rho}_i\right) = H(\mathbf{p}) + \sum_i p_i S(\tilde{\rho}_i). \quad (12.22)$$

Proof. Since the states have orthogonal supports, we can write the joint spectral decomposition of the normalised states $\tilde{\rho}_i$ as

$$\tilde{\rho}_i |e_j^{(i)}\rangle = \lambda_j^{(i)} |e_j^{(i)}\rangle.$$

Consequently, we have

$$\sum_i p_i \tilde{\rho}_i |e_j^{(k)}\rangle = p_k \lambda_j^{(k)} |e_j^{(k)}\rangle, \quad \sum_j \lambda_j^{(k)} = 1 \forall k.$$

We therefore have

$$\begin{aligned} S\left(\sum_i p_i \tilde{\rho}_i\right) &= -\sum_{jk} p_k \lambda_j^{(k)} \log p_k \lambda_j^{(k)} \\ &= -\sum_k p_k \log p_k - \sum_k p_k \sum_j \lambda_j^{(k)} \log \lambda_j^{(k)} \\ &= H(\mathbf{p}) + \sum_i p_i S(\tilde{\rho}_i). \end{aligned} \quad \square$$

As a corollary of this case we have that the quantum von Neumann entropy of a joint classical-quantum system AX in the *quantum-classical state*

$$\rho_{AX} = \sum_j p_j \rho_A^j \otimes |j\rangle\langle j|_X,$$

associated with the ensemble of states $E_A = \{p_j \rho_A^j\}$, is given by

$$S(AX) = \sum_j p_j S(\rho_A^j) + H(\mathbf{p}), \quad (12.23)$$

and the quantum-classical mutual information is rewritten as

$$S(A : X) = S\left(\sum_j p_j \rho_A^j\right) - \sum_j p_j S(\rho_A^j). \quad (12.24)$$

8. **Concavity.** For $\mathbf{p} := (p_1, p_2, \dots, p_k)$ a general probability vector and $\tilde{\rho}_1, \tilde{\rho}_2, \dots, \tilde{\rho}_k \in \text{St}_1(A)$, one has

$$S\left(\sum_{i=1}^k p_i \tilde{\rho}_i\right) \geq \sum_{i=1}^k p_i S(\tilde{\rho}_i). \quad (12.25)$$

Proof. Introduce an auxiliary system B such that the state of AB is

$$\rho_{AB} = \sum_{n=1}^k p_n \tilde{\rho}_n \otimes |n\rangle\langle n|.$$

Clearly we have $S(A) = S\left(\sum_{n=1}^k p_n \tilde{\rho}_n\right)$, $S(B) = H(\mathbf{p})$, hence, using subadditivity and identity (12.23), we obtain

$$H(\mathbf{p}) + \sum_{i=1}^k p_i S(\tilde{\rho}_i) = S(AB) \leq S(A) + S(B) = S\left(\sum_{n=1}^k p_n \tilde{\rho}_n\right) + H(\mathbf{p}),$$

namely the statement. \square

9. **Non-selective projective measurements.** Let $\mathcal{P} : A \rightarrow A$ be a projective measurement, i.e.

$$\mathcal{P}(\rho) = \sum_{i=1}^k P_i \rho P_i, \quad P_i P_j = \delta_{ij} P_i, \quad \sum_{i=1}^k P_i = I_A. \quad (12.26)$$

Then one has

$$S(\mathcal{P}(\rho)) \geq S(\rho), \quad (12.27)$$

with equality holding iff $\mathcal{P}(\rho) = \rho$.

Proof. Notice that \mathcal{P} is an idempotent projective map, i.e. $\mathcal{P}^2 = \mathcal{P}$, and that $\mathcal{P}(\rho)$ commutes with P_i for every i . Now, reminding that $\sum_{i=1}^k P_i = I_A$, one has

$$\begin{aligned}\text{Tr}[\rho \log_2 \mathcal{P}(\rho)] &= \text{Tr}\left[\sum_{i=1}^k P_i \rho \log_2 \mathcal{P}(\rho)\right] \\ &= \text{Tr}\left[\sum_{i=1}^k P_i \rho \log_2 \mathcal{P}(\rho) P_i\right] = \text{Tr}\left[\sum_i P_i \rho P_i \log \mathcal{P}(\rho)\right]\end{aligned}$$

namely

$$\text{Tr}[\rho \log \mathcal{P}(\rho)] = \text{Tr}[\mathcal{P}(\rho) \log \mathcal{P}(\rho)].$$

This implies that

$$\begin{aligned}S(\rho \| \mathcal{P}(\rho)) &= -S(\rho) - \text{Tr}[\rho \log \mathcal{P}(\rho)] \\ &= -S(\rho) + S(\mathcal{P}(\rho)),\end{aligned}$$

and the statement then follows from Klein's inequality. Consequently, equality holds iff $\mathcal{P}(\rho) = \rho$. \square

10. Upper bound. The von Neumann entropy is bounded as

$$S\left(\sum_{n=1}^k p_n \tilde{\rho}_n\right) \leq \sum_{n=1}^k p_n S(\tilde{\rho}_n) + H(\mathbf{p}), \quad (\text{with equality iff } \tilde{\rho}_i \tilde{\rho}_j = 0 \text{ for } i \neq j). \quad (12.28)$$

Proof. Let us first consider the pure case $\tilde{\rho}_n = |\psi_n\rangle\langle\psi_n| \in \text{St}_1(A)$. Take the following purification of $\rho := \sum_{n=1}^k p_n \rho_n$ on $\mathcal{H}_A \otimes \mathcal{H}_B$

$$|\Psi\rangle\rangle = \sum_{n=1}^k \sqrt{p_n} |\psi_n\rangle \otimes |n\rangle.$$

Since the system AB is in a pure state, we have

$$S(\rho) = S(A) = S(B).$$

Suppose now we perform a projective measurement over B in the orthonormal basis $\{|n\rangle\}$. After the measurement we have

$$\rho_{AB'} = \sum_n p_n |\psi_n\rangle\langle\psi_n| \otimes |n\rangle\langle n|, \quad \rho_{B'} = \sum_n p_n |n\rangle\langle n|.$$

Now, since non-selective projective measurements never decrease the entropy [see identity (12.27)], we have $S(AB') \geq S(AB) = 0$, and at the same time $S(B') \geq S(B) = S(\rho)$. Therefore,

$$S(\rho) \leq S(B') = H(\mathbf{p}), \quad \text{with equality iff } \{|\psi_n\rangle\} \text{ are orthogonal},$$

and since $S(\tilde{\rho}_n) = 0$, we have trivially

$$S(\rho) \leq \sum_{n=1}^k p_n S(\tilde{\rho}_n) + H(\mathbf{p}).$$

The case of mixed ρ_n is proved as follows. Consider the orthonormal decomposition

$$\tilde{\rho}_n = \sum_{j=1}^{d_A} p_j^{(n)} |e_j^{(n)}\rangle\langle e_j^{(n)}|.$$

This provides a pure ensemble $\{p_n p_j^{(n)} |e_j^{(n)}\rangle\langle e_j^{(n)}|\}$ for ρ , since

$$\rho = \sum_{nj} p_n p_j^{(n)} |e_j^{(n)}\rangle\langle e_j^{(n)}|,$$

and using our theorem for the pure case, we have

$$\begin{aligned} S(\rho) &\leq H(\{p_n p_j^{(n)}\}) = - \sum_{jn} p_n p_j^{(n)} \log p_n p_j^{(n)} \\ &= - \sum_n p_n \log p_n - \sum_n p_n \sum_j p_j^{(n)} \log p_j^{(n)} \\ &= H(\mathbf{p}) + \sum_n p_n S(\tilde{\rho}_n), \end{aligned}$$

which is the statement. Notice that the bound is achieved iff all vectors $|e_j^{(n)}\rangle$ are mutually orthogonal, i.e. the states $\tilde{\rho}_n$ have orthogonal supports. \square

11. **Entropy of state preparation.** If a pure state is drawn randomly from the ensemble $\mathsf{E} = \{p_n, |\psi_n\rangle\langle\psi_n|\}$, one has

$$S\left(\sum_n p_n |\psi_n\rangle\langle\psi_n|\right) \leq H(\mathbf{p}), \quad (12.29)$$

with equality iff the states are mutually orthogonal.

Proof. Indeed, the statement is just the upper bound (12.28) of item 10 for an ensemble of pure states. \square

Notice that this result implies that the von Neumann entropy could be defined as

$$S(\rho) := \min_{\mathsf{E}_\rho^P} H(\mathbf{p}_\mathsf{E}),$$

where E_ρ^P is the set of pure ensembles $\mathsf{E} = \{p_i |\psi_i\rangle\langle\psi_i|\}$ such that $\sum_i p_i |\psi\rangle\langle\psi|_i = \rho$, and $\mathbf{p}_\mathsf{E} := \{p_i\}$.

- 12. Entropy of measurement.** Let $\rho \in \text{St}_1(\mathcal{A})$, and let \mathbf{P} be an *atomic* POVM $\{|\psi_i\rangle\langle\psi_i|\}$. Then

$$S(\rho) \leq H(X), \quad (12.30)$$

where X is the random variable

$$X = \begin{cases} \text{Rng}(X) &= \{x_n\}, \\ \mathbb{P}_X[X = x_n] &= \{\langle\psi_n|\rho|\psi_n\rangle\}. \end{cases} \quad (12.31)$$

Proof. Let us define the following isometric operator

$$V := \sum_{i=1}^{|\text{Rng}(X)|} |i\rangle\langle\psi_i| : \mathcal{H}_A \rightarrow \mathcal{H}_B = \mathbb{C}^{|\text{Rng}(X)|},$$

where $\langle i|j\rangle = \delta_{i,j}$. Indeed, it is easy to check that $V^\dagger V = I_A$. Now, let us define $\mathcal{P} : B \rightarrow B$ as

$$\mathcal{P}(\sigma) := \sum_{i=1}^{d_B} |i\rangle\langle i|\sigma|i\rangle\langle i|.$$

By items 9 and 3, we now have

$$S(\mathcal{P}(V\rho V^\dagger)) \geq S(V\rho V^\dagger) = S(\rho).$$

Observing that

$$\begin{aligned} \mathcal{P}(V\rho V^\dagger) &= \sum_{i=1}^{d_B} |i\rangle\langle\psi_i|\rho|\psi_i\rangle\langle i| \\ &= \sum_{i=1}^{d_B} |i\rangle\langle i|\mathbb{P}_X[X = x_i], \end{aligned}$$

we conclude that $S(\mathcal{P}(V\rho V^\dagger)) = H(X)$. Notice that, having used item 9, equality holds if and only if $\mathcal{P}(V\rho V^\dagger) = V\rho V^\dagger$, i.e.

$$\begin{aligned} V\rho V^\dagger &= \sum_{i,i'=1}^{d_B} |i\rangle\langle\psi_i|\rho|\psi_{i'}\rangle\langle i'| \\ &= \sum_{i=1}^{d_B} |i\rangle\langle\psi_i|\rho|\psi_i\rangle\langle i|, \\ \Leftrightarrow \langle\psi_i|\rho|\psi_{i'}\rangle &= \delta_{i,i'}\langle\psi_i|\rho|\psi_i\rangle, \end{aligned}$$

namely if and only if $|\psi_i\rangle\langle\psi_i|$ are projections on the eigenvectors of ρ . \square

Notice that also this result provides an alternate definition of von Neumann entropy:

$$S(\rho) := \min_{\mathbf{P} \in P_A} H(\mathbf{p}_{\mathbf{P}, \rho}),$$

where P_A is the set of *atomic* POVMs, i.e. POVMs made of rank one operators $\mathbf{P} = \{|\psi_i\rangle\langle\psi_i|\}$, and $\mathbf{p}_{\mathbf{P}, \rho}$ is the probability vector $\mathbf{p} = (\langle\psi_i|\rho|\psi_i\rangle)$.

13. Araki-Lieb triangle inequality. The following bound holds

$$S(AB) \geq |S(A) - S(B)|. \quad (12.32)$$

Proof. Let AB be prepared in the state ρ_{AB} . We introduce a system R which purifies ρ_{AB} . From subadditivity, we have

$$S(AR) \leq S(A) + S(R). \quad (12.33)$$

However, since ABR is pure, we have

$$S(AR) = S(B), \quad S(AB) = S(R),$$

namely, using inequality (12.33), we obtain

$$S(B) = S(AR) \leq S(A) + S(R) = S(A) + S(AB),$$

and finally

$$S(AB) \geq S(B) - S(A).$$

Due to symmetry between A and B in the same way we find also

$$S(AB) \geq S(A) - S(B),$$

and the two bounds together are equivalent to the statement. \square

Chapter 13

Lecture 15: Properties of quantum relative entropy and Lieb's theorem

13.1 Mathematical properties of the quantum relative entropy

1. **Isometric invariance.** The quantum relative entropy is invariant under isometric transformations, namely

$$S(V\rho V^\dagger \| V\sigma V^\dagger) = S(\rho\|\sigma), \quad V^\dagger V = I. \quad (13.1)$$

Proof. In the case where $\text{Ker}(\sigma) \not\subseteq \text{Ker}(\rho)$, clearly $\text{Ker}(V\sigma V^\dagger) \not\subseteq \text{Ker}(V\rho V^\dagger)$. Let indeed $|\varphi\rangle \in \text{Ker}(\sigma)$ while $|\varphi\rangle \notin \text{Ker}(\rho)$. Then $V|\varphi\rangle \in \text{Ker}(V\sigma V^\dagger)$ while $V|\varphi\rangle \notin \text{Ker}(V\rho V^\dagger)$. In the complementary case the thesis is an immediate consequence of the identity $f(VAV^\dagger) = Vf(A)V^\dagger$ for any selfadjoint operator A , and the invariance of the trace under cyclic permutations of its argument. \square

2. **Invariance under extension.** The quantum relative entropy is invariant under extension with a fixed state, namely

$$S(\rho \otimes \nu \| \sigma \otimes \nu) = S(\rho\|\sigma). \quad (13.2)$$

Proof. If $\text{Ker}(\sigma) \not\subseteq \text{Ker}(\rho)$, then also $\text{Ker}(\sigma \otimes \nu) \not\subseteq \text{Ker}(\rho \otimes \nu)$. When $\text{Ker}(\sigma) \subseteq \text{Ker}(\rho)$, the thesis is an easy consequence of the definition. One has

$$\begin{aligned} S(\rho \otimes \nu \| \sigma \otimes \nu) &= \text{Tr}[(\rho \otimes \nu) \log_2 (\rho \otimes \nu)] - \text{Tr}[(\rho \otimes \nu) \log_2 (\sigma \otimes \nu)] \\ &= \text{Tr}[(\rho \otimes \nu)(\log_2 \rho \otimes I_B + I_A \otimes \log_2 \nu)] \\ &\quad - \text{Tr}[(\rho \otimes \nu)(\log_2 \sigma \otimes I_B + I_A \otimes \log_2 \nu)] \\ &= \text{Tr}[\rho \log_2 \rho] + \text{Tr}[\nu \log_2 \nu] - \text{Tr}[\rho \log_2 \sigma] - \text{Tr}[\nu \log_2 \nu] \\ &= S(\rho\|\sigma). \end{aligned} \quad \square$$

3. **Monotonicity under marginalisation.** Let ρ and σ be two states of a composite system AB. The quantum relative entropy is monotonically decreasing when discarding one system, namely

$$S(\text{Tr}_2(\rho) \| \text{Tr}_2(\sigma)) \leq S(\rho\|\sigma). \quad (13.3)$$

This result requires joint convexity of relative entropy, which we prove later, using Lieb's theorem.

4. **Uhlmann monotonicity theorem.** The quantum relative entropy is monotonically non-increasing under the action of a channel, namely one has

$$S(\mathcal{E}(\rho)\|\mathcal{E}(\sigma)) \leq S(\rho\|\sigma). \quad (13.4)$$

Proof. This result easily follows from items 1, 2, and 3. \square

13.2 Preliminary results for Lieb's theorem

The aim of this chapter is to provide a proof of Lieb's theorem, a technical result whose proof was provided by E. Lieb, and used for strong subadditivity by Lieb and Ruskai. The theorem is preliminary to many crucial results about properties of von Neumann entropy. The proof in this chapter is based on J. Watrous's lecture notes [3]. We start by introducing some definitions and lemmas that will be used to prove the main result.

Definition 13.1 (Spectral radius). Given an operator $A : \mathcal{H}_A \rightarrow \mathcal{H}_A$, we define its spectral radius $\rho(A)$ as the maximum among absolute values of its eigenvalues. In formula

$$\rho(A) := \max_{\lambda \in \text{Spec}(A)} |\lambda|, \quad (13.5)$$

$\text{Spec}(A)$ denoting the spectrum of A .

Remark 18. In the infinite-dimensional case the spectral radius $\rho(A)$ is defined as the supremum of $|\lambda|$ for $\lambda \in \text{Spec}(A)$.

The spectral radius is related to the *uniform norm*, also called *sup-norm*, inducing the *uniform operator topology* in $\mathcal{L}(\mathcal{H}_A)$.

Definition 13.2 (Uniform norm). Given an operator $A : \mathcal{H}_A \rightarrow \mathcal{H}_A$, we define its uniform norm (sup-norm) $\|A\|_\infty$ as follows

$$\|A\|_\infty := \sup_{\psi \in \mathcal{H}_A} \frac{\|A\psi\|_A}{\|\psi\|_A} = \sup_{\substack{\psi \in \mathcal{H}_A \\ \|\psi\|_A=1}} \|A\psi\|_A \quad (13.6)$$

where $\|\cdot\|_A$ denotes the norm induced by the inner product in \mathcal{H}_A .

The following property (the proof is left as an exercise) holds for the uniform norm

$$\|AB\|_\infty \leq \|A\|_\infty \|B\|_\infty.$$

Now, the spectral radius is bounded by the sup-norm as follows.

Lemma 13.3. *For an operator $A : \mathcal{H}_A \rightarrow \mathcal{H}_A$, the spectral radius is upper-bounded by the uniform norm.*

Proof. Let ψ be any normalised eigenvector of A . Then, by definition, the sup norm satisfies

$$\|A\|_\infty \geq \|A\psi\|_A = |\lambda|.$$

This implies that the maximum (supremum) of $|\lambda|$ over $\text{Spec}(A)$ is also bounded from above by $\|A\|_\infty$. \square

For normal operators $A : \mathcal{H}_A \rightarrow \mathcal{H}_A$ a stronger result holds.

Lemma 13.4. *Let $A : \mathcal{H}_A \rightarrow \mathcal{H}_A$ be normal. Then $\|A\|_\infty = \rho(A)$.*

Proof. If A is normal, it is diagonalizable, and can then be written as

$$A = \sum_{\lambda_i \in \text{Spec}(A)} \lambda_i \Pi_i, \quad \Pi_i \Pi_j = \delta_{ij} \Pi_i, \quad \sum_{\lambda_i \in \text{Spec}(A)} \Pi_i = I_A \quad (13.7)$$

where Π_i denotes the projection on the eigenspace corresponding to the eigenvalue λ_i . For $\|\psi\|_A = 1$ one then has

$$\begin{aligned} \|A\psi\|_A &= \left(\sum_{\lambda_i \in \text{Spec}(A)} |\lambda_i|^2 \langle \psi | \Pi_i | \psi \rangle \right)^{\frac{1}{2}} \\ &\leq \rho(A) \left(\sum_{\lambda_i \in \text{Spec}(A)} \langle \psi | \Pi_i | \psi \rangle \right)^{\frac{1}{2}} \\ &= \rho(A) \|\psi\|_A. \end{aligned}$$

Taking the sup on l.h.s. we have $\|A\|_\infty \leq \rho(A)$. Along with lemma 13.3, this implies $\|A\|_\infty = \rho(A)$. \square

In the case of selfadjoint operators, we have also the following useful result.

Lemma 13.5. *Let $A = A^\dagger : \mathcal{H}_A \rightarrow \mathcal{H}_A$ be a selfadjoint operator. Then*

$$-\|A\|_\infty I_A \leq A \leq \|A\|_\infty I_A. \quad (13.8)$$

Proof. Since A is normal, by lemma 13.4 we have $\rho(A) = \|A\|_\infty$. Moreover, the eigenvalues of A are real: $\text{Spec}(A) \subseteq \mathbb{R}$. Then we have $-\rho(A) \leq -|\lambda_i| \leq \lambda_i \leq |\lambda_i| \leq \rho(A)$ for all $\lambda_i \in \text{Spec}(A)$. Upon diagonalizing A as $A = \sum_{\lambda_i \in \text{Spec}(A)} \lambda_i \Pi_i$, this finally implies

$$\begin{aligned} -\|A\|_\infty I_A &= \sum_{\lambda_i \in \text{Spec}(A)} [-\rho(A)] \Pi_i \\ &\leq \sum_{\lambda_i \in \text{Spec}(A)} \lambda_i \Pi_i \\ &\leq \sum_{\lambda_i \in \text{Spec}(A)} \rho(A) \Pi_i \\ &= \|A\|_\infty I_A. \end{aligned}$$
 \square

The last result has a relevant consequence on the limit of Cauchy sequences of operators in the uniform norm.

Definition 13.6. Let $\{A_n\}_{n \in \mathbb{N}} \subseteq \mathcal{L}(\mathcal{H}_A)$ be a sequence of operators. We say that the sequence uniformly converges to A , in formula $\lim_{n \rightarrow \infty} A_n = A$, if, $\forall \varepsilon > 0$, $\exists n_0$ such that $\forall n > n_0$

$$\|A_n - A\|_\infty < \varepsilon. \quad (13.9)$$

We now prove that if a sequence $\{A_n\}_{n \in \mathbb{N}}$ of selfadjoint operators uniformly converges to A , then also A is selfadjoint.

Lemma 13.7. Let $\{A_n\}_{n \in \mathbb{N}}$ be a sequence of selfadjoint operators $A_n^\dagger = A_n$, that uniformly converges to A . Then $A = A^\dagger$.

Proof. Let us consider the polar decomposition $B = U|B|$ of a general operator B , with $U : \mathcal{H}_A \rightarrow \mathcal{H}_A$ unitary. Then one has $\|B^\dagger\|_\infty = \|U^\dagger B U^\dagger\|_\infty = \|B\|_\infty$. Thus,

$$\begin{aligned} \|A_n - A^\dagger\|_\infty &= \|A_n^\dagger - A^\dagger\|_\infty \\ &= \|A_n - A\|_\infty, \end{aligned}$$

and thus A_n converges to A^\dagger . By uniqueness of the limit we finally have $A = A^\dagger$. \square

Lemma 13.5 allows us now to prove the following theorem.

Theorem 13.8. If the sequence $\{A_n\}_{n \in \mathbb{N}} \subseteq \mathcal{L}(\mathcal{H}_A)$ of selfadjoint operators uniformly converges to A , and for all $n > \bar{n}$ we have $A_n \geq 0$, then also $A \geq 0$.

Proof. By definition of uniform convergence, for every $\varepsilon > 0$ there exists n_0 such that for $n > n_0$ we have $\|A_n - A\|_\infty < \varepsilon$. This implies that

$$-\varepsilon I_A \leq A_n - A \leq \varepsilon I_A,$$

and finally, taking the expectation value of all members on the arbitrary vector $\psi \in \mathcal{H}_A$, we have

$$|\langle \psi | A_n | \psi \rangle - \langle \psi | A | \psi \rangle| \leq \varepsilon.$$

Finally, since a sequence of positive real numbers $\{a_n\}_{n \in \mathbb{N}}$ such that for $n > \bar{n}$ $a_n \geq 0$ converges to a positive $a \geq 0$, we conclude that also $\langle \psi | A | \psi \rangle \geq 0$. \square

13.3 Lieb's theorem

In this lecture we prove Lieb's theorem and subsequently we will use it to prove some key properties of von Neumann entropy. For this purpose, we need some preliminary lemmas. The next few results regard partitioned matrices. Let the Hilbert space \mathcal{H}_A be partitioned into orthogonal subspaces $\mathcal{H}_A = \bigoplus_{r=1}^b \mathcal{H}_r$. If we define the projections $\{P_r\}_{r=1}^b$ on the b subspaces, then any operator $A \in \mathcal{L}(\mathcal{H}_A)$ can be partitioned into blocks $A^{(p,q)} := P_p A P_q$. The matrix representing A in a basis that is a collection made of one orthonormal basis for each subspace \mathcal{H}_r is then a partitioned matrix, with blocks $A^{(p,q)}$. In the following, when we write *partitioned matrix* the partitioning of the Hilbert space \mathcal{H}_A will be implicit.

Lemma 13.9. Let $A, B : \mathcal{H}_A \rightarrow \mathcal{H}_A$ be partitioned matrices. Then the following identities hold

$$(AB)^{(p,q)} = \sum_{r=1}^b A^{(p,r)} B^{(r,q)}, \quad (A^\dagger)^{(p,q)} = A^{(q,p)\dagger}. \quad (13.10)$$

Proof. We can write $(AB)^{(p,q)} = P_p A B P_q$ as

$$P_p A I_A B P_q = \sum_{r=1}^b P_p A P_r B P_q = \sum_{r=1}^b A^{(p,r)} B^{(r,q)}.$$

Moreover, $(A^\dagger)^{(p,q)} = P_p A^\dagger P_q = (P_q A P_p)^\dagger = A^{(q,p)\dagger}$. \square

For positive partitioned matrices the following result holds.

Lemma 13.10. Let $A : \mathcal{H}_A \rightarrow \mathcal{H}_A$ be a partitioned matrix. Then A is positive if and only if there exists a Hilbert space \mathcal{H}_B and a set $\{T^{(r)}\}_{r=1}^b \subseteq \mathcal{L}(\mathcal{H}_A, \mathcal{H}_B)$ such that $\text{Supp}(T^{(r)}) \subseteq \mathcal{H}_r$ and $A^{(p,q)} = T^{(p)\dagger} T^{(q)}$.

Proof. First suppose that $A^{(p,q)} = T^{(p)\dagger} T^{(q)}$. Define $X := \sum_{p=1}^b T^{(p)}$. It is immediate to verify that $A = \sum_{p,q=1}^b A^{(p,q)}$ can be written as $A = X^\dagger X \geq 0$. Conversely, if $A \geq 0$ then $A = X^\dagger X$ for some $X : \mathcal{H}_A \rightarrow \mathcal{H}_A$, and by lemma 13.9 we have

$$A^{(p,q)} = \sum_{r=1}^b X^{(r,p)\dagger} X^{(r,q)}.$$

Notice that $X^{(r,p)\dagger} X^{(s,q)} = (X^\dagger)^{(p,r)} X^{(s,q)} = \delta_{rs} X^{(r,p)\dagger} X^{(s,q)}$. If we now define $T^{(p)} := \sum_{r=1}^b X^{(r,p)}$, we then have $A^{(p,q)} = T^{(p)\dagger} T^{(q)}$. \square

Lemma 13.11. For a partitioned matrix $A : \mathcal{H}_A \rightarrow \mathcal{H}_A$, A is positive iff $A^{(p,p)} \geq 0$ and $A^{(p,q)} = (A^{(p,p)})^{\frac{1}{2}} U_p^\dagger U_q (A^{(q,q)})^{\frac{1}{2}}$, where $U_p : \mathcal{H}_A \rightarrow \mathcal{H}_A$ is unitary for every p .

Proof. By lemma 13.10, if $A \geq 0$ one has $A^{(p,p)} = T^{(p)\dagger} T^{(p)}$ with $T^{(p)} \in \mathcal{L}(\mathcal{H}_A)$, and then by the polar decomposition $T^{(p)} = U_p (A^{(p,p)})^{\frac{1}{2}}$, where U_p is unitary. Then $A^{(p,q)} = T^{(p)\dagger} T^{(q)} = (A^{(p,p)})^{\frac{1}{2}} U_p^\dagger U_q (A^{(q,q)})^{\frac{1}{2}}$. On the other hand, if $A^{(p,q)} = (A^{(p,p)})^{\frac{1}{2}} U_p^\dagger U_q (A^{(q,q)})^{\frac{1}{2}}$, one can define $T^{(p)} := U_p (A^{(p,p)})^{\frac{1}{2}}$, and clearly

$$A = \sum_{p,q} A^{(p,q)} = \sum_{p,q} T^{(p)\dagger} T^{(q)}. \quad \square$$

Lemma 13.12. Let $M : \mathcal{H} \rightarrow \mathcal{H}$ with $\mathcal{H} = \mathcal{H}_0 \oplus \mathcal{H}_1$ and $\mathcal{H}_0 \sim \mathcal{H}_1 \sim \mathcal{K}$. Let M correspond to a matrix of the form

$$M = \left(\begin{array}{c|c} X & H \\ \hline H & Y \end{array} \right), \quad (13.11)$$

with $H, X, Y \in \mathcal{L}(\mathcal{K})$. Let $H = H^\dagger$ and $[X, Y] = 0$. If $M \geq 0$ then $X, Y \geq 0$ and $H \leq X^{\frac{1}{2}} Y^{\frac{1}{2}}$.

Proof. Let $W_i : \mathcal{K} \rightarrow \mathcal{H}$ be the isometry such that $W_i \mathcal{K} = \mathcal{H}_i$. Thus, $W_0 X W_0^\dagger = M^{(0,0)}$, $W_1 Y W_1^\dagger = M^{(1,1)}$, $W_1 H W_0^\dagger = M^{(1,0)}$, and $W_0 H W_1^\dagger = M^{(0,1)}$. By lemma 13.11, M is positive if and only if $X, Y \geq 0$ and

$$\begin{aligned} H &= W_0^\dagger M^{(0,1)} W_1 \\ &= W_0^\dagger M^{(0,0)\frac{1}{2}} U_0^\dagger U_1 M^{(1,1)\frac{1}{2}} W_1 \\ &= X^{\frac{1}{2}} V_0^\dagger V_1 Y^{\frac{1}{2}}, \end{aligned}$$

where $V_i := U_i W_i$, and $\{U_i\}$ are the unitary operators of lemma 13.11. Moreover, since $H = H^\dagger$, we have $H = X^{\frac{1}{2}} V_0^\dagger V_1 Y^{\frac{1}{2}} = Y^{\frac{1}{2}} V_1^\dagger V_0 X^{\frac{1}{2}}$. Now, denoting by Z^{-1} the generalized inverse of Z , and defining $Q := X^{-\frac{1}{2}} H Y^{-\frac{1}{2}}$ we have

$$\rho(Q) \leq \|Q\|_\infty = \|P_X V_0^\dagger V_1 P_Y\|_\infty \leq \|P_X\|_\infty \|V_0^\dagger\|_\infty \|V_1\|_\infty \|P_Y\|_\infty = 1,$$

where P_Z denotes the projection on the support $\text{Supp}(Z)$ of Z . Notice that $P_X H P_Y = P_Y H P_X = H$, and then $P_H \leq P_X$ and $P_H \leq P_Y$. This implies that both X and Y are invertible on $\text{Supp}(H)$. As a consequence, the selfadjoint operator $R := X^{-\frac{1}{4}} Y^{-\frac{1}{4}} H Y^{-\frac{1}{4}} X^{-\frac{1}{4}}$ can be written as $R = Y^{-\frac{1}{4}} X^{\frac{1}{4}} Q X^{-\frac{1}{4}} Y^{\frac{1}{4}} = T Q T^{-1}$, where

$$T := Y^{-\frac{1}{4}} X^{\frac{1}{4}} + W(I - P_{Y^{-\frac{1}{4}} X^{\frac{1}{4}}}),$$

and $W \text{Ker}(Y^{-\frac{1}{4}} X^{\frac{1}{4}}) = \text{Co} - \text{Ker}(Y^{-\frac{1}{4}} X^{\frac{1}{4}})$. The operator T is invertible. Since $R = T Q T^{-1}$ for an invertible T , it holds that $\text{Spec}(R) = \text{Spec}(Q)$, and then $\rho(R) = \rho(Q)$. Finally, since R is selfadjoint, we have

$$\|R\|_\infty = \rho(R) \leq 1.$$

Now, by lemma 13.5 we conclude $R \leq I_A$, and finally

$$H = X^{\frac{1}{4}} Y^{\frac{1}{4}} R Y^{\frac{1}{4}} X^{\frac{1}{4}} \leq X^{\frac{1}{4}} Y^{\frac{1}{2}} X^{\frac{1}{4}} = X^{\frac{1}{2}} Y^{\frac{1}{2}}. \quad \square$$

Lemma 13.13. *Let $\mathcal{L}(\mathcal{H}_A) \ni A_0, A_1 \geq 0$, $\mathcal{L}(\mathcal{H}_B) \ni B_0, B_1 \geq 0$. Then for every $p \in [0, 1]$ we have*

$$(A_0 + A_1)^p \otimes (B_0 + B_1)^{1-p} \geq A_0^p \otimes B_0^{1-p} + A_1^p \otimes B_1^{1-p}. \quad (13.12)$$

Proof. Let $X_0(p) := A_0^p \otimes B_0^{1-p}$, $X_1(p) := A_1^p \otimes B_1^{1-p}$, and $X_2(p) := (A_0 + A_1)^p \otimes (B_0 + B_1)^{1-p}$. It is clear that

$$[X_i(p), X_i(q)] = 0, \quad \forall i = 0, 1, 2, \quad \forall p, q \in [0, 1].$$

Moreover,

$$(X_i(p) X_i(q))^{\frac{1}{2}} = X_i(\frac{p+q}{2}), \quad \forall i = 0, 1, 2.$$

Let us now consider the partitioned matrices $Y_0(p, q)$ and $Y_1(p, q)$ defined as

$$Y_i(p, q) := \left(\begin{array}{c|c} X_i(p) & X_i(\frac{p+q}{2}) \\ \hline X_i(\frac{p+q}{2}) & X_i(q) \end{array} \right) = \left(\begin{array}{c|c} X_i(p)^{\frac{1}{2}} & 0 \\ \hline X_i(q)^{\frac{1}{2}} & 0 \end{array} \right) \left(\begin{array}{c|c} X_i(p)^{\frac{1}{2}} & X_i(q)^{\frac{1}{2}} \\ \hline 0 & 0 \end{array} \right).$$

It is clear that $Y_i(p, q) \geq 0$ for all $i = 0, 1$ and for all $p, q \in [0, 1]$. Moreover, if $X_2(p) \geq X_0(p) + X_1(p)$ and $X_2(q) \geq X_0(q) + X_1(q)$, then

$$\begin{aligned} & \left(\begin{array}{c|c} X_2(p) - X_0(p) - X_1(p) & 0 \\ \hline 0 & X_2(q) - X_0(q) - X_1(q) \end{array} \right) + Y_0(p, q) + Y_1(p, q) = \\ & \left(\begin{array}{c|c} X_2(p) & X_0(\frac{p+q}{2}) + X_1(\frac{p+q}{2}) \\ \hline X_0(\frac{p+q}{2}) + X_1(\frac{p+q}{2}) & X_2(q) \end{array} \right) \geq 0. \end{aligned}$$

By lemma 13.12 we then have $X_0(\frac{p+q}{2}) + X_1(\frac{p+q}{2}) \leq (X_2(p)X_2(q))^{\frac{1}{2}} = X_2(\frac{p+q}{2})$. It is easy to check that for $p = 0$ and $p = 1$ we have $X_2(p) \geq X_0(p) + X_1(p)$, since $A_i^0 = P_{A_i} \leq P_{A_0+A_1} = (A_0 + A_1)^0$, and $B_i^0 = P_{B_i} \leq P_{B_0+B_1} = (B_0 + B_1)^0$. Iterating this argument, we find that the thesis is true for all $p \in J := \{\frac{n}{2^k} \mid k \in \mathbb{N}, 0 \leq n \leq 2^k\}$, which is dense in $[0, 1]$. In order to complete the proof, it is then sufficient to prove continuity of the function $p \mapsto \Delta(p) := X_2(p) - X_0(p) - X_1(p)$ in the uniform topology. This can be proved considering that for every $\varepsilon > 0$ there exists δ such that if $|p - q| < \delta$ we have

$$\begin{aligned} \|\Delta(p) - \Delta(q)\|_\infty & \leq \sum_{i=0,1,2} \|X_i(p) - X_i(q)\|_\infty \\ & \leq \sum_{i=0,1,2} \|X_i(p)\|_\infty \|I - A_i^{q-p} \otimes B_i^{p-q}\|_\infty \\ & \leq \sum_{i=0,1,2} \rho(X_i(p)) |1 - k_i^{q-p}| < \varepsilon, \end{aligned}$$

where we defined $A_2 := A_0 + A_1$ and $B_2 := B_0 + B_1$. Now, since for any sequence $\{p_n\}_{n \in \mathbb{N}} \subseteq J$ converging to $p \in [0, 1]$ we have $\Delta(p_n) \geq 0$, by theorem 13.8 we also have $\Delta(p) \geq 0$ for all $p \in [0, 1]$. \square

Using the last result, we can prove Lieb's theorem. We introduce the set $\mathcal{P}(A) := \{X \in \mathcal{L}(A) \mid X \geq 0\} = \{\lambda\rho \mid 0 \leq \lambda \in \mathbb{R}, \rho \in \text{St}_1(A)\}$

Theorem 13.14 (Lieb's theorem). *Let $F_p^{(K)} : \mathcal{P}(A) \times \mathcal{P}(B) \rightarrow \mathbb{R}$ be defined as $F_p^{(K)}(A, B) := \text{Tr}[A^p K B^{1-p} K^\dagger]$ for an arbitrary $K : \mathcal{H}_B \rightarrow \mathcal{H}_A$. The function F_p is jointly concave in A, B for any K and any $p \in [0, 1]$, i.e.*

$$F_p^{(K)}(\lambda A_1 + (1-\lambda)A_2, \lambda B_1 + (1-\lambda)B_2) \geq \lambda F_p^{(K)}(A_1, B_1) + (1-\lambda)F_p^{(K)}(A_2, B_2), \quad (13.13)$$

for every $0 \leq \lambda \leq 1$.

Proof. It is sufficient to notice that

$$F_p^{(K)}(A, B) = \langle\langle K | A^p \otimes (B^T)^{1-p} | K \rangle\rangle. \quad (13.14)$$

If we then apply Eq. (13.12) with λA_1 and $(1-\lambda)A_2$ instead of A_1 and A_2 , and λB_1^T and $(1-\lambda)B_2^T$ instead of B_1 and B_2 , by lemma 13.13 we have

$$(\lambda A_1 + (1-\lambda)A_2)^p \otimes (\lambda B_1^T + (1-\lambda)B_2^T)^{1-p} \geq \lambda A_1^p \otimes (B_1^T)^{1-p} + (1-\lambda)A_2^p \otimes (B_2^T)^{1-p}.$$

The thesis then follows reminding Eq. (13.14). \square

Chapter 14

Lecture 16: Monotonicity and Holevo bound

In the following we use Lieb's theorem to prove four crucial results: Joint convexity of quantum relative entropy, Uhlmann's monotonicity, concavity of quantum conditional entropy and strong subadditivity. The logical scheme shown in fig. 14.1 is useful to clarify the importance of Lieb's theorem, as well as the relation between the above mentioned theorems.

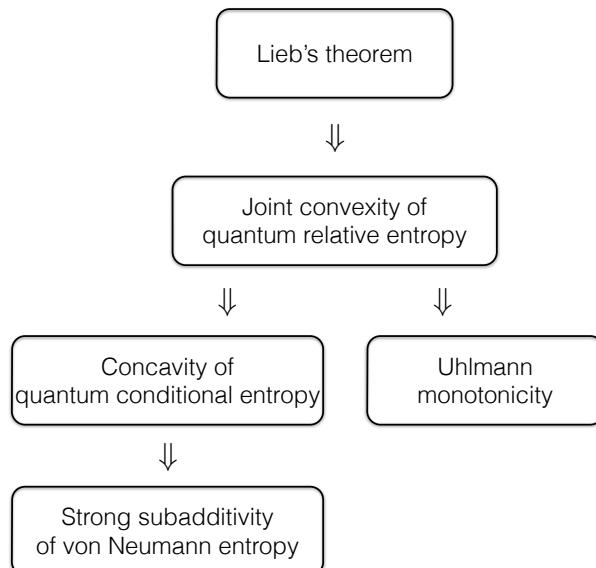


Figure 14.1 Logical dependence of the main results in the chapter

14.1 Joint convexity of quantum relative entropy

Now that we have proved the Lieb's theorem, we are in position to prove the joint convexity of relative entropy. First, we prove an equivalent way of expressing $S(\rho\|\sigma)$.

Definition 14.1. Let us define the function $G_t : \mathcal{P}(A) \times \mathcal{P}(B) \rightarrow \mathbb{R}$

$$G_t(A, B) := \text{Tr}[A^{1-t}B^t] - \text{Tr}[A].$$

Let us also define

$$G(A, B) := \lim_{t \rightarrow 0^+} \frac{G_t(A, B)}{t},$$

where we recall that $\lim_{t \rightarrow 0^+} A^t := P_{\text{Supp}(A)}$ for $A \geq 0$.

Lemma 14.2. *The quantum relative entropy $S(\rho\|\sigma)$ can be expressed as*

$$S(\rho\|\sigma) = -\frac{1}{\ln 2} G(\rho, \sigma). \quad (14.1)$$

Proof. Let us remind definition 12.2 of $S(\rho\|\sigma)$, in particular equation (12.10):

$$S(\rho\|\sigma) := \begin{cases} \text{Tr}[\rho \log_2 \rho] - \text{Tr}[\rho \log_2 \sigma] & \text{Ker}(\sigma) \subseteq \text{Ker}(\rho) \\ \infty & \text{otherwise.} \end{cases}$$

Let us first consider the case $\text{Ker}(\sigma) \not\subseteq \text{Ker}(\rho)$. Since by definition one has $\lim_{t \rightarrow 0^+} A^t := P_{\text{Supp}(A)}$, we have

$$\lim_{t \rightarrow 0^+} G_t(\rho, \sigma) = -\text{Tr}[\rho(I - P_{\text{Supp}(\sigma)})] = -\text{Tr}[\rho P_{\text{Ker}(\sigma)}] < 0.$$

Thus,

$$-\frac{1}{\ln 2} \lim_{t \rightarrow 0^+} \frac{G_t(\rho, \sigma)}{t} = +\infty.$$

Let us now consider the case $\text{Ker}(\sigma) \subseteq \text{Ker}(\rho)$. In this case we have

$$\begin{aligned} G_0(\rho, \sigma) &= -\text{Tr}[\rho(I - P_{\text{Supp}(\sigma)})] \\ &= -\text{Tr}[\rho P_{\text{Ker}(\sigma)}] \\ &= 0, \end{aligned}$$

and thus

$$\begin{aligned} \lim_{t \rightarrow 0^+} \frac{G_t(\rho, \sigma)}{t} &= \lim_{t \rightarrow 0^+} \frac{G_t(\rho, \sigma) - G_0(\rho, \sigma)}{t} \\ &= \lim_{t \rightarrow 0^+} \frac{d}{dt} \text{Tr}[\rho^{1-t} \sigma^t] \\ &= -\lim_{t \rightarrow 0^+} \{\text{Tr}[\ln \rho \rho^{1-t} \sigma^t] - \text{Tr}[\rho^{1-t} \ln \sigma \sigma^t]\}. \end{aligned}$$

Finally, this implies

$$-\frac{1}{\ln 2} G(\rho, \sigma) = \text{Tr}[\rho \log_2 \rho] - \text{Tr}[\rho \log_2 \sigma]. \quad \square$$

Theorem 14.3 (Joint convexity of quantum relative entropy). *The relative entropy $S(\rho\|\sigma)$ is jointly convex in ρ and σ , that is*

$$S(p\rho_1 + (1-p)\rho_2\|p\sigma_1 + (1-p)\sigma_2) \leq pS(\rho_1\|\sigma_1) + (1-p)S(\rho_2\|\sigma_2), \quad (14.2)$$

for all $p \in [0, 1]$.

Proof. If we use the expression 14.1, we have

$$S(\rho\|\sigma) = -\frac{1}{\ln 2} \lim_{t \rightarrow 0^+} \frac{\text{Tr}[\rho^{1-t}\sigma^t] - \text{Tr}[\rho]}{t}.$$

We can write

$$G_t(\rho, \sigma) = \text{Tr}[\rho^{1-t}\sigma^t] - \text{Tr}[\rho] = \langle\langle I | \rho^{1-t} \otimes (\sigma^T)^t | I \rangle\rangle - \text{Tr}[\rho],$$

and using now Lieb's theorem 13.14, we have

$$G_t(p\rho_1 + (1-p)\rho_2, p\sigma_1 + (1-p)\sigma_2) \geq pG_t(\rho_1, \sigma_1) + (1-p)G_t(\rho_2, \sigma_2).$$

Then, we have

$$\begin{aligned} S(p\rho_1 + (1-p)\rho_2\|p\sigma_1 + (1-p)\sigma_2) &= -\frac{1}{\ln 2} \lim_{t \rightarrow 0^+} \frac{G_t(p\rho_1 + (1-p)\rho_2, p\sigma_1 + (1-p)\sigma_2)}{t} \\ &\leq -\frac{1}{\ln 2} \lim_{t \rightarrow 0^+} \frac{pG_t(\rho_1, \sigma_1) + (1-p)G_t(\rho_2, \sigma_2)}{t} \\ &= p \left\{ -\frac{1}{\ln 2} \lim_{t \rightarrow 0^+} \frac{G_t(\rho_1, \sigma_1)}{t} \right\} + (1-p) \left\{ -\frac{1}{\ln 2} \lim_{t \rightarrow 0^+} \frac{G_t(\rho_2, \sigma_2)}{t} \right\} \\ &= pS(\rho_1\|\sigma_1) + (1-p)S(\rho_2\|\sigma_2). \end{aligned} \quad \square$$

14.2 Concavity of quantum conditional entropy

From joint convexity of relative entropy, concavity of the conditional entropy follows.

Theorem 14.4 (Concavity of quantum conditional entropy). *The quantum conditional entropy is concave versus the joint state ρ_{AB} , that is*

$$S_{p\rho^{(1)} + (1-p)\rho^{(2)}}(A|B) \geq pS_{\rho^{(1)}}(A|B) + (1-p)S_{\rho^{(2)}}(A|B). \quad (14.3)$$

Proof. Let us use the following identity

$$\begin{aligned} S(\rho_{AB} \| \frac{I_A}{d_A} \otimes \rho_B) &= -S(AB) - \text{Tr}[\rho_{AB} \log_2(\frac{I_A}{d_A} \otimes \rho_B)] \\ &= -S(AB) - \text{Tr}[\rho_{AB} \{ \log_2(\frac{I_A}{d_A}) \otimes I_B \}] - \text{Tr}[\rho_{AB} \{ I_A \otimes \log_2(\rho_B) \}] \\ &= -S(AB) + S(B) + \log_2 d_A \\ &= -S(A|B) + \log_2 d_A. \end{aligned}$$

Since the partial trace Tr_A is linear, we have $\text{Tr}_B[p\rho_{AB}^{(1)} + (1-p)\rho_{AB}^{(2)}] = p\rho_B^{(1)} + (1-p)\rho_B^{(2)}$. Then, using joint convexity of quantum relative entropy, one has

$$\begin{aligned} S_{p\rho^{(1)} + (1-p)\rho^{(2)}}(A|B) &= \log_2 d_A - S(p\rho_{AB}^{(1)} + (1-p)\rho_{AB}^{(2)} \| p(\frac{I_A}{d_A} \otimes \rho_B^{(1)}) + (1-p)(\frac{I_A}{d_A} \otimes \rho_B^{(2)})) \\ &\geq \log_2 d_A - pS(\rho_{AB}^{(1)} \| \frac{I_A}{d_A} \otimes \rho_B^{(1)}) - (1-p)S(\rho_{AB}^{(2)} \| \frac{I_A}{d_A} \otimes \rho_B^{(2)}) \\ &= p(\log_2 d_A - S(\rho_{AB}^{(1)} \| \frac{I_A}{d_A} \otimes \rho_B^{(1)})) + (1-p)(\log_2 d_A - S(\rho_{AB}^{(2)} \| \frac{I_A}{d_A} \otimes \rho_B^{(2)})) \\ &= pS_{\rho^{(1)}}(A|B) + (1-p)S_{\rho^{(2)}}(A|B) \end{aligned} \quad \square$$

14.3 Strong subadditivity of von Neumann entropy

Using concavity of conditional entropy we can now prove the strong subadditivity property of the von Neumann entropy.

Theorem 14.5 (Strong subadditivity of von Neumann entropy). *For any three quantum systems A, B, and C, one has*

$$S(A) + S(B) \leq S(AC) + S(BC), \quad (14.4)$$

or equivalently

$$S(ABC) + S(B) \leq S(AB) + S(BC). \quad (14.5)$$

Proof. The two inequalities are equivalent, as we will see. Define the function

$$T(\rho_{ABC}) := -S(C|A) - S(C|B) = S(A) + S(B) - S(AC) - S(BC),$$

Since conditional entropies are concave versus their joint state, the function T is convex versus ρ_{ABC} . Then, it follows that for $\rho_{ABC} = \sum_i p_i |i\rangle\langle i|$ one has

$$T(\rho_{ABC}) \leq \sum_i p_i T(|i\rangle\langle i|).$$

Now, T is zero on pure states, since in this case $S(AC) = S(B)$ and $S(BC) = S(A)$. Thus, one has $T(\rho_{ABC}) \leq 0$, namely

$$S(A) + S(B) - S(AC) - S(BC) \leq 0, \quad (14.6)$$

which is the first inequality of the statement. To obtain the second inequality, let us purify ρ_{ABC} by adding a reference system R. Inequality (14.6) for systems R, B, C is

$$S(R) + S(B) \leq S(RC) + S(BC),$$

but since ρ_{RABC} is now pure, one has $S(R) = S(ABC)$ and $S(RC) = S(AB)$, and the last inequality becomes

$$S(ABC) + S(B) \leq S(AB) + S(BC),$$

which is the second inequality of the statement. \square

Theorem 14.6 (Conditioning reduces the quantum entropy).

$$S(A|BC) \leq S(A|B). \quad (14.7)$$

Proof. Using the definition of conditional entropy, the statement of the theorem is equivalent to

$$S(ABC) - S(BC) \leq S(AB) - S(B),$$

which is just strong subadditivity. \square

Theorem 14.7 (Discarding systems and quantum mutual information). *Discarding quantum systems never increases the quantum mutual information, namely*

$$S(A : B) \leq S(A : BC). \quad (14.8)$$

Proof. Using the definition of mutual information, we see that the statement is equivalent to

$$S(A) + S(B) - S(AB) \leq S(A) + S(BC) - S(ABC),$$

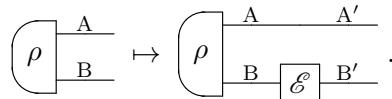
namely

$$S(ABC) + S(B) \leq S(AB) + S(BC),$$

which is strong subadditivity. \square

Theorem 14.8 (Local channels and quantum mutual information). *Local quantum channels never increase the quantum mutual information.*

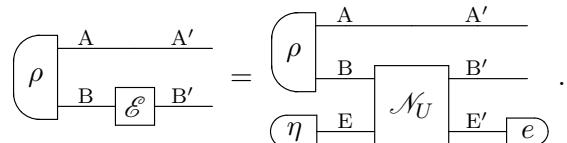
Proof. Denote by A' and B' the two quantum systems A and B after the action of the local channel \mathcal{E} on B , as in the following diagram



We want to prove that

$$S(A' : B') \leq S(A : B).$$

By the unitary dilation theorem, the action of \mathcal{E} on B is equivalent to a unitary interaction \mathcal{N}_U with another system E which is then discarded, as represented in the following diagram



Before the action of \mathcal{N}_U , B and E are uncorrelated, then

$$\begin{aligned} S(A : BE) &= S(A) + S(BE) - S(ABE) \\ &= S(A) + S(B) + S(E) - S(AB) - S(E) = S(A : B). \end{aligned}$$

Invariance of von Neumann entropy under unitary transformations leads to

$$S(A' : B'E') = S(A : BE).$$

Finally, discarding a quantum system cannot increase the mutual information, thus

$$S(A' : B') \leq S(A' : B'E') = S(A : BE) = S(A : B). \quad \square$$

14.4 Uhlmann monotonicity theorem

In this section we derive the last of the main consequences of Lieb's theorem: Uhlmann monotonicity theorem for joint entropy. In particular, the proof uses joint convexity of quantum relative entropy.

The first lemma that we need involves the so-called *shift-and-multiply group* for a system A. This is a group of d_A^2 unitary operators $\{U_{p,q}\}_{p,q=0}^{d_A-1}$ that are given in the canonical basis $\{\varphi_i\}_{i=0}^{d_A-1}$ as $U_{p,q} := Z^p W^q$ in terms of the phase operator W and shift operator Z , defined as

$$\begin{aligned} W &:= \sum_{i=0}^{d_A-1} \omega_A^i |\varphi_i\rangle\langle\varphi_i|, & \omega_A &:= \exp(-i\frac{2\pi}{d_A}), \\ Z &:= \sum_{i=0}^{d_A-1} |\varphi_{i\oplus 1}\rangle\langle\varphi_i|, & a \oplus b &:= a + b \pmod{d_A}. \end{aligned} \quad (14.9)$$

The set $U_{p,q}$ provides a *projective unitary representation* of the group $\mathbb{Z}_{d_A} \times \mathbb{Z}_{d_A}$. The property we are interested in is given by the following lemma.

Lemma 14.9. *Let $\{U_{p,q}\}_{p,q=0}^{d_A-1}$ be the shift-and-multiply group for a system A with dimension d_A . then the following identity holds*

$$\frac{1}{d_A^2} \sum_{p,q=0}^{d_A-1} U_{p,q} \otimes U_{p,q}^* = \frac{1}{d_A} |I_A\rangle\langle I_A|. \quad (14.10)$$

Proof. By straightforward algebra we have

$$\begin{aligned} \frac{1}{d_A^2} \sum_{p,q=0}^{d_A-1} U_{p,q} \otimes U_{p,q}^* &= \frac{1}{d_A^2} \sum_{p,q=0}^{d_A-1} Z^p W^q \otimes Z^p W^{-q} \\ &= \frac{1}{d_A^2} \sum_{p,q=0}^{d_A-1} \sum_{j,k=0}^{d_A-1} \omega_A^{(j-k)q} |\varphi_{j\oplus p}\rangle\langle\varphi_j| \otimes |\varphi_{k\oplus p}\rangle\langle\varphi_k|. \end{aligned}$$

Now, since

$$\frac{1}{d_A} \sum_{q=0}^{d_A-1} \omega_A^{(j-k)q} = \begin{cases} 1 & j = k, \\ 0 & j \neq k, \end{cases}$$

we conclude

$$\begin{aligned} \frac{1}{d_A^2} \sum_{p,q=0}^{d_A-1} U_{p,q} \otimes U_{p,q}^* &= \frac{1}{d_A} \sum_{p=0}^{d_A-1} \sum_{j=0}^{d_A-1} |\varphi_{j\oplus p}\rangle\langle\varphi_j| \otimes |\varphi_{j\oplus p}\rangle\langle\varphi_j| \\ &= \frac{1}{d_A} \sum_{p=0}^{d_A-1} |\varphi_p\rangle\langle\varphi_p| \sum_{j=0}^{d_A-1} \langle\varphi_j|\langle\varphi_j| \\ &= \frac{1}{d_A} |I_A\rangle\langle I_A|. \end{aligned} \quad \square$$

Corollary 14.10. *For any system A and every operator $X : \mathcal{H}_A \rightarrow \mathcal{H}_A$ the following identity holds*

$$\frac{1}{d_A^2} \sum_{p,q=0}^{d_A-1} U_{p,q} X U_{p,q}^\dagger = \text{Tr}[X] \frac{I_A}{d_A} \quad (14.11)$$

Proof. Let us use the Vec isomorphism. We have

$$\begin{aligned} \left| \frac{1}{d_A^2} \sum_{p,q=0}^{d_A-1} U_{p,q} X U_{p,q}^\dagger \right\rangle &= \frac{1}{d_A^2} \sum_{p,q=0}^{d_A-1} U_{p,q} \otimes U_{p,q}^* |X\rangle \\ &= \frac{1}{d_A} |I_A\rangle \langle I_A| X \rangle \\ &= \frac{\text{Tr}[X]}{d_A} |I_A\rangle \langle I_A|. \end{aligned} \quad \square$$

We can now prove the following result.

Lemma 14.11 (Monotonicity of $S(\rho\|\sigma)$ under marginalisation). *Let ρ and σ be states of a composite system AB. The quantum relative entropy is monotonically decreasing when ignoring one system, namely*

$$S(\rho_A\|\sigma_A) \leq S(\rho_{AB}\|\sigma_{AB}). \quad (14.12)$$

Proof. Using the result (14.11) of corollary 14.10, we can write the partial trace Tr_B as follows

$$\tau_A \otimes \frac{I}{d_B} = \frac{1}{d_B^2} \sum_{p,q=0}^{d_B-1} (I_A \otimes U_{p,q}) \tau_{AB} (I_A \otimes U_{p,q}^\dagger).$$

Using properties (13.1) and (13.2), along with the convexity of the relative entropy we have

$$\begin{aligned} S(\rho_A\|\sigma_A) &= S\left(\rho_A \otimes \frac{I_B}{d_B} \|\sigma_A \otimes \frac{I_B}{d_B}\right) \\ &\leq \frac{1}{d_B^2} \sum_{p,q=0}^{d_B-1} S((I_A \otimes U_{p,q}) \rho_{AB} (I_A \otimes U_{p,q}^\dagger) \| (I_A \otimes U_{p,q}) \sigma_{AB} (I_A \otimes U_{p,q}^\dagger)) \\ &= S(\rho_{AB}\|\sigma_{AB}). \end{aligned} \quad \square$$

We can finally prove Uhlmann's monotonicity theorem.

Theorem 14.12 (Uhlmann monotonicity of relative entropy). *The quantum relative entropy is monotonically non-increasing under the action of a channel, namely for every $\mathcal{E} : A \rightarrow B$ one has*

$$S(\mathcal{E}(\rho)\|\mathcal{E}(\sigma)) \leq S(\rho\|\sigma). \quad (14.13)$$

Proof. The inequality in equation 14.13 can be simply proved using properties (13.1), (13.2), and (14.12). From the unitary dilation theorem for channels we know indeed that there exist systems C and D such that $\mathcal{H}_A \otimes \mathcal{H}_C \simeq \mathcal{H}_B \otimes \mathcal{H}_D$, a unitary operator $U : \mathcal{H}_A \otimes \mathcal{H}_C \rightarrow \mathcal{H}_B \otimes \mathcal{H}_D$, and a (pure) state $\nu \in \text{St}_1(C)$, such that the channel \mathcal{E} can be written as $\mathcal{E}(\rho) = \text{Tr}_D[U(\rho \otimes \nu)U^\dagger]$. Therefore, one has

$$\begin{aligned} S(\mathcal{E}(\rho) \| \mathcal{E}(\sigma)) &= S(\text{Tr}_D[U(\rho \otimes \nu)U^\dagger] \| \text{Tr}_D[U(\sigma \otimes \nu)U^\dagger]) \\ &\leq S(U(\rho \otimes \nu)U^\dagger \| U(\sigma \otimes \nu)U^\dagger) \\ &= S(\rho \otimes \nu \| \sigma \otimes \nu) \\ &= S(\rho \| \sigma). \end{aligned}$$

□

14.5 The Holevo bound

We now use Uhlmann's monotonicity theorem to prove the Holevo bound. This proof exploits the channel associated with a POVM. This notion is a particular case of channel associated with a quantum instrument. This kind of channel describes a transformation

$$\text{quantum} \rightarrow \text{quantum} + \text{classical},$$

namely a transformation that consists in a measurement with corresponding state reduction along with a classical record of the outcome. The precise mathematical definition is the following.

Definition 14.13 (Channel associated with an instrument). Let $\{\mathcal{A}_i\}_{i \in X} \subseteq \text{QO}(A \rightarrow B)$ be a quantum instrument. Let C be a system with $d_C = |X|$, and $\{\varphi_i\}_{i \in X}$ an orthonormal basis for $\mathcal{H}_C \simeq \mathbb{C}^{|X|}$. The channel associated with the instrument $\{\mathcal{A}_i\}_{i \in X}$ is the channel $\mathcal{C}_{\mathcal{A}} : A \rightarrow BC$

$$\mathcal{C}_{\mathcal{A}}(\rho) := \sum_{i \in X} \mathcal{A}_i(\rho) \otimes |\varphi_i\rangle\langle\varphi_i|. \quad (14.14)$$

When we discard the reduced state of system B, and only consider the outcomes, the instrument becomes a POVM, and the associated channel becomes

$$\text{quantum} \rightarrow \text{classical}.$$

The precise definition of a channel associated with a POVM is the following.

Definition 14.14 (Channel associated with a POVM). Let $\mathbf{P} := \{P_i\}_{i \in X} \subseteq \text{Eff}(A)$ be a POVM. Let C be a system with $d_C = |X|$, and $\{\varphi_i\}_{i \in X}$ an orthonormal basis for $\mathcal{H}_C \simeq \mathbb{C}^{|X|}$. The channel associated with the POVM \mathbf{P} is the channel $\mathcal{C}_{\mathbf{P}} : A \rightarrow C$

$$\mathcal{C}_{\mathbf{P}}(\rho) := \sum_{i \in X} \text{Tr}[\rho P_i] |\varphi_i\rangle\langle\varphi_i|. \quad (14.15)$$

Remark 19. As one can write the channel $\mathcal{C}_{\mathbf{P}}$ as

$$\mathcal{C}_{\mathbf{P}}(\rho) = \text{Tr}_A[(I_C \otimes \rho^T) \sum_{i \in X} |\varphi_i\rangle\langle\varphi_i| \otimes P_i^T],$$

it is immediate to deduce the Choi representative of $\mathcal{C}_{\mathbf{P}}$ as

$$\text{Ch}(\mathcal{C}_{\mathbf{P}}) = \sum_{i \in X} |\varphi_i\rangle\langle\varphi_i| \otimes P_i^T.$$

Let now $E = \{\rho_i\}_{i \in X}$ be an ensemble. We define $\rho^E := \sum_{i \in X} \rho_i = \sum_{i \in X} p_i \tilde{\rho}_i$, with $p_i := \text{Tr}[\rho_i]$. The *Holevo χ -quantity* of the ensemble E is

$$\chi(E) := \sum_{j \in X} p_j S(\tilde{\rho}_j \| \rho^E) = S(\rho^E) - \sum_{j \in X} p_j S(\tilde{\rho}_j). \quad (14.16)$$

Notice that the equality is obtained from the definition of the quantum relative entropy, considering that since $\rho^E = \sum_{i \in X} p_i \tilde{\rho}_i$ one has $\text{Ker}(\rho^E) \subseteq \text{Ker}(\rho_i)$ for every $i \in X$. The Holevo bound sets an upper bound on the accessible information $\text{Acc}(E)$ in terms of the χ -quantity.

Theorem 14.15 (Holevo bound). *For every POVM $\mathbf{P} = \{P_k\}_{k \in Y}$ and any ensemble of states $E = \{\rho_j\}_{j \in X}$ the mutual information is bounded as follows*

$$I(E, \mathbf{P}) \leq \chi(E). \quad (14.17)$$

Proof. By applying Uhlmann's monotonicity Theorem 14.12 with the channel $\mathcal{C}_{\mathbf{P}}$ associated with the POVM \mathbf{P} of Eq. (14.15) and the states $\tilde{\rho}_j$ and $\rho^E = \sum_{j \in X} p_j \tilde{\rho}_j$, and averaging on j , one obtains

$$\sum_j p_j S(\mathcal{C}_{\mathbf{P}}(\tilde{\rho}_j) \| \mathcal{C}_{\mathbf{P}}(\rho^E)) \leq \sum_{j \in X} p_j S(\tilde{\rho}_j \| \rho^E) = \chi(E).$$

We now prove that the left hand side of the bound is just the mutual information $I(E, \mathbf{P})$. Indeed, defining $p_{k|j} := \text{Tr}[P_k \tilde{\rho}_j]$, $p_{k,j} := \text{Tr}[\rho_j P_k] = p_j p_{k|j}$ and $p_k = \text{Tr}[P_k \rho^E]$ one has

$$S(\mathcal{C}_{\mathbf{P}}(\tilde{\rho}_j) \| \mathcal{C}_{\mathbf{P}}(\rho)) = S \left(\sum_{k \in Y} p_{k|j} |\varphi_k\rangle\langle\varphi_k| \middle\| \sum_{k \in Y} p_k |\varphi_k\rangle\langle\varphi_k| \right) = \sum_{k \in Y} p_{k|j} \log_2 \frac{p_{k|j}}{p_k},$$

and finally, averaging over X this implies

$$\sum_{j \in X} p_j S(\mathcal{C}_{\mathbf{P}}(\tilde{\rho}_j) \| \mathcal{C}_{\mathbf{P}}(\rho)) = \sum_{j \in X} \sum_{k \in Y} p_{k,j} \log_2 \frac{p_{k,j}}{p_j p_k} = I(E, \mathbf{P}). \quad \square$$

It should be emphasised that the Holevo bound is not generally saturated, namely there are ensembles of states for which the bound is not achieved for any POVM. The bound can indeed be achieved if and only if the states of the ensemble commute, namely

$$[\rho_i, \rho_j] = 0. \quad (14.18)$$

In particular, for orthogonal states ρ_j one has $I(X : Y) = H(X)$ and the POVM achieving the bound is $P_i = |\psi_i\rangle\langle\psi_i|$, where $\{\psi_i\}_{i=1}^{d_A}$ is the orthonormal basis of common eigenvectors of the states ρ_j . Generally, one has the following chain of bounds

$$I(E, \mathbf{P}) \leq \text{Acc}(E) \leq \chi(E) \leq S(\rho^E) \leq \log_2 d_A. \quad (14.19)$$

The third bound is achieved with equality for ensembles made of pure states.

The Holevo theorem is one of the most relevant results in quantum information theory. Indeed, it establishes the maximal amount of information that can be extracted by a measurement on a quantum system. The last bound in Eq. (14.19) gives the maximum accessible information from a quantum system in any state and for any measurement, which is just the \log_2 of the Hilbert space dimension.

Notice that the rightmost bound tells us that the maximum amount of classical information that can be encoded in a quantum system A coincides with its dimension (measured in qubits): $\log_2 d_A$. From this point of view, quantum systems do not provide any advantage versus classical systems.

Exercise 14.1

Show that the Holevo theorem is saturated for Abelian ensembles, and that for orthogonal states ones one has $I(X : Y) = H(X)$. Show that the optimal POVM is the von Neumann measurement on the orthonormal basis on which the states ρ_j are jointly diagonal. [It can be shown that condition (14.18) is also necessary for achieving the Holevo bound].

Answer of exercise 14.1

Consider the POVM $P_n = |n\rangle\langle n|$ where $\{|n\rangle\}$ denote the orthonormal basis jointly diagonalizing the states ρ_j , namely

$$\tilde{\rho}_j = \sum_n \lambda_n^{(j)} |n\rangle\langle n| \quad (14.20)$$

One has

$$\mathcal{C}_{\mathbf{P}}(\tilde{\rho}_j) = \sum_n \text{Tr}[P_n \tilde{\rho}_j] |n\rangle\langle n| = \sum_n \lambda_n^{(j)} |n\rangle\langle n| = \tilde{\rho}_j, \quad (14.21)$$

which implies also $\mathcal{C}_{\mathbf{P}}(\rho) = \rho$, and then

$$S(\mathcal{C}_{\mathbf{P}}(\tilde{\rho}_j) \| \mathcal{C}_{\mathbf{P}}(\rho)) = S(\tilde{\rho}_j \| \rho). \quad (14.22)$$

Finally, this implies

$$I(\mathbf{E}, \mathbf{P}) = \sum_j p_j S(\mathcal{C}_{\mathbf{P}}(\tilde{\rho}_j) \| \mathcal{C}_{\mathbf{P}}(\rho)) = \sum_j p_j S(\tilde{\rho}_j \| \rho) = \chi(\mathbf{E}). \quad (14.23)$$

If the states are orthogonal, one has $\text{Tr}[P_n \tilde{\rho}_j] = \lambda_n^{(j)}$ if $|n\rangle \in \text{Supp}(\rho_j)$, and $\text{Tr}[P_n \tilde{\rho}_j] = 0$ otherwise. Then, the random variable X with values j is a deterministic function of Y with values n , and thus $H(X|Y) = 0$. Finally, this implies that $I(X : Y) = H(X) - H(X|Y) = H(X)$.

Example 14.16 (Two equiprobable nonorthogonal pure states). To get an idea of the bound at work, it is interesting to consider the mutual information when the states are non-orthogonal. For simplicity, let's consider just the case of two equiprobable pure states: $|\psi_0\rangle = |0\rangle$ and $|\psi_1\rangle = \cos\theta|0\rangle + \sin\theta|1\rangle$. The *a priori* density matrix of the ensemble is

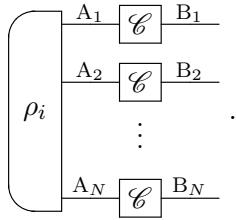
$$\rho^{\mathbf{E}} = \frac{1}{2} \begin{pmatrix} 1 + \cos^2 \theta & \cos \theta \sin \theta \\ \cos \theta \sin \theta & \sin^2 \theta \end{pmatrix}, \quad (14.24)$$

with eigenvalues $\lambda_{\pm} = \frac{1}{2}(1 \pm \cos \theta)$. Since the two states are pure, the Holevo theorem gives $\chi(\mathcal{E}) = S(\rho) = H_2(\lambda_+)$. The bound is maximized for $\theta = \frac{\pi}{2}$, namely when the states are orthogonal.

14.6 Holevo-Schumacher-Westmoreland theorem

Asymptotically for large n the Holevo bound for accessible information on output states of a channel \mathcal{C} can be achieved as the classical transmission rate of \mathcal{C} . In particular, this is the case if one uses separable encodings, i.e. without exploiting entanglement. This is the content of Holevo-Schumacher-Westmoreland theorem. The theorem thus provides a lower bound on the classical capacity of \mathcal{C} , which could be in principle overcome by using entangled encodings. Let us see this result in more detail.

We are considering the following scenario



Let us consider the ensemble $\mathcal{E}_o^N := \{\mathcal{C}^{\otimes N}(\rho_i)\}_{i \in X}$ of output states provided an input ensemble $\mathcal{E}^N := \{\rho_i\}_{i \in X}$. The classical capacity of the quantum channel \mathcal{C} can be defined as

$$C_C(\mathcal{C}) := \lim_{N \rightarrow \infty} \frac{\max_{\mathcal{E}^N, \mathbf{P}^N} I(\mathcal{E}_o^N, \mathbf{P}^N)}{N}.$$

This quantity is generally difficult to calculate. One can try to make the problem simpler by considering separable encodings. The Holevo-Schumacher-Westmoreland theorem indeed provides a nice expression for the classical capacity

$$C_C^P(\mathcal{C}) := \lim_{N \rightarrow \infty} \frac{\max_{\mathcal{E}^{\otimes N}, \mathbf{P}^N} I(\mathcal{E}_o^{\otimes N}, \mathbf{P}^N)}{N} \quad (14.25)$$

under the restrictive hypothesis of product encoding, in terms of the single-use Holevo quantity $\max_{\mathcal{E}} \chi(\mathcal{C}(\mathcal{E}))$. This result then shifts the question about classical capacity to whether the assumption of separable encodings is really restrictive—which means that an entangled encoding provides better performances—or the use of entanglement does not provide an advantage in communication of classical information.

In order to prove that $C_C^P(\mathcal{C}) \leq \max_{\mathcal{E}} \chi(\mathcal{C}(\mathcal{E}))$, let us define the ensemble $\mathcal{E}^{\otimes N}$ as follows

$$\mathcal{E}^{\otimes N} := \{p_{i_1, i_2, \dots, i_N} \tilde{\rho}_{i_1} \otimes \tilde{\rho}_{i_2} \otimes \dots \otimes \tilde{\rho}_{i_N}\}, \quad (14.26)$$

for arbitrary probability distribution p_{i_1, i_2, \dots, i_N} . We define

$$\begin{aligned} \tilde{\rho}_{\mathbf{i}} &:= \tilde{\rho}_{i_1} \otimes \tilde{\rho}_{i_2} \otimes \dots \otimes \tilde{\rho}_{i_N}, \\ p_{\mathbf{i}} &:= p_{i_1, i_2, \dots, i_N}. \end{aligned}$$

If Alice is allowed to encode over tensor product states, by choosing any prior probability distribution, then she can use any ensemble of the form $\mathsf{E}^{\otimes N}$ defined in equation (14.26), obtaining an output ensemble of the same form

$$\begin{aligned}\mathsf{E}_o^{\otimes N} &= \{p_{i_1, i_2, \dots, i_N} \mathcal{C}(\tilde{\rho}_{i_1}) \otimes \mathcal{C}(\tilde{\rho}_{i_2}) \otimes \dots \otimes \mathcal{C}(\tilde{\rho}_{i_N})\} \\ &= \{p_{i_1, i_2, \dots, i_N} \tilde{\rho}'_{i_1} \otimes \tilde{\rho}'_{i_2} \otimes \dots \otimes \tilde{\rho}'_{i_N}\},\end{aligned}$$

where $\tilde{\rho}'_i := \mathcal{C}(\tilde{\rho}_i)$. We will also denote by E_l the l -th marginal ensemble of $\mathsf{E}_o^{\otimes N}$, namely

$$\mathsf{E}_l := \{p_{i_l}^{(l)} \tilde{\rho}'_{i_l}\}, \quad p_{i_l}^{(l)} := \sum_{i_1, i_2, \dots, i_{l-1}, i_{l+1}, \dots, i_N} p_{i_1, i_2, \dots, i_N}.$$

Notice that

$$\begin{aligned}&\sum_{i_1, i_2, \dots, i_N} p_{i_1, i_2, \dots, i_N} S(\tilde{\rho}'_{i_1} \otimes \tilde{\rho}'_{i_2} \otimes \dots \otimes \tilde{\rho}'_{i_N}) \\ &= \sum_{i_1} p_{i_1}^{(1)} S(\tilde{\rho}'_{i_1}) + \sum_{i_2} p_{i_2}^{(2)} S(\tilde{\rho}'_{i_2}) + \dots + \sum_{i_N} p_{i_N}^{(N)} S(\tilde{\rho}'_{i_N}),\end{aligned}$$

or, more compactly

$$\sum_{\mathbf{i}} p_{\mathbf{i}} S(\tilde{\rho}'_{\mathbf{i}}) = \sum_{l=1}^N \sum_{i_l} p_{i_l}^{(l)} S(\tilde{\rho}'_{i_l}).$$

The last identity, along with subadditivity of the von Neumann entropy

$$S(\rho^{\mathsf{E}_o^{\otimes N}}) \leq \sum_{l=1}^N S(\rho'^{\mathsf{E}_l}),$$

gives the bound

$$\chi(\mathsf{E}_o^{\otimes N}) \leq \sum_{l=1}^N \chi(\mathsf{E}_l) \leq N \max_{\mathsf{E}} \chi(\mathcal{C}(\mathsf{E})). \quad (14.27)$$

If we now define single-shot channel capacity as

$$C_C^{(1)}(\mathcal{C}) := \max_{\mathsf{E}} \chi(\mathcal{C}(\mathsf{E})), \quad (14.28)$$

it is clear that from the Holevo bound and from equation 14.27 one has

$$\begin{aligned}C_C^P(\mathcal{C}) &\leq \lim_{N \rightarrow \infty} \frac{1}{N} C_C^{(1)}(\mathcal{C}^{\otimes N}) \\ &\leq C_C^{(1)}(\mathcal{C}).\end{aligned}$$

The Holevo-Schumacher-Westmoreland theorem assures that the capacity in terms of transmission rate is actually given by the Holevo quantity, and that the last bound is achievable asymptotically for large N , namely

$$\lim_{N \rightarrow \infty} \frac{1}{N} \chi(\mathsf{E}_o^{\otimes N}) = C_C^{(1)}(\mathcal{C}) = \max_{\mathsf{E}} \chi(\mathcal{C}(\mathsf{E})), \quad (14.29)$$

for suitable encoding and decoding on ensembles of factorised states.

The question about classical capacity without constraints then boils down to whether $\frac{1}{k}C_C^{(1)}(\mathcal{C}^{\otimes k}) > C_C^{(1)}(\mathcal{C})$. A weaker version of the question—i.e. whether $C_C^{(1)}(\mathcal{C}_1 \otimes \mathcal{C}_2) > C_C^{(1)}(\mathcal{C}_1) + C_C^{(1)}(\mathcal{C}_2)$ —is known as the *superadditivity question*, and was solved for the positive by an example provided in 2009 by Hastings. This implies that entanglement can make the classical capacity better than the Holevo-Schumacher-Westmoreland rate. However, this also implies that in general one has no simple formula to express the classical capacity.

Chapter 15

Lecture 17: Quantum information, compression, and Uhlmann fidelity

What is quantum information? In order to answer this question, let us analyse the task of transmitting an unknown quantum state ρ over a quantum channel \mathcal{C} , as in the following diagram

$$(\rho) \xrightarrow{\text{A}} \mathcal{C} \xrightarrow{\text{B}} .$$

Let us consider, for example, the channel $\mathcal{C} : A \rightarrow A$ given by the following Kraus form, in terms of an orthonormal basis $\{\varphi_i\}_{i=1}^{d_A}$

$$\mathcal{C}(\rho) = \sum_{i=1}^{d_A} |\varphi_i\rangle\langle\varphi_i| \rho |\varphi_i\rangle\langle\varphi_i|.$$

One can easily verify that this channel has a very large classical capacity, essentially behaving as an ideal classical channel. However, it is also clear that \mathcal{C} completely destroys superpositions of the states φ_i . Let us also consider the following scenario, in which the systems A and B of two agents Alice and Bob are entangled, and Alice sends her quantum system to Charlie C:

$$\begin{array}{c} P \\ \swarrow \quad \searrow \\ \text{A} \quad \mathcal{C} \quad \text{C} \\ | \qquad \qquad \qquad | \\ \text{B} \quad \text{B}' \end{array} .$$

The question is: are the systems C and B' still entangled? As we will see in the following, the answer to this question is strictly related to the possibility of sending superpositions. In other words, the ability to send quantum superpositions is equivalent to the ability to send entanglement: quantum information is equivalent to entanglement.

In quantitative terms, the amount of quantum information sent by a quantum channel represents the dimension of the largest Hilbert space whose superpositions can be reliably transmitted, or equivalently the amount of entanglement that can be reliably transmitted. This amount of quantum information is normally evaluated in *qubits*—the units of quantum information—and it corresponds then to the base two logarithm of the dimension.

The task that we consider in this chapter is *compression*.

We start from an example. Let $E = \{p_i|\varphi_i\rangle\langle\varphi_i|\} \subseteq \text{St}(A)$. Let us consider the state

$$(\rho^E)^{\otimes N} = \rho^{E^{\otimes N}}.$$

The encoding $\mathcal{C}(\rho^{E^{\otimes N}})$ generally provides a compressed state with smaller support than $(\rho^E)^{\otimes N}$ —let us say $\dim \text{Supp}(\mathcal{C}(\rho^{E^{\otimes N}})) = d_E(N)$. This will introduce some error probability in the decoding, but the optimal encoding will make the error arbitrarily small in the asymptotic limit $N \rightarrow \infty$. The compression ratio is given by

$$R = \lim_{N \rightarrow \infty} \frac{\log_2 d_E(N)}{N}$$

We will prove Schumacher's theorem, that provides the optimal compression ratio for a reliable compression scheme as $d_E(N) = 2^{NS(\rho^E)}$, i.e.

$$R = S(\rho^E).$$

Therefore, the von Neumann entropy can be interpreted as the number of qubits of quantum information per use that are needed in average to store or transmit a quantum message described by an ensemble with average state ρ^E .

15.1 Quantum compression

As an example, suppose that the ensemble is binary, corresponding to the two pure non-orthogonal states

$$|\psi_0\rangle = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad p_0 = \frac{1}{2}, \quad |\psi_1\rangle = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad p_1 = \frac{1}{2},$$

so that the *a priori* state is given by

$$\rho = \frac{1}{4} \begin{pmatrix} 3 & 1 \\ 1 & 1 \end{pmatrix}.$$

The eigenstates of ρ are given by

$$|0'\rangle = \begin{pmatrix} \cos \frac{\pi}{8} \\ \sin \frac{\pi}{8} \end{pmatrix}, \quad |1'\rangle = \begin{pmatrix} \sin \frac{\pi}{8} \\ -\cos \frac{\pi}{8} \end{pmatrix},$$

corresponding to the eigenvalues $\lambda_0 = \cos^2 \frac{\pi}{8}$ and $\lambda_1 = \sin^2 \frac{\pi}{8}$. The eigenstate $|0'\rangle$ has large overlap with both input states

$$|\langle 0'|\psi_0\rangle|^2 = |\langle 0'|\psi_1\rangle|^2 = \cos^2 \frac{\pi}{8} = .8535,$$

whereas $|1'\rangle$ have small overlap with both

$$|\langle 1'|\psi_0\rangle|^2 = |\langle 1'|\psi_1\rangle|^2 = \sin^2 \frac{\pi}{8} = .1465.$$

Therefore, if Alice compresses to $\mathcal{H} = \mathbb{C}$ (0 qubits), Bob can only guess the input state. The best guess that Bob can make is $|\psi\rangle = |0'\rangle$. This observation is based on the overlap with the input states.

Notice that just guessing the state would have given an average overlap

$$\frac{1}{2}|\langle\psi_0|\psi_0\rangle|^2 + \frac{1}{2}|\langle\psi_0|\psi_1\rangle|^2 = 3/4. \quad (15.1)$$

Now, suppose that Alice needs to send a message of 3 letters to Bob, but she can afford using only two qubits. Can she have Bob reconstruct her state, and with which overlap? If she sends to Bob only two of the three letters and asks Bob to guess the third one, then Bob will receive two letters with overlap 1, and one letter with overlap .8535, hence with overall overlap .8535.

Is there a better procedure? Indeed, there is. By diagonalising ρ we decomposed the Hilbert space into a “high overlap subspace” $\text{Span}\{|0'\rangle\}$ and an “low overlap subspace” $\text{Span}\{|1'\rangle\}$. Similarly we can decompose the Hilbert space of three qubits. Indeed, for any input string of states $|\Psi\rangle = |\psi_{i_1}\rangle|\psi_{i_2}\rangle|\psi_{i_3}\rangle$ we have

$$\begin{aligned} |\langle 0'0'0'|\Psi\rangle|^2 &= \cos^6 \frac{\pi}{8} = .6219, \\ |\langle 1'0'0'|\Psi\rangle|^2 &= |\langle 0'1'0'|\Psi\rangle|^2 = |\langle 0'0'1'|\Psi\rangle|^2 = \cos^4 \frac{\pi}{8} \sin^2 \frac{\pi}{8} = .1067, \\ |\langle 1'1'0'|\Psi\rangle|^2 &= |\langle 0'1'1'|\Psi\rangle|^2 = |\langle 1'0'1'|\Psi\rangle|^2 = \cos^2 \frac{\pi}{8} \sin^4 \frac{\pi}{8} = .0183, \\ |\langle 1'1'1'|\Psi\rangle|^2 &= \sin^6 \frac{\pi}{8} = .0031. \end{aligned}$$

Thus, the “high overlap subspace” L will be the span of the four states

$$L = \text{Span}\{|0'0'0'\rangle, |1'0'0'\rangle, |0'1'0'\rangle, |0'0'1'\rangle\},$$

and the “low overlap subspace” will be the orthogonal complement L^\perp . Now, if we want to project the state $|\Psi\rangle$ over L , the probability of success is

$$p = \text{Tr}[P_L|\Psi\rangle\langle\Psi|] = .6219 + 3(.1067) = .9419,$$

whereas the probability of projecting over the low overlap space is $1 - p = .0581$. The compression technique of Alice is precisely based on such projection. This can be achieved by the following measuring procedure. First she makes a unitary transformation that maps L and L^\perp into two-qubit spaces as follows

$$\begin{array}{ll} U|0'0'0'\rangle = |0\rangle|0\rangle|0\rangle & U|1'1'1'\rangle = |0\rangle|0\rangle|1\rangle \\ U|1'0'0'\rangle = |1\rangle|0\rangle|0\rangle & U|1'1'0'\rangle = |1\rangle|0\rangle|1\rangle \\ U|0'1'0'\rangle = |0\rangle|1\rangle|0\rangle & U|0'1'1'\rangle = |0\rangle|1\rangle|1\rangle \\ U|0'0'1'\rangle = |1\rangle|1\rangle|0\rangle & U|1'0'1'\rangle = |1\rangle|1\rangle|1\rangle \end{array}$$

that is to say

$$|\Psi\rangle \in L, U|\Psi\rangle = |\cdot\rangle|\cdot\rangle|0\rangle, \quad |\Psi\rangle \in L^\perp, U|\Psi\rangle = |\cdot\rangle|\cdot\rangle|1\rangle.$$

Then she measures the third qubit, and if she finds 0 (i.e. the input state has been successfully projected) she sends the remaining two qubits to Bob. Let's call the state of the two qubits that she sends $|\Psi_{\text{comp}}\rangle$. When Bob receives the two qubits in the state $|\Psi_{\text{comp}}\rangle$, he decompresses them by appending $|0\rangle$ and performing the inverse unitary transformation U^{-1} , obtaining

$$|\Psi'\rangle = U^{-1}|\Psi_{\text{comp}}\rangle|0\rangle.$$

If Alice's measurement of the third qubit yields 1, this means that the state has been projected to L^\perp , in which case, the best thing she can do is to send the state that Bob will decompress to the highest overlap state $|0'0'0'\rangle$, namely she sends the state $|\Psi_{\text{comp}}\rangle = |0\rangle|0\rangle$ such that

$$|\Psi'\rangle = U^{-1}|\Psi_{\text{comp}}\rangle|0\rangle = |0'0'0'\rangle.$$

The overall compression channel from Alice to Bob is then given by

$$\mathcal{C}(\rho) = P_L \rho P_L + \text{Tr}[(I - P_L)\rho]|0'0'0'\rangle\langle 0'0'0'|,$$

with ρ being any 3-qubit state, and with $P_L = U^\dagger(I \otimes I \otimes |0\rangle\langle 0|)U$. Thus one has

$$\mathcal{C}(|\Psi\rangle\langle\Psi|) = P_L|\Psi\rangle\langle\Psi|P_L + |0'0'0'\rangle\langle 0'0'0'|,$$

with overlap

$$\begin{aligned} \langle\Psi|\mathcal{C}(|\Psi\rangle\langle\Psi|)|\Psi\rangle &= \langle\Psi|P_L|\Psi\rangle^2 + \langle\Psi|(I - P_L)|\Psi\rangle|\langle 0'0'0'|\Psi\rangle|^2 \\ &= .9419^2 + .0581 \times .6219 \\ &= .9234, \end{aligned}$$

which is much better than the value .8535 corresponding to the naive strategy of sending only two qubits and guessing the third one. Clearly, as we make the string of qubits longer, the overlap can be improved. According to the Schumacher compression theorem, we can shorten a long message by a factor

$$S(\rho) = H_2(\cos^2 \pi/8) = .6009$$

and asymptotically get the decoding close to the input state at will. Up to now, we considered the overlap as an intuitive measure of quality of the compression. However, we now need to provide a justification for this, through the introduction of *fidelity* as a figure of merit.

15.2 Fidelity

We start with the definition of fidelity, due to Uhlmann.

Definition 15.1 (Uhlmann's fidelity). For any two normalized states ρ and σ in $\text{St}_1(A)$ the Uhlmann's fidelity is defined as

$$\begin{aligned} F(\rho, \sigma) &:= \text{Tr}[(\sigma^{\frac{1}{2}}\rho\sigma^{\frac{1}{2}})^{\frac{1}{2}}] \\ &= \text{Tr}[|\sqrt{\rho}\sqrt{\sigma}|] \\ &= \|\sqrt{\rho}\sqrt{\sigma}\|_1. \end{aligned} \tag{15.2}$$

We remind that the trace-norm $\|X\|_1$ is defined as follows.

Definition 15.2 (Trace-norm). Let $X : \mathcal{H}_A \rightarrow \mathcal{H}_A$. The trace-norm $\|X\|_1$ of X is defined as

$$\|X\|_1 := \text{Tr}[(X^\dagger X)^{\frac{1}{2}}]. \quad (15.3)$$

Lemma 15.3. *Let $A, U \in \mathcal{L}(\mathcal{H}_A)$, and U be unitary. The following bound holds*

$$|\text{Tr}[AU]| \leq \text{Tr}[|A|] = \|A\|_1, \quad (15.4)$$

with equality iff $A = U^\dagger |A|$.

Proof. Using the polar decomposition $A = V|A|$, one can write

$$|\text{Tr}[AU]| = |\text{Tr}[V|A|U]| = |\text{Tr}[|A|^{\frac{1}{2}}|A|^{\frac{1}{2}}UV]|, \quad (15.5)$$

and using the Cauchy-Schwarz inequality for the Frobenius product $|\text{Tr}[A^\dagger B]|^2 \leq \text{Tr}[A^\dagger A] \text{Tr}[B^\dagger B]$ one has

$$|\text{Tr}[AU]| \leq \sqrt{\text{Tr}[|A|] \text{Tr}[V^\dagger U^\dagger |A| UV]} = \text{Tr}[|A|], \quad (15.6)$$

with equality attained iff $V^\dagger U^\dagger |A| = |A|$, i.e. $U^\dagger |A| = V|A| = A$. \square

We can now prove the following important theorems, providing two useful ways of calculating the fidelity

Theorem 15.4 (Uhlmann's theorem). *For any two normalised states ρ and σ in $\text{St}_1(A)$ and any fixed system B, the Uhlmann's fidelity is given by*

$$F(\rho, \sigma) = \max_{\Psi_\rho, \Psi_\sigma} |\langle\langle \Psi_\rho | \Psi_\sigma \rangle\rangle|, \quad (15.7)$$

where $\Psi_\rho, \Psi_\sigma \in \text{St}_1(AB)$ denote any two purifications of ρ and σ , respectively.

Proof. First of all, in order to have a purification of σ and ρ in $\mathcal{H}_A \otimes \mathcal{H}_B$, it must be $d_B \geq \max\{\text{rank}(\rho), \text{rank}(\sigma)\}$. By the purification theorem, a purification $|\Psi_\rho\rangle\rangle \langle\langle \Psi_\rho|$ of ρ can be obtained by $|\Psi_\rho\rangle\rangle = |\rho^{\frac{1}{2}}U\rangle\rangle$, where $UU^\dagger \geq P_{\text{Supp}(\rho)}$. Analogously for σ . Having fixed the system B, the purification theorem 11.2 then states that every other purification of ρ can be obtained from a given one $|\Psi_\rho\rangle\rangle = |\rho^{\frac{1}{2}}U\rangle\rangle$ by a unitary map W on system B, i.e. $|\Psi'_\rho\rangle\rangle = |\rho^{\frac{1}{2}}UW\rangle\rangle$ (the same for σ , with V and Z in place of U and W). Let us then consider the expression

$$|\langle\langle \Psi'_\rho | \Psi'_\sigma \rangle\rangle| = |\langle\langle \rho^{\frac{1}{2}}UW | \sigma^{\frac{1}{2}}VZ \rangle\rangle| = |\text{Tr}[W^\dagger U^\dagger \rho^{\frac{1}{2}} \sigma^{\frac{1}{2}} VZ]| = |\text{Tr}[ZW^\dagger U^\dagger \rho^{\frac{1}{2}} \sigma^{\frac{1}{2}} V]|.$$

By lemma 15.3, considering that ZW^\dagger is unitary, one has

$$|\langle\langle \Psi'_\rho | \Psi'_\sigma \rangle\rangle| \leq \|U^\dagger \rho^{\frac{1}{2}} \sigma^{\frac{1}{2}} V\|_1 = \text{Tr}[(V^\dagger \sigma^{\frac{1}{2}} \rho^{\frac{1}{2}} U U^\dagger \rho^{\frac{1}{2}} \sigma^{\frac{1}{2}} V)^{\frac{1}{2}}].$$

Now, since $UU^\dagger \geq P_{\text{Supp}(\rho)}$, the above inequality becomes

$$|\langle\langle \Psi'_\rho | \Psi'_\sigma \rangle\rangle| \leq \text{Tr}[(V^\dagger \sigma^{\frac{1}{2}} \rho \sigma^{\frac{1}{2}} V)^{\frac{1}{2}}].$$

Moreover, since $VV^\dagger \geq P_{\text{Supp}(\sigma)}$, and $\text{Supp}[(\sigma^{\frac{1}{2}} \rho \sigma^{\frac{1}{2}})^{\frac{1}{2}}] \subseteq \text{Supp}(\sigma)$, one has

$$V^\dagger (\sigma^{\frac{1}{2}} \rho \sigma^{\frac{1}{2}})^{\frac{1}{2}} V V^\dagger (\sigma^{\frac{1}{2}} \rho \sigma^{\frac{1}{2}})^{\frac{1}{2}} V = V^\dagger \sigma^{\frac{1}{2}} \rho \sigma^{\frac{1}{2}} V,$$

namely $(V^\dagger \sigma^{\frac{1}{2}} \rho \sigma^{\frac{1}{2}} V)^{\frac{1}{2}} = V^\dagger (\sigma^{\frac{1}{2}} \rho \sigma^{\frac{1}{2}})^{\frac{1}{2}} V$. Finally, this gives

$$|\langle\langle \Psi'_\rho | \Psi'_\sigma \rangle\rangle| \leq \text{Tr}[VV^\dagger (\sigma^{\frac{1}{2}} \rho \sigma^{\frac{1}{2}})^{\frac{1}{2}}] = \text{Tr}[(\sigma^{\frac{1}{2}} \rho \sigma^{\frac{1}{2}})^{\frac{1}{2}}] = F(\rho, \sigma).$$

By lemma 15.3, the bound can be achieved by a suitable choice of ZW^\dagger . \square

Theorem 15.5. *The Uhlmann's fidelity is given by the minimum over all possible measurements of the classical fidelity between the probability distributions of the two states, namely*

$$F(\rho, \sigma) = \min_{\mathbf{P}} F_C(\mathbf{p}, \mathbf{q}), \quad \mathbf{p} := \text{Tr}[\mathbf{P}\rho], \quad \mathbf{q} := \text{Tr}[\mathbf{P}\sigma], \quad (15.8)$$

where

$$F_C(\mathbf{p}, \mathbf{q}) := \sum_n \sqrt{p_n q_n}. \quad (15.9)$$

Proof. From lemma 15.3 we learned that there exists a unitary U on \mathcal{H}_A such that $F(\rho, \sigma) = \text{Tr}(\rho^{\frac{1}{2}} \sigma^{\frac{1}{2}} U)$. Then

$$\begin{aligned} F(\rho, \sigma) &= \text{Tr}(U \rho^{\frac{1}{2}} \sigma^{\frac{1}{2}}) = \sum_n \text{Tr}(U \rho^{\frac{1}{2}} \sqrt{P_n} \sqrt{P_n} \sigma^{\frac{1}{2}}) \\ &\leq \sum_n \sqrt{\text{Tr}(\rho P_n) \text{Tr}(\sigma P_n)} =: F_C(\mathbf{p}, \mathbf{q}) \end{aligned} \quad (15.10)$$

where we used the Cauchy-Schwarz inequality, along with the normalisation of the arbitrary POVM \mathbf{P} , $\sum_n P_n = I_A$. The Cauchy-Schwarz inequality is saturated when, for every n , either

$$U \rho^{\frac{1}{2}} \sqrt{P_n} - \alpha_n \sigma^{\frac{1}{2}} \sqrt{P_n} = 0, \quad (15.11)$$

or $\sigma^{\frac{1}{2}} \sqrt{P_n} = 0$. Multiplying by $\sqrt{P_n}$ on the right hand side and by $P_{\text{Supp}(\sigma)}$ on the left, one obtains the following necessary and sufficient condition

$$(P_{\text{Supp}(\sigma)} U \rho^{\frac{1}{2}} - \alpha_n \sigma^{\frac{1}{2}}) P_n = 0. \quad (15.12)$$

Now, reminding that U is such that $U \rho^{\frac{1}{2}} \sigma^{\frac{1}{2}} = \sqrt{\sigma^{\frac{1}{2}} \rho \sigma^{\frac{1}{2}}}$, we have

$$\sigma^{-\frac{1}{2}} U \rho^{\frac{1}{2}} P_{\text{Supp}(\sigma)} = \sigma^{-\frac{1}{2}} \sqrt{\sigma^{\frac{1}{2}} \rho \sigma^{\frac{1}{2}}} \sigma^{-\frac{1}{2}} =: M. \quad (15.13)$$

The condition of Eq. (15.12) can be satisfied by taking the POVM $P_n = |n\rangle\langle n|$, with $M|n\rangle = \beta_n|n\rangle$, since by definition of M it is $\text{Supp}(M) \subseteq \text{Supp}(\sigma)$. Thus one obtains $\sigma^{\frac{1}{2}} M P_n = P_{\text{Supp}(\sigma)} U \rho^{\frac{1}{2}} P_n = \sigma^{\frac{1}{2}} \beta_n P_n$, which means that condition (15.12) is satisfied taking $\alpha_n = \beta_n$. For $|n\rangle \in \text{Ker}(\sigma)$, clearly $\text{Tr}[\sigma P_n] = 0$, then the n -th term does not contribute to the sum in Eq. (15.10). \square

Lemma 15.6. *The Uhlmann's fidelity satisfies the following properties:*

1. $0 \leq F(\rho, \sigma) \leq 1$, and $F(\rho, \sigma) = 1$ iff $\rho = \sigma$, and $F(\rho, \sigma) = 0$ iff $\text{Supp}(\rho) \perp \text{Supp}(\sigma)$;
2. $F(\rho, \sigma) = F(\sigma, \rho)$;
3. $F(\rho, |\psi\rangle\langle\psi|) = \langle\psi|\rho|\psi\rangle^{\frac{1}{2}}$.

Proof. Item 1 can be proved considering that $0 \leq (\rho^{\frac{1}{2}}\sigma\rho^{\frac{1}{2}})^{\frac{1}{2}}$, and thus $F(\rho, \sigma) \geq 0$. Moreover, by theorem 15.4 one clearly has $F(\rho, \sigma) \leq 1$. Now, the case $F(\rho, \sigma) = 0$ corresponds to the case $\|\rho^{\frac{1}{2}}\sigma^{\frac{1}{2}}\|_1 = 0$, i.e. $\rho^{\frac{1}{2}}\sigma^{\frac{1}{2}} = 0$. It is then easy to obtain $\rho\sigma = \sigma\rho = 0$. Moreover, again by theorem 15.4, if $F(\rho, \sigma) = 1$ then there exist purifications Ψ_ρ, Ψ_σ such that $\langle\langle\Psi_\rho|\Psi_\sigma\rangle\rangle = 1$, and thus $\Psi_\rho = \Psi_\sigma$, which in turn implies $\rho = \sigma$. Item 2 is an immediate consequence of theorem 15.4. Finally, for item 3 it is sufficient to observe that $(|\psi\rangle\langle\psi|)^{\frac{1}{2}} = |\psi\rangle\langle\psi|$ for every $|\psi\rangle\langle\psi| \in \text{St}_1(\mathbf{A})$. \square

Theorem 15.7 (Monotonicity of Uhlmann's fidelity). *For any quantum channel \mathcal{E} , one has*

$$F[\mathcal{E}(\rho), \mathcal{E}(\sigma)] \geq F(\rho, \sigma). \quad (15.14)$$

Proof. Denote by $|\Psi_\rho\rangle\rangle$ and $|\Psi_\sigma\rangle\rangle$ any two purifications achieving the fidelity $F(\rho, \sigma)$. The reversible dilation theorem 11.14 then provides a purification for the output of the channel, and one has

$$F[\mathcal{E}(\rho), \mathcal{E}(\sigma)] \geq |\langle\langle\Psi_\rho|\langle 0|(I_A \otimes U^\dagger U)|\Psi_\sigma\rangle\rangle| = |\langle\langle\Psi_\rho|\Psi_\sigma\rangle\rangle| = F(\rho, \sigma). \quad \square$$

Since the (partial) trace is a special channel, fidelity is monotonic under partial trace.

Corollary 15.8 (Monotonicity of Uhlmann's fidelity under partial trace). *For any two bipartite states ρ_{AB}, σ_{AB} one has*

$$F(\rho_A, \sigma_A) \geq F(\rho_{AB}, \sigma_{AB}). \quad (15.15)$$

Lemma 15.9 (Strong concavity of Uhlmann's fidelity). *For any two ensembles of states $\{r_i\tilde{\rho}_i\}$ and $\{s_i\tilde{\sigma}_i\}$, one has*

$$F\left(\sum_i r_i \tilde{\rho}_i, \sum_i s_i \tilde{\sigma}_i\right) \geq \sum_i \sqrt{r_i s_i} F(\tilde{\rho}_i, \tilde{\sigma}_i). \quad (15.16)$$

Proof. Denote by $|\psi_i\rangle$ and $|\phi_i\rangle$ the vectors on $\mathcal{H}_A \otimes \mathcal{H}_A$ corresponding to purifications of ρ_i and σ_i such that $\langle\psi_i|\phi_i\rangle = F(\tilde{\rho}_i, \tilde{\sigma}_i)$. Then $|\Psi\rangle = \sum_i \sqrt{r_i} |\psi_i\rangle \otimes |i\rangle$ and $|\Phi\rangle = \sum_i \sqrt{s_i} |\phi_i\rangle \otimes |i\rangle$ are purifications of $\sum_i r_i \tilde{\rho}_i$ and $\sum_i s_i \tilde{\sigma}_i$, respectively. By Uhlmann's theorem one has

$$\begin{aligned} F\left(\sum_i r_i \tilde{\rho}_i, \sum_i s_i \tilde{\sigma}_i\right) &\geq |\langle\langle\Psi|\Phi\rangle\rangle| = \sum_i \sqrt{r_i s_i} \langle\psi_i|\phi_i\rangle \\ &= \sum_i \sqrt{r_i s_i} F(\tilde{\rho}_i, \tilde{\sigma}_i). \end{aligned} \quad \square$$

Corollary 15.10. *The Uhlmann's fidelity is jointly concave in both entries.*

Proof. It is sufficient to take $r_i = s_i$ in equation (15.16), obtaining

$$F\left(\sum_i r_i \tilde{\rho}_i, \sum_i r_i \tilde{\sigma}_i\right) \geq \sum_i r_i F(\tilde{\rho}_i, \tilde{\sigma}_i). \quad \square$$

The last result that we prove relates the Uhlmann fidelity to the *trace-norm distance* defined as

Definition 15.11 (Trace-norm distance). Let $X, Y : \mathcal{H}_A \rightarrow \mathcal{H}_A$. The trace-norm distance $D(X, Y)$ is defined as

$$D(X, Y) := \frac{1}{2} \|X - Y\|_1. \quad (15.17)$$

We now show that the trace norm distance for states has an immediate operational interpretation in terms of discrimination probability. The scenario that we consider is the following. An agent is provided with a system A prepared in a state, which is either ρ_0 or ρ_1 , both cases having prior probability 1/2. The agent then performs a binary measurement $\{P_0, P_1\}$ with outcomes 0 and 1, summarising the discrimination strategy: upon reading 0 the agent will declare " ρ_0 ", and upon reading 1 they declare " ρ_1 ". The total error probability for the above strategy is thus given by

$$\begin{aligned} p_{\text{err}} &= \frac{1}{2} (\text{Tr}[\rho_0 P_1] + \text{Tr}[\rho_1 P_0]) \\ &= \frac{1}{2} (\text{Tr}[\rho_0 P_1] + \text{Tr}[\rho_1 (I_A - P_1)]) \\ &= \frac{1}{2} (1 - \text{Tr}[(\rho_1 - \rho_0) P_1]) \\ &= \frac{1}{2} (\text{Tr}[\rho_0 (I_A - P_0)] + \text{Tr}[\rho_1 P_0]) \\ &= \frac{1}{2} (1 + \text{Tr}[(\rho_1 - \rho_0) P_0]) \\ &= \frac{1}{4} (2 - \text{Tr}[(\rho_1 - \rho_0)(P_1 - P_0)]). \end{aligned}$$

Notice that $\text{Tr}[(\rho_1 - \rho_0)(P_1 + P_0)] = \text{Tr}[\rho_1 - \rho_0] = 0$, and thus $\text{Tr}[(\rho_1 - \rho_0)(P_1 - P_0)] = 2 \text{Tr}[(\rho_1 - \rho_0)P_1]$. If we minimise the error probability for the above discrimination strategy, we obtain

$$p_{\text{err}} = \frac{1}{2} (1 - \max_{0 \leq P \leq I_A} \text{Tr}[(\rho_1 - \rho_0)P]). \quad (15.18)$$

The next lemma shows that the above expression is related to the trace-norm distance as follows

$$p_{\text{err}} = \frac{1}{4} (2 - \|\rho_1 - \rho_0\|_1).$$

Lemma 15.12. *Let $\rho_0, \rho_1 \in \text{St}_1(\mathcal{A})$. Then*

$$\frac{1}{2} \|\rho_0 - \rho_1\|_1 = \max_{0 \leq P \leq I_A} \text{Tr}[(\rho_0 - \rho_1)P]. \quad (15.19)$$

Proof. First, we remark that $\rho_0 - \rho_1$ is selfadjoint, and thus it can be expressed as $N_+ - N_-$, with $N_+N_- = N_-N_+ = 0$, and $N_\pm \geq 0$. Now, clearly $|\rho_0 - \rho_1| = \{(N_+ - N_-)^2\}^{\frac{1}{2}} = (N_+^2 + N_-^2)^{\frac{1}{2}} = N_+ + N_-$, since $(N_+ + N_-)^2 = N_+^2 + N_-^2$. Moreover, since $\text{Tr}[N_+ - N_-] = \text{Tr}[\rho_0 - \rho_1] = 0$, we have $\text{Tr}[N_+] = \text{Tr}[N_-] = \frac{1}{2} \text{Tr}[N_+ + N_-] = \frac{1}{2} \text{Tr}[|\rho_0 - \rho_1|] = \frac{1}{2} \|\rho_0 - \rho_1\|_1$. Then, one has

$$\begin{aligned} \text{Tr}[(\rho_0 - \rho_1)P] &= \text{Tr}[N_+P] - \text{Tr}[N_-P] \\ &\leq \text{Tr}[N_+P] \\ &\leq \text{Tr}[N_+] \\ &= \frac{1}{2} \|\rho_0 - \rho_1\|_1. \end{aligned}$$

Moreover, one can clearly achieve the upper bound by taking $P = P_+$, where P_+ is the projection on the support of N_+ , i.e. the projection on the direct sum of all the eigenspaces of $\rho_0 - \rho_1$ with non-negative eigenvalues. Thus,

$$\max_{0 \leq P \leq I_A} \text{Tr}[(\rho_0 - \rho_1)P] = \frac{1}{2} \|\rho_0 - \rho_1\|_1. \quad \square$$

The above result can be generalised as follows.

Theorem 15.13. *The trace-norm distance of two states $\rho, \sigma \in \text{St}_1(\mathcal{A})$ is equal to*

$$\|\rho - \sigma\|_1 = \max_{\mathbf{P}} \sum_{i=1}^{l(\mathbf{P})} |\text{Tr}[P_i(\rho - \sigma)]|, \quad (15.20)$$

where \mathbf{P} is a generic POVM with finitely many elements, and $l(\mathbf{P})$ is the number of elements of \mathbf{P} .

Proof. Let us evaluate the argument of the maximum on r.h.s. in Eq. (15.20) separating the indices i of terms for which $\text{Tr}[P_i(\rho - \sigma)] \geq 0$ from those for which $\text{Tr}[P_i(\rho - \sigma)] < 0$, thus forming the sets S_+ and S_- :

$$\begin{aligned} \sum_{i=1}^{l(\mathbf{P})} |\text{Tr}[P_i(\rho - \sigma)]| &= \sum_{i \in S_+} \text{Tr}[P_i(\rho - \sigma)] - \sum_{i \in S_-} \text{Tr}[P_i(\rho - \sigma)] \\ &= \text{Tr}[P_+(\rho - \sigma)] - \text{Tr}[P_-(\rho - \sigma)], \end{aligned}$$

where $P_\pm = \sum_{i \in S_\pm} P_i$. Thus, it is not restrictive to evaluate the maximum on the set of binary POVMs. As a result, one obtains the thesis by virtue of lemma 15.12. \square

For generally mixed states, we have the following relevant theorem

Theorem 15.14 (Equivalence of trace-norm and Uhlmann's distances). *The following double bound holds*

$$1 - F(\rho, \sigma) \leq \frac{1}{2} \|\rho - \sigma\|_1 \leq \sqrt{1 - F(\rho, \sigma)^2}, \quad (15.21)$$

with both equalities iff $\rho = \sigma$.

Proof. By theorem 15.5, one has

$$1 - F(\rho, \sigma) = 1 - \min_{\mathbf{P}} F_C(\mathbf{p}, \mathbf{q}) = 1 - F_C(\mathrm{Tr}[\rho \mathbf{P}], \mathrm{Tr}[\sigma \mathbf{P}]),$$

where \mathbf{P} is the POVM minimising the classical fidelity $F_C(\mathrm{Tr}[\rho \mathbf{P}], \mathrm{Tr}[\sigma \mathbf{P}])$. We then have

$$\begin{aligned} 1 - F(\rho, \sigma) &= 1 - \sum_{i=1}^{l(\mathbf{P})} \sqrt{\mathrm{Tr}[\rho P_i] \mathrm{Tr}[\sigma P_i]} \\ &= \frac{1}{2} \sum_{i=1}^{l(\mathbf{P})} (\sqrt{\mathrm{Tr}[\rho P_i]} - \sqrt{\mathrm{Tr}[\sigma P_i]})^2 \\ &\leq \frac{1}{2} \sum_{i=1}^{l(\mathbf{P})} \{ |(\sqrt{\mathrm{Tr}[\rho P_i]} - \sqrt{\mathrm{Tr}[\sigma P_i]})| (\sqrt{\mathrm{Tr}[\rho P_i]} + \sqrt{\mathrm{Tr}[\sigma P_i]}) \} \\ &= \frac{1}{2} \sum_{i=1}^{l(\mathbf{P})} |(\mathrm{Tr}[\rho P_i] - \mathrm{Tr}[\sigma P_i])| \\ &\leq \frac{1}{2} \|\rho - \sigma\|_1, \end{aligned}$$

where I^+ and I^- in the fifth line are the sets of values of i such that $(\mathrm{Tr}[\rho P_i] - \mathrm{Tr}[\sigma P_i]) \geq 0$ and $(\mathrm{Tr}[\rho P_i] - \mathrm{Tr}[\sigma P_i]) < 0$, respectively, and $P'_0 = \sum_{i \in I^+} P_i$. In the second to last step, we used the observation that $\mathrm{Tr}[(\rho - \sigma)(P'_0 - P'_1)] = \mathrm{Tr}[(\rho - \sigma)(2P'_0 - I)] = 2 \mathrm{Tr}[(\rho - \sigma)P'_0]$. Now, by theorem 15.13 we have

$$\begin{aligned} \frac{1}{2} \|\rho - \sigma\|_1 &= \frac{1}{2} \max_{\mathbf{P}} \sum_{i=1}^{l(\mathbf{P})} |\mathrm{Tr}[P_i \rho] - \mathrm{Tr}[P_i \sigma]| \\ &\leq \frac{1}{2} \max_{\mathbf{P}} \sum_{i=1}^{l(\mathbf{P})} \{ |\sqrt{\mathrm{Tr}[P_i \rho]} - \sqrt{\mathrm{Tr}[P_i \sigma]}| |\sqrt{\mathrm{Tr}[P_i \rho]} + \sqrt{\mathrm{Tr}[P_i \sigma]}| \} \\ &\leq \frac{1}{2} \max_{\mathbf{P}} \left\{ \sum_{i=1}^{l(\mathbf{P})} |\sqrt{\mathrm{Tr}[P_i \rho]} - \sqrt{\mathrm{Tr}[P_i \sigma]}|^2 \right\}^{\frac{1}{2}} \left\{ \sum_{i'=1}^{l(\mathbf{P})} |\sqrt{\mathrm{Tr}[P_{i'} \rho]} + \sqrt{\mathrm{Tr}[P_{i'} \sigma]}|^2 \right\}^{\frac{1}{2}} \\ &= \frac{1}{2} \max_{\mathbf{P}} \{ [2 - 2F_C(\mathbf{p}, \mathbf{q})][2 + 2F_C(\mathbf{p}, \mathbf{q})] \}^{\frac{1}{2}} \\ &= \max_{\mathbf{P}} [1 - F_C(\mathbf{p}, \mathbf{q})]^2^{\frac{1}{2}} \\ &= \{1 - [\min_{\mathbf{P}} F_C(\mathbf{p}, \mathbf{q})]^2\}^{\frac{1}{2}} \\ &= \sqrt{1 - F(\rho, \sigma)^2}. \end{aligned}$$

□

Chapter 16

Lecture 19: Schumacher quantum source coding theorem

16.1 Entanglement fidelity

The square of the fidelity $\langle\langle \Phi | \mathcal{E} \otimes \mathcal{I}_B(|\Phi\rangle\langle\Phi|) |\Phi \rangle\rangle$ between input and output of a channel acting locally on a pure bipartite state $|\Phi\rangle\langle\Phi| \in \text{St}_1(AB)$ only depends on the channel $\mathcal{E} : A \rightarrow A$ and on the local state $\rho_A := \text{Tr}_B[|\Phi\rangle\langle\Phi|]$. Indeed, for every purification $|\Phi'\rangle\langle\Phi'| \in \text{St}_1(AC)$ of the same state ρ_A one has

$$\begin{aligned}\rho_A &= \text{Tr}_C[|\Phi'\rangle\langle\Phi'|] \\ &= \Phi'\Phi'^{\dagger} \\ &= \Phi\Phi^{\dagger},\end{aligned}$$

and using the identity

$$\begin{aligned}|\Phi'\rangle\langle\Phi'| &= (I_A \otimes \Phi'^T)|I_A\rangle\langle I_A|(I_A \otimes \Phi'^*) \\ &= \mathcal{I}_A \otimes \mathcal{N}_{\Phi'^T}(|I_A\rangle\langle I_A|),\end{aligned}$$

one can easily verify that

$$\begin{aligned}\langle\langle \Phi' | \mathcal{E} \otimes \mathcal{I}_C(|\Phi'\rangle\langle\Phi'|) |\Phi' \rangle\rangle &= \langle\langle \Phi' | (\mathcal{E} \otimes \mathcal{I}_C)(\mathcal{I}_A \otimes \mathcal{N}_{\Phi'^T})(|I_A\rangle\langle I_A|) |\Phi' \rangle\rangle \\ &= \langle\langle \Phi' | (\mathcal{I}_A \otimes \mathcal{N}_{\Phi'^T})(\mathcal{E} \otimes \mathcal{I}_{A'})(|I_A\rangle\langle I_A|) |\Phi' \rangle\rangle \\ &= \langle\langle \Phi' | (I_A \otimes \Phi'^T) \text{Ch}(\mathcal{E})(I \otimes \Phi'^*) |\Phi' \rangle\rangle \\ &= \langle\langle \Phi' \Phi'^{\dagger} | \text{Ch}(\mathcal{E}) | \Phi' \Phi'^{\dagger} \rangle\rangle \\ &= \langle\langle \rho | \text{Ch}(\mathcal{E}) | \rho \rangle\rangle.\end{aligned}$$

This observation motivates the following definition

Definition 16.1 (Entanglement fidelity of a channel). The entanglement fidelity of a channel \mathcal{E} for a state ρ is defined as the square of the fidelity between input and output of a channel acting locally on any purification $|\Phi_\rho\rangle\langle\Phi_\rho|$ of ρ , i.e. $\rho = \Phi_\rho\Phi_\rho^\dagger$

$$\begin{aligned}\mathcal{F}(\rho, \mathcal{E}) &:= \langle\langle \Phi_\rho | \mathcal{E} \otimes \mathcal{I}(|\Phi_\rho\rangle\langle\Phi_\rho|) |\Phi_\rho \rangle\rangle \\ &= \langle\langle \rho | \text{Ch}(\mathcal{E}) | \rho \rangle\rangle.\end{aligned}\tag{16.1}$$

By the above observation, the entanglement fidelity $\mathcal{F}(\rho, \mathcal{E})$ is well defined, since it does not depend on the particular purification $|\Phi_\rho\rangle\langle\Phi_\rho|$ of ρ .

Lemma 16.2. *The following identity holds*

$$\mathcal{F}(\rho, \mathcal{E}) = \sum_i |\text{Tr}[\rho E_i]|^2, \quad (16.2)$$

for an arbitrary Kraus decomposition $\mathcal{E}(\rho) = \sum_i E_i \rho E_i^\dagger$ of \mathcal{E} .

Proof. We remind that if $\mathcal{E}(\rho) = \sum_i E_i \rho E_i^\dagger$ of \mathcal{E} is a Kraus decomposition of \mathcal{E} then one has

$$\text{Ch}(\mathcal{E}) = \sum_i |E_i\rangle\langle E_i|,$$

and by the second line of the definition (16.1), one has

$$\mathcal{F}(\rho, \mathcal{E}) = \sum_i \langle\langle \rho | E_i \rangle\rangle \langle E_i | \rho \rangle = \sum_i |\text{Tr}[\rho E_i]|^2. \quad \square$$

Lemma 16.3. *The entanglement fidelity satisfies the following properties*

1. $0 \leq \mathcal{F}(\rho, \mathcal{E}) \leq 1$;
2. $\mathcal{F}(\rho, \mathcal{E})$ is linear in \mathcal{E} ;
3. For pure states one has

$$\mathcal{F}(|\psi\rangle\langle\psi|, \mathcal{E}) = F(|\psi\rangle\langle\psi|, \mathcal{E}(|\psi\rangle\langle\psi|))^2. \quad (16.3)$$

Proof. Item 1 can be proved reminding that $\mathcal{F}(\rho, \mathcal{E})$ is the square of the fidelity of two states, and the Uhlmann fidelity is positive and bounded by one. As to item 2, it is easy to observe that the second line in the definition (16.1) is linear in $\text{Ch}(\mathcal{E})$, and the Choi map Ch is linear, too. Finally, let us consider item 3. An arbitrary purification of a pure state $|\psi\rangle\langle\psi|$ is a factorised state

$$|\Phi_\psi\rangle\langle\Phi_\psi| = |\psi\rangle\langle\psi| \otimes |\eta\rangle\langle\eta|,$$

for some pure state $|\eta\rangle\langle\eta| \in \text{St}_1(C)$. Thus,

$$\begin{aligned} \mathcal{F}(|\psi\rangle\langle\psi|, \mathcal{E}) &= \langle\psi|\langle\eta|(\mathcal{E} \otimes \mathcal{I}_C)[|\psi\rangle\langle\psi| \otimes |\eta\rangle\langle\eta|]|\psi\rangle|\eta\rangle \\ &= \langle\psi|\mathcal{E}(|\psi\rangle\langle\psi|)|\psi\rangle|\langle\eta|\eta\rangle|^2 \\ &= \langle\psi|\mathcal{E}(|\psi\rangle\langle\psi|)|\psi\rangle \\ &= F(|\psi\rangle\langle\psi|, \mathcal{E}(|\psi\rangle\langle\psi|))^2. \end{aligned} \quad \square$$

Theorem 16.4. *The entanglement fidelity $\mathcal{F}(\rho, \mathcal{E})$ is convex in ρ and for every ensemble $\{p_i \tilde{\rho}_i\}$ satisfies the following bound*

$$\begin{aligned}\mathcal{F}\left(\sum_j p_j \tilde{\rho}_j, \mathcal{E}\right) &\leq \sum_j p_j \mathcal{F}(\tilde{\rho}_j, \mathcal{E}) \\ &\leq \sum_j p_j F(\tilde{\rho}_j, \mathcal{E}(\rho_j))^2.\end{aligned}\quad (16.4)$$

Proof. From the definition of entanglement fidelity and the monotonicity of the Uhlmann fidelity under partial trace, one has

$$\mathcal{F}(\rho, \mathcal{E}) = F(|\Phi_\rho\rangle\langle|\Phi_\rho|, \mathcal{E} \otimes \mathcal{I}(|\Phi_\rho\rangle\langle|\Phi_\rho|))^2 \leq F(\rho, \mathcal{E}(\rho))^2.$$

Moreover, the entanglement fidelity is a convex function of ρ , since one has

$$\begin{aligned}\mathcal{F}(p\rho_1 + (1-p)\rho_2, \mathcal{E}) &= p^2 \langle\langle \rho_1 | \text{Ch}(\mathcal{E}) | \rho_1 \rangle\rangle + (1-p)^2 \langle\langle \rho_2 | \text{Ch}(\mathcal{E}) | \rho_2 \rangle\rangle \\ &\quad + p(1-p) \langle\langle \rho_1 | \text{Ch}(\mathcal{E}) | \rho_2 \rangle\rangle + p(1-p) \langle\langle \rho_2 | \text{Ch}(\mathcal{E}) | \rho_1 \rangle\rangle\end{aligned}$$

for all $0 \leq p \leq 1$. By the Cauchy-Schwartz inequality, one has $|\langle\langle \rho_1 | \text{Ch}(\mathcal{E}) | \rho_2 \rangle\rangle| \leq \sqrt{\langle\langle \rho_1 | \text{Ch}(\mathcal{E}) | \rho_1 \rangle\rangle \langle\langle \rho_2 | \text{Ch}(\mathcal{E}) | \rho_2 \rangle\rangle}$, and finally

$$\begin{aligned}\mathcal{F}(p\rho_1 + (1-p)\rho_2, \mathcal{E}) &\leq (p\sqrt{\langle\langle \rho_1 | \text{Ch}(\mathcal{E}) | \rho_1 \rangle\rangle} + (1-p)\sqrt{\langle\langle \rho_2 | \text{Ch}(\mathcal{E}) | \rho_2 \rangle\rangle})^2 \\ &\leq p\langle\langle \rho_1 | \text{Ch}(\mathcal{E}) | \rho_1 \rangle\rangle + (1-p)\langle\langle \rho_2 | \text{Ch}(\mathcal{E}) | \rho_2 \rangle\rangle \\ &= p\mathcal{F}(\rho_1, \mathcal{E}) + (1-p)\mathcal{F}(\rho_2, \mathcal{E}),\end{aligned}$$

where the second inequality is a straightforward consequence of convexity of $f(x) = x^2$ in the variable x . Therefore, one has

$$F\left(\sum_j p_j \rho_j, \mathcal{E}\right) \leq \sum_j p_j \mathcal{F}(\rho_j, \mathcal{E}) \leq \sum_j p_j F(\rho_j, \mathcal{E}(\rho_j))^2. \quad \square$$

16.2 Refinement set

Let us say that the state σ *refines* ρ if the preparation of the state ρ can be achieved by an algorithm that involves preparation in state σ with some non-null probability. The remainder of the section is dedicated to the relation between entanglement fidelity of a channel \mathcal{E} for state ρ and the behaviour of \mathcal{E} on the set of every possible state σ that refines ρ . This set is called the *refinement set* of ρ .

Definition 16.5 (Refinement set). The *refinement set* of a state $\rho \in \text{St}_1(A)$ is the set $\text{Ref}_1(\rho)$ of all the deterministic states σ that refine ρ , namely those states $\sigma \in \text{St}_1(A)$ for which there exist $0 < p \leq 1$ and $\tau \in \text{St}_1(A)$ such that

$$\rho = p\sigma + (1-p)\tau. \quad (16.5)$$

We now prove a theorem that regards the structure of the refinement set of a state ρ , and will be useful for the following results.

Theorem 16.6. *Let $\rho, \sigma \in \text{St}_1(\mathcal{A})$. Then one has $\text{Supp}(\sigma) \subseteq \text{Supp}(\rho)$ iff $\sigma \in \text{Ref}_1(\rho)$.*

Proof. Let us prove that the condition is necessary. Indeed, if $\rho = \lambda\sigma + (1 - \lambda)\tau$ with $0 < \lambda \leq 1$ then $0 \leq \lambda\sigma \leq \rho$, and consequently

$$\langle \psi | \rho | \psi \rangle = 0 \Rightarrow \langle \psi | \sigma | \psi \rangle = 0, \quad \forall \psi \in \mathcal{H}_{\mathcal{A}}.$$

This implies that $\text{Ker}(\rho) \subseteq \text{Ker}(\sigma)$, i.e. $\text{Supp}(\sigma) \subseteq \text{Supp}(\rho)$. On the other hand, the condition is sufficient. Indeed, if $\text{Ker}(\rho) \subseteq \text{Ker}(\sigma)$, let ρ^{-1} denote the Moore-Penrose generalised inverse of ρ , namely $\rho\rho^{-1} = \rho^{-1}\rho = \Pi_{\text{Supp}(\rho)}$, and $\text{Supp}(\rho^{-1}) = \text{Supp}(\rho)$. In this case, let us define

$$\lambda_M := \rho(\rho^{-\frac{1}{2}}\sigma\rho^{-\frac{1}{2}}),$$

where $\rho(X)$ denotes the spectral radius of X . Since $\rho^{-\frac{1}{2}}\sigma\rho^{-\frac{1}{2}}$ is selfadjoint, by lemma 13.4 it is $\lambda_M = \|\rho^{-\frac{1}{2}}\sigma\rho^{-\frac{1}{2}}\|_{\infty}$. Then by lemma 13.5 we have

$$\rho^{-\frac{1}{2}}\sigma\rho^{-\frac{1}{2}} \leq \lambda_M I_{\mathcal{A}}. \quad (16.6)$$

Now, multiplying on both sides by $\rho^{\frac{1}{2}}$, and reminding that $\text{Supp}(\sigma) \subseteq \text{Supp}(\rho)$, we have

$$\sigma \leq \lambda_M \rho.$$

There must then exist a positive operator $\tilde{\tau}$ such that

$$\rho = \frac{1}{\lambda_M} \sigma + \tilde{\tau},$$

and clearly $0 \leq \text{Tr}[\tilde{\tau}] = 1 - \frac{1}{\lambda_M}$, thus $\lambda_M \geq 1$. Finally

$$\rho = \mu\sigma + (1 - \mu)\tau,$$

with $\mu = \frac{1}{\lambda_M} \leq 1$ and $(1 - \mu)\tau = \tilde{\tau}$. \square

Then, the above theorem states that σ refines ρ if and only if its support is contained in the support of ρ . In other words,

$$\text{Ref}_1(\rho) = \{\sigma \in \text{St}_1(\mathcal{A}) | \text{Supp}(\sigma) \subseteq \text{Supp}(\rho)\}. \quad (16.7)$$

We now introduce the notion of a *refinement* of the state ρ , as follows

Definition 16.7. Let $\rho \in \text{St}_1(\mathcal{A})$. The ensemble $\{\rho_i\}_{i \in I}$ is a *refinement* of ρ if

$$\sum_{i \in I} \rho_i = \rho. \quad (16.8)$$

Clearly, for every ρ_i in a refinement of ρ , the deterministic state $\tilde{\rho}_i := \rho_i / \text{Tr}[\rho_i]$ refines ρ . The next result deals with refinements of a given state ρ .

Theorem 16.8. *Let $\rho \in \text{St}_1(\text{A})$. Then $\{\rho_i\}_{i \in I}$ is a refinement of ρ iff for every purification $\Psi \in \text{St}_1(\text{AB})$ of ρ , there exists a POVM $\{Q_i\}_{i \in I}$ such that*

$$\langle \rho_i | A \rangle = \left(\begin{array}{c} \Psi \\ \eta \end{array} \right) \quad (16.9)$$

The diagram shows a box labeled $\langle \rho_i |$ with an input line labeled A . This is equal to a box labeled Ψ with an output line labeled A and a line labeled B leading to a box labeled Q_i .

Proof. Eq. (16.9) is clearly sufficient for $\{\rho_i\}_{i \in I}$ to be an ensemble. Let us now prove that it is also necessary. First, let us construct the state $P = \sum_{i \in I} \rho_i \otimes |\varphi_i\rangle\langle\varphi_i| \in \text{St}_1(\text{AC})$, then consider its purification $\Sigma \in \text{St}_1(\text{ACD})$. Clearly, the POVM $\{P_i \otimes I_D\}_{i \in I}$ with $P_i := |\varphi_i\rangle\langle\varphi_i|$ is such that

$$\langle \rho_i | A \rangle = \left(\begin{array}{c} \Sigma \\ \eta \end{array} \right) \quad .$$

The diagram shows a box labeled $\langle \rho_i |$ with an input line labeled A . This is equal to a box labeled Σ with an output line labeled A and a line labeled CD leading to a box labeled $P_i \otimes I_D$.

Now, let $\Psi \in \text{St}_1(\text{AB})$ be any purification of ρ . Clearly, given two pure states $\eta \in \text{St}_1(\text{CD})$ and $\phi \in \text{St}_1(\text{B})$, the pure states $\Sigma \otimes \phi$ and $\Psi \otimes \eta$ are purifications of ρ with the same purifying system BCD, thus by theorem 11.3 there exists a unitary operator U on \mathcal{H}_{BCD} such that

$$\left(\begin{array}{c} \Psi \\ \eta \end{array} \right) = \left(\begin{array}{c} \Sigma \\ \eta \end{array} \right) \quad .$$

The diagram shows two boxes. The left box has an input line labeled B and an output line labeled A . The right box has an input line labeled CD and an output line labeled A . Both boxes have a line labeled B leading to a box labeled \mathcal{N}_U , which in turn has a line labeled CD leading to the right box's CD input.

Thus, upon defining the POVM $\{Q_i\}_{i \in I}$ as

$$\langle B | Q_i \rangle := \left(\begin{array}{c} B \\ \eta \end{array} \right) \quad ,$$

The diagram shows a box labeled $\langle B | Q_i \rangle$ with an input line labeled B . This is equal to a box labeled η with an input line labeled CD and an output line labeled B , followed by a box labeled \mathcal{N}_U , followed by a box labeled $P_i \otimes I_D$ with an input line labeled CD and an output line labeled B , followed by a box labeled I_B .

one obtains

$$\langle \rho_i | A \rangle = \left(\begin{array}{c} \Psi \\ \eta \end{array} \right) \quad . \quad \square$$

The diagram shows a box labeled $\langle \rho_i |$ with an input line labeled A . This is equal to a box labeled Ψ with an output line labeled A and a line labeled B leading to a box labeled Q_i .

Corollary 16.9 (Conditional marginal). *Let $\sigma, \rho \in \text{St}_1(\text{A})$. Then $\sigma \in \text{Ref}_1(\rho)$ iff there exists $0 < \lambda \leq 1$, and for any purification $\Phi_\rho \in \text{St}_1(\text{AB})$ there exists an effect $P \in \text{Eff}(\text{B})$, such that*

$$\lambda \langle \sigma | A \rangle = \left(\begin{array}{c} \Phi_\rho \\ \eta \end{array} \right) \quad . \quad (16.10)$$

The diagram shows a box labeled $\lambda \langle \sigma |$ with an input line labeled A . This is equal to a box labeled Φ_ρ with an output line labeled A and a line labeled B leading to a box labeled P .

Proof. By definition $\sigma \in \text{Ref}_1(\rho)$ iff there exists $0 < \lambda \leq 1$ and $\tau \in \text{St}_1(A)$ such that the ensemble $\{\lambda\sigma, (1 - \lambda)\tau\}$ refines ρ . By theorem 16.8, this is true iff for every purification $\Phi_\rho \in \text{St}_1(AB)$ there exists a binary POVM $\{P, I - P\}$ such that one has

$$\lambda (\sigma) = \begin{array}{c} \text{---} \\ | \\ \Phi_\rho \\ | \\ \text{---} \end{array} \quad , \quad (\tau) = \begin{array}{c} \text{---} \\ | \\ \Phi_\rho \\ | \\ \text{---} \end{array} \quad , \quad \square$$

16.3 Equality upon input

The notion that we introduce now is that of *equality upon input*, and determines when two channels behave in the same way on some subspace, that is typically identified with the support of a state ρ .

Definition 16.10 (Equality upon input). Let \mathcal{E}, \mathcal{F} be two quantum operations in $\text{QO}(A \rightarrow B)$. We say that \mathcal{E} is equal to \mathcal{F} upon input of ρ , and write

$$\mathcal{E} =_\rho \mathcal{F}, \tag{16.11}$$

if $\mathcal{E}(\sigma) = \mathcal{F}(\sigma)$ for all $\sigma \in \text{Ref}_1(\rho)$.

Lemma 16.11. Let $\mathcal{E}, \mathcal{F} \in \text{QO}(A \rightarrow B)$. Then one has $\mathcal{E} =_\rho \mathcal{F}$ if and only if $\mathcal{E} \otimes \mathcal{I}_C(\Phi_\rho) = \mathcal{F} \otimes \mathcal{I}_C(\Phi_\rho)$ for any purification $\Phi_\rho \in \text{St}_1(AC)$ of ρ .

Proof. Indeed, $\mathcal{E} \otimes \mathcal{I}_B(\Phi_\rho) = \mathcal{F} \otimes \mathcal{I}_B(\Phi_\rho)$ if and only if

$$\begin{array}{c} \text{---} \\ | \\ \Phi_\rho \\ | \\ \text{---} \end{array} \quad \begin{array}{c} \text{---} \\ | \\ \mathcal{E} \\ | \\ \text{---} \\ \text{---} \\ | \\ \mathcal{F} \\ | \\ \text{---} \\ \text{---} \end{array} \quad .$$

Then, equivalently we can write

$$\begin{array}{c} \text{---} \\ | \\ \Phi_\rho \\ | \\ \text{---} \end{array} \quad \begin{array}{c} \text{---} \\ | \\ \mathcal{E} \\ | \\ Q \\ | \\ \text{---} \\ \text{---} \\ | \\ \mathcal{F} \\ | \\ Q \\ | \\ \text{---} \\ \text{---} \end{array} \quad = \quad \begin{array}{c} \text{---} \\ | \\ \Phi_\rho \\ | \\ \text{---} \end{array} \quad \begin{array}{c} \text{---} \\ | \\ \mathcal{F} \\ | \\ Q \\ | \\ \text{---} \\ \text{---} \\ | \\ \mathcal{E} \\ | \\ P \\ | \\ \text{---} \\ \text{---} \end{array} \quad ,$$

for every $P \in \text{Eff}(A)$ and $Q \in \text{Eff}(B)$, which in turn is equivalent to

$$\begin{array}{c} \text{---} \\ | \\ \Phi_\rho \\ | \\ \text{---} \end{array} \quad \begin{array}{c} \text{---} \\ | \\ \mathcal{E} \\ | \\ \text{---} \\ \text{---} \\ | \\ \mathcal{F} \\ | \\ \text{---} \\ \text{---} \end{array} \quad = \quad \begin{array}{c} \text{---} \\ | \\ \Phi_\rho \\ | \\ \text{---} \end{array} \quad \begin{array}{c} \text{---} \\ | \\ \mathcal{F} \\ | \\ \text{---} \\ \text{---} \\ | \\ \mathcal{E} \\ | \\ P \\ | \\ \text{---} \\ \text{---} \end{array} \quad ,$$

for every $P \in \text{Eff}(A)$. Finally, by corollary 16.9 the latter is equivalent to $\lambda \mathcal{E}(\sigma) = \lambda \mathcal{F}(\sigma)$ for every $\sigma \in \text{Ref}_1(\rho)$, and some $0 < \lambda \leq 1$ depending on σ , i.e. $\mathcal{E}(\sigma) = \mathcal{F}(\sigma)$ for every $\sigma \in \text{Ref}_1(\rho)$. \square

Now, the following result holds

Theorem 16.12. Let $\mathcal{E} \in \text{QC}(A \rightarrow A)$ be a quantum channel. Then

$$\mathcal{F}(\rho, \mathcal{E}) = 1 \Leftrightarrow \mathcal{E} =_\rho \mathcal{I}_A \tag{16.12}$$

Proof. By definition, $\mathcal{F}(\rho, \mathcal{E}) = F[\Phi_\rho, (\mathcal{E} \otimes \mathcal{I}_B)(\Phi_\rho)]^2$, and thus $\mathcal{F}(\rho, \mathcal{E}) = 1$ if and only if

$$\begin{array}{c} \text{Circuit Diagram} \\ \Phi_\rho \xrightarrow{\mathcal{E}} A \\ \text{Circuit Diagram} \\ \Phi_\rho \xrightarrow{\mathcal{I}} A \end{array} = \quad (16.13)$$

for any purification Φ_ρ of ρ . Finally, by lemma 16.11, this is equivalent to $\mathcal{E} =_\rho \mathcal{I}_A$. \square

What can we say when the entanglement fidelity is *approximately* unit? One expects that, in some sense, in that case \mathcal{E} is approximately equal to \mathcal{I} upon input of ρ . Is that true? The answer to this question is provided by the following result.

Theorem 16.13. *Approximate equality upon input* For every $\varepsilon > 0$ there exists $\delta > 0$ such that if $\mathcal{F}(\rho, \mathcal{E}) \geq 1 - \delta$, then for every refinement $\{\rho_i\}_{i \in I} = \{p_i \tilde{\rho}_i\}$ of ρ one has

$$\frac{1}{2} \sum_{i \in I} p_i \|\mathcal{E}(\tilde{\rho}_i) - \rho_i\|_1 \leq \varepsilon. \quad (16.14)$$

Proof. Let Φ_ρ be any purification of ρ . Then

$$\frac{1}{2} \sum_{i \in I} p_i \|\mathcal{E}(\tilde{\rho}_i) - \rho_i\|_1 = \sum_{i \in I} p_i \text{Tr}[\{\mathcal{E}(\tilde{\rho}_i) - \rho_i\} Q_i],$$

for suitably chosen effects Q_i . By theorem 16.8, this is equivalent to

$$\begin{aligned} \frac{1}{2} \sum_{i \in I} p_i \|\mathcal{E}(\tilde{\rho}_i) - \rho_i\|_1 &= \sum_{i \in I} \left\{ \begin{array}{c} \text{Circuit Diagram} \\ \Phi_\rho \xrightarrow{\mathcal{E}} A \\ \text{Circuit Diagram} \\ \Phi_\rho \xrightarrow{Q_i} A \\ \text{Circuit Diagram} \\ \Phi_\rho \xrightarrow{P_i} A \end{array} - \begin{array}{c} \text{Circuit Diagram} \\ \Phi_\rho \xrightarrow{Q_i} A \\ \text{Circuit Diagram} \\ \Phi_\rho \xrightarrow{P_i} A \end{array} \right\} \\ &= \left\{ \begin{array}{c} \text{Circuit Diagram} \\ \Phi_\rho \xrightarrow{\mathcal{E}} A \\ \text{Circuit Diagram} \\ \Phi_\rho \xrightarrow{A} A \end{array} - \begin{array}{c} \text{Circuit Diagram} \\ \Phi_\rho \xrightarrow{A} A \\ \text{Circuit Diagram} \\ \Phi_\rho \xrightarrow{A} A \end{array} \right\}, \end{aligned}$$

where $A = \sum_{i \in I} P_i \otimes Q_i$ is an effect, corresponding to a coarse graining of the POVM $\{P_i \otimes Q_i, P_i \otimes (I - Q_i)\}_{i \in I}$. Then, by equation (15.19), this implies that

$$\frac{1}{2} \sum_{i \in I} p_i \|\mathcal{E}(\tilde{\rho}_i) - \rho_i\|_1 \leq \frac{1}{2} \|\mathcal{E} \otimes \mathcal{I}(\Phi_\rho) - \Phi_\rho\|_1.$$

Now, by the Fuchs-van de Graaf inequality we have

$$\begin{aligned} \frac{1}{2} \sum_{i \in I} p_i \|\mathcal{E}(\tilde{\rho}_i) - \rho_i\|_1 &\leq \sqrt{1 - F^2[\mathcal{E} \otimes \mathcal{I}(\Phi_\rho), \Phi_\rho]} \\ &= \sqrt{1 - \mathcal{F}(\rho, \mathcal{E})} \leq \sqrt{\delta}. \quad \square \end{aligned}$$

16.4 Schumacher's quantum source coding theorem

The main idea at the basis of the quantum source coding theorem is a quantum version of the notions of typical sequence and typical set.

Definition 16.14 (Typical subspace). Given a quantum state ρ with orthonormal decomposition $\rho = \sum_{x_i \in \text{Rng}(X)} p_i |x_i\rangle\langle x_i|$, the ε -typical subspace $H_{N,\varepsilon}(\rho)$ of $\mathcal{H}^{\otimes N}$ is defined as

$$H_{N,\varepsilon}(\rho) := \text{Span}\{|x_i\rangle \mid x_i \in T_{N,\varepsilon}(X)\}. \quad (16.15)$$

where $|x_i\rangle := |x_{i_1}\rangle|x_{i_2}\rangle\dots|x_{i_N}\rangle$, and X is the random variable with $\text{Rng}(X) = \{x_i\}$ and $\mathbb{P}_X(x_i) := p_i$.

We remind that N i.i.d. copies of the state ρ have the following form

$$\rho^{\otimes N} = \sum_{x_i \in \text{Rng}(X)^N} p_i |x_i\rangle\langle x_i|,$$

where $p_i = p_{i_1}p_{i_2}\dots p_{i_N}$.

It is an immediate consequence of the definition of typical subspace that

$$H_{N,\varepsilon}(\rho) := \text{Span}\left\{|x_i\rangle \mid \left|\frac{1}{N} \log_2 \frac{1}{\mathbb{P}_{X^N}(x_i)} - S(\rho)\right| \leq \varepsilon\right\}.$$

We will also denote the projector on the typical subspace as

$$\begin{aligned} P_{N,\varepsilon}(\rho) &:= \sum_{x_i \in T_{N,\varepsilon}(X)} |x_i\rangle\langle x_i| \\ &= \sum_{x_i \in T_{N,\varepsilon}(X)} |x_{i_1}\rangle\langle x_{i_1}| \otimes |x_{i_2}\rangle\langle x_{i_2}| \otimes \dots |x_{i_N}\rangle\langle x_{i_N}|, \end{aligned} \quad (16.16)$$

and we have that

$$\dim(H_{N,\varepsilon}(\rho)) = \text{Tr}[P_{N,\varepsilon}(\rho)] = |T_{N,\varepsilon}(X)| \quad (16.17)$$

It is immediate to see that

$$\text{Tr}[P_{N,\varepsilon}(\rho)\rho^{\otimes N}] = \sum_{x_i \in T_{N,\varepsilon}(X)} \mathbb{P}_{X^N}(x_i) = \mathbb{P}_{X^N}[x_i \in T_{N,\varepsilon}(X)]. \quad (16.18)$$

Theorem 16.15 (Typical subspace). *The following statements hold:*

1. For every $\varepsilon > 0$ and $\delta > 0$ there exists N_0 such that for every $N \geq N_0$

$$\text{Tr}[P_{N,\varepsilon}(\rho)\rho^{\otimes N}] \geq 1 - \delta. \quad (16.19)$$

2. For every $\epsilon > 0$ and $\delta > 0$ there exists N_0 such that for every $N \geq N_0$ the dimension of the typical subspace $H_{N,\varepsilon}(\rho)$ is bounded as

$$(1 - \delta)2^{N(S(\rho) - \varepsilon)} \leq \dim(H_{N,\varepsilon}(\rho)) \leq 2^{N(S(\rho) + \varepsilon)} \quad (16.20)$$

3. Let S_N denote for every N an arbitrary orthogonal projection on a subspace of $\mathcal{H}^{\otimes N}$ with dimension $\text{Tr}(S_N) < 2^{NR}$, with $R < S(\rho)$ fixed. Then for every $\delta > 0$ there exists N_0 such that for every $N \geq N_0$ and every choice of S_N

$$\text{Tr}[S_N \rho^{\otimes N}] \leq \delta. \quad (16.21)$$

Proof. 1. This is an immediate consequence of identity (16.18) and of item 2 of the equipartition theorem 4.15.

2. Using identity (16.17), this is just item 3 of the equipartition theorem.
3. This is the only item that is not an immediate consequence of the asymptotic equipartition theorem: let us choose ε such that $R < S(\rho) - \varepsilon$, and write

$$\text{Tr}[S_N \rho^{\otimes N}] = \text{Tr}[S_N \rho^{\otimes N} P_{N,\varepsilon}(\rho)] + \text{Tr}[S_N \rho^{\otimes N} (I - P_{N,\varepsilon}(\rho))].$$

Reminding that $\rho^{\otimes N}$ commutes with $P_{N,\varepsilon}(\rho)$, we have

$$\text{Tr}[S_N \rho^{\otimes N} P_{N,\varepsilon}(\rho)] = \text{Tr}[S_N P_{N,\varepsilon}(\rho) \rho^{\otimes N} P_{N,\varepsilon}(\rho)],$$

and since $P_{N,\varepsilon}(\rho) \rho^{\otimes N} P_{N,\varepsilon}(\rho) = \sum_{x_i \in \mathcal{T}_{N,\varepsilon}(X)} \mathbb{P}_{X^N}(x_i) |x_i\rangle\langle x_i|$, using item 2 of the asymptotic equipartition theorem one has

$$P_{N,\varepsilon}(\rho) \rho^{\otimes N} P_{N,\varepsilon}(\rho) \leq 2^{-N(S(\rho)-\varepsilon)} I,$$

which implies

$$\begin{aligned} \text{Tr}[S_N \rho^{\otimes N} P_{N,\varepsilon}(\rho)] &\leq \text{Tr}[S_N] 2^{-N(S(\rho)-\varepsilon)} \\ &\leq 2^{N[R-S(\rho)+\varepsilon]} \\ &\leq \frac{\delta}{2}, \end{aligned}$$

since the dimension of the subspace on which S_N projects is smaller than 2^{NR} . Notice that the bound above holds with the same δ for every S_N satisfying condition (16.21). On the other hand, using item 2 of this theorem, for every $\alpha > 0$ there is N_0 such that for $N \geq N_0$ we can bound the second term as

$$\text{Tr}[S_N \rho^{\otimes N} (I - P_{N,\varepsilon})] \leq \text{Tr}[\rho^{\otimes N} (I - P_{N,\varepsilon})] \leq 1 - (1 - \alpha) = \alpha.$$

Notice that also the second bound above holds with the same α independently of S_N . Therefore, in conclusion for sufficiently large N we have

$$\text{Tr}[S_N \rho^{\otimes N}] \leq \alpha + 2^{N(R-S(\rho)+\varepsilon)}. \quad (16.22)$$

Since $R < S(\rho)$, for every $\delta > 0$ one can find N_0 such that for every $N \geq N_0$ the last bound can be made smaller than δ , by making $\alpha < \delta/2$ and $2^{N(R-S(\rho)+\varepsilon)} \leq \delta/2$.

□

Theorem 16.16 (Schumacher). *For every i.i.d. quantum source with Hilbert space \mathcal{H}_A and prior state $\rho \in \text{St}_1(\mathbf{A})$, for every $\delta > 0$ and $R > S(\rho)$ there exists N_0 such that for every $N \geq N_0$ one has a compression scheme $\{\mathcal{E}^N, \mathcal{D}^N\}$ with rate R , and $\mathcal{F}(\rho^{\otimes N}, \mathcal{D}^N \mathcal{E}^N) \geq 1 - \delta$. Conversely, for every $R < S(\rho)$ there is $\delta \geq 0$ such that for every compression scheme $\{\mathcal{E}^N, \mathcal{D}^N\}$ with rate R one has $\mathcal{F}(\rho^{\otimes N}, \mathcal{D}^N \mathcal{E}^N) \leq \delta$.*

Proof. Let $R > S(\rho)$. Take $\varepsilon > 0$ such that $R \geq S(\rho) + \varepsilon$. Then, by the typical subspace theorem one has that for every $\delta > 0$ there exists N_0 such that for every $N \geq N_0$

$$\text{Tr}[P_{N,\varepsilon}(\rho)\rho^{\otimes N}] \geq 1 - \delta, \quad \dim(\mathcal{H}_{N,\varepsilon}(\rho)) \leq 2^{N(S(\rho)+\varepsilon)} \leq 2^{NR}.$$

Now, consider a subspace $\mathcal{H}_{\text{comp}}^N$ of $\mathcal{H}_A^{\otimes N}$ with dimension 2^{NR} , containing $\mathcal{H}_{N,\varepsilon}$, and use the following compression scheme.

1. Perform the von Neumann-Lüders measurement $\{P_{N,\varepsilon}(\rho), I - P_{N,\varepsilon}(\rho)\}$. If the outcome corresponding to $P_{N,\varepsilon}(\rho)$ occurs, then leave the state unchanged. Otherwise, if the outcome corresponding to $I - P_{N,\varepsilon}(\rho)$ occurs, replace the state by a standard state $|S\rangle\langle S|$, with $|S\rangle \in \mathcal{H}_{N,\varepsilon}(\rho)$. Finally, isometrically embed $\mathcal{H}_{\text{comp}}^N$ into $\mathcal{H}_{B^N} \simeq \mathcal{H}_{\text{comp}}^N$ by the co-isometry V^\dagger [with $VV^\dagger = P_{\mathcal{H}_{\text{comp}}^N}$], complemented by an arbitrary channel from $(\mathcal{H}_{\text{comp}}^N)^\perp$ to \mathcal{H}_{B^N} , e.g. $\mathcal{A}(\sigma) = |\psi_0\rangle\langle\psi_0|$. Such an encoding is described by the quantum channel

$$\mathcal{E}^N : \mathbf{A}^{\otimes N} \rightarrow \mathbf{B}^N, \quad \mathcal{E}^N(\sigma) := V^\dagger P_{N,\varepsilon}(\rho) \sigma P_{N,\varepsilon}(\rho) V + \sum_i A_i \sigma A_i^\dagger,$$

where $A_i := V^\dagger |S\rangle\langle i|, \{|i\rangle\}$ denoting an orthonormal basis for $\mathcal{H}_{N,\varepsilon}(\rho)^\perp$.

2. The decoding is the channel $\mathcal{D}^N : \mathbf{B}^N \rightarrow \mathbf{A}^{\otimes N}$ which is just the isometry V inverting the embedding V^\dagger .

With the above compression scheme, we have according to the typical subspace theorem

$$\begin{aligned} \mathcal{F}(\rho^{\otimes N}, \mathcal{D}^N \mathcal{E}^N) &= |\text{Tr}[\rho^{\otimes N} P_{N,\varepsilon}(\rho)]|^2 + \sum_i |\text{Tr}(\rho^{\otimes N} V A_i)|^2 \\ &\geq |\text{Tr}[\rho^{\otimes N} P_{N,\varepsilon}(\rho)]|^2 \\ &\geq |1 - \delta|^2 \\ &\geq 1 - 2\delta, \end{aligned}$$

where we used identity (16.2) in the first equality. Hence there exists a reliable compression scheme with rate R for $R > S(\rho)$.

To prove the converse, suppose now that $R < S(\rho)$, and consider any compression scheme $\{\mathcal{E}^N, \mathcal{D}^N\}$ which compresses states on $\mathcal{H}_A^{\otimes N}$ to states on a system \mathbf{B}^N with $\dim(\mathcal{H}_{B^N}) = 2^{NR}$. In terms of the Kraus decompositions $\mathcal{E}^N(\sigma) = \sum_k C_k \sigma C_k^\dagger$ and $\mathcal{D}^N(\tau) = \sum_k D_k \tau D_k^\dagger$ we have

$$\mathcal{F}(\rho^{\otimes N}, \mathcal{D}^N \mathcal{E}^N) = \sum_{ij} |\text{Tr}(D_i C_j \rho^{\otimes N})|^2.$$

In order to compress to \mathcal{H}_B^N , the map \mathcal{E} clearly has Kraus operators with $\text{Rng}(C_j) \subseteq \mathcal{H}_{B^N}$, and the map \mathcal{D} has Kraus operators with support $\text{Supp}(D_i) \subseteq \mathcal{H}_{B^N}$. Denote by P_i the projection on $\text{Rng}(D_i) \simeq \text{Supp}(D_i)$, namely the space to which the compressed state is mapped by D_i . Thus $D_i C_j = P_i D_i C_j$, and

$$\begin{aligned}\mathcal{F}(\rho^{\otimes N}, \mathcal{D}^N \mathcal{E}^N) &= \sum_{ij} |\text{Tr}[D_i C_j \rho^{\otimes N} P_i]|^2 \\ &\leq \sum_{ij} \text{Tr}[D_i C_j \rho^{\otimes N} C_j^\dagger D_i^\dagger] \text{Tr}[P_i \rho^{\otimes N}],\end{aligned}$$

where we applied the Cauchy-Schwarz inequality

$$|\text{Tr}[A^\dagger B]|^2 \leq \text{Tr}[A^\dagger A] \text{Tr}[B^\dagger B],$$

with $A = \sqrt{\rho^{\otimes N}} C_j^\dagger D_i^\dagger$, and $B = \sqrt{\rho^{\otimes N}} P_i$. Since for every value of the index i one has $\text{Tr}[P_i] = \dim(\text{Supp}(P_i)) \leq \dim(\mathcal{H}_{B^N}) = 2^{NR}$, applying item 3 of the typical subspace theorem we have that for every i and for any $\delta > 0$ there is N_0 such that for $N \geq N_0$ one has $\text{Tr}[P_i \rho^{\otimes N}] \leq \delta$. Moreover, N_0 does not depend on the particular projections P_i as they all project on subspaces whose dimension is $\text{Tr}[P_i] \leq 2^{NR}$. Thus the following bound holds for all schemes $(\mathcal{E}^N, \mathcal{D}^N)$ with rate R and $N \geq N_0$

$$\mathcal{F}(\rho^{\otimes N}, \mathcal{D}^N \mathcal{E}^N) \leq \delta \sum_{ij} \text{Tr}(D_i C_j \rho^{\otimes N} C_j^\dagger D_i^\dagger) = \delta,$$

since both maps $\mathcal{D}^N, \mathcal{E}^N$ are trace-preserving. \square

Chapter 17

Lecture 20: Entropy exchange and coherent information

17.1 Entropy exchange

17.1.1 Quantum data-processing theorem

Is there a quantum version of the data-processing theorem? If so, what is the quantity that plays the role of mutual information? Indeed, as we will see in this chapter, there exists a quantum version of the data-processing theorem, in which the role of mutual information is played by the *coherent information*. Before introducing the coherent information, we need to introduce the notion of *entropy exchange*.

Definition 17.1 (Entropy exchange). Let $\mathcal{E} : Q \rightarrow Q'$ be a quantum channel and $\rho \in St_1(Q)$ be a state of the input system. We define the entropy exchange of \mathcal{E} upon input of ρ as the following quantity

$$S(\rho, \mathcal{E}) := S(\mathcal{E} \otimes \mathcal{I}_R(|\Phi_\rho\rangle\langle\Phi_\rho|)), \quad (17.1)$$

where $|\Phi_\rho\rangle\langle\Phi_\rho| \in St_1(QR)$ is an arbitrary purification of ρ .

We remind that any purification $|\Phi_\rho\rangle\langle\Phi_\rho|$ of ρ can be written as

$$\begin{aligned} |\Phi_\rho\rangle\langle\Phi_\rho| &= \mathcal{I}_Q \otimes \mathcal{N}_V(|\Phi\rangle\langle\Phi|), \\ |\Phi\rangle &= \sum_{j=1}^r \sqrt{p_i} |\psi_i\rangle |i\rangle, \\ r &:= \text{rank}(\rho), \quad \{|i\rangle\}_{i=1}^r \text{ o.n.b. in } \mathbb{C}^r =: \mathcal{H}_C, \end{aligned} \quad (17.2)$$

where $|\Phi\rangle\langle\Phi| \in St_1(QC)$ is the *minimal purification* of ρ , and $V : \mathcal{H}_C \rightarrow \mathcal{H}_R$ is an isometry. Then, due to isometric invariance of von Neumann entropy, the entropy exchange is invariant under local isometries on the purifying system, namely

$$\begin{aligned} S(\mathcal{E} \otimes \mathcal{I}_R(|\Phi_\rho\rangle\langle\Phi_\rho|)) &= S(\mathcal{E} \otimes \mathcal{N}_V(|\Phi\rangle\langle\Phi|)) \\ &= S[\mathcal{I}_{Q'} \otimes \mathcal{N}_V(\mathcal{E} \otimes \mathcal{I}_C(|\Phi\rangle\langle\Phi|))] \\ &= S(\mathcal{E} \otimes \mathcal{I}_C(|\Phi\rangle\langle\Phi|)). \end{aligned} \quad (17.3)$$

The definition of entropy exchange is then independent of the particular purification $|\Phi_\rho\rangle\langle\Phi_\rho|$.

Why the name entropy exchange? Because, if one considers that the channel \mathcal{E} has a unitary dilation in terms of an interaction between the quantum system Q and an environment E in a pure state η , and one purifies the initial state by the reference system R in the pure state $\Phi_\rho := |\Phi_\rho\rangle\langle\Phi_\rho|$ as in the following diagram



then the state of $Q'E'R'$ after the interaction is still pure, and one has

$$S(E') = S(Q'R') = S(\rho, \mathcal{E}). \quad (17.5)$$

In other words, the entropy exchange is exactly the amount of entropy introduced (“exchanged”) in the environment. Due to the Holevo bound, the classical mutual information between Alice preparing the system Q and any user Eve controlling the environment E is bounded by

$$I(E, P) \leq S(\rho^{E'}) - \sum_{i \in X} p_i S(\rho_i^{E'}) \leq S(\rho^{E'}) = S(\rho, \mathcal{E}).$$

A small entropy exchange is thus not only a signature of a reliable communication, but also of a reliably *private* one. This feature is distinctive of quantum communication theory compared to the classical one.

Notice that this amount of entropy depends neither on the particular unitary interaction with the environment used to model the channel \mathcal{E} , nor on the particular purification Φ_ρ of ρ , but only on the channel itself and the initial state ρ of the system Q .

It is interesting to write down explicitly the state of the environment after the application of the unitary transformation which dilates the channel on the system Q . Let $\{E_n\}_{n \in X} \subseteq \mathcal{L}(\mathcal{H}_Q \rightarrow \mathcal{H}_{Q'})$ be the canonical Kraus of \mathcal{E} . Then the operator $V : \mathcal{H}_Q \rightarrow \mathcal{H}_{Q'E'}$ defined as

$$V := \sum_{n \in X} E_n \otimes |n\rangle,$$

for an orthonormal set $\{|n\rangle\}_{n=1}^{d_E}$ on $\mathcal{H}_{E'}$ with $d_{E'} \geq |X|$ is isometric, since

$$V^\dagger V = \sum_{n \in X} E_n^\dagger E_n = I_Q.$$

The isometry V can be extended to a unitary transformation over system Q and environment E provided that $d_E d_Q = d_{E'} d_{Q'}$. Such a unitary must then satisfy

$$U|\psi\rangle \otimes |0\rangle = \sum_{n \in X} E_n |\psi\rangle \otimes |n\rangle, \quad (17.6)$$

where $|\psi\rangle$ is an arbitrary vector in \mathcal{H}_Q . For mixed input state we have

$$U(\rho \otimes |0\rangle\langle 0|)U^\dagger = \sum_{nm} E_n \rho E_m^\dagger \otimes |n\rangle\langle m|,$$

and the state of the environment E' after the interaction is given by

$$\sigma := \rho^{E'} = \sum_{nm} \sigma_{nm} |n\rangle\langle m|, \quad \sigma_{nm} := \text{Tr}[E_n \rho E_m^\dagger].$$

Therefore, according to Eq. (17.5) the entropy exchange can be written as

$$S(\rho, \mathcal{E}) = S(\sigma) = -\text{Tr}(\sigma \log_2 \sigma). \quad (17.7)$$

Exercise 17.1

Prove that the entropy exchange has the property that $S(d_Q^{-1} I_Q, \mathcal{E}) = 0$ if and only if the channel $\mathcal{E} : Q \rightarrow Q'$ is an isometry.

Answer of exercise 17.1

The state of the environment E' after the interaction is given by

$$\sigma_{m,n} = \text{Tr}[E_m \rho E_n^\dagger], \quad (17.8)$$

and if we diagonalise it we obtain

$$\delta_{i,j} q_i = \sum_{m,n=1}^{d_{E'}} W_{i,m} \sigma_{m,n} W_{n,j}^\dagger = \text{Tr}[F_i \rho F_j^\dagger],$$

where $F_i := \sum_{m=1}^{d_{E'}} W_{i,m} E_m$ is an alternate Kraus for \mathcal{E} . The entropy exchange in terms of the canonical environment density matrix is given by

$$-\text{Tr}(\sigma \log_2 \sigma) = -\sum_i q_i \log_2 q_i, \quad (17.9)$$

where $\{q_i\} := \text{Spec}(\sigma)$, and

$$\text{Tr}[F_i \rho F_i^\dagger] = q_i, \quad \text{Tr}[F_i \rho F_j^\dagger] = 0 \quad i \neq j. \quad (17.10)$$

Now, $S(d_Q^{-1} I_Q, \mathcal{E}) = 0$ if and only if $q_i = \delta_{ii_0}$, namely

$$\frac{1}{d_Q} \text{Tr}[F_i F_i^\dagger] = \delta_{i,i_0}, \quad \text{Tr}[F_i F_j^\dagger] = 0 \quad i \neq j. \quad (17.11)$$

Since $0 \leq F_{i_0}^\dagger F_{i_0} \leq I_Q$, the only way to satisfy the first identity in 17.11 is to have $F_i^\dagger F_i = \delta_{i,i_0} I_Q$, which implies that F_{i_0} is isometric and, since the map \mathcal{E} is trace preserving, all the remaining Kraus operators are equal to zero.

Exercise 17.2

Prove concavity of $S(\rho, \mathcal{E})$ versus the channel \mathcal{E} .

Answer of exercise 17.2

The property follows immediately from concavity of the von Neumann entropy. Intuitively this means that mixing two different maps produces to an increase of the noise, and consequently to a larger inflation of entropy into the environment. Indeed, one has

$$\begin{aligned} S(\rho, p\mathcal{E}_1 + (1-p)\mathcal{E}_2) &= S[(p\mathcal{E}_1 + (1-p)\mathcal{E}_2) \otimes \mathcal{I}_R(|\Phi_\rho\rangle\langle\Phi_\rho|)] \\ &\leq pS(\mathcal{E}_1 \otimes \mathcal{I}_Q(|\Phi_\rho\rangle\langle\Phi_\rho|)) + (1-p)S(\mathcal{E}_2 \otimes \mathcal{I}_Q(|\Phi_\rho\rangle\langle\Phi_\rho|)) \\ &= pS(\rho, \mathcal{E}_1) + (1-p)S(\rho, \mathcal{E}_2). \end{aligned}$$

17.1.2 Quantum Fano inequality

A very interesting property of the entropy exchange is the quantum version of Fano's inequality

Theorem 17.2 (Quantum Fano's inequality). *The following inequality holds*

$$S(\rho, \mathcal{E}) \leq H_2(F(\rho, \mathcal{E})) + (1 - F(\rho, \mathcal{E})) \log_2(d_Q r - 1), \quad (17.12)$$

for any quantum channel $\mathcal{E} : Q \rightarrow Q$ and any state $\rho \in \text{St}_1(Q)$ with rank $r := \text{rank}(\rho)$.

Proof. Consider an orthonormal basis $\{|B_i\rangle\}_{i=1}^{d_Q r}$ for system and a reference R with $d_R = r := \text{rank}(\rho)$. Let $B_1 = \Phi_\rho$, and consider the probability distribution

$$p_i = \langle\langle B_i | \mathcal{E} \otimes \mathcal{I}_R(|\Phi_\rho\rangle\langle\Phi_\rho|) | B_i \rangle\rangle.$$

Using bound (12.27) we have

$$\begin{aligned} S(\rho, \mathcal{E}) &= S(\mathcal{E} \otimes \mathcal{I}_R(|\Phi_\rho\rangle\langle\Phi_\rho|)) \\ &\leq S\left(\sum_{i=1}^{d_Q r} |B_i\rangle\langle B_i| \mathcal{E} \otimes \mathcal{I}_R(|\Phi_\rho\rangle\langle\Phi_\rho|) |B_i\rangle\langle B_i|\right) \\ &= S\left(\sum_{i=1}^{d_Q r} p_i |B_i\rangle\langle B_i|\right) \\ &= H(\mathbf{p}). \end{aligned}$$

Then, it is easy to see that

$$\begin{aligned} H(\mathbf{p}) &= H_2(p_1) + (1 - p_1)H(\mathbf{q}) \\ &\leq H_2(p_1) + (1 - p_1) \log_2(d_Q r - 1), \\ q_j &:= \frac{p_j}{1 - p_1}, \quad 2 \leq j \leq d_Q r \end{aligned}$$

On the other hand, since we have chosen $B_1 = \Phi_\rho$, we have

$$p_1 = \langle\langle \Phi_\rho | \mathcal{E} \otimes \mathcal{I}_R(|\Phi_\rho\rangle\rangle \langle\langle \Phi_\rho|) |\Phi_\rho\rangle\rangle = F(\rho, \mathcal{E}),$$

Finally, by substitution we get the statement. \square

The quantum Fano inequality allows us to conclude that a necessary condition for a reliable quantum communication, corresponding to $F(\rho, \mathcal{E}) > 1 - \delta$, is granted by a small entropy exchange, namely a small leakage of information in the environment.

17.2 Reversibility upon input

In this section we introduce the notion of reversibility upon input of a channel, which lies at the basis of the theory of quantum error correction.

Definition 17.3 (Channel reversible upon input). A channel $\mathcal{E} : Q \rightarrow Q'$ is reversible upon input of $\rho \in \text{St}_1(Q)$ if there exists another channel $\mathcal{R} : Q' \rightarrow Q$ such that $\mathcal{R}\mathcal{E} =_\rho \mathcal{I}_Q$, or equivalently $F(\rho, \mathcal{R}\mathcal{E}) = 1$.

According to corollary 16.11, the above definition is equivalent to the condition

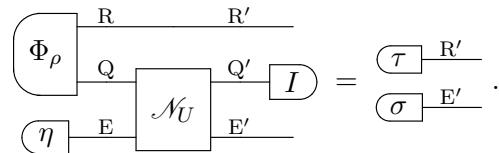
$$\mathcal{R}\mathcal{E} \otimes \mathcal{I}_R(|\Phi_\rho\rangle\rangle \langle\langle \Phi_\rho|) = |\Phi_\rho\rangle\rangle \langle\langle \Phi_\rho|, \quad (17.13)$$

for any purification $|\Phi_\rho\rangle\rangle \langle\langle \Phi_\rho| \in \text{St}_1(QR)$ of ρ .

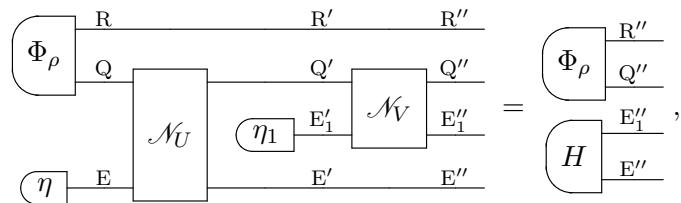
Theorem 17.4. *The channel $\mathcal{E} : Q \rightarrow Q'$ is reversible upon input of $\rho \in \text{St}_1(Q)$ if and only if the reference and the environment remain uncorrelated after the action of the channel. In formula, if we consider a general purification $\Phi_\rho := |\Phi_\rho\rangle\rangle \langle\langle \Phi_\rho| \in \text{St}_1(QR)$ of ρ , and a general unitary dilation (E, η, U) of \mathcal{E} , the equivalent condition for reversibility of \mathcal{E} upon input of ρ is*

$$\rho_{R'E'} = \text{Tr}_{Q'}[(I_R \otimes U)(\Phi_\rho \otimes \eta)(I_R \otimes U^\dagger)] = \rho_{R'} \otimes \rho_{E'} = \tau \otimes \sigma. \quad (17.14)$$

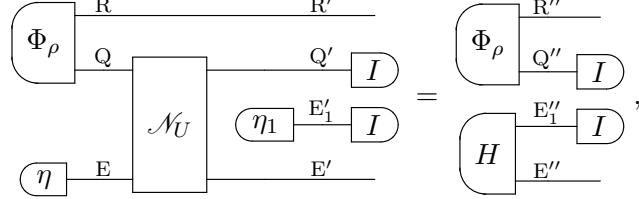
In diagrams



Proof. The converse statement can be proved by considering a reversible dilation (E'_1, η_1, V) of the correcting channel \mathcal{R} . In this case, we have that $\mathcal{I}_R \otimes \mathcal{R}\mathcal{E}(|\Phi_\rho\rangle\rangle \langle\langle \Phi_\rho|) = |\Phi_\rho\rangle\rangle \langle\langle \Phi_\rho|$, and then



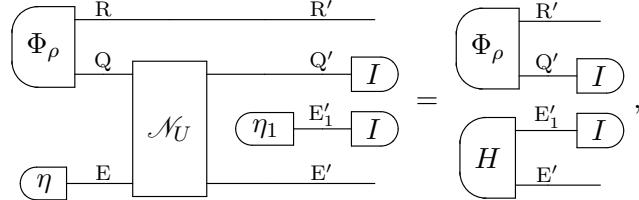
for some pure state $H \in \text{St}_1(E''_1 E'')$, being the state of $R'' Q'' E''_1 E''$ a purification of Φ_ρ . It is now easy to check that, applying the deterministic effect on $Q'' E''_1$ we obtain



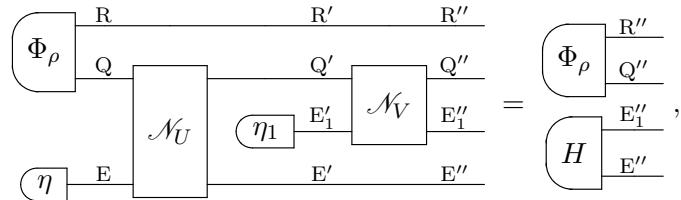
namely

$$\begin{array}{ccc} \text{Quantum circuit diagram showing the equivalence of two states. On the left, a state } \Phi_\rho \text{ is processed by a unitary } \mathcal{N}_U, \text{ followed by a measurement } M_{\eta_1} \text{ on line } E'_1, \text{ resulting in state } H. \text{ On the right, state } \Phi_\rho \text{ is processed by a unitary } \mathcal{N}_U, \text{ followed by a measurement } M_H \text{ on line } E'', \text{ resulting in state } H. \text{ The lines are labeled } R, Q, E, R', Q', E', R'', Q'', E'', E''_1. & & (17.15) \end{array}$$

Starting from equation (17.15) one can now prove sufficiency. Indeed, if equation (17.15) holds, then we can introduce a pure state $\eta_1 \in \text{St}_1(E'_1)$ such that



where $H \in \text{St}_1(E'E'_1)$ is a purification of σ . Thus, by uniqueness of purification there must exist a unitary $V \in \mathcal{L}(\mathcal{H}_{Q'E'_1})$ such that



and then discarding $E'' E''_1$ we obtain

$$\begin{array}{ccc} \text{Quantum circuit diagram showing the equivalence of two states. On the left, a state } \Phi_\rho \text{ is processed by a unitary } \mathcal{E}, \text{ followed by a unitary } \mathcal{R}, \text{ resulting in state } H. \text{ On the right, state } \Phi_\rho \text{ is processed by a unitary } \mathcal{R}, \text{ resulting in state } H. & & \end{array}$$

namely $\mathcal{R}\mathcal{E} \otimes \mathcal{I}_R(|\Phi_\rho\rangle\langle\Phi_\rho|) = |\Phi_\rho\rangle\langle\Phi_\rho|$. \square

We now provide a lemma that justifies the subject of the next section.

Lemma 17.5. *The condition $\rho_{R'E'} = \rho_{R'}\rho_{E'} = \tau \otimes \sigma$ is equivalent to*

$$S(\mathcal{E}(\rho)) = S(\rho) + S(\rho, \mathcal{E}). \quad (17.16)$$

Proof. Notice that since the state $\rho_{R'Q'E'}$ is pure one has

$$S(\mathcal{E}(\rho)) = S(Q') = S(R'E').$$

Moreover, by equation 17.5 one has

$$S(\rho, \mathcal{E}) = S(E').$$

Finally,

$$S(\rho) = S(Q) = S(R) = S(R'),$$

since QR is in a pure state, and $\rho_{R'} = \rho_R$. Now, identity (17.16) is equivalent to

$$S(R'E') = S(R') + S(E'),$$

which is satisfied if and only if R' and E' after the interaction remain in an uncorrelated state. \square

17.3 Coherent information and quantum data processing theorem

Let us consider again the identity (17.16), which is equivalent to $S(R'E') = S(R') + S(E')$. In general, one has $S(R'E') \leq S(R') + S(E')$, namely $S(\mathcal{E}(\rho)) \leq S(\rho) + S(\rho, \mathcal{E})$. This observation suggests the next definition.

Definition 17.6 (Coherent information). Let $\mathcal{E} : Q \rightarrow Q'$ be a quantum channel, and $\rho \in \text{St}_1(Q)$ a state of the input system. The coherent information of \mathcal{E} for state ρ is defined as

$$I(\rho, \mathcal{E}) := S(\mathcal{E}(\rho)) - S(\rho, \mathcal{E}). \quad (17.17)$$

Remark 20. Notice that in the quantum Fano inequality, the role of $H(X|Y)$ is played by $S(\rho, \mathcal{E})$. In a similar way, we will now see that there is a quantum version of the data processing theorem, and the role of $I(X : Y)$ in this case is played by the coherent information $I(\rho, \mathcal{E})$.

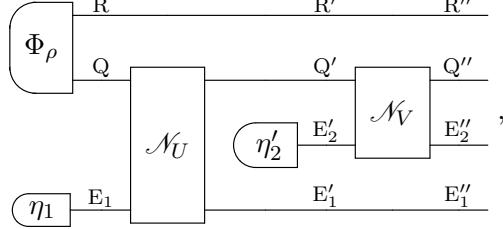
Theorem 17.7 (Quantum data-processing theorem). *Let $\rho \in \text{St}_1(A)$, $\mathcal{E}_1 : A \rightarrow B$ and $\mathcal{E}_2 : B \rightarrow C$ two quantum channels. The following chain of inequalities holds.*

$$I(\rho, \mathcal{E}_2 \mathcal{E}_1) \leq I(\rho, \mathcal{E}_1) \leq S(\rho). \quad (17.18)$$

Moreover, $I(\rho, \mathcal{E}_1) = S(\rho)$ if and only if \mathcal{E}_1 is reversible upon input of ρ .

Proof. The second bound is just a special case of the first one for $\mathcal{E}_1 = \mathcal{I}_A$. The fact that the last inequality is an identity if and only if \mathcal{E} is reversible upon input of ρ is already proved in lemma 17.5. Therefore, it is now sufficient to prove the first inequality. We

consider a dilation (E_1, η_1, U) and (η'_2, E'_2, V) for \mathcal{E}_1 and \mathcal{E}_2 , respectively. Including also the reference system R purifying ρ , we have



after the application of both channels. Using the strong subadditivity inequality for the von Neumann entropy

$$S(R''E''_1E''_2) + S(E''_1) \leq S(R''E''_1) + S(E''_1E''_2), \quad (17.19)$$

and since the total state of the four systems is pure, one has

$$S(R''E''_1E''_2) = S(Q'') = S(\mathcal{E}_2\mathcal{E}_1(\rho)).$$

Since R and E_1 are not involved in the second map, one has

$$S(R''E''_1) = S(R'E'_1) = S(Q') = S(\mathcal{E}_1(\rho)),$$

where we used also the fact that the state of $R'E'_1Q'$ is pure. The environmental entropies are, by definition, entropy exchanges, namely

$$S(E''_1) = S(E'_1) = S(\rho, \mathcal{E}_1), \quad S(E''_1E''_2) = S(\rho, \mathcal{E}_2\mathcal{E}_1),$$

and substitution into (17.19) gives

$$S(\mathcal{E}_2\mathcal{E}_1(\rho)) + S(\rho, \mathcal{E}_1) \leq S(\mathcal{E}_1(\rho)) + S(\rho, \mathcal{E}_2\mathcal{E}_1),$$

which is the inequality

$$I(\rho, \mathcal{E}_2\mathcal{E}_1) \leq I(\rho, \mathcal{E}_1). \quad \square$$

Exercise 17.3

Show that the quantum data-processing theorem reduces to the classical one for classical states and channels.

Answer of exercise 17.3

Corresponding to a random variable X with probability distribution $p(x)$ we have the quantum state $\rho_X = \sum_x p(x)|x\rangle\langle x|$, with $\{|x\rangle\}$ orthonormal set, and with the identification of entropies

$$S(\rho_X) = H(X).$$

Corresponding to the classical channel with conditional probabilities $p(y|x)$ one has a quantum channel that prepares the state $|y\rangle\langle y|$ with probability $p(y|x)$ for input state $|x\rangle\langle x|$, namely

$$\mathcal{E}_1(\rho) = \sum_{xy} K_{xy}\rho K_{xy}^\dagger, \quad K_{xy} = \sqrt{p(y|x)}|y\rangle\langle x|.$$

The output state is then

$$\mathcal{E}_1(\rho) = \sum_{xy} p(xy)|y\rangle\langle y| \quad (17.20)$$

$$= \sum_y p(y)|y\rangle\langle y|, \quad (17.21)$$

and

$$S(\mathcal{E}_1(\rho)) = H(Y).$$

The corresponding entropy exchange is given by

$$S(\rho, \mathcal{E}_1) = -\text{Tr}[\sigma \log_2 \sigma],$$

where

$$\sigma_{(xy), (x'y')} = \text{Tr}[K_{xy}\rho_X K_{x'y'}^\dagger] = \delta_{yy'}\delta_{xx'}p(x, y),$$

where $p(x, y) = p(y|x)p(x)$ is the joint probability of X, Y . Therefore, one obtains

$$S(\rho, \mathcal{E}_1) = H(X, Y).$$

Let now $X \leftrightarrow Y \leftrightarrow Z$ be a Markov chain, and let us define the channel \mathcal{E}_2 as

$$\mathcal{E}_2(\sigma) = \sum_{yz} K_{yz}\rho K_{yz}^\dagger, \quad K_{yz} = \sqrt{p(z|y)}|z\rangle\langle y|.$$

We then have

$$\mathcal{E}_2\mathcal{E}_1(\tau) = \sum_{xyz} p(z|y)p(y|x)|z\rangle\langle x|\tau|x\rangle\langle z|.$$

By the Markov property, we have $p(z|y) = p(z|xy)$, then $p(z|y)p(y|x) = p(z|xy)p(y|x) = p(yz|x)$. This implies that

$$\begin{aligned} \mathcal{E}_2\mathcal{E}_1(\tau) &= \sum_{xyz} p(yz|x)|z\rangle\langle x|\tau|x\rangle\langle z| \\ &= \sum_{xz} p(z|x)|z\rangle\langle x|\tau|x\rangle\langle z|. \end{aligned}$$

A Kraus form for $\mathcal{E}_2\mathcal{E}_1$ is then given by $H_{zx} = \sqrt{p(z|x)}|z\rangle\langle x|$, and

$$\mathcal{E}_2\mathcal{E}_1(\rho) = \sum_z p(z)|z\rangle\langle z|.$$

Repeating the argument that we used for \mathcal{E}_1 , we obtain $S(\mathcal{E}_2\mathcal{E}_1(\rho)) = H(Z)$ and $S(\rho, \mathcal{E}_2\mathcal{E}_1) = H(X, Z)$. Upon substituting the above identifications in the quantum data-processing theorem

$$S(\mathcal{E}_2\mathcal{E}_1(\rho)) - S(\rho, \mathcal{E}_2\mathcal{E}_1) \leq S(\mathcal{E}_1(\rho)) - S(\rho, \mathcal{E}_1).$$

one obtains

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

and by adding $H(X)$ to both members we have

$$I(X : Z) \leq I(X : Y),$$

which is the classical data-processing theorem.

Chapter 18

Lecture 21: The theory of entanglement

18.1 Entanglement

We defined entanglement as the property of a state $\rho \in \text{St}_1(\text{AB})$ of being not separable.

Definition 18.1 (Separable state). Let $\rho \in \text{St}_1(\text{AB})$. We say that ρ is *separable* if there exist $\{p_i\}_{i \in I}$ $0 \leq p_i \forall i$, $\sum_{i \in I} p_i = 1$, and $\{\sigma_i\}_{i \in I} \subseteq \text{St}_1(\text{A})$, $\{\tau_i\}_{i \in I} \subseteq \text{St}_1(\text{B})$, such that

$$\rho = \sum_{i \in I} p_i (\sigma_i \otimes \tau_i).$$

Definition 18.2 (Entangled state). Let $\rho \in \text{St}_1(\text{AB})$. We say that ρ is *entangled* if ρ is not separable.

The questions that we consider in these lectures is then: what do the properties of separability and entanglement mean operationally? Can we evaluate how much entanglement is there in a given state, and how? Let us start from the first question. Suppose Alice and Bob want to prepare a separable state $\rho = \sum_{i \in I} p_i (\sigma_i \otimes \tau_i)$. Then Alice can use a random source to sample i with probability distribution $\{p_i\}_{i \in I}$, and prepare σ_{i_0} upon reading $i = i_0$. Alice can then send a classical message to Bob, saying "I read $i = i_0$ ". Bob then prepares τ_{i_0} . At the end of the protocol, they prepared $\sigma_{i_0} \otimes \tau_{i_0}$, and this happens with probability p_{i_0} . Overall, they then prepare

$$\rho = \sum_{i \in I} p_i (\sigma_i \otimes \tau_i)$$

Separable states are precisely those states that can be prepared by Local Operations and Classical Communication (LOCC). Entanglement is the main resource in quantum information—though not the only one. It enables:

- Teleportation
- Dense coding
- EPR - QKD

- Most quantum algorithms

Thus, we are not satisfied with a definition, we would also like to assess the *amount* of entanglement in a given state $\rho \in \text{St}_1(\text{AB})$. For this purpose, measures of entanglement are introduced. Every reasonable entanglement measure must obey two rules:

1. separable states have 0 entanglement
2. LOCC cannot increase entanglement

18.2 What does LOCC mean, precisely?

An instrument can be achieved by a LOCC (bipartite) protocol if it corresponds to a coarse graining of an instrument corresponding to:

- An instrument $\{\mathcal{A}_{i_1}\}_{i_1 \in I_1}$ of Alice, who keeps the output and sends the outcome i_1 (classical information) to Bob;
- An instrument $\{\mathcal{B}_{j_1}^{(i_1)}\}_{j_1 \in J_1}$ of Bob that is chosen depending on the information i_1 . Bob keeps the output and sends the outcome j_1 to Alice;
- An instrument $\{\mathcal{A}_{i_2}^{(i_1, j_1)}\}_{i_2 \in I_2}$ of Alice chosen depending on i_1, j_1 . Alice keeps the output and sends i_2 to Bob; ...

The protocol can last n rounds, with $n \in \mathbb{N}$. Notice that n can be arbitrarily large.

Definition 18.3. An instrument $(\mathcal{M}_i)_{i \in I}$ is *one-way* local with respect to Alice [Bob] if

$$\mathcal{M}_i = \mathcal{A}_i \otimes \mathcal{B}^{(i)} \quad [\mathcal{M}_i = \mathcal{A}^{(i)} \otimes \mathcal{B}_i]$$

where $(\mathcal{A}_i)_{i \in I}$ is an instrument $[(\mathcal{B}_i)_{i \in I}]$ while $\forall i \in I, \mathcal{B}^{(i)}$ is a channel $[\mathcal{A}^{(i)}]$

The interpretation in this case is easy: Alice [Bob] does the measurement described by the instrument $(\mathcal{A}_i)_{i \in I}$ $[(\mathcal{B}_i)_{i \in I}]$, sends i to Bob [Alice], who performs the channel $\mathcal{B}^{(i)}$ $[\mathcal{A}^{(i)}]$.

Definition 18.4. (LOCC-linked instrument). Let $(\mathcal{M}_i)_{i \in I}$ be an instrument $\text{AB} \rightarrow \text{AB}$. We say that $(\mathcal{M}'_j)_{j \in J}$ is LOCC-linked to $(\mathcal{M}_i)_{i \in I}$ if there exists a collection of one-way local instruments $(\mathcal{F}_k^{(i)})_{k \in K}$ such that $(\mathcal{M}'_j)_{j \in J}$ is a coarse-graining of $(\mathcal{F}_k^{(i)} \circ \mathcal{M}_i)_{i \in I, k \in K}$.

We can now give the following definitions:

- $(\mathcal{M}_i)_{i \in I}$ is in LOCC_1 if $(\mathcal{M}_i)_{i \in I}$ is one-way local followed by coarse-graining;
- $(\mathcal{M}_i)_{i \in I}$ is in LOCC_n ($n \geq 2$) if it is LOCC-linked to $(\mathcal{F}_j^{(i)})_{j \in J} \in \text{LOCC}_{n-1}$ (it is implicitly assumed that if LOCC_1 is one-way local with respect to Alice [Bob], then the link with LOCC_2 is via a one-way local instrument with respect to Bob [Alice], etc.);

- $(\mathcal{M}_i)_{i \in I}$ is in $\text{LOCC}_{\mathbb{N}}$ if it belongs to LOCC_n for some $n \in \mathbb{N}$.

We then introduce a topology on instruments and close $\text{LOCC}_{\mathbb{N}}$ with respect to such topology.

Definition 18.5 (Diamond norm). We define the *diamond norm* of a linear map $\mathcal{M} : A \rightarrow A$ as follows

$$\|\mathcal{M}\|_{\diamond} := \max_{B; \rho \in \text{St}_1(AB)} \|(\mathcal{M} \otimes \mathcal{I}_B)(\rho)\|_1$$

Definition 18.6 (Diamond norm distance). We define the diamond norm distance of two instruments $(\mathcal{M}_i)_{i \in I}$ and $(\mathcal{F}_i)_{i \in I}$ as

$$D_{\diamond}(\underline{\mathcal{M}}, \underline{\mathcal{F}}) := \max_{B; \rho \in \text{St}_1(AB)} \sum_{i \in I} \|(\mathcal{M}_i \otimes \mathcal{I}_B - \mathcal{F}_i \otimes \mathcal{I}_B)(\rho)\|_1$$

We can then first close $\text{LOCC}_{\mathbb{N}}$ including protocols with infinitely many rounds, as follows:

- $(\mathcal{M}_i)_{i \in I} \in \text{LOCC}$ if there exists a sequence $(\mathcal{F}_i^{(n)})_{i \in I_n} \subseteq \text{LOCC}_{\mathbb{N}}$ such that
 1. $(\mathcal{F}_i^{(n)})_{i \in I_n}$ is LOCC-linked to $(\mathcal{F}_j^{(n-1)})_{j \in I_{n-1}}$
 2. $\forall n \in \mathbb{N}$ there exists a coarse-graining $(\mathcal{M}_i^{(n)})_{i \in I}$ of $(\mathcal{F}_i^{(n)})_{i \in I_n}$ such that $\forall \varepsilon > 0 \ \exists n_0 : \forall n \geq n_0 D_{\diamond}(\underline{\mathcal{M}}^{(n)}, \underline{\mathcal{M}}) < \varepsilon$.

Finally, we can topologically close LOCC as follows:

- $(\mathcal{M}_i)_{i \in I} \in \overline{\text{LOCC}}$ if there exists a sequence $(\mathcal{M}_i^{(n)})_{i \in I}$ in LOCC such that $\underline{\mathcal{M}}^{(n)}$ converges to $\underline{\mathcal{M}}$.

It holds that

$$\text{LOCC}_1 \subset \text{LOCC}_r \subset \text{LOCC}_{r+1} \subset \text{LOCC}_{\mathbb{N}} \subset \text{LOCC} \subset \overline{\text{LOCC}} \subset \text{SEP},$$

where **SEP** is the set of all separable operations, that we define below.

Definition 18.7. An instrument $\underline{\mathcal{M}}$ belongs to **SEP** if $\underline{\mathcal{M}}$ is a coarse graining of $\underline{\mathcal{F}}$ with $\mathcal{F}_i = \mathcal{A}_i \otimes \mathcal{B}_i$ for a set of quantum operations $\mathcal{A}_i : A \rightarrow A$ and $\mathcal{B}_i : B \rightarrow B$.

While in the case of states $\text{LOCC}_1 \equiv \text{SEP}$, for channels and instruments the hierarchy is non-trivial. More precisely, every inclusion is proved to be strict.

An interesting example of a **SEP** instrument that is not $\overline{\text{LOCC}}$ is provided by the famous result “Quantum nonlocality without entanglement” [1]. The result states that in

$\mathbb{C}^3 \otimes \mathbb{C}^3$, the orthonormal basis:

$$\begin{aligned} |\psi_{11}\rangle &= |1\rangle \otimes |1\rangle \\ |\psi_{01}^\pm\rangle &= |0\rangle \otimes \frac{1}{\sqrt{2}}(|0\rangle \pm |1\rangle) \\ |\phi_{01}^\pm\rangle &= \frac{1}{\sqrt{2}}(|0\rangle \pm |1\rangle) \otimes |2\rangle \\ |\psi_{12}^\pm\rangle &= |2\rangle \otimes \frac{1}{\sqrt{2}}(|1\rangle \pm |2\rangle) \\ |\phi_{12}^\pm\rangle &= \frac{1}{\sqrt{2}}(|1\rangle \pm |2\rangle) \otimes |0\rangle \end{aligned}$$

cannot be discriminated by $\overline{\text{LOCC}}$ protocols [1].

We omit the proof. However, notice that the result implies that the instrument $\{\mathcal{P}_i\}_{i \in I}$ with $\mathcal{P}_i(\tau) = P_i \tau P_i$, where $I = \{\phi_{01}^\pm, \psi_{01}^\pm, \phi_{12}^\pm, \psi_{12}^\pm, \psi_{11}\}$ and $P_i = |i\rangle \langle i|$, is not in $\overline{\text{LOCC}}$, while it is in SEP .

For example:

$$P_{\phi_{01}^\pm} = P_{01}^\pm \otimes P_2, \text{ where } P_{ij}^\pm := |\eta_{ij}^\pm\rangle \langle \eta_{ij}^\pm|, |\eta_{ij}^\pm\rangle = \frac{1}{\sqrt{2}}(|i\rangle \pm |j\rangle) \text{ and } P_k = |k\rangle \langle k|.$$

If \mathcal{P} were in $\overline{\text{LOCC}}$, this would indeed contradict the mentioned result of Ref. [1].

18.2.1 Entanglement and LOCC

Since $\overline{\text{LOCC}} \subset \text{SEP}$, it is clear that any $\overline{\text{LOCC}}$ protocol will map separable states to separable states. Indeed, let $\rho = \sum_{i \in I} p_i (\sigma_i \otimes \tau_i)$ and $\mathcal{C} = \sum_{j \in J} \mathcal{A}_j \otimes \mathcal{B}_j$. Then

$$\mathcal{C}(\rho) = \sum_{i \in I, j \in J} p_i [\mathcal{A}_j(\sigma_i) \otimes \mathcal{B}_j(\tau_i)] = \sum_{i,j} \tilde{p}_{ij} (\tilde{\sigma}_{ij} \otimes \tilde{\tau}_{ij}),$$

where $\tilde{\sigma}_{ij} = \frac{\mathcal{A}_j(\sigma_i)}{\text{Tr}[\mathcal{A}_j(\sigma_i)]}$, $\tilde{\tau}_{ij} = \frac{\mathcal{B}_j(\tau_i)}{\text{Tr}[\mathcal{B}_j(\tau_i)]}$ and $\tilde{p}_{ij} = p_i \text{Tr}[\mathcal{A}_j(\sigma_i)] \text{Tr}[\mathcal{B}_j(\tau_i)]$.

Moreover, we expect that $\overline{\text{LOCC}}$ protocols cannot increase the amount of entanglement of a given state. However, the last statement is not provable. Actually, the situation is exactly the opposite. We *define* measures of entanglement as to be non increasing under $\overline{\text{LOCC}}$ maps. This is the paradigm of *resource theories*.

18.2.2 Resource theories

In any resource theory one starts defining *free operations*. Then set of free states must be invariant under free operations.

Normally, free operations form an OPT (operational probabilistic theory), as they have the following properties:

1. the sequence $\mathcal{A} \circ \mathcal{B}$ of any two free operations is free;
2. the parallel compositions $\mathcal{A} \otimes \mathcal{B}$ of free operations is free;

3. the identity map \mathcal{I}_A is considered as free for all systems A

One can define a relation between states as follows:

Definition 18.8. We say that $\rho \succ \sigma$ if $\exists \mathcal{E}$ free such that $\mathcal{E}(\rho) = \sigma$.

Since \mathcal{I}_A is free $\forall A$, clearly, the binary relation denoted as “ \succ ” is *reflexive*. Moreover, since $\forall \mathcal{A}, \mathcal{B}$ free, $\mathcal{B} \circ \mathcal{A}$ is free, \succ is *transitive*. These two properties classify \succ as a pseudo-order (or pre-order) relation¹. One can then introduce the relation

Definition 18.9. $\rho \sim \sigma$ if $\rho \succ \sigma$ and $\sigma \succ \rho$.

$\rho \sim \sigma$ is an equivalence relation. One can then order equivalence classes ρ by a (partial) order: $[\rho] \succeq [\sigma]$ if $\rho \succeq \sigma$. One can easily check that \succeq is well defined: the condition $[\rho] \succeq [\sigma]$ if $\rho \succeq \sigma$ holds independently of the elements $\rho \in [\rho]$ and $\sigma \in [\sigma]$. It is also easy to check that \succeq is a *partial ordering* relation². Thus, resource theories allow one to introduce an order relation on the set of states according to the “amount” of resource that the states contain. In the resource theory of entanglement, free operations are LOCC operations. In this way, entanglement is essentially defined as that resource that consists in those quantum correlations that cannot be produced by just exchanging classical information.

Any measure of the resource of interest must be monotonic under the above defined ordering.

18.2.3 Maximally entangled states

The ordering of states under the relation \succeq naturally poses a question: are there states that represent “maximal” resources? Let us start considering two-qubit states. In this case we can suspect that maximally entangled states are maximal resources. Maximally entangled states are all pure states of the form

$$\begin{aligned}\rho &= \frac{1}{2}|U\rangle\langle U| \\ U^\dagger U &= UU^\dagger = I_A\end{aligned}$$

Notice that given $\frac{1}{\sqrt{2}}|U\rangle\rangle, \frac{1}{\sqrt{2}}|W\rangle\rangle$, one can always find V such that $V^\dagger V = VV^\dagger = I_A$, and

$$(V \otimes I_{A'}) \frac{1}{\sqrt{2}}|U\rangle\rangle = \frac{1}{\sqrt{2}}|W\rangle\rangle$$

Indeed, the choice of $V := WU^\dagger$ will do the job. This implies that there exist LOCC maps that convert any two-qubit maximally entangled state into any other. In the following we will consider the singlet state $\frac{1}{\sqrt{2}}|i\sigma_y\rangle\rangle$ as our reference maximally entangled state.

¹We remind that an ordering relation has the properties of being reflexive, antisymmetric and transitive.

²Reflexive, antisymmetric, transitive; however, differently from total ordering, given two states ρ, ρ' it can happen that neither $[\rho]_\sim \succeq [\rho']_\sim$, nor $[\rho']_\sim \succeq [\rho]_\sim$.

However, we now need to show that every two qubit state can be obtained via LOCC from a maximally entangled one. First, let $\rho = |\psi\rangle\langle\psi| \otimes |\phi\rangle\langle\phi|$. The LOCC map $\mathcal{Q}_\psi \circ \mathcal{Q}_\phi$, with

$$\mathcal{Q}_\eta(\sigma) := \text{Tr}[\sigma] |\eta\rangle\langle\eta|$$

will clearly turn $\frac{1}{\sqrt{2}}|i\sigma_y\rangle\rangle$ into $|\psi\rangle\langle\psi| \otimes |\phi\rangle\langle\phi|$. Let now $\rho = |F\rangle\rangle\langle\langle F|$ with $|F\rangle\rangle$ pure, entangled but not maximally entangled, i.e.

$$\begin{aligned}\text{rank}(FF^\dagger) &= 2 \\ FF^\dagger &\neq \frac{I_A}{2}\end{aligned}$$

In general one can reduce $|F\rangle\rangle$ in the Schmidt form

$$|F\rangle\rangle = \sum_{i=0}^1 f_i |\eta_i\rangle |\theta_i\rangle = f_0 |\eta_0\rangle |\theta_0\rangle + f_1 |\eta_1\rangle |\theta_1\rangle.$$

Via local unitary maps (that are LOCC) one can then turn $|F\rangle\rangle$ into $|\tilde{F}\rangle\rangle$, where

$$|\tilde{F}\rangle\rangle = f_0 |0\rangle |1\rangle + f_1 |1\rangle |0\rangle$$

and vice-versa. We then just need to prove that $\rho = \frac{1}{2}|U\rangle\rangle\langle\langle U| \succ |F\rangle\rangle\langle\langle F|$. We remind that $\frac{1}{\sqrt{2}}|i\sigma_y\rangle\rangle = \frac{1}{\sqrt{2}}(|0\rangle|1\rangle - |1\rangle|0\rangle)$. Let us then consider the following protocol: Alice prepares an ancillary system E in the pure state $|0\rangle\langle 0|$. The state $\rho = \frac{1}{2}|i\sigma_y\rangle\rangle\langle\langle i\sigma_y|$ is then mapped to $|0\rangle\langle 0| \otimes \frac{1}{2}|i\sigma_y\rangle\rangle\langle\langle i\sigma_y| = |\phi\rangle\langle\phi|$, where

$$|\phi\rangle = \frac{1}{\sqrt{2}}(|0\rangle_E|0\rangle_A|1\rangle_B - |0\rangle_E|1\rangle_A|0\rangle_B)$$

Alice controls the system EA. She now performs the unitary map \mathcal{N}_U defined by the following identities

$$\begin{aligned}U|00\rangle &= f_0|00\rangle + f_1|11\rangle, \\ U|01\rangle &= -f_1|01\rangle - f_0|10\rangle, \\ U|10\rangle &= f_0^*|01\rangle - f_1^*|10\rangle, \\ U|11\rangle &= -f_1^*|00\rangle + f_0^*|11\rangle.\end{aligned}$$

The matrix for U in the basis $\{|ij\rangle\}_{i,j=0}^1$ is

$$U = \begin{pmatrix} f_0 & 0 & 0 & -f_1^* \\ 0 & -f_1 & f_0^* & 0 \\ 0 & -f_0 & -f_1^* & 0 \\ f_1 & 0 & 0 & f_0^* \end{pmatrix}.$$

One has $\mathcal{N}_U \otimes \mathcal{I}_B (|\phi\rangle\langle\phi|) = |\phi'\rangle\langle\phi'|$, with

$$\begin{aligned}|\phi'\rangle &= \frac{1}{\sqrt{2}}(f_0|0\rangle_E|0\rangle_A|1\rangle_B + f_1|1\rangle_E|1\rangle_A|1\rangle_B + f_1|0\rangle_E|1\rangle_A|0\rangle_B + f_0|1\rangle_E|0\rangle_A|0\rangle_B) = \\ &= \frac{1}{\sqrt{2}}[|0\rangle_E(f_0|0\rangle_A|1\rangle_B + f_1|1\rangle_A|0\rangle_B) + |1\rangle_E(f_0|0\rangle_A|0\rangle_B + f_1|1\rangle_A|1\rangle_B)].\end{aligned}$$

Now Alice can measure the ancillary system E with the POVM $\{|0\rangle\langle 0|, |1\rangle\langle 1|\}$, obtaining for the system AB the following conditional states:

$$\rho'_x := \frac{\text{Tr}_E[(P_x \otimes I_{AB}) |\phi'\rangle\langle\phi'|]}{\text{Tr}[(P_x \otimes I_{AB}) |\phi'\rangle\langle\phi'|]}, \quad x = 0, 1.$$

Notice that

$$(P_0 \otimes I_{AB}) |\phi'\rangle = |0\rangle \otimes \frac{1}{\sqrt{2}} (f_0|0\rangle|1\rangle + f_1|1\rangle|0\rangle),$$

$$(P_1 \otimes I_{AB}) |\phi'\rangle = |1\rangle \otimes \frac{1}{\sqrt{2}} (f_0|0\rangle|0\rangle + f_1|1\rangle|1\rangle).$$

Thus $\rho_0 = |\tilde{F}\rangle\langle\tilde{F}|$, $\rho_1 = |G\rangle\langle G|$, where $|G\rangle\langle G| := \frac{1}{\sqrt{2}} (f_0|0\rangle|0\rangle + f_1|1\rangle|1\rangle)$. If Alice sends the outcome x to Bob, then Bob does:

- Nothing if $x = 0 \implies \rho''_{AB} = |\tilde{F}\rangle\langle\tilde{F}|$
 - Applies σ_x if $x = 1 \implies \rho''_{AB} = \mathcal{I}_A \otimes \mathcal{N}_{\sigma_x}(|G\rangle\langle G|) = |\tilde{F}\rangle\langle\tilde{F}|$.
- Indeed, $I_A \otimes \sigma_x |G\rangle\langle G| = \frac{1}{\sqrt{2}} (f_0|0\rangle|1\rangle + f_1|1\rangle|0\rangle) = |\tilde{F}\rangle\langle\tilde{F}|$.

Every maximally entangled two-qubit state can be mapped via LOCC to every pure two-qubit state. In order to produce mixed states

$$\rho = \sum_i p_i |\phi_i\rangle\langle\phi_i|,$$

it is then sufficient that Alice produces a classical variable x_i with $p(x_i) = p_i$, then send x_i to Bob, and finally they start the LOCC protocol to turn $\frac{1}{2}|i\sigma_y\rangle\langle i\sigma_y|$ into $|\phi_i\rangle\langle\phi_i|$.

However, in the case of d-dimensional systems, there are *incomparable pairs* of pure maximal states ρ, σ . This means that, despite the fact that it both holds that $\tau \succ \rho$ only if $\tau \sim \rho$, and $\tau \succ \sigma$ only if $\tau \sim \sigma$, neither $\rho \succ \sigma$ nor $\sigma \succ \rho$. It is then impossible to establish a notion of *maximally entangled state* in the same way as we did for the singlet in the case of qubits.

Moreover, it can happen that $\rho \not\succ \sigma$, but $\rho \succ \tau_\varepsilon$ for a family of states $\{\tau_\varepsilon\}$ with $\|\sigma - \tau_\varepsilon\|_1 < \varepsilon$. The way in which equivalence classes are ordered in our resource theory of entanglement seems then to be too restrictive: if it is possible to produce an arbitrarily good approximation of σ via LOCC from ρ , it seems that our order relation should be extended so that $\rho \succ \sigma$.

To overcome these issues, one can introduce a scenario where, instead of changing the order relation between equivalence classes, one uses it in order to introduce a way of assessing an “exchange rate” between states, that allows one to compare also states of different systems. For example, let $\rho \in \text{St}_1(A)$ and $\sigma \in \text{St}_1(B)$. First of all, $\forall \varepsilon > 0$ we seek an LOCC protocol and a state τ_ε of suitably many copies of system B (say m), such that $\exists n: \rho^{\otimes n} \succ \tau_\varepsilon$ with $\|\sigma^{\otimes m} - \tau_\varepsilon\|_1 < \varepsilon$. The largest possible ratio $\frac{m}{n}$ for which the above LOCC transformation can be achieved is an “exchange rate” between ρ and σ , i.e. the relative entanglement content of σ with respect to ρ .

18.3 Entanglement cost and distillable entanglement

With the above observations in mind, we can now keep the singlet as our *standard unit* of entanglement, and define an entanglement measure for $\rho \in \text{St}_1(\text{AB})$ that accounts for the entanglement needed to produce ρ via LOCC. This measure is the *entanglement cost* of ρ . In the following, the symbol Q will denote a qubit system.

Definition 18.10 (Feasible production ratio). Let $\rho \in \text{St}_1(\text{AB})$. We will say that the triple (p, q, ε) , is *feasible* if there exists a protocol $\mathcal{L} \in \overline{\text{LOCC}}$, with $\mathcal{L} : Q^{\otimes 2p} \rightarrow (\text{AB})^{\otimes q}$, such that

$$\frac{1}{2} \|\mathcal{L}(\Sigma^{\otimes p}) - \rho^{\otimes q}\|_1 < \varepsilon,$$

where $\Sigma := 1/2|i\sigma_y\rangle\langle i\sigma_y|$ is the singlet state. The number p/q , in particular, is a *feasible production ratio* to tolerance ε . The set of feasible triples (p, q, ε) is denoted as F_C .

Let F_ε be the set of values of (m, n) such that $(m, n, \varepsilon) \in F_C$ is feasible. The number m/n evaluates the number of singlets per copy of ρ needed to produce n copies of ρ via LOCC starting from m singlets. Now, the value $r(\varepsilon) = \inf_{(m, n) \in F_\varepsilon} m/n$ denotes the ratio that can be achieved if the degree of approximation ε is tolerated. Notice that $r(\varepsilon)$ is a monotonically increasing function, since the ratios $r = m/n$ achievable for approximation ε_1 are straightforwardly achievable for every $\varepsilon_2 > \varepsilon_1$. The achievable ratio such that the approximation ε can be made arbitrarily small is then the limit $r(\varepsilon)$ for $\varepsilon \rightarrow 0$. Indeed, considering larger values of $r = m/n$ would mean overestimating the amount of singlets required to produce each copy of ρ . The definition of entanglement cost is then the following.

Definition 18.11 (Entanglement cost). The entanglement cost of $\rho \in \text{St}_1(\text{AB})$ is

$$E_C(\rho) := \inf\{r \in \mathbb{R} \mid \forall \varepsilon > 0, \delta > 0, \exists (m, n, \varepsilon) \in F_C, |r - m/n| < \delta\}.$$

$E_C(\rho)$ is then the optimal ratio $r = \frac{m}{n}$ of the number m of starting copies of a singlet state Σ and the number n of copies of ρ that can be obtained via LOCC, asymptotically, and with arbitrarily small approximation.

A second measure of entanglement accounts, on the contrary, for amount of entanglement that can be extracted from ρ via LOCC. This quantity is called *distillable entanglement* of ρ , and is defined as follows.

Definition 18.12 (Feasible extraction ratio). Let $\rho \in \text{St}_1(\text{AB})$. We will say that the triple (q, p, ε) , is *feasible* if there exists a protocol $\mathcal{L} \in \overline{\text{LOCC}}$, with $\mathcal{L} : (\text{AB})^{\otimes q} \rightarrow Q^{\otimes 2p}$, such that

$$\frac{1}{2} \|\mathcal{L}(\rho^{\otimes q}) - \Sigma^{\otimes p}\|_1 \leq \varepsilon.$$

The number p/q , in particular, is a *feasible extraction ratio* to tolerance ε . The set of feasible triples (p, q, ε) is denoted as F_D .

Let now \mathcal{G}_ε be the set of values of (m, n) such that $(m, n, \varepsilon) \in \mathcal{F}_D$ is feasible. The number n/m evaluates the number of singlets that can be produced out of each copy of ρ , starting from m copies of ρ , via LOCC and producing n singlets. Similarly to the case of the entanglement cost, the value $r(\varepsilon) = \sup_{(m,n) \in \mathcal{G}_\varepsilon} n/m$ denotes the ratio that can be achieved if the degree of approximation ε is tolerated. Notice that $r(\varepsilon)$ is a monotonically increasing function, since the ratios $r = n/m$ achievable for approximation ε_1 are straightforwardly achievable for every $\varepsilon_2 > \varepsilon_1$. The achievable ratio such that the approximation ε can be made arbitrarily small is then the limit $r(\varepsilon)$ for $\varepsilon \rightarrow 0$. Indeed, considering smaller values of $r = n/m$ would mean underestimating the number of singlets that can be produced from each copy of ρ . The definition of distillable entanglement is then the following.

Definition 18.13. (Distillable entanglement) The distillable entanglement $E_D(\rho)$ of a state $\rho \in \text{St}_1(AB)$ is

$$E_D(\rho) := \sup\{r \in \mathbb{R} \mid \forall \varepsilon > 0, \delta > 0, \exists (m, n, \varepsilon) \in \mathcal{F}_D, |r - n/m| < \delta\}.$$

$E_C(\rho)$ and $E_D(\rho)$ then measure the amount of entanglement measured in *singlets*, that is needed to produce one copy of ρ or that can be produced out of one copy of ρ , respectively, in the asymptotic scenario and with arbitrary accuracy. As one might expect, $E_C(\rho) \geq E_D(\rho)$. One might think that actually $E_C(\rho) = E_D(\rho)$, always. This is not the case, however. While $E_C(\rho) > 0$ for every entangled state, it can happen that $E_C(\rho) = 0$ even if ρ is entangled. In this case we speak of “bound entanglement”, that can be interpreted as a number of singlets that must be spent to produce ρ but cannot be extracted anymore out of ρ .

18.4 Criteria for entanglement

We thus gave a thorough definition of entanglement through two measures with a clear operational interpretation. However, calculating $E_D(\rho)$ or $E_C(\rho)$ is hard. It is even hard just to establish whether a given ρ_{AB} is separable or not. For this reason it is useful to introduce criteria for entanglement that are easy to calculate, even though they are just sufficient, and not quantitative.

18.4.1 The PPT criterion - Positive partial transpose

Also known as Peres-Horodecki criterion, this criterion is a sufficient condition for ρ_{AB} to be entangled. It states that:

Lemma 18.14. *If $(\mathcal{I}_A \otimes \Theta_B)(\rho_{AB}) \not\succeq 0$, then ρ_{AB} is entangled. (Alternatively if ρ_{AB} is separable, then $(\mathcal{I}_A \otimes \Theta_B)(\rho_{AB}) \succeq 0$.)*

Proof. The proof is very simple. Let

$$\rho_{AB} = \sum_i p_i \sigma_i \otimes \tau_i \quad \text{with } p_i \geq 0, \sigma_i \geq 0, \tau_i \geq 0.$$

Then

$$(\mathcal{J}_A \otimes \Theta_B)(\rho_{AB}) = \sum_i p_i \sigma_i \otimes \tau_i^T \geq 0.$$

□

Remark 21. In the case where $\mathcal{H}_A \otimes \mathcal{H}_B = \mathbb{C}^2 \otimes \mathbb{C}^2$, $\mathcal{H}_A \otimes \mathcal{H}_B = \mathbb{C}^3 \otimes \mathbb{C}^2$, or $\mathcal{H}_A \otimes \mathcal{H}_B = \mathbb{C}^2 \otimes \mathbb{C}^3$, the PPT criterion becomes necessary and sufficient.

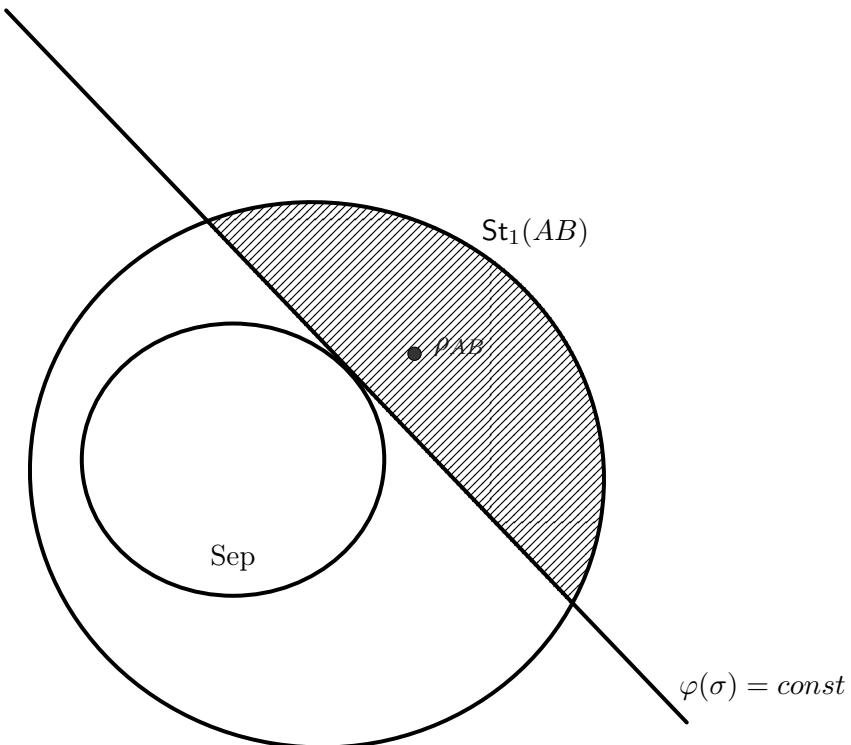
Remark 22. For every positive non-CP map \mathcal{L} , a similar criterion can be given: if ρ_{AB} is separable, then $(\mathcal{J}_A \otimes \mathcal{L}_B)(\rho_{AB}) \geq 0$.

18.4.2 Entanglement witnesses

Since separable states form a convex set, the separating hyperplane theorem ensures that given ρ_{AB} entangled, there is a functional $\phi : \text{St}_1(AB) \rightarrow \mathbb{R}$ such that

$$\begin{aligned}\phi(\sigma) &\geq 0 \quad \forall \sigma \text{ separable;} \\ \phi(\rho_{AB}) &< 0.\end{aligned}$$

A famous example is given by Bell's inequalities (or CHSH inequality).



Remark 23. The PPT criterion is sufficient for *bound entanglement*. Precisely, if ρ_{AB} is entangled but $(\mathcal{J}_A \otimes \Theta_B)(\rho_{AB}) \geq 0$, then $E_D(\rho) = 0$ while $(E_C(\rho) > 0)$.

18.4.3 Bound entanglement

Indeed, let $E_D(\rho) > 0$. Then $\exists n, \exists \mathcal{E} \in \overline{\text{LOCC}}$, and r s.t. $|r - E_D(\rho)| < \delta$ and

$$\frac{1}{2} \|\mathcal{E}(\rho^{\otimes n}) - \Sigma^{\otimes rn}\|_1 < \varepsilon. \quad (18.1)$$

Now, since $\frac{1}{2} \|\tau - \eta\|_1 = \max_{0 \leq P \leq I} \text{Tr}[(\tau - \eta)P]$,

$$\frac{1}{2} \|\mathcal{C}(\tau) - \mathcal{C}(\eta)\|_1 = \max_{0 \leq P \leq I} \text{Tr}[(\tau - \eta)\mathcal{C}^\dagger(P)] \leq \max_{0 \leq Q \leq I} \text{Tr}[(\tau - \eta)Q] \leq \frac{1}{2} \|\tau - \eta\|_1.$$

Thus, discarding all the rn systems but one in Eq. (18.1), one has

$$\frac{1}{2} \|\text{Tr}_{rn-1} [\mathcal{E}(\rho^{\otimes n}) - \Sigma^{\otimes rn}]\|_1 = \frac{1}{2} \|\mathcal{E}'(\rho^{\otimes n}) - \Sigma\|_1 < \varepsilon,$$

where we considered the channel $\mathcal{C} := \text{Tr}_{rn-1} : (\text{AB})^{\otimes rn} \rightarrow \text{AB}$. Thus, $\mathcal{E}'(\rho^{\otimes n})$ must be entangled. Indeed $F(\Sigma, \tau) = \frac{1}{2} \langle i\sigma_y | \tau | i\sigma_y \rangle$ and for separable $\tau = \sum_i p_i \theta_i \otimes \gamma_i$

$$F(\Sigma, \tau) = \frac{1}{2} \sum_i p_i \text{Tr}[\theta_i \sigma_y \gamma_i^T \sigma_y] \leq \frac{1}{2} \sum_i p_i \text{Tr}[\theta_i] = \frac{1}{2}.$$

By the Fuchs-Van der Graaf inequalities one has

$$\frac{1}{2} \|\Sigma - \tau\|_1 \geq 1 - F(\Sigma, \tau) \geq \frac{1}{2}$$

for separable τ . Thus, for sufficiently small ε ,

$$\frac{1}{2} \|\mathcal{E}'(\rho^{\otimes n}) - \Sigma\|_1 < \varepsilon \Rightarrow \mathcal{E}'(\rho^{\otimes n}) \text{ entangled.}$$

The map \mathcal{E}' is in $\overline{\text{LOCC}}$, thus it is in SEP :

$$\mathcal{E}'(\rho^{\otimes n}) = \sum_i (A_i \otimes B_i) \rho^{\otimes n} (A_i^\dagger \otimes B_i^\dagger).$$

There must then exist i_0 such that $(A_{i_0} \otimes B_{i_0}) \rho^{\otimes n} (A_{i_0}^\dagger \otimes B_{i_0}^\dagger)$ is entangled. Since $A_{i_0} : \mathcal{H}_A^{\otimes n} \rightarrow \mathbb{C}^2$, and similarly for $B_{i_0} : \mathcal{H}_B^{\otimes n} \rightarrow \mathbb{C}^2$, by the singular values decomposition we have

$$\begin{aligned} A_{i_0} &= |0\rangle \langle \psi_0| + |1\rangle \langle \psi_1|, \\ B_{i_0} &= |0\rangle \langle \varphi_0| + |1\rangle \langle \varphi_1|. \end{aligned}$$

Defining P_A and P_B as the projections on $\text{Span}(|\psi_0\rangle, |\psi_1\rangle)$ and $\text{Span}(|\varphi_0\rangle, |\varphi_1\rangle)$ respectively, we have

$$(A_{i_0} \otimes B_{i_0}) \rho^{\otimes n} (A_{i_0}^\dagger \otimes B_{i_0}^\dagger) = (A_{i_0} \otimes B_{i_0})(P_A \otimes P_B) \rho^{\otimes n} (P_A \otimes P_B) (A_{i_0}^\dagger \otimes B_{i_0}^\dagger),$$

and since LOCC maps cannot create entanglement, $\rho' = (P_A \otimes P_B)\rho^{\otimes n}(P_A \otimes P_B)$ must be entangled. Finally since ρ' has support in a $\mathbb{C}^2 \otimes \mathbb{C}^2$ space, ρ' has a non-positive partial transpose.

$$\rho'^{\tau_B} = (I_A \otimes \Theta_B)[(P_A \otimes P_B)\rho^{\otimes n}(P_A \otimes P_B)] = (P_A \otimes P_B^T)(\rho^{\tau_B})^{\otimes n}(P_A \otimes P_B^T),$$

where we used the shorthand notation $X^{\tau_B} := \mathcal{J}_A \otimes \Theta_B(X)$. Since ρ'^{τ_B} has a negative eigenvalue, it must be:

$$\langle\langle \Phi | (P_A \otimes P_B^T)(\rho^{\tau_B})^{\otimes n}(P_A \otimes P_B^T) | \Phi \rangle\rangle < 0,$$

for some $|\Phi\rangle\rangle$ —e.g. by taking $|\Phi\rangle\rangle$ as the corresponding eigenvector—i.e.

$$\langle\langle \Psi | (\rho^{\tau_B})^{\otimes n} | \Psi \rangle\rangle < 0, \quad |\Psi\rangle\rangle := P_A \otimes P_B |\Phi\rangle\rangle.$$

Finally, if $\rho^{\tau_B} \geq 0$ also $(\rho^{\tau_B})^{\otimes n} \geq 0$, thus it must be $\rho^{\tau_B} \not\geq 0$. In other words we proved that if $E_D(\rho) > 0$ it must be $\rho^{\tau_B} \not\geq 0$, i.e. if $\rho^{\tau_B} \geq 0$, then $E_D(\rho) = 0$.

As a final remark on this point, we observe that it can be proved that entangled states with positive partial transpose exist, thus there exist states exhibiting bound entanglement.

18.4.4 Entanglement of formation

For pure states $\rho_{AB} = |\Psi\rangle\rangle\langle\langle \Psi|$, a measure of entanglement is the *entropy of entanglement*:

$$S(\rho_A) = S(\rho_B) = S(\Psi\Psi^\dagger) =: E(\rho)$$

Every entanglement measure must reduce to $E(\rho)$ for pure states ρ .

Definition 18.15 (Entanglement of formation).

$$E_F(\rho) := \inf \left\{ \sum_i p_i E(|\psi_i\rangle\rangle\langle\langle \psi_i|) \middle| \rho = \sum_i p_i |\psi_i\rangle\rangle\langle\langle \psi_i| \right\}$$

The *regularized* entanglement of formation is:

$$E_F^\infty(\rho) := \lim_{n \rightarrow \infty} \frac{E_F(\rho^{\otimes n})}{n}$$

It holds that $E_F^\infty(\rho) = E_C(\rho)$.

Acknowledgements

These notes are largely based on the lecture notes of G. M. D'Ariano who was formerly the teacher of the present course. Some parts are inspired by J. Watrous's lecture notes. I wish to express my gratitude to Maddalena Ghio, Francesco Tacchino, Andrea Olivo and Francesco Zagaria for their careful reading, and for pointing out various errors in earlier drafts of these lecture notes. A special thank to Gianmarco Ricciardi, Edoardo Centofanti, Daniele Gilio, Federico Quetti, Lia Rapella, Lorenzo Rossi and Alessandro Summer for turning my (awfully) handwritten notes for the last lecture into a L^AT_EX file. Finally, I acknowledge the thoughtful and detailed review of the above notes made by Leonardo Vaglini.

Bibliography

- [1] Charles H. Bennett, David P. DiVincenzo, Christopher A. Fuchs, Tal Mor, Eric Rains, Peter W. Shor, John A. Smolin, and William K. Wootters. Quantum nonlocality without entanglement. *Phys. Rev. A*, 59:1070–1091, Feb 1999.
- [2] C. E. Shannon. A mathematical theory of communication. *SIGMOBILE Mob. Comput. Commun. Rev.*, 5(1):3–55, January 2001.
- [3] J. Watrous. *Quantum entropy and source coding*. lecture notes available at <https://cs.uwaterloo.ca/~watrous/CS766/DraftChapters/5.QuantumEntropy.pdf>.