

Marco Tomamichel

EE5139/EE6139:
Information Theory and its Applications

(Semester I, 2021–2022)

Disclaimer: These notes are not yet free of typos or always presented in the most clear way. Any comments that help reduce these deficiencies are very much appreciated. Parts of Chapters 0 and Chapter 7 are based on notes by Vincent Y. F. Tan. Most figures were contributed by Michael X. Cao.

Contents

0	Review of mathematical notation and foundations	2
0.1	Notation	2
0.2	Probability theory	3
0.2.1	Probability space	3
0.2.2	Random variables	5
0.2.3	Expectation and variance	7
0.2.4	Markov chains	8
0.3	Tail bounds	8
0.3.1	Basic bounds	9
0.3.2	Central limit theorem	10
0.4	Vector norms and Cauchy-Schwarz inequality	11
0.5	Convexity and Jensen's inequality	12
0.6	Finite field arithmetic	14
1	Information measures	16
1.1	Surprisal and entropy	16
1.1.1	Surprisal	16
1.1.2	Entropy	17
1.2	Conditional entropy, and mutual information	20
1.2.1	Joint entropy	20
1.2.2	Conditional entropy	20
1.2.3	Mutual information	23
1.3	Relative entropy	24
2	Source coding	28
2.1	Problem setup and definitions	28
2.1.1	Data source	28
2.1.2	Source codes	29
2.2	Variable-length codes	31
2.2.1	Optimal codeword lengths	32
2.2.2	Optimal expected codeword length	32
2.2.3	Shannon code	33
2.2.4	Huffman codes	34

2.3	Fixed-length block codes	37
2.3.1	Setup for block coding	38
2.3.2	Proof of converse and Fano's inequality	40
2.3.3	Proof of achievability and typical sets	41
2.3.4	Strong converse via typical sets	44
3	Cryptography: randomness extraction	47
3.1	Problem setup	47
3.2	Guessing probability and min-entropy	49
3.3	Achievability via two-universal hash functions	50
3.4	Converse via an entropy inequality	53
4	Information theory in statistics: hypothesis testing	56
4.1	Binary hypothesis testing	56
4.2	Symmetric hypothesis testing	57
4.2.1	Total variation distance	58
4.2.2	Chernoff exponent	59
4.3	Asymmetric hypothesis testing and the information spectrum method	60
5	Error correcting codes	64
5.1	Definitions and bounds on codebook size	64
5.2	Linear codes	66
5.3	Reed-Solomon codes	68
5.4	Low density parity check (LDPC) codes	69
5.4.1	Decoding with belief propagation	69
6	Noisy channel coding	70
6.1	Channel mutual information	70
6.2	The channel coding theorem	74
6.2.1	The meta-converse	76
6.2.2	Proof of converse and types	77
6.2.3	Proof of achievability and random codes	80
6.2.4	Maximum probability of error	83
6.3	Source-channel separation theorem	83
6.4	Gaussian channels	85
6.4.1	Differential entropy and mutual information	86
6.4.2	Channel coding theorem for the AWGN channel	88
7	Learning theory: Multiarmed stochastic bandits	90
7.1	Problem setup and objective	90
7.2	A lower-bound on minimax regret	92
7.2.1	Decomposing the regret	92
7.2.2	Constructing worst-case environments	93

7.2.3	Lower-bounding the regret	94
-------	-------------------------------------	----

\emptyset	empty set $\{\}$
$[M]$	the set $\{1, 2, \dots, M\}$
$\mathcal{P}(\mathcal{X})$	the power set of \mathcal{X} , i.e. $\{A : A \subseteq \mathcal{X}\}$
$\mathcal{X} \times \mathcal{Y}$	the set of tuples $\{(x, y) : x \in \mathcal{X}, y \in \mathcal{Y}\}$
\mathcal{X}^n	the set of n -tuples with each element taking values in \mathcal{X} , e.g., $\mathcal{X}^2 = \mathcal{X} \times \mathcal{X}$
$\{0, 1\}^n$	the set of n -bit strings
$\{0, 1\}^*$	the set of bit strings of arbitrary length
$\max \mathcal{X}$	largest $x^* \in \mathcal{X}$, might not always exist
$\sup \mathcal{X}$	smallest $x^* \in \mathbb{R}$ such that $x \leq x^*$ for all $x \in \mathcal{X}$; equals the maximum, $\max \mathcal{X}$, if it exists
$\min \mathcal{X}$	smallest $x^* \in \mathcal{X}$, might not always exist
$\inf \mathcal{X}$	largest $x^* \in \mathbb{R}$ such that $x \geq x^*$ for all $x \in \mathcal{X}$; equals the minimum, $\min \mathcal{X}$, if it exists
$\mathbf{1}\{x = y\}$	indicator function, evaluates to 1 if the condition is true and 0 otherwise, so that, for example, $\mathbf{1}\{x = y\} + \mathbf{1}\{x \neq y\} = 1$
δ_{xy}	shorthand for $\mathbf{1}\{x = y\}$
$P_X(x)$	probability mass function (pmf), $P_X(x) = P[X = x]$
$p_X(x)$	probability density function (pdf), i.e. $P[X \in (1, 2)] = \int_1^2 p_X(x) dx$
$P[X \in \mathcal{A}]$	probability of a random variable X being in some set \mathcal{A} , i.e. $P[X \in \mathcal{A}] = \mathbb{P}(\{\omega : X(\omega) \in \mathcal{A}\}) = \sum_{x \in \mathcal{A}} P_X(x)$
$P[5 \leq X < 6]$	another way of writing $P[X \in [5, 6)]$
\log	logarithm; in these notes we take the logarithm to base 2, i.e. $\log = \log_2$

Table 1: Some basic notation used in this module.

pmf	probability mass function
pdf	probability density function
cdf	cumulative density function
rv	random variable
DMS	discrete memoryless source
DMC	discrete memoryless channel

Table 2: Some abbreviations used in this module.

Chapter 0

Review of mathematical notation and foundations

[Week 1]

Intended learning outcomes:

- You are familiar with common notation used throughout the lecture.
- You are comfortable with the main mathematical concepts needed in this module, namely basic probability theory including random variables, conditional probabilities and Markov chains.
- You can apply basic bounds on tail probabilities, and can prove the weak law of large numbers.
- You can compute vector norms and apply the Cauchy-Schwarz inequality.
- You know what convex and concave functions are and can apply Jensen's inequality.
- You know what finite fields are and how to come up with the multiplication table for simple examples.

0.1 Notation

We will use standard notation and abbreviations that you should be familiar with from other modules. Some of the less frequently encountered mathematical expressions are summarised in Tables 1 and 2 on the previous page.

0.2 Probability theory

We will not directly need the framework of probability theory in its most abstract formulation as presented in the following, but it is good to know that both discrete and continuous random variables can be seen as emanating from a shared mathematical framework.

0.2.1 Probability space

A probability space is represented by a triple $(\Omega, \Sigma, \mathbb{P})$. Here Ω is a set that is called the *sample space*. Moreover, Σ is a σ -algebra, i.e. a collection of subsets of Ω , called events, with the following properties:

- $\Omega \in \Sigma$
- If $A \in \Sigma$, then its *complement*, $A^c = \Omega \setminus A$ is also in Σ , i.e. $A^c \in \Sigma$.
- If $A_1, A_2, \dots, A_n, \dots \in \Sigma$, then $\bigcup_{i=1}^{\infty} A_i \in \Sigma$

Question 0.1. *Show that the above also implies that $\emptyset \in \Sigma$ and $\bigcap_{i=1}^{\infty} A_i \in \Sigma$.*

For example, let $\Omega = [0, 1]$, and we are interested in the probability of subsets of Ω that are intervals of the form $[a, b]$ where $0 \leq a < b \leq 1$, but not individual points in Ω . Then we should also be able to say something about the probability of the union, intersections, complement and so on of such intervals. This is captured by the definition of a σ -algebra. Think of Σ as the properties of Ω that can actually be observed.

Example 0.2. *If your random variable is the location an athlete lands after a long jump then it makes sense to take Ω to be positive real numbers, \mathbb{R}_+ indicating the distance jumped (say, in meters). However, even with arbitrarily good equipment we cannot actually measure a real number, we can only ever say that he landed in some interval, the size of which is given by our measurement precision. Thus, Σ , comprised of the events we can actually observe, is built up by including all (arbitrarily small) intervals in \mathbb{R}_+ and their unions and complements. Or another way of looking at this is that the probability of the jumper landing exactly at 9m is always zero — it is simply the wrong question to ask. But the probability of landing within 1cm or some arbitrarily small interval around 9m might very well be nonzero.*

Question 0.3. *For the advanced reader: Note however that in the above example $\{x\} \in \Sigma$ for any $x \in \mathbb{R}^+$, that is, single points are also elements of the σ -algebra. Can you see why? Use an infinite intersection to construct it.*

Finally, the probability measure \mathbb{P} is a function $\mathbb{P} : \Sigma \rightarrow [0, 1]$ defined on the measurable space (Ω, Σ) , and represents your “belief” about the events in Σ . In order for \mathbb{P} to be called a probability measure, it must satisfy the following two properties:

1. $\mathbb{P}(\Omega) = 1$

2. For A_1, A_2, \dots such that $A_i \cap A_j = \emptyset$ for all $i \neq j$, i.e. for mutually *disjoint* sets, we have

$$\mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mathbb{P}(A_i) \quad (1)$$

Some basic and very useful properties that can be derived from the above definition. The union bound in particular is very often used when analysing problems in information theory.

Proposition 0.4. *Let $(\Omega, \Sigma, \mathbb{P})$ be a probability space. The following holds true:*

1. $\mathbb{P}(A^c) = 1 - \mathbb{P}(A)$
2. If $A \subset B$, then $\mathbb{P}(A) \leq \mathbb{P}(B)$
3. $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B) \leq \mathbb{P}(A) + \mathbb{P}(B)$, which is called the union bound. Clearly, by induction, the union bound works for finitely many sets $A_i, i = 1, \dots, k$, namely

$$\mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right) \leq \sum_{i=1}^{\infty} \mathbb{P}(A_i). \quad (2)$$

Proof. Property 1 follows since $A^c \cap A = \emptyset$, and thus $\mathbb{P}(A) + \mathbb{P}(A^c) = \mathbb{P}(\Omega) = 1$ by (1). For Property 2, note that $B \setminus A = B \cap A^c \in \Sigma$ and since $A \cap (B \setminus A) = \emptyset$ we again argue that $\mathbb{P}(A) + \mathbb{P}(B \setminus A) = \mathbb{P}(B)$, from which the desired inequality follows.

For Property 3 note that $A \cup B$ can be decomposed in three different ways into mutually disjoint sets:

$$A \cup B = A \cup (B \setminus A) = B \cup (A \setminus B) = (A \setminus B) \cup (B \setminus A) \cup (A \cap B). \quad (3)$$

Again using (1) for each of these decompositions we have

$$2\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) + \mathbb{P}(B \setminus A) + \mathbb{P}(A \setminus B) \quad (4)$$

$$= \mathbb{P}(A) + \mathbb{P}(B) + \mathbb{P}(A \cup B) - \mathbb{P}(A \cap B), \quad (5)$$

which implies the desired equality. \square

Question 0.5. *Show that $0 \leq \mathbb{P}(A) \leq 1$ for every $A \in \Sigma$.*

Sometimes we have two conflicting beliefs, or models, about the underlying probability distribution, and so we will consider two compatible probability spaces $(\Omega, \Sigma, \mathbb{P})$ and $(\Omega, \Sigma, \mathbb{Q})$. They offer different predictions about the probability with which the events in Σ occur, and one fundamental task in statistics is to find out which model is the correct one from the frequency with which certain events occur. We will cover this later in the module.

0.2.2 Random variables

We will usually not deal directly with the probability space but with random variables. A *random variable* (rv) $X : \Omega \rightarrow \mathcal{X}$ is a function from the space (Ω, Σ) to a measurable space (\mathcal{X}, Σ_X) . In order for X to make any sense, the mapping has to ensure that $\{\omega \in \Omega : X(\omega) \in \mathcal{B}\} \in \Sigma$ for all $\mathcal{B} \in \Sigma_X$, because we are restricted to observing events in Σ and our random variable can thus not be more fine-grained than what Σ allows. Functions satisfying this property are called a *measurable function*. A random variable is then more formally defined as a measurable mapping from (Ω, Σ) to (\mathcal{X}, Σ_X) .

The only two examples of interest for us in the following are discrete and continuous random variables:

discrete rv: \mathcal{X} is a discrete set and Σ_X is the power set $\mathcal{P}(\mathcal{X})$ of \mathcal{X} , i.e. the set of all subsets of \mathcal{X} .

continuous rv: $\mathcal{X} = \mathbb{R}$ and $\Sigma_X = \mathcal{B}$, the Borel σ -algebra. This is the smallest σ -algebra containing all open intervals in \mathbb{R} .

The probability measure \mathbb{P} induces a probability measure P_X on (\mathcal{X}, Σ_X) , given by

$$P_X(B) = P[X \in B] = \mathbb{P}(\{\omega \in \Omega : X(\omega) \in B\}) \quad (6)$$

for all $B \in \Sigma'$. P_X is called the *distribution* of the random variable X .

If $\mathcal{X} = \{a_1, \dots, a_d\}$ is discrete (and Σ_X the power set of \mathcal{X}), then we say that X is a *discrete random variable*. The distribution of X is then also known as the *probability mass function* (pmf) of X and is fully characterised by all the events consisting of a single value, i.e. the values $P_X(a_1), P_X(a_2), \dots, P_X(a_d)$.

Question 0.6. *Can you give a formal argument why the values at these points are sufficient?*

Some random variables are not random at all. If there is an a_i with $P_X(a_i) = 1$ (and thus $P_X(a_j) = 0$ for all $j \neq i$), then we call this random variable *deterministic*. On the other extreme we have *uniformly distributed* random variables, where $P_X(a_i) = \frac{1}{d}$ for all $i \in [d]$.

Example 0.7. *The simplest example is the Bernoulli random variable. It is defined on a binary alphabet $\mathcal{X} = \{0, 1\}$ and we write $X \sim \text{Bern}(\epsilon)$ to denote the rv with $P[X = 1] = \epsilon$ and $P[X = 0] = 1 - \epsilon$.*

Let us now consider a real-valued random variable X . If there exists a function $p_X : \mathbb{R} \rightarrow [0, \infty)$ such that for all $A \in \Sigma_X$, we have

$$P[X \in A] = \int_A p_X(x) dx \quad (7)$$

then we say that X is a *continuous random variable*. The function p_X is called the *probability density function* (pdf) of X . We also define the *cumulative distribution function* (cdf) by integrating $p_X(x)$, that is, the cdf is given by $F_X(a) = \mathbb{P}[X \leq a] = \int_{-\infty}^a p_X(x) dx$.

Question 0.8. Show that $\int_{\mathcal{X}} p_X(x) = 1$. Moreover, if p_X is continuous at some point x , it must satisfy $p_X(x) \geq 0$. Can $p_X(x)$ ever be larger than 1?

In this class, we deal mainly with discrete rvs, although we will also encounter Gaussian random variables, which are continuous, later on.

Example 0.9. We denote the pdf of a Gaussian random variable X as

$$\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad (8)$$

where μ is the mean and σ the standard deviation of X . The variance of X is σ^2 . A normal Gaussian random variable has $\mu = 0$ and $\sigma = 1$. The corresponding cdf is denoted as

$$\Phi(y) = \int_{-\infty}^y \mathcal{N}(x; 0, 1) dx. \quad (9)$$

Some additional notations and definitions for discrete random variables are given below. The counterparts for continuous random variables can be obtained by simply replacing pmfs with pdfs. Thus assume now that X and Y are discrete random variables taking on values in \mathcal{X} and \mathcal{Y} respectively. The joint pmf of X and Y is defined as

$$P_{X,Y}(x, y) = P[X = x \wedge Y = y] = \mathbb{P}(\{\omega \in \Omega : X(\omega) = x \wedge Y(\omega) = y\}). \quad (10)$$

Question 0.10. Verify that $P_Y(y) = \sum_{x' \in \mathcal{X}} P_{X,Y}(x', y)$.

With this in hand we can define conditional pmf's and a notion of independence of random variables.

- The conditional pmf is given by

$$P_{X|Y}(x|y) = \frac{P_{X,Y}(x, y)}{P_Y(y)} = \frac{P_{Y|X}(y|x)P_X(x)}{P_Y(y)}, \quad \text{for } P_Y(y) > 0, \quad (11)$$

where the second expression is often referred to as Bayes' rule. If $P_Y(y) = 0$ then the conditional pmf is simply not defined.

- X and Y are *independent random variables*, if and only if, for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$,

$$P_{X,Y}(x, y) = P_X(x)P_Y(y) \quad (12)$$

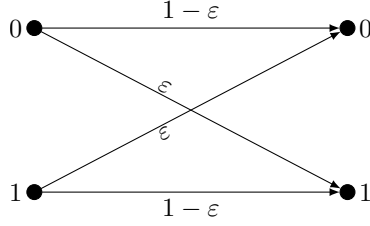
or equivalently $P_{X|Y}(x|y) = P_X(x)$. The latter condition simply states that the conditional distribution $P_{X|Y}(x|y)$ does not depend on y .

Example 0.11 (Binary symmetric channel). $X \sim \text{Bern}(p)$ is a bit that is sent over channel and is corrupted by additive noise $Z \sim \text{Bern}(\epsilon)$, where X and Z are independent. The output of the channel is $Y = X \oplus Z$. The channel is fully defined by the conditional distribution $P_{Y|X}$, which we can compute as follows:

$$P_{Y|X}(y|x) = P[X \oplus Z = y | X = x] = P[Z = y \oplus x | X = x] = P[Z = y \oplus x] = P_Z(y \oplus x) \quad (13)$$

Hence, the channel can be given as a matrix or pictorially as follows:

x	y	$P_{Y X}$
0	0	$1 - \epsilon$
1	0	ϵ
0	1	ϵ
1	1	$1 - \epsilon$



0.2.3 Expectation and variance

The expectation of a random variable X is defined to be

$$\mathbb{E}[X] = \int_{\Omega} X(\omega) d\mathbb{P}(\omega). \quad (14)$$

This definition has a very precise mathematical meaning in measure theory, but here we are only interested in two special cases. If X is a discrete random variable this reduces to the familiar formula

$$\mathbb{E}[X] = \sum_{x \in \mathcal{X}} x P_X(x). \quad (15)$$

If X is a continuous random variable with pdf $f_X(x)$, we have

$$\mathbb{E}[X] = \int_{\mathbb{R}} x p_X(x) dx. \quad (16)$$

Note that the expectation is a statistical summary of the distribution of X , rather than depending on the realised value of X . If there are two different models \mathbb{P} and \mathbb{Q} we need to specify which probability measure we are using. We only do this when necessary (because the model is not obvious from context) by adding a subscript \mathbb{E}_P or \mathbb{E}_Q .

If g is a function, the expectation of $g(X)$ is given by

$$\mathbb{E}[g(X)] = \int_{\mathbb{R}} g(x) p_X(x) dx. \quad (17)$$

Question 0.12. Show that the expectation is linear, i.e. $\mathbb{E}[aX + b] = a\mathbb{E}[X] + b$.

The variance of X is the expectation of $g(X) = (X - \mathbb{E}[X])^2$. Thus,

$$\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2] = \int_{\mathbb{R}} (x - \mathbb{E}[X])^2 p_X(x) dx. \quad (18)$$

Question 0.13. Check from the above definition that the variance can also be expressed as

$$\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2. \quad (19)$$

Question 0.14. Verify that $\mathcal{N}(x; \mu, \sigma^2)$ indeed has expectation μ and variance σ^2 .

0.2.4 Markov chains

Markov chains describe a notion of conditional independence. Let's start with the three random variables X, Y and Z . They are said to form a *Markov chain in the order*

$$X - Y - Z$$

if their joint distribution P_{XYZ} satisfies

$$P_{XYZ}(x, y, z) = P_X(x)P_{Y|X}(y|x)P_{Z|Y}(z|y) \quad \text{for all} \quad (x, y, z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Z}. \quad (20)$$

This is the same as saying that X and Z are *conditionally independent given Y* .

Question 0.15. Assume $X - Y - Z$. Show that it is also true that $Z - Y - X$.

Notice that if we do not assume anything about the joint distribution P_{XYZ} , then it factorizes (by repeated applications of Bayes rule) as

$$P_{XYZ}(x, y, z) = P_X(x)P_{Y|X}(y|x)P_{Z|XY}(z|x, y) \quad \text{for all} \quad (x, y, z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Z} \quad (21)$$

so what Markovianity in the order $X - Y - Z$ buys us is that $P_{Z|XY}(z|x, y) = P_{Z|Y}(z|y)$ (i.e., we can drop the conditioning on X). In essence all the information that we can learn about Z is already contained in Y . No other information about Z can be gleaned from knowing X if we already know Y . Another way of saying this is that the conditional distribution of X and Z given $Y = y$ can be factorised as

$$P_{XZ|Y}(x, z|y) = P_{X|Y}(x|y)P_{Z|Y}(z|y) \quad \text{for all} \quad (x, y, z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Z}. \quad (22)$$

Notice that this is in direct analogy to the situation where X and Z are (marginally) independent. Simply set Y to be a deterministic random variable (with only one possible outcome) to recover the definition of independence.

Question 0.16. If Z is a deterministic function of Y , show that $X - Y - Z$ is true.

Question 0.17. If X and Z are conditionally independent given Y , this does not imply that X and Z are marginally independent (in general). Construct a counterexample.

0.3 Tail bounds

In this section, we summarise some bounds on probabilities that we use extensively in the sequel. More precisely, we are interested in showing that the probability of a random variable deviating too far from its expectation value is small.

0.3.1 Basic bounds

We start with the familiar Markov and Chebyshev inequalities.

Proposition 0.18 (Markov's inequality). *Let X be a real-valued non-negative random variable with pdf p_X . Then for any $a > 0$, we have*

$$P[X > a] \leq \frac{\mathbb{E}[X]}{a}. \quad (23)$$

Proof. By the definition of the expectation, we have

$$\mathbb{E}[X] = \int_0^\infty xp_X(x) dx \geq \int_a^\infty xp_X(x) dx \geq a \int_a^\infty p_X(x) dx = aP[X > a]. \quad (24)$$

and we are done. \square

Note that this bound only becomes nontrivial if a exceeds the expectation value $\mathbb{E}[X]$.

Question 0.19. *In which step is non-negativity of X used?*

Question 0.20. *Can you do the proof also for discrete random variables?*

If we let X above be the non-negative random variable $(X - \mathbb{E}[X])^2$, we obtain Chebyshev's inequality.

Proposition 0.21 (Chebyshev's inequality). *Let X be a real-valued random variable with mean μ and variance σ^2 . Then for any $a > 0$, we have*

$$P[|X - \mu| > a\sigma] \leq \frac{1}{a^2}. \quad (25)$$

Proof. Let X in Markov's inequality be the random variable $g(X) = (X - \mathbb{E}[X])^2$. This is clearly non-negative and the expectation of $g(X)$ is $\text{Var}(X) = \sigma^2$. Thus, by Markov's inequality, we have

$$P[g(X) > a^2\sigma^2] \leq \frac{\sigma^2}{a^2\sigma^2} = \frac{1}{a^2}. \quad (26)$$

Now, $g(X) > a^2\sigma^2$ if and only if $|X - \mu| > a\sigma$ so the claim is proved. \square

We now consider a collection of real-valued random variables that are independent and identically distributed (i.i.d.). In particular, let $X^n = (X_1, \dots, X_n)$ be a collection of independent random variables where each X_i has distribution P with zero mean and finite variance σ^2 .

Proposition 0.22 (Weak Law of Large Numbers). *For every $\epsilon > 0$, we have*

$$\lim_{n \rightarrow \infty} P \left[\left| \frac{1}{n} \sum_{i=1}^n X_i \right| > \epsilon \right] = 0. \quad (27)$$

Consequently, the average $\frac{1}{n} \sum_{i=1}^n X_i$ converges to 0 in probability.

Note that for a sequence of random variables $\{S_n\}_{n=1}^\infty$, we say that this sequence *converges to a number $b \in \mathbb{R}$ in probability* if for all $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P[|S_n - b| > \epsilon] = 0. \quad (28)$$

We also write this as $S_n \xrightarrow{p} b$. Contrast this to convergence of numbers: We say that a sequence of numbers $\{s_n\}_{n=1}^\infty$ *converges to a number $b \in \mathbb{R}$* if we have $\lim_{n \rightarrow \infty} |s_n - b| = 0$.

Proof. Let $\frac{1}{n} \sum_{i=1}^n X_i$ take the role of X in Chebyshev's inequality. Clearly, the mean is zero. The variance of X is

$$\text{Var} \left(\frac{1}{n} \sum_{i=1}^n X_i \right) = \frac{1}{n^2} \text{Var} \left(\sum_{i=1}^n X_i \right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{\sigma^2}{n}. \quad (29)$$

Thus, we have

$$P \left(\left| \frac{1}{n} \sum_{i=1}^n X_i \right| > \epsilon \right) \leq \frac{\sigma^2}{n\epsilon^2} \rightarrow 0 \quad (30)$$

as $n \rightarrow \infty$, which proves the claim. \square

Some further useful bounds are derived in the homework.

0.3.2 Central limit theorem

We can actually say quite a bit more than the weak law of large numbers dictates. If the scaling in front of the sum in the statement of the law of large numbers Proposition 0.22 is $1/\sqrt{n}$ instead of $1/n$, the resultant random variable $\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i$ converges in distribution to a Gaussian random variable. As in Proposition 0.22, let X^n be a collection of i.i.d. random variables where each X_i is zero mean with finite variance σ^2 .

Proposition 0.23 (Central limit theorem). *For any $a \in \mathbb{R}$, we have*

$$\lim_{n \rightarrow \infty} P \left(\frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i < a \right) = \Phi(a). \quad (31)$$

In other words,

$$\frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i \xrightarrow{d} Z \quad (32)$$

where \xrightarrow{d} means convergence in distribution and Z is a standard Gaussian random variable.

For a sequence of random variables $\{S_n\}_{n=1}^\infty$, we say that this sequence of random variables *converges in distribution* to another random variable \bar{S} if

$$\lim_{n \rightarrow \infty} P(S_n < a) = P(\bar{S} < a)$$

for all $a \in \mathbb{R}$. The proof of this statement requires tools that are outside the scope of these notes, but can be found in any textbook on probability theory.

0.4 Vector norms and Cauchy-Schwarz inequality

We can naturally interpret pmf's on an alphabet with d symbols as row vectors in a d -dimensional inner-product space. Without loss of generality we take the alphabet to be $\mathcal{X} = \{1, 2, \dots, d\} = [d]$ and define the vector $p \in \mathbb{R}^d$ by its elements $p_x = P_X(x)$ for $x \in [d]$. The inner product is denoted by $\langle \cdot, \cdot \rangle$. For two general vectors $u, v \in \mathbb{R}^d$, it evaluates to

$$\langle u, v \rangle = uv^T = \sum_{i=1}^d u_i v_i, \quad (33)$$

where v^T denotes the transpose of the vector v , and is a column vector. The Cauchy-Schwarz inequality then states that for any two vectors we have

$$|\langle u, v \rangle|^2 \leq \langle u, u \rangle \langle v, v \rangle. \quad (34)$$

On these vector spaces we can also define the p -norms for $p \geq 1$ as

$$\|u\|_p = \left(\sum_{x=1}^d |u_x|^p \right)^{\frac{1}{p}} \quad (35)$$

We will mostly encounter the 1-norm and the 2-norm, the latter being the usual Euclidian norm of the vector. The following special case of the Cauchy-Schwarz inequality will be encountered later.

Lemma 0.24. *Let $u, v \in \mathbb{R}^d$. Then,*

$$\|u \cdot v\|_1 \leq \|u\|_2 \|v\|_2, \quad (36)$$

where \cdot denotes the element-wise product of the vectors, i.e. $(u \cdot v)_i = u_i v_i$.

Proof. Define $k \in \mathbb{R}^d$ using $k_i = \text{sgn}^*(u_i v_i)$, where sgn^* is the modified sign function, i.e.

$$\text{sgn}^*(x) = \begin{cases} -1 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}. \quad (37)$$

Then since $\text{sgn}^*(x)^2 = 1$ for all $x \in \mathbb{R}$, the Cauchy-Schwarz inequality yields

$$|\langle k \cdot u, v \rangle| \leq \langle u, u \rangle \langle v, v \rangle = \|k \cdot u\|_2 \|v\|_2 = \|u\|_2 \|v\|_2. \quad (38)$$

Moreover, we have

$$\langle k \cdot u, v \rangle = \sum_{x=1}^d k_x u_x v_x = \sum_{x=1}^d |u_x v_x| = \|u \cdot v\|_1. \quad (39)$$

□

Question 0.25. *Using the above, can you show that $\|u\|_1 \leq \sqrt{d} \|u\|_2$?*

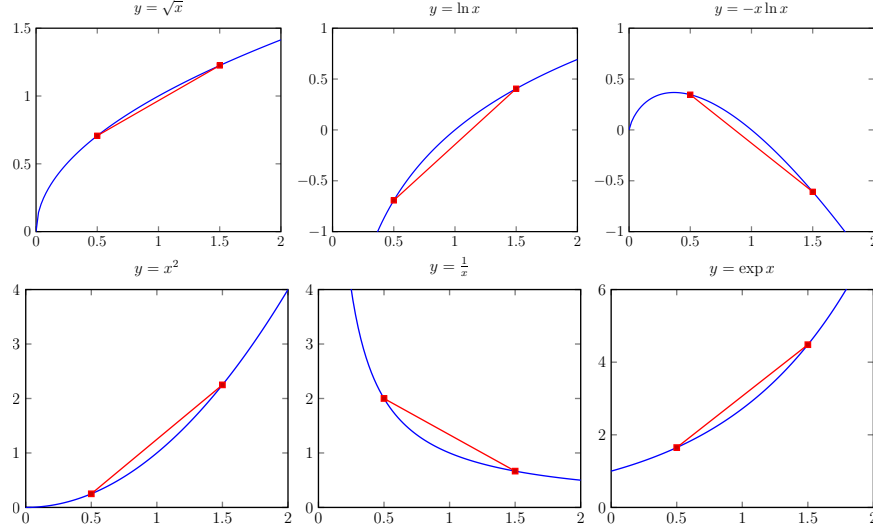


Figure 0.1: Examples of concave (upper ones) and convex (lower ones) functions. The straight line between two points of the curve is either below or above the plot of the function, which is exactly what the definition requires.

0.5 Convexity and Jensen's inequality

A function $f(x)$ is said to be *convex* on $[a, b]$ if for all $x, y \in [a, b]$ and $\lambda \in [0, 1]$, we have

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y). \quad (40)$$

If we do not mention any interval then we mean that the function is convex on its full domain, i.e. the statement $\log(x)$ is concave should be understood as $\log(x)$ is concave on $(0, \infty)$.

The function f is *strictly convex* if equality in (40) holds only if $\lambda = 0$ or 1 , or $x = y$. The function f is *concave* if $-f$ is convex, and *strictly concave* if $-f$ is strictly convex.

In the homework you will show the following lemma:

Lemma 0.26. *If f is convex on $[a, b]$, then for any $a \leq x_1 < x_2 \leq x_3 < x_4 \leq b$, we have*

$$\frac{f(x_2) - f(x_1)}{x_2 - x_1} \leq \frac{f(x_4) - f(x_3)}{x_4 - x_3} \quad (41)$$

Proposition 0.27 (Jensen's inequality). *If $f(x)$ is convex and X is a random variable on \mathbb{R} , then*

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)] \quad (42)$$

We only give a proof for discrete distributions here.

Proof. We give a proof by induction. Due to convexity, we have

$$p_1 f(x_1) + p_2 f(x_2) \geq f(p_1 x_1 + p_2 x_2), \quad (43)$$

which proves the statement if $|\mathcal{X}| = 2$.

Suppose the statement $\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$ is true when $|\mathcal{X}| = k - 1$. Then consider a pmf with k mass points $\{p_1, p_2, \dots, p_k\}$. Define another pmf on $k - 1$ points given by the probabilities

$$p'_i = \frac{p_i}{1 - p_k}, \quad i = 1, \dots, k - 1. \quad (44)$$

We then have

$$\sum_{i=1}^k p_i f(x_i) = p_k f(x_k) + (1 - p_k) \sum_{i=1}^{k-1} p'_i f(x_i) \quad (45)$$

$$\geq p_k f(x_k) + (1 - p_k) f\left(\sum_{i=1}^{k-1} p'_i x_i\right) \quad (46)$$

$$\geq f\left(p_k x_k + (1 - p_k) \sum_{i=1}^{k-1} p'_i x_i\right) \quad (47)$$

$$= f\left(\sum_{i=1}^k p_i x_i\right) \quad (48)$$

where the first inequality is from the induction hypothesis and the second by convexity (of two points). By the definition of expectation we have $\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$. \square

Often it is hard to check convexity directly. But for twice differentiable functions, this is easy.

Proposition 0.28. *Let $f : [a, b] \rightarrow \mathbb{R}$ be twice differentiable. The function f is convex if and only if $f''(x) \geq 0$ for all $x \in (a, b)$, and strictly convex if $f''(x) > 0$ for all $x \in (a, b)$.*

Proof. Assume $f''(x) > 0$ for all $x \in [a, b]$. By Taylor expansion of f around $x_0 \in (a, b)$, we have

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x^*)}{2}(x - x_0)^2 \quad (49)$$

where $x^* \in [x_0, x]$. By assumption $f''(x^*) > 0$ so the quadratic term is strictly positive unless $x = x_0$, in which case it is still non-negative. Now let $x_0 = \lambda x_1 + (1 - \lambda)x_2$. Further let $x = x_1$. Then we have

$$f(x_1) \geq f(x_0) + f'(x_0)((1 - \lambda)(x_1 - x_2)). \quad (50)$$

Now let $x = x_2$. Then we have

$$f(x_2) \geq f(x_0) + f'(x_0)(\lambda(x_2 - x_1)). \quad (51)$$

Both of these inequalities are strict unless $\lambda \in \{0, 1\}$ or $x_1 = x_2$. Multiplying the first inequality by λ and the second by $1 - \lambda$ and adding them up, we recover the definition of strict convexity. If we instead had assumed only $f''(x) \geq 0$ the same argument would ensure convexity (but no longer strict convexity).

For the other direction, choose $a < x_1 < x_2 < x_3 < x_4 < b$. By the property shown in Lemma 0.26,

$$\frac{f(x_2) - f(x_1)}{x_2 - x_1} \leq \frac{f(x_4) - f(x_3)}{x_4 - x_3} \quad (52)$$

Now let $x_2 \searrow x_1$ and $x_3 \nearrow x_4$. We see that $f'(x_1) \leq f'(x_4)$, and since these were arbitrary points, f' is increasing on (a, b) . So $f''(x) \geq 0$ for all $x \in (a, b)$. \square

0.6 Finite field arithmetic

This is a rather informal discussion, but it is sufficient for our purposes.

A finite field is a field (on which addition, subtraction, multiplication and division are defined) with a finite number of elements. Such fields are denoted by F_q where q is the number of elements in the field, or its *dimension*. For each dimension, the field is unique up to a relabelling of the elements. The idea is that such fields behave like \mathbb{Q} , \mathbb{R} or \mathbb{C} , with the usual rules for addition and multiplication.

A bit more formally, we have two binary operations on F_q denoted by $+$ and \cdot and the following properties (here a, b and c are any elements of F_q):

Associativity: $a + (b + c) = (a + b) + c$ and $a \cdot (b \cdot c) = (a \cdot b) \cdot c$.

Commutativity: $a + b = b + a$ and $a \cdot b = b \cdot a$.

Identities: There exist two different elements $0, 1$ such that $a + 0 = a$ and $a \cdot 1 = a$.

Additive inverse: Every a has an additive inverse, denoted $-a$, such that $a + (-a) = 0$.

Multiplicative inverse: Every $a \neq 0$ has a multiplicative inverse, denoted a^{-1} , such that $a \cdot a^{-1} = 1$.

Distributivity: $a \cdot (b + c) = (a \cdot b) + (a \cdot c)$.

Such fields exist only for particular numbers of elements, namely when $q = p^\ell$ for some prime p and $\ell \in \mathbb{N}$.

For F_q where q is prime we can always simply denote the elements of F_q by the integers $\{0, 1, \dots, q-1\}$ and use integer addition and multiplication modulo q as our operations.

Question 0.29. Verify the above properties for F_2 and F_3 . Can you find the inverse of 2 for a general F_q with prime q ? Use that $q + 1$ is even...

This strategy fails when $q = 4$ is not a prime. The problem is that using multiplication modulo 4 we would, for example, get

$$2 \times 0 = 0 \quad (53)$$

$$2 \times 1 = 2 \quad (54)$$

$$2 \times 2 = 4 \pmod{4} = 0 \quad (55)$$

$$2 \times 3 = 6 \pmod{4} = 2, \quad (56)$$

and hence 2 does not have a multiplicative inverse.

When q is a prime power $q = p^\ell$ we can derive the arithmetic using a polynomial ring. We give the construction here; but we do not attempt to show that it actually works or that it is unique. First, we denote the elements of F_q by strings of length ℓ with elements in F_p . In particular, if the underlying prime is 2, these are simply binary strings, e.g., $F_4 = \{00, 01, 10, 11\}$. We can then interpret these elements as polynomials of degree $\ell - 1$ with coefficients in F_p . Again, for F_4 the polynomials corresponding to the four elements would be $00 \rightarrow 0$, $01 \rightarrow 1$, $10 \rightarrow x$ and $11 \rightarrow x + 1$. We can add these polynomials modulo p for each coefficient individually, so in particular for the binary case the negation of each number is just the number itself. For multiplication, we simply do this modulo an irreducible polynomial. (An irreducible polynomial is one that has no roots in F_p .) The choice of irreducible polynomial turns out not to matter—the resulting fields are equivalent up to relabelling of elements. For F_4 we can take the irreducible polynomial to be $x^2 + x + 1$.

Question 0.30. Verify that $x^2 + x + 1$ is indeed irreducible over F_2 ? Is it also irreducible over F_3 ?

So for the above labelings $\{00, 01, 10, 11\}$ of elements, we get

$$10 \times 00 \rightarrow x \times 0 = 0 \rightarrow 00 \quad \implies 10 \times 00 = 00 \quad (57)$$

$$10 \times 01 \rightarrow x \times 1 = x \rightarrow 10 \quad \implies 10 \times 01 = 10 \quad (58)$$

$$10 \times 10 \rightarrow x \times x = x^2 \pmod{x^2 + x + 1} = x + 1 \rightarrow 11 \quad \implies 10 \times 10 = 11 \quad (59)$$

$$10 \times 11 \rightarrow x \times (x + 1) = x^2 + x \pmod{x^2 + x + 1} = 1 \rightarrow 01 \quad \implies 10 \times 11 = 01. \quad (60)$$

Hence, 10 and 11 are multiplicative inverses of each other. The full addition and multiplication tables can then be written down as follows:

+	00	01	10	11	×	00	01	10	11
00	00	01	10	11	00	00	00	00	00
01	01	00	11	10	01	00	01	10	11
10	10	11	00	01	10	00	10	11	01
11	11	10	01	00	11	00	11	01	10

Similar constructions can be done for every prime power, and, quite importantly for practical applications, all of this arithmetic can be implemented highly efficiently in computer programs.

Chapter 1

Information measures

[Week 2]

Intended learning outcomes:

- You can compute the entropy and conditional entropy for any discrete random variable and understand the basic properties of these two quantities, e.g. you can apply the chain rule or sub-additivity.
- You can compute mutual information and now how it relates to entropy and conditional entropy. You can apply the data-processing inequality for mutual information.
- You can compute the relative entropy and understand how entropy and mutual information can be expressed in terms of the relative entropy.

Book reference: Chapter 2 in Cover & Thomas [1], but we are not following it too closely.

1.1 Surprisal and entropy

It is not immediately clear how to model our intuitive notion of “information” in a mathematical language. In this chapter we take a somewhat axiomatic approach to information measures, i.e. we try to build them up from our intuitive understanding of what entropy and information “should” be. But we will only really be able to justify the choices we make here once we start analysing practical problems in information theory, and see that the quantities we derive here pop up again and again.

1.1.1 Surprisal

It turns out to be fruitful to start not by finding an expression for the information contained in a random variable, but rather the lack of information, or uncertainty inherent in a random experiment. Let us consider a discrete random variable X taking values in \mathcal{X} following the pmf $P_X(x) = p_x$. How surprised are we to see a particular outcome $x \in \mathcal{X}$ of this random

experiment? Clearly this depends on the probability p_x and not the value of x itself. In fact, we do not even need to know what \mathcal{X} really is. On the one hand, if $p_x = 1$ we are not surprised at all since we already knew that we would see x . On the other hand, the smaller p_x is the more surprised we are to see this particular outcome. If $p_x = 0$ we will never see x , so our surprise when seeing it anyway would be literally off the scale. Furthermore—and this turns out to be a very convenient choice—if we do a random experiment twice independently and both times observe x , we say that we will be twice as surprised as if we had seen x once in a single random experiment.

The above notions can be formalised, and that is essentially what Shannon did when he introduced the notion of *surprisal*. Let us denote the surprisal of x as $s(p_x)$. We want this function to satisfy the following three conditions:

1. **Monotonicity:** $s(p_x) = 0$ if $p_x = 1$ and $s(p_x)$ increases monotonically as p_x decreases.
2. **Additivity:** The surprisal of seeing a pair of outcomes of independent random experiments is simply the sum of the individual surprisals, i.e. $s(p_x p_y) = s(p_x) + s(p_y)$.
3. **Normalisation:** $s(\frac{1}{2}) = 1$

Question 1.1. *We do not really need the condition $s(p_x) = 0$ if $p_x = 1$ under Point 1 as it follows directly from additivity. Can you see how?*

It turns out that the only positive function that satisfies these three properties is the logarithm. To show this one uses a result by Erdős that characterises additive functions, but that is beyond the scope here. We therefore pick

$$s(p_x) = \log \frac{1}{p_x} . \quad (1.1)$$

where the logarithm is taken to base 2 (as everywhere in these notes) so that the normalisation requirement is satisfied.

We can see the surprisal as another random variable, say S , that takes the value $s(p_x) = \log \frac{1}{p_x}$ with probability p_x . Since $S = S(X)$ is a function of X we usually simply write this new random variable as

$$S(X) = \log \frac{1}{P_X(X)} . \quad (1.2)$$

1.1.2 Entropy

Entropy measures how much we can learn by looking at the outcome of a random experiment, or, in other words, how much uncertainty there is about the outcome. It is simply the expected surprisal of X .

Definition 1.2. *Given a discrete random variable X , the entropy of X is defined as*

$$H(X) := \mathbb{E}[S(X)] = \mathbb{E} \left[\log \frac{1}{P_X(X)} \right] = \sum_{x \in \mathcal{X}} p_x \log \frac{1}{p_x} \quad (1.3)$$

Here and throughout we use the convention that $0 \log 0 = 0$. This is reasonable since $\lim_{\epsilon \rightarrow 0} \epsilon \log \epsilon = 0$, and thus we simply continuously extend the function to the point 0.

Question 1.3. *Can you verify $\epsilon \log \epsilon \rightarrow 0$ as $\epsilon \rightarrow 0$?*

Note again that the entropy $H(X)$ is really only a function of the pmf of X , and in particular independent of the alphabet \mathcal{X} , in contrast to potential alternative uncertainty measures like the variance of X .

Sometimes we are interested in more than just the expected surprisal. The minimum surprisal, or min-entropy, for example, has applications in cryptography (see Chapter 3) and the variance of $S(X)$ has itself operational meaning in many information-theoretic problems when we go beyond first order asymptotics.

Question 1.4. *Can you find an expression for $\text{Var}[S(X)]$ in terms of the probabilities p_x ?*

Now let us explore the entropy a bit. First we want to show the following basic property.

Proposition 1.5. *Let X be a discrete random variable taking values in \mathcal{X} . We have*

$$H(X) \geq 0, \tag{1.4}$$

with equality if and only if X is deterministic.

Proof. Since $p_x \leq 1$, we have $\log \frac{1}{p_x} \geq 0$ for every $x \in \mathcal{X}$, so the expectation of this quantity over x must be non-negative too. In fact, $\log \frac{1}{p_x}$ equals 0 if and only if $p_x = 1$ and hence $H(X) = 0$ only if there exists an $x \in \mathcal{X}$ for which $p_x = 1$, which is the hallmark of a deterministic rv. \square

The entropy is a strictly concave function of the probability mass function P_X . To see this, we first verify that $f(t) = t \log \frac{1}{t} = -t \log t$ is concave on $(0, 1)$ by taking its second derivative:

$$f'(t) = -\log t - \log e, \quad f''(t) = -\frac{\log e}{t}. \tag{1.5}$$

Since the latter is always negative for $t \in (0, 1)$, the function is indeed strictly concave according to Lemma 0.28. Now since the entropy is simply the sum $\sum_{x \in \mathcal{X}} f(p_x)$ it is indeed a concave function of the pmf. This simple property, together with Jensen's inequality, has profound implications. The first one is that the entropy has a unique maximum. Intuitively we would want that entropy is maximal when uncertainty about the outcome is greatest, namely when the rv is uniformly distributed. And this is indeed the case.

Proposition 1.6. *Let X be a discrete random variable taking values in \mathcal{X} . We have*

$$H(X) \leq \log |\mathcal{X}|, \tag{1.6}$$

with equality if and only if X is uniformly distributed.

The general case will be covered in the homework but here we give a proof for the case when the set \mathcal{X} is a bit, i.e. when the random variable is binary.

Proof for $\mathcal{X} = \{0, 1\}$. It is easy to verify by a simple computation that $H(X) = 1$ for a uniformly distributed random variable, so the difficulty is only in showing that this is the maximum and only achieved for the uniform distribution.

Let now $\{p, 1 - p\}$ for $p \in [0, 1]$ be a general pmf for the random variable X . We use the function $f(t) = -t \log t$ to simplify notation. Then we can write

$$H(X) = f(p) + f(1 - p) \quad (1.7)$$

$$= \frac{1}{2} (f(p) + f(1 - p)) + \frac{1}{2} (f(p) + f(1 - p)) \quad (1.8)$$

$$\leq f\left(\frac{1}{2}p + \frac{1}{2}(1 - p)\right) + f\left(\frac{1}{2}p + \frac{1}{2}(1 - p)\right) \quad (1.9)$$

$$= f\left(\frac{1}{2}\right) + f\left(\frac{1}{2}\right) \quad (1.10)$$

$$= 1. \quad (1.11)$$

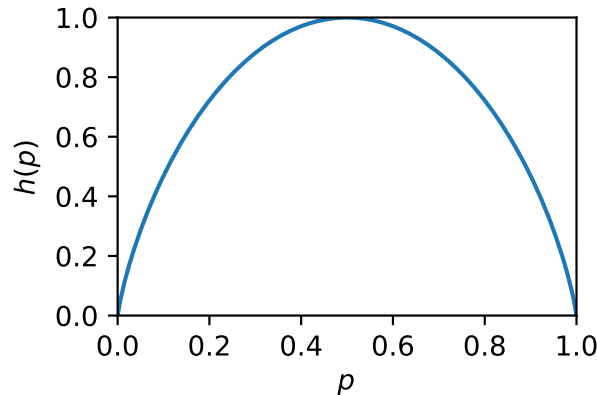
The inequality is due to Jensen's inequality and the strict concavity of f , and equality holds only if $p = 1 - p = \frac{1}{2}$, i.e. when the random variable X follows the uniform distribution. \square

Concavity in fact has even stronger consequences, and we will show a few additional properties of entropy later on using it.

Example 1.7. *The simplest example of a random variable is the Bernoulli random variable X with $\mathcal{X} = \{0, 1\}$ and $P_X(0) = p$ for $p \in [0, 1]$. The entropy of the Bernoulli random variable is called the binary entropy,*

$$H(X) = p \log \frac{1}{p} + (1 - p) \log \frac{1}{1 - p} =: h(p). \quad (1.12)$$

From the plot we can easily verify all the properties we discussed above.



1.2 Conditional entropy, and mutual information

1.2.1 Joint entropy

For two discrete random variables X and Y with joint pmf $P_{XY}(x, y) = p_{xy}$ we can simply consider (X, Y) as one single random variable and use the same construction to define the surprisal of a tuple (X, Y) as $S(X, Y) = -\log P_{XY}(X, Y)$. Its expectation is the *joint entropy*, $H(XY)$, given by

$$H(XY) := \mathbb{E}[S(X, Y)] = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{xy} \log \frac{1}{p_{xy}} \quad (1.13)$$

The first thing to note is that —if X and Y are independent— then $p_{xy} = p_x \cdot p_y$ and thus the expression simplifies to

$$H(XY) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{xy} \log \frac{1}{p_x p_y} \quad (1.14)$$

$$= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{xy} \log \frac{1}{p_x} + \log \frac{1}{p_y} \quad (1.15)$$

$$= \sum_{x \in \mathcal{X}} p_x \log \frac{1}{p_x} + \sum_{y \in \mathcal{Y}} p_y \log \frac{1}{p_y} \quad (1.16)$$

$$= H(X) + H(Y). \quad (1.17)$$

This is not true in general though if the two random variables are correlated.

Question 1.8. Find an example for which $H(XY) = H(X) = H(Y) = 1$.

1.2.2 Conditional entropy

So why do these entropies not just add up? Fundamentally, this is because once we learn X we might not be so surprised seeing some particular outcomes of the random variable Y anymore. In fact, in the most extreme case, we have $Y = f(X)$ for some function f ; hence, once we know that X takes on the value x , we can immediately deduce that Y will take on the value $f(x)$ with probability one, and thus there is no surprisal anymore! We model this “conditional surprisal” using the conditional pmfs, $P_{Y|X}(y|x) = p_{y|x}$, which leads us to conditional entropy.

Definition 1.9. The conditional entropy of Y given X is defined as

$$H(Y|X) = \mathbb{E} \left[\log \frac{1}{P_{Y|X}(Y|X)} \right] = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{xy} \log \frac{1}{p_{y|x}}. \quad (1.18)$$

This can be interpreted as the expectation of the entropy of Y over all outcomes X . We sometimes use the notation $H(Y|X = x) = H(Y_x)$ to denote the entropy of the random variable Y_x that follows the pmf $\{p_{y|x}\}_{y \in \mathcal{Y}}$, i.e., the pmf of Y when we already know that $X = x$. Using this and the expression in (1.18) we can write the conditional entropy as

$$H(Y|X) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{xy} \log \frac{1}{p_{y|x}} \quad (1.19)$$

$$= \sum_{x \in \mathcal{X}} p_x \sum_{y \in \mathcal{Y}} p_{y|x} \log \frac{1}{p_{y|x}} \quad (1.20)$$

$$= \sum_x p_x H(Y|X = x). \quad (1.21)$$

The last line which expresses the conditional entropy in terms of an average of (unconditional) entropies is particularly useful since it allows us to immediately conclude that the conditional entropy is also bounded from below and above, like the entropy. We thus have

$$0 \leq H(Y|X) \leq \log |\mathcal{Y}|. \quad (1.22)$$

Moreover, our definition of conditional entropy also allows us to establish a *chain rule* for the conditional entropy, which sometimes is in fact used as the definition of conditional entropy itself. This rule is very useful because it allows us to write the joint entropy as a sum of its parts, even if the two random variables are not independent.

Proposition 1.10. *We have $H(XY) = H(X) + H(Y|X)$.*

Proof. We take advantage of $p_{xy} = p_x p_{y|x}$ to write

$$H(XY) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{xy} \log \frac{1}{p_{xy}} \quad (1.23)$$

$$= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{xy} \log \frac{1}{p_x} + \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{xy} \log \frac{1}{p_{y|x}} \quad (1.24)$$

$$= \sum_{x \in \mathcal{X}} p_x \log \frac{1}{p_x} + \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{xy} \log \frac{1}{p_{y|x}} \quad (1.25)$$

$$= H(X) + H(Y|X). \quad (1.26)$$

□

Question 1.11. *Show that $H(Y|X) = H(Y)$ for independent random variables. Using the chain rule, find a different proof that $H(XY) = H(X) + H(Y)$ in this case.*

Now we have put everything in place to show our first entropic inequality, which relates the entropy of two random variables with their joint entropy. This result shows the *subadditivity* of entropy.

Proposition 1.12. *Let X and Y be two discrete random variables. Then*

$$H(XY) \leq H(X) + H(Y) \quad \text{or, equivalently,} \quad H(X|Y) \leq H(X). \quad (1.27)$$

Equality holds in either statement only if X and Y are independent.

Proof. The equivalence of the two relations follows directly from the chain rule, we thus only need to show the second statement.

We start with Eq. (1.21), which states that

$$H(Y|X) = \sum_x p_x H(Y|X = x) \quad (1.28)$$

$$= \sum_x p_x \sum_y p_{y|x} \log \frac{1}{p_{y|x}} \quad (1.29)$$

$$= \mathbb{E} \left[\sum_y p_{y|X} \log \frac{1}{p_{y|X}} \right] \quad (1.30)$$

Note that sum inside the expectation is simply another expectation, as in the definition of entropy—but since we only want to apply Jensen’s inequality on the outer expectation we spell this one out explicitly. Moreover, by definition of the conditional pmf we have $\mathbb{E}[p_{y|X}] = \sum_x p_x p_{y|x} = \sum_x p_{xy} = p_y$. Hence, using concavity of the entropy as a function of the pmf and Jensen’s inequality for the outer expectation, we find

$$H(Y|X) = \mathbb{E} \left[\sum_y p_{y|X} \log \frac{1}{p_{y|X}} \right] \quad (1.31)$$

$$\leq \sum_y (\mathbb{E}[p_{y|X}]) \log \frac{1}{\mathbb{E}[p_{y|X}]} \quad (1.32)$$

$$= \sum_y p_y \log \frac{1}{p_y} \quad (1.33)$$

$$= H(Y). \quad (1.34)$$

Equality in Jensen’s inequality only holds if either X is deterministic or if $p_{y|x} = p_y$ for all x and y , but this only holds if X and Y are in fact independent. \square

The second relation in Eq. (1.27) can be strengthened by considering three random variables X , Y and Z . In that case, we have

$$H(X|YZ) \leq H(X|Z). \quad (1.35)$$

This is sometimes referred to as *strong sub-additivity*. The proof follows from (regular) sub-additivity, applied to the entropies $H(X|Y, Z = z)$ and $H(X|Z = z)$, and averaging the resulting inequalities.

Question 1.13. *Can you construct a formal proof out of the above sketch?*

1.2.3 Mutual information

We have already established that $H(XY) \neq H(X) + H(Y)$ in general, and hence also $H(Y|X) \neq H(Y)$ by the chain rule. The difference between these two quantities clearly tells us something about how much the uncertainty about Y changes when we learn X , or in other words, about how much information X contains about Y . This leads us to the definition of mutual information,

Definition 1.14. *The mutual information between X and Y is defined as*

$$I(X : Y) := H(Y) - H(Y|X) \quad (1.36)$$

It is not evident immediately from the way we defined it here but this expression is symmetric between X and Y . Namely, using the chain rule for conditional entropy (recall Proposition 1.10) twice, we can write

$$I(X : Y) = H(Y) - H(Y|X) = H(Y) + H(X) - H(XY) = H(X) - H(X|Y). \quad (1.37)$$

The mutual information is thus a symmetric measure of the correlation between the two random variables.

Using these various equivalent expressions it is then easy to derive some bounds on the mutual information. First, sub-additivity of the entropy directly implies that $I(X : Y) \geq 0$, so the mutual information is non-negative, and it vanishes only if the two random variables are independent (a consequence of Proposition 1.12). This is consistent with our intuitive notion of information—we cannot know less than nothing after all! We also cannot know more than everything, i.e. the mutual information can never exceed the minimal entropy of its constituent parts.

Question 1.15. *Using the bounds on entropies established in the previous sections, show that $I(X : Y) \leq \log \min\{|\mathcal{X}|, |\mathcal{Y}|\}$. Give an example that saturates the bound.*

Example 1.16. *Consider two binary random variables X and Y with joint pmf*

$$P_{XY}(0,0) = P_{XY}(1,1) = \frac{1}{4}(1+r), \quad P_{XY}(0,1) = P_{XY}(1,0) = \frac{1}{4}(1-r) \quad (1.38)$$

for $r \in [-1, 1]$. We can compute the mutual information between X and Y as follows:

$$I(X : Y) = H(X) - H(X|Y) = 1 - h\left(\frac{1+r}{2}\right) \quad (1.39)$$

So this function takes its maximum at $r = -1$ and $r = 1$ and drops to zero for $r = 0$.

Question 1.17. *You might have heard of the correlation coefficient in statistics:*

$$\rho = \frac{\mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]}{\sqrt{\text{Var}(X)\text{Var}(Y)}} \quad (1.40)$$

Can you determine ρ as a function of r ?

If we have three random variables X , Y and Z we can ask for the mutual information between X and Y conditioned on knowing Z , the *conditional mutual information*. It is defined as

$$I(X : Y|Z) := \sum_z P_Z(z) I(X : Y|Z = z). \quad (1.41)$$

Various equivalent expressions can then be readily derived, e.g.,

$$I(X : Y|Z) = H(Y|Z) - H(Y|XZ) = H(X|Z) - H(X|YZ). \quad (1.42)$$

Moreover, the *chain rule* for the mutual information states that

$$I(X : YZ) = I(X : Y) + I(X : Z|Y), \quad (1.43)$$

which can be verified by a close inspection of the definition of both conditional and unconditional mutual information.

Consider now the special case where $X - Z - Y$ form a Markov chain. In this case $P_{X|YZ} = P_{X|Z}$ and thus $H(X|YZ) = H(X|Z)$. As a consequence, the conditional mutual information $I(X : Y|Z)$ as written in (1.42) vanishes.

One of the most intriguing properties of the mutual information is the *data-processing inequality* for mutual information. It states that the mutual information can never increase when we apply an operation that only acts on one of the parts. Intuitively this tells us that by manipulating one of the random variables without looking at the other we cannot increase the correlations between the pair.

We can formalise this using the notion of Markov chains.

Proposition 1.18. *Let $X - Y - Z$ form a Markov chain. Then, $I(X : Y) \geq I(X : Z)$.*

Proof. Since $I(X : Z|Y) = 0$, the chain rule for mutual information implies that $I(X : Y) = I(X : YZ)$. It thus remains to show that

$$I(X : Z) \leq I(X : YZ). \quad (1.44)$$

But, since $I(X : Z) = H(X) - H(X|Z)$ and $I(X : YZ) = H(X) - H(X|YZ)$, the relation in Eq. (1.44) is equivalent to the condition $H(X|Z) \geq H(X|YZ)$, which is in turn ensured by the strong sub-additivity of entropy. \square

Question 1.19. *Can you also show that $I(Y : Z) \geq I(X : Z)$ under the same assumption?*

1.3 Relative entropy

The relative entropy appears when we want to compare two different probability distributions. We define it here only for discrete random variables (or rather the respective pmfs), but this can be readily generalised to other probability measures.

Definition 1.20. Let P and Q be two pmfs on an alphabet \mathcal{X} . The relative entropy of P with regards to Q is defined as

$$D(P\|Q) := \sum_{\substack{x \in \mathcal{X} \\ P(x) > 0}} P(x) \log \frac{P(x)}{Q(x)}. \quad (1.45)$$

if $P(x) > 0 \implies Q(x) > 0$ for all $x \in \mathcal{X}$, and $D(P\|Q) = +\infty$ otherwise.

In the following, instead of restricting the sum, we will use the convention that $0 \log \frac{0}{0} = 0$.

We can alternatively see the relative entropy as the expectation value of the *log-likelihood ratio*, namely we can write

$$D(P\|Q) = \mathbb{E}[Z(X)], \quad \text{where} \quad Z(X) = \log \frac{P(X)}{Q(X)} \quad (1.46)$$

and X is distributed according to P . The random variable $Z(X)$ is called the log-likelihood ratio. It takes on the role of the surprisal in the definition of entropy. We will explore this random variable and its distribution much more when we discuss the information spectrum method and hypothesis testing later on.

Just by manipulating the definition, we are able to show the following equivalences.

Proposition 1.21. Let X and Y be random variables on alphabets \mathcal{X} and \mathcal{Y} . Moreover, let U be a uniform random variable on \mathcal{X} . Then the following relations are true:

$$H(X) = \log |\mathcal{X}| - D(P_X\|U_X) \quad (1.47)$$

$$H(X|Y) = \log |\mathcal{X}| - D(P_{XY}\|U_X \times P_Y) \quad (1.48)$$

$$I(X : Y) = D(P_{XY}\|P_X \times P_Y). \quad (1.49)$$

You will prove these equivalences in the homework. They turn out to be very useful because they essentially tell us that once we established properties of the relative entropy this has immediate consequences also for the derived quantities

We will need two important properties of the relative entropy. The first proposition establishes that the relative entropy is always positive.

Proposition 1.22. For any two pmfs P and Q , we have $D(P\|Q) \geq 0$ with equality if and only if $P = Q$.

Proof. We can assume without loss of generality that the quantity is finite, as otherwise the

statement is trivially true. We first note that $x \mapsto -\log x$ is strictly convex. Hence,

$$D(P\|Q) = \sum_{x:P(x)>0} P(x) \log \frac{P(x)}{Q(x)} \quad (1.50)$$

$$= \sum_{x:P(x)>0} P(x) \left(-\log \frac{Q(x)}{P(x)} \right) \quad (1.51)$$

$$\geq -\log \left(\sum_{x:P(x)>0} P(x) \frac{Q(x)}{P(x)} \right) \quad (1.52)$$

$$= -\log \left(\sum_{x:P(x)>0} Q(x) \right) \quad (1.53)$$

$$\geq -\log \left(\sum_x Q(x) \right) = 0. \quad (1.54)$$

Equality in the second inequality only holds if P and Q have the same support. Moreover, equality in the first inequality holds if $\frac{Q(x)}{P(x)}$ is independent of x for any x in the support of P . These two statements are both true only if $P(x) = Q(x)$ for all $x \in \mathcal{X}$, and thus $P = Q$. \square

An immediate corollary of Propositions 1.21 and 1.22 is that $I(X : Y)$ is positive and zero only if X and Y are independent.

Question 1.23. *Can you see why?*

Finally, there is one property of the relative entropy that implies all other properties of both entropy and mutual information. It states that applying a noisy operation, i.e. a stochastic map or channel, on both arguments of the relative entropy will never increase the relative entropy. Together with the positivity of relative entropy this justifies that we think of it as a measure of similarity or distinguishability. If the relative entropy is small the two pmfs are similar and hard to distinguish by observing the outcomes of a random experiment. Observing the outcomes after further noise has been applied should make distinguishing them even harder, and that is exactly what the *data-processing inequality* for relative entropy tells us.

Proposition 1.24. *Let P_X and Q_X be two pmfs on an alphabet \mathcal{X} (the input distributions), and let $P_{Y|X}$ be a conditional pmf (the channel). Define the marginals (the output distributions)*

$$P_Y(y) = \sum_{x \in \mathcal{X}} P_{Y|X}(y|x) P_X(x) \quad \text{and} \quad Q_Y(y) = \sum_{x \in \mathcal{X}} P_{Y|X}(y|x) Q_X(x). \quad (1.55)$$

Then, the data-processing inequality (DPI) states that

$$D(P_X\|Q_X) \geq D(P_Y\|Q_Y). \quad (1.56)$$

Proof. Consider now the joint distributions $P_{XY}(x, y) = P_{Y|X}(y|x)P_X(x)$ and $Q_{XY}(x, y) = P_{Y|X}(y|x)Q_X(x)$, using the usual shorthand notation for conditional and marginal distributions. We first show that

$$D(P_{XY}||Q_{XY}) - D(P_Y||Q_Y) = \left(\sum_{x,y} p_{xy} \log \frac{p_{xy}}{q_{xy}} \right) - \left(\sum_y p_y \log \frac{p_y}{q_y} \right) \quad (1.57)$$

$$= \sum_{x,y} p_{xy} \left(\log \frac{p_{xy}}{q_{xy}} - \log \frac{p_y}{q_y} \right) \quad (1.58)$$

$$= \sum_y p_y \sum_x p_{x|y} \log \frac{p_{x|y}}{q_{x|y}} \quad (1.59)$$

$$= \sum_y p_y D(P_{X|Y=y}||Q_{X|Y=y}) \geq 0, \quad (1.60)$$

where we have used the positivity of relative entropy in the last step. Similarly, we have

$$D(P_{XY}||Q_{XY}) - D(P_X||Q_X) = \sum_x p_x D(P_{Y|X=x}||Q_{Y|X=x}) = 0 \quad (1.61)$$

since $Q_{Y|X} = P_{Y|X}$ by construction of the joint distribution. Combining Eqs. (1.57)–(1.60) and (1.61) yields the desired inequality. \square

It turns out that all the properties of entropy, conditional entropy and mutual information we discussed previously can be derived from the DPI. As an example we give here a strengthening of the strong sub-additivity, which we call the data-processing inequality for conditional entropy. It intuitively states that any processing of the side information can at most increase the conditional entropy.

Corollary 1.25. *Let P_{XY} be a joint pmf and $P_{Z|Y}$ a conditional pmf. Define now the pmf*

$$P_{XZ}(x, z) = \sum_y P_{XY}(x, y) P_{Z|Y}(z|y). \quad (1.62)$$

Then, we have $H(X|Y) \leq H(X|Z)$.

Proof. Let us express the inequality in terms of relative entropies using Proposition 1.21. This reads

$$\log |\mathcal{X}| - D(P_{XY}||U_X \times P_Y) \leq \log |\mathcal{X}| - D(P_{XZ}||U_X \times P_Z). \quad (1.63)$$

or simply $D(P_{XY}||U_X \times P_Y) \geq D(P_{XZ}||U_X \times P_Z)$. But this is imply the DPI applied to the channel $P_{Z|Y}$ that happens to leave X untouched. \square

Chapter 2

Source coding

[Week 3–5]

Intended learning outcomes:

- You can determine if a code is instantaneous and if the codeword lengths are optimal using the McMillan–Kraft inequality.
- You can evaluate the quality of a variable-length code by comparing it to the fundamental limits.
- You can construct a Huffman code for any discrete source, and understand the algorithm and its properties.
- You understand the mathematical model used to study block codes asymptotically, and can compute the code rate.
- You understand Fano’s inequality and typical sets and how they can be used to derive the fundamental limits of compression.

Book reference: Chapter 5 in Cover & Thomas [1].

2.1 Problem setup and definitions

In this chapter we are concerned with removing redundancy. In the first section we will introduce the formal mathematical model we use to investigate source coding, or compression.

2.1.1 Data source

We are given data as a sequence of symbols, for example these could be numbers, letters, colours of pixels, etc., and we would like to store that data in a (preferably short) sequence of bits. (We could generalise to larger alphabets, but conceptually nothing changes so we restrict ourselves to bits here to simplify presentation.) We start by formally defining what we mean by a *data source*, or simply *source* in the remainder of this chapter.

Definition 2.1. A data source is an infinite sequence of random variables

$$\mathbf{X} = X_1, X_2, \dots, X_k, \dots \quad (2.1)$$

- A source is called *discrete* if the random variables are discrete, i.e. if the source outputs in each iteration $i \in \mathbb{N}$ values from a finite set \mathcal{X} .
- A source is furthermore called *memoryless* if the X_i are independent and identically distributed (i.i.d.) according to the same pmf P_X , i.e., if we have

$$P_{X_1 X_2 \dots X_k \dots}(x_1, x_2, \dots, x_k, \dots) = P_X(x_1) P_X(x_2) \dots P_X(x_k) \dots \quad (2.2)$$

Memoryless here refers to the fact that the distribution of X_i does not depend on the value of X_{i-1} or any other symbol in the sequence, or, formally, $P_{X_i | X_1 X_2 \dots X_{i-1}} = P_{X_i} = P_X$. We will not consider sources that output continuous values in this module.

Question 2.2. Can you see why we cannot expect to store and perfectly recover a continuous variable using finite (digital) memory?

For most of our theoretical analysis, we will consider a *discrete memoryless source (DMS)*. An example of such a source is the sequence of face values one gets by throwing the same (fair or unfair) die repeatedly. Generally, the assumption that a source is memoryless is simplifying the analysis but in fact most sources do not satisfy this exactly. For example, think of a book (written in English) as a sequence of letters and a source reproducing them one by one. If $X_{i-1} = \text{'q'}$, then $X_i = \text{'u'}$ with much higher probability than the frequency of 'u' in English text would otherwise suggest. Hence, this source is far from memoryless and the corresponding distribution of the random variable is not i.i.d.. Nonetheless, understanding the simple case of discrete memoryless sources properly will allow us to get an intuition for more loosely structured sources as well. Various more complicated models of sources have been analysed in the literature.

2.1.2 Source codes

Next we introduce *source codes*, or simply *codes* for the remainder of this chapter.

Definition 2.3. A source code is a map e that maps outputs of the source $x \in \mathcal{X}$ to bit strings of variable length, $C(x) \in \{0, 1\}^*$. We denote by $\ell(x)$ the length of the codeword $C(x)$.

- A code is called a *fixed-length code* if $\ell(x) = \ell$ is constant, otherwise it is called a *variable-length code*.

(NS) A code is called *non-singular* if C is injective, i.e. if every $x \in \mathcal{X}$ is mapped to a unique bit string.

- (P) A code is called a *prefix code* if for any pair $x, x' \in \mathcal{X}$ with $x \neq x'$, the codeword $C(x)$ is not a prefix of the codeword $C(x')$.¹
- (U) A code is *uniquely decodable* if there exists a decoder that, for any $n \in \mathbb{N}$ and any sequence $x^n \in \mathcal{X}^n$, can uniquely recover x^n from the bit string $C(x_1)C(x_2) \dots C(x_n)$.
- (I) A code is *instantaneous* if it is uniquely decodable and if the decoder can deduce the k -th symbol x_k as soon it has seen the bit string $C(x_1)C(x_2) \dots C(x_k)C(x_{k+1}) \dots$ up to and including all of $C(x_k)$, even if there is no guarantee that the string is complete.

Let us note that the codes we consider here are not as general as they could be. In fact, we could also consider codes that take a variable length sequence of input symbols to a codeword (of either fixed or variable length). Such codes are in fact often used in practical applications. A prominent example is the Lempel–Ziv–Markov algorithm for lossless compression, which uses a dictionary to replace often reoccurring variable-length sequences with shorter codewords.

Let us now discuss some of the interrelations between all these code properties. Clearly, any uniquely decodable code is non-singular by definition. However, we observe that not every non-singular code is uniquely decodable. Consider a code on $\mathcal{X} = \{0, 1, 2, 3\}$ that yields the binary representation

$$C(0) = 0, \quad C(1) = 1, \quad C(2) = 10, \quad C(3) = 11. \quad (2.3)$$

This code is non-singular but not a prefix code. The codeword string 110 could either be produced by the source sequence (3, 0), by (1, 2) or even by (1, 1, 0), so there is no way for a decoder to distinguish between these three cases.

Proposition 2.4. *A code is instantaneous if and only if it is a prefix code*

Proof. We first show that a prefix code is instantaneous by constructing a decoder. The decoder will read the sequence $C(x_1)C(x_2) \dots$ bit by bit. Once $C(x_1)$ is fully read we can immediately deduce that the first source symbol was x_1 since $C(x_1)$ cannot be a prefix for a longer codeword. Similarly, with this rule in mind, since there is no other codeword that is a prefix to $C(x_1)$ we can be assured that this is indeed the first symbol we will decode. The same procedure continues for the remainder of the string with $C(x_2)C(x_3) \dots$. We know where every codeword ends and can decode them individually and instantaneously.

To verify the other direction, simply note that if a decoder can decide instantly once it has seen the codeword $C(x)$ this implies that $C(x)$ cannot be a prefix to any other codeword $C(x')$. Since this is true for any $x \neq x'$, the code must be a prefix code. \square

Thus, any prefix code is uniquely decodable. However, the converse is not true in general, i.e. not every uniquely decodable code is a prefix code. Consider as an example the code

$$C(\text{'a'}) = 1, \quad C(\text{'b'}) = 10, \quad C(\text{'c'}) = 00. \quad (2.4)$$

¹We say that a bit string is a prefix of another bit string if the latter starts with the former, e.g. '01' is a prefix to '0100'.

After seeing 10 we cannot decide if the first symbol was ‘a’ or ‘b’ even though we have seen the full codeword of the first symbol, hence this code is not instantaneous. However, once we have seen a full sequence of codewords, for example 100, we can decode uniquely by looking at the parity of the number of 0’s between two 1’s.

Question 2.5. *Can you come up with a formal decoder for this code?*

Codes can be conveniently represented by binary trees. Binary trees are connected graphs without cycles (trees) where each node (except the root) has exactly one parent and either zero (in which case it is called a leaf) or two children. The two branches emanating from the root and each node are assigned values ‘0’ or ‘1’ and codewords are composed by following a path from the root to a node.

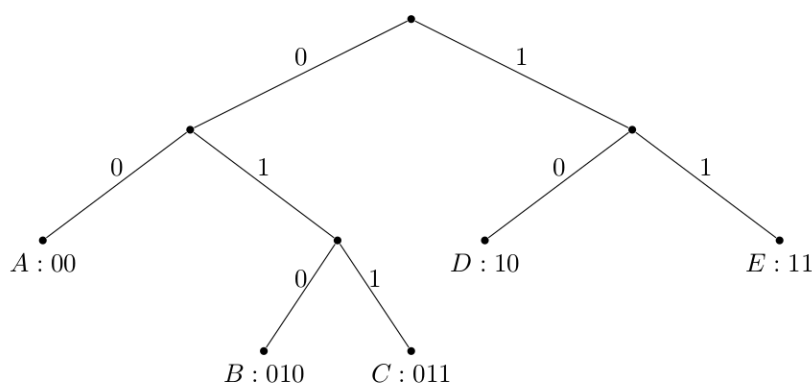


Figure 2.1: Example of a code tree.

The codeword length is equivalent to the depth (i.e. the distance from the root) of the node in the tree. For a fixed-length code all the codewords are at the same depth (or level) of the tree. A code is a prefix code if and only if all the codewords are on leaves of the tree.

Question 2.6. *Can you see why this is the case?*

The next two sections will be devoted to variable-length and fixed-length block codes (defined later), respectively.

2.2 Variable-length codes

Before we discuss particular codes, we first want to establish some fundamental limits all codes need to satisfy.

2.2.1 Optimal codeword lengths

The Kraft inequality gives a lower bound on the lengths of codewords in any instantaneous code. It is the first fundamental limit we will establish, it shows us that no code with shorter codeword lengths can exist, and thus if a code achieves equality in (2.5) we know it is optimal in this regard. We present a slightly more general result, the MacMillan–Kraft inequality, which applies for any uniquely decodable code (and not just prefix codes).

Proposition 2.7 (MacMillan–Kraft inequality). *Any uniquely decodable code must satisfy the inequality*

$$\sum_{x \in \mathcal{X}} 2^{-\ell(x)} \leq 1. \quad (2.5)$$

Conversely, given a set of codeword lengths satisfying Eq. (2.5), it is possible to construct a prefix code with these lengths.

We show the inequality only for instantaneous codes, and the general proof for uniquely decodable (not necessarily instantaneous) codes will be covered in the homework.

Proof. We use the one-to-one correspondence of prefix codes with binary trees for which the codewords are leaves. For any tree we may assign to every node a value 2^{-d} where d is the depth in the tree. The root thus gets the value 1. To show the Kraft inequality we simply need to show that the sum of all the leaf values in a tree cannot exceed 1. To see that this is correct, simply note that by our construction any parent node’s value is simply the sum of its children’s values, so as we build up the binary tree from the root the sum of the values on all leaves is always exactly 1.

Conversely, given a set of codeword lengths satisfying the inequality, we can always create a binary tree with leaves at the corresponding depths and populate the leaves with codewords. If the inequality is strict there will be unused leaves in the tree. \square

2.2.2 Optimal expected codeword length

Finding codewords with short lengths is however only half of the problem—we also need to assign those codewords to elements of \mathcal{X} . And we want to do this in such a way as to minimise the expected length of the codeword. That is, for a code $C(x)$ with codeword lengths $\ell(x)$, we define the *expected codeword length* as

$$\bar{\ell}(C) := \sum_x P_X(x) \ell(x) = \mathbb{E}[\ell(X)] \quad (2.6)$$

Using the MacMillan–Kraft inequality from Proposition 2.7, we can lower bound $\bar{\ell}(C)$ with the entropy $H(X)$ for any uniquely decodable code.

Proposition 2.8. *For any uniquely decodable code C for a discrete source X with distribution P_X , we have*

$$\bar{\ell}(C) \geq H(X). \quad (2.7)$$

Moreover, the equality is saturated if only if the codeword lengths saturate the McMillan-Kraft inequality and $P_X(x) = 2^{-\ell_x}$ for some set of numbers $\ell_x \in \mathbb{N}$.

Proof. We evaluate

$$\bar{\ell}(C) - H(X) = \mathbb{E} \left[\ell(X) - \log \frac{1}{P_X(X)} \right] \quad (2.8)$$

$$= \mathbb{E} \left[\log \frac{P_X(X)}{2^{-\ell(X)}} \right] \quad (2.9)$$

$$\geq \mathbb{E} \left[\log \frac{t \cdot P_X(X)}{2^{-\ell(X)}} \right] \quad (2.10)$$

$$= D(P_X \| Q_X) \quad (2.11)$$

$$\geq 0, \quad (2.12)$$

where we introduced the constant $t = \sum_x 2^{-\ell(x)} \leq 1$ (by the McMillan-Kraft inequality) and the pmf $Q(x) = \frac{1}{t} \cdot 2^{-\ell(x)}$. The final inequality is simply due to the non-negativity of relative entropy.

If the conditions for saturation are met we can see that the two inequalities become equalities as $t = 1$ and $P_X = Q_X$ if we choose $\ell(x) = \ell_x$. Conversely, using the positive definiteness of the relative entropy, we see that these conditions are in fact necessary to achieve equality. \square

2.2.3 Shannon code

The above result allows us to show that certain codes have optimal expected codeword lengths. For example for the source X that outputs symbols ‘a’, ‘b’ and ‘c’ with probabilities $\frac{1}{2}$, $\frac{1}{4}$ and $\frac{1}{4}$, respectively, the code

$$C(\text{‘a’}) = 0, \quad C(\text{‘b’}) = 10, \quad C(\text{‘c’}) = 11 \quad (2.13)$$

satisfies $\bar{\ell}(C) = \frac{1}{2} + 2 \cdot \frac{1}{4} \cdot 2 = \frac{3}{2}$ and $H(X) = \frac{1}{2} \log 2 + 2 \cdot \frac{1}{4} \log 4 = \frac{3}{2}$, and thus, we know that it is optimal thanks to Proposition 2.8.

The above example is constructed in such a way that all the probabilities are negative powers of 2 and the codeword lengths satisfy the Kraft inequality with equality, and in this particular case it is easy to see from the proof of Proposition 2.8 that $\bar{\ell}(C) = H(X)$. If the probabilities do not have this form the same construction does not generally work. However, we can show the following.

Proposition 2.9. *For a discrete source X with distribution P_X there always exists a prefix code C with $\bar{\ell}(C) \leq H(X) + 1$.*

The code we construct to prove this is called the Shannon code, and it not optimal in general. However, the this bound, together with Proposition 2.8, shows that we lose at most 1 bit per symbol using this code.

Proof. We can choose codeword lengths $\ell(x) = \lceil \log \frac{1}{P_X(x)} \rceil$. These satisfy the Kraft inequality since

$$\sum_x 2^{-\ell(x)} = \sum_x 2^{-\lceil \log \frac{1}{P_X(x)} \rceil} \leq \sum_x 2^{-\log \frac{1}{P_X(x)}} = \sum_x P_X(x) = 1. \quad (2.14)$$

Hence, using Proposition 2.7 we may construct a prefix codes using these lengths. Moreover, for this code we have

$$\bar{\ell}(C) = \sum_x P_X(x) \left\lceil \log \frac{1}{P_X(x)} \right\rceil \quad (2.15)$$

$$\leq \sum_x P_X(x) \left(\log \frac{1}{P_X(x)} + 1 \right) \quad (2.16)$$

$$= H(X) + 1. \quad (2.17)$$

□

2.2.4 Huffman codes

In this section we will construct prefix codes with optimal expected codeword length, so-called Huffman codes. We will first learn how to construct the codes and then use this construction to show optimality. Interestingly, the codes were invented by a student that was in the same situation as you are right now!

In 1951, David Huffman and his MIT information theory classmates were given the choice of a term paper or a final exam. The professor, Robert Fano, assigned a term paper on the problem of finding the most efficient binary code. Huffman, unable to prove any codes were the most efficient, was about to give up and start studying for the final when he hit upon the idea of using a frequency-sorted binary tree and quickly proved this method the most efficient. In doing so, Huffman outdid Fano, who had worked with information theory inventor Claude Shannon to develop a similar code.

The code is constructed using Algorithm 2.1, which step-by-step merges a forest of trivial binary trees into a single binary tree.

The construction is not unique because in each step we can assign the labels ‘0’ and ‘1’ in either way to the two children. Moreover, we are asked to select the two trees with smallest probabilities, but there might be different such pairs, e.g. if we start with a source X with

Input: List of symbols $x \in \mathcal{X}$ with probabilities $p_x = P_X(x)$

Output: Binary tree for the Huffman code

```
% initialise forest
for each  $x \in \mathcal{X}$  do
    Create a tree with a root node labelled by the probability  $p_x$  (and the symbol  $x$ ) and no
    other nodes;
    Add this tree to the forest;
end
% condense forest into a single tree
while number of trees in the forest are larger than 1 do
    select the two trees whose roots have the smallest probabilities, say  $p$  and  $p'$ ;
    join the two trees by adding a new root with probability  $p + p'$  with the two trees as
    children, the edges are labelled with '0' and '1';
end
return last remaining tree in the forest;
```

Algorithm 2.1: Construction of a Huffman code tree.

symbols and probabilities ('a', $\frac{1}{3}$), ('b', $\frac{1}{3}$), ('c', $\frac{1}{6}$), ('d', $\frac{1}{6}$) then both of these codes, C_1 and C_2 , are possible Huffman codes:

$$C_1('a') = 00, \quad C_1('b') = 01, \quad C_1('c') = 10, \quad C_1('d') = 11 \quad (2.18)$$

$$C_2('a') = 0, \quad C_2('b') = 10, \quad C_2('c') = 110, \quad C_2('d') = 111 \quad (2.19)$$

$$(2.20)$$

Question 2.10. *Can you retrace how they are created step-by-step?*

The codeword lengths for both codes are optimal according to the Kraft inequality, that is, we have

$$\sum_x 2^{-\ell(x)} = 4 \cdot 2^{-2} = 1 \quad \text{and} \quad (2.21)$$

$$\sum_x 2^{-\ell(x)} = 2^{-1} + 2^{-2} + 2 \cdot 2^{-3} = 1 \quad (2.22)$$

for C_1 and C_2 , respectively. We can further compute their respective expected codeword lengths. This yields

$$\bar{\ell}(C_1) = 2 \quad (2.23)$$

$$\bar{\ell}(C_2) = \frac{1}{3} \cdot 1 + \frac{1}{3} \cdot 2 + \frac{1}{6} \cdot 3 + \frac{1}{6} \cdot 3 = 2 \quad (2.24)$$

Let us compare this to the entropy $H(X) = 2 \cdot \frac{1}{3} \log 3 + 2 \cdot \frac{1}{6} \log 6 \approx 1.92$. So from the entropy bound alone we cannot deduce that these codes are optimal in terms of the expected codeword length—but they in fact are!

Proposition 2.11. *Given a source X with probability distribution P_X , any code constructed using Algorithm 2.1 achieves the minimal expected codeword length for any prefix code.*

We call such a code with minimal expected length an *optimal prefix code*. The proof relies on the following lemma, which we show first.

Lemma 2.12. *There exists an optimal prefix code with the following property:*

- *The two longest codewords are siblings and their respective source symbols have the two smallest probabilities (this is not always unique).*

Proof. An optimal code always corresponds to a binary tree with no unused leaves—if not we can compress the tree by removing the parent of the unused leaf, reducing the expected codeword length. There is always at least one pair of leaves at maximum depth, and those are thus occupied with codewords. If those codewords would not correspond to the two symbols with smallest probability we could exchange symbols to move them there, a process which clearly cannot increase the expected codeword length. \square

RecursiveHuffman:

Input: Forest f_{in}

Output: Forest f_{out}

if number of trees in $f_{\text{in}} = 1$ **then**

 return $f_{\text{out}} = f_{\text{in}}$;

else

 select the two trees in f_{in} whose roots have the smallest probabilities, say p and p' ;

 create new forest f' from f_{in} by joining the selected trees as in Algorithm 2.1;

 return $f_{\text{out}} = \text{RecursiveHuffman}(f')$;

end

Algorithm 2.2: Recursive formulation of the Huffman algorithm.

Proof of Proposition 2.11. The recursive formulation of the Huffman algorithm in Algorithm 2.2 is useful here. Clearly the algorithm produces an optimal code when $|\mathcal{X}| = 2$ since in this case the optimal code simply assigns the codewords 0 and 1 to the two symbols, and this is exactly what the output of the Huffman algorithm will be.

We will thus prove optimality by induction as follows. Let us, without loss of generality, label elements such that our source symbols have probabilities $p_1 \geq p_2 \geq \dots \geq p_n$ and ordered such that the Huffman algorithm will pick p_{n-1} and p_n as the smallest elements if there are ambiguities. By induction we may assume that the recursive Huffman algorithm provides us with an optimal tree when we call it for $n - 1$ symbols with the probabilities $p_1, \dots, p_2, \dots, p_{n-2}, p_{n-1} + p_n$. We denote the expected codeword length of this optimal tree by $L_{n-1}^*(p_1, \dots, p_2, \dots, p_{n-2}, p_{n-1} + p_n)$.

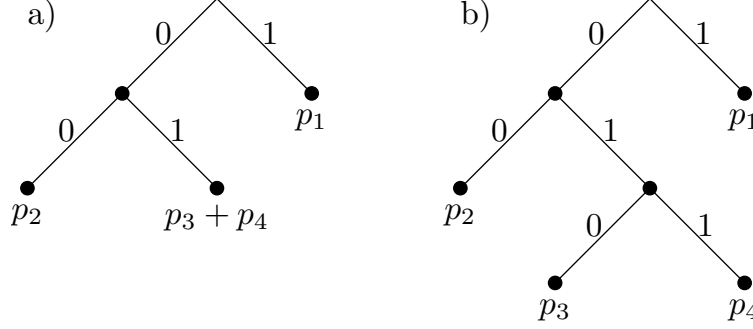


Figure 2.2: **Optimality of Huffman coding, recursive step.** a) Example of a construction with an optimal code tree for $(p_1, p_2, p_3 + p_4)$ with $L_3^* = p_1 + 2p_2 + 2(p_3 + p_4)$. (Note that we must have $p_3 + p_4 \leq p_1$ since otherwise this tree would not be optimal for L_3^* .) b) The Huffman tree for (p_1, p_2, p_3, p_4) with $L_4 = p_1 + 2p_2 + 3p_3 + 3p_4 = L_3^* + p_3 + p_4$.

The Huffman algorithm for n symbols, by definition in its recursive form, produces exactly this tree but with the $(n - 1)$ -th node split into two siblings with probabilities p_{n-1} and p_n (see Figure 2.2 for an example). Its expected codeword length, L_n , thus satisfies

$$L_n = L_{n-1}^*(p_1, \dots, p_2, \dots, p_{n-2}, p_{n-1} + p_n) + p_{n-1} + p_n. \quad (2.25)$$

Note that we added $p_{n-1} + p_n$ as compared to the tree for $n - 1$ symbols those two leaves are now one level deeper, which increases the codeword length by 1 with probability $p_{n-1} + p_n$.

On the other hand, we have

$$L_{n-1}^*(p_1, \dots, p_2, \dots, p_{n-2}, p_{n-1} + p_n) \leq L_n^*(p_1, \dots, p_2, \dots, p_{n-1}, p_n) - p_{n-1} - p_n \quad (2.26)$$

since a valid (although not necessarily optimal) prefix code for the $n - 1$ probabilities can be constructed from an optimal prefix code for n symbols by merging the two leaves at maximum depth with minimal probability (which exist due to Lemma 2.12) into a single leaf. Such a code has expected codeword length $L_n^* - p_{n-1} - p_n$, and thus in particular the optimal length of such a tree for $n - 1$ symbols must satisfy $L_{n-1}^* \leq L_n^* - p_{n-1} - p_n$.

Combining Eqs. (2.25) and (2.26) yields $L_n \leq L_n^*(p_1, \dots, p_2, \dots, p_{n-1}, p_n)$, and since L_n can never be smaller than the optimal codeword length (by definition of the latter), these two quantities must in fact be equal. This proves that the Huffman code construction is optimal. \square

2.3 Fixed-length block codes

Fixed-length codes have the property that all codewords are equally long. If we require error-free compression, then there is not much flexibility: the expected codeword length has to be equal to $\lceil \log |\mathcal{X}| \rceil$ (assuming that every source symbol appears with strictly positive probability). The picture gets dramatically more interesting if we encode a whole block of n

source symbols and only require that the probability of a decoding error vanishes as $n \rightarrow \infty$. A *block code* of length n takes a sequence of n source outputs $x^n \in \mathcal{X}^n$ as input and outputs a binary string $C(x^n)$.

2.3.1 Setup for block coding

As we are observing a long sequence of symbols $X_1, X_2, \dots, X_k, \dots$, one thing we can do is to treat a block of, say n , symbols as a single symbol (with a much larger alphabet of size $|\mathcal{X}|^n$) and then try to find efficient codes for such blocks. Obviously then we can no longer encode and decode instantaneously as we will need to wait for the full block to perform the encoding, and the decoding operation will in turn yield a full block as well.

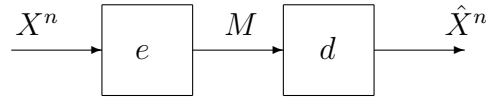


Figure 2.3: Illustration of the fixed-to-fixed length source coding problem.

So for a block code, we observe a sequence of random variables $X^n = (X_1, \dots, X_n)$ from a discrete memoryless source and we would like to compress it into a random variable $M \in \{0, 1\}^L$, the codeword, using an encoder, a function e from X^n to M . Later on, the decoder d will produce an estimate \hat{X}^n of X^n from M . See Fig. 2.3 for an illustration. If we demand that

$$P(\hat{X}^n \neq X^n) = 0 \quad (2.27)$$

then M needs to take on at least

$$|\{x \in \mathcal{X} : P_X(x) > 0\}|^n \quad (2.28)$$

different values, and thus we need to choose $L \geq \lceil n \log |\{x \in \mathcal{X} : P_X(x) > 0\}| \rceil$.

Question 2.13. *Argue why this is sufficient and optimal.*

Comparing this to the bound $L \geq n \lceil \log |\{x \in \mathcal{X} : P_X(x) > 0\}| \rceil$, which we would have arrived at by encoding each source symbol separately using a fixed length code, we see that there is some improvement. However, it is at most n bits. Can we do better if we relax the stringent condition in (2.27) to be such that

$$P(\hat{X}^n \neq X^n) \leq \epsilon \quad (2.29)$$

for any $\epsilon > 0$ positive but arbitrarily small? Let us formalise this.

Definition 2.14 (Block code). *An $(n, 2^L)$ -fixed-length source code (or simply an $(n, 2^L)$ -code) consists of an encoder, e , and a decoder, d , where*

- $e : \mathcal{X}^n \rightarrow \{0, 1\}^L$ and

- $d : \{0, 1\}^L \rightarrow \mathcal{X}^n$

The number n is called the *block length* of the code; L is the length of the codeword; and $R = \frac{L}{n}$ is called the rate of the code. The rate simply evaluates how many bits of codeword this codes uses per symbol of the source to store its output.

Definition 2.15 (Achievable rate). *A rate R is achievable for a DMS \mathbf{X} if there exists a sequence of $(n, 2^{\lfloor nR \rfloor})$ -codes for $n \in \mathbb{N}$ with encoder e_n and decoder d_n such that*

$$\lim_{n \rightarrow \infty} P(\hat{X}^n \neq X^n) = 0 \quad (2.30)$$

where

$$\hat{X}^n = d_n(M_n), \quad \text{and} \quad M_n = e_n(X^n) \quad (2.31)$$

are the reconstructed source and the codeword, respectively.

Question 2.16. *Show that if R is achievable so is any $R' \geq R$.*

Thus, what we really are interested in is the smallest R that is still achievable.

Definition 2.17 (Optimal source coding rate). *The optimal source coding rate for the DMS \mathbf{X} , denoted as $R^*(\mathbf{X})$, is defined to be the infimum of all achievable rates, i.e.,*

$$R^*(\mathbf{X}) = \inf\{R : R \text{ is achievable}\}. \quad (2.32)$$

Finding the optimal source coding rate looks like a formidable problem to solve at first sight. Note that the notion of achievable rate is asymptotic, namely we need to consider a sequence of codes, a code for each $n \in \mathbb{N}$, so that it is not even obvious that $R^*(\mathbf{X})$ is computable in finite time from a complexity-theoretic perspective. However, in his seminal work Shannon [6] showed that $R^*(X)$ has a simple form for a DMS.

Theorem 2.18 (Fixed-length data compression). *For any DMS \mathbf{X} with pmf P_X , we have*

$$R^*(\mathbf{X}) = H(X) \quad (2.33)$$

To prove that $R^*(\mathbf{X}) = H(X)$, we must prove the *achievability part*, $R^*(\mathbf{X}) \leq H(X)$, and the *converse part*, $R^*(\mathbf{X}) \geq H(X)$.

- Achievability means that, for every $R > H(X)$, we must exhibit a sequence of $(n, 2^{\lfloor nR \rfloor})$ -codes such that the error probability vanishes as $n \rightarrow \infty$.
- The converse implies that we cannot do better than this, i.e., there is no sequence of $(n, 2^{\lfloor nR \rfloor})$ -codes where $R < H(X)$ such that we have a vanishing error probability.

In most problems in information theory the two proofs (achievability and converse) use quite different techniques and we thus treat them separately.

2.3.2 Proof of converse and Fano's inequality

For the converse we will use Fano's inequality, which is a fundamental tool in the analysis of information processing tasks. We formulate it here as a statement about conditional entropies of strongly correlated random variables.

Lemma 2.19 (Fano's inequality). *Let X, Y be random variables with joint pmf P_{XY} and let $\epsilon := P[X \neq Y]$. Then*

$$H(X|Y)_P \leq h(\epsilon) + \epsilon \log(|X| - 1) \leq 1 + \epsilon \log |X|, \quad (2.34)$$

where $h(\epsilon) = -\epsilon \log \epsilon - (1 - \epsilon) \log(1 - \epsilon)$ is the binary entropy.

This essentially tells us that if the probability that the two random variables differ is small, then so is the conditional entropy. We should expect this, since the conditional entropy measures the uncertainty of X given side information Y : if ϵ is small then knowing the value y of Y allows us to guess that $X = y$ as well, which is correct with high probability.

Proof. First, we recall that $H(X|Y)_P = \sum_y P_Y(y) H(X|Y = y)$. We will bound these terms individually first. Define $\epsilon_y = 1 - P_{X=y|Y=y}$ so that $\sum_y P_Y(y) \epsilon_y = \epsilon$. Then, we find

$$H(X|Y = y) = - \sum_x P_{X|Y=y}(x) \log P_{X|Y=y}(x) \quad (2.35)$$

$$= -(1 - \epsilon_y) \log(1 - \epsilon_y) - \sum_{x \neq y} P_{X|Y=y}(x) \log P_{X|Y=y}(x) \quad (2.36)$$

$$= -(1 - \epsilon_y) \log(1 - \epsilon_y) - \epsilon_y \log \epsilon_y - \epsilon_y \sum_{x \neq y} \frac{P_{X|Y=y}(x)}{\epsilon_y} \log \frac{P_{X|Y=y}(x)}{\epsilon_y} \quad (2.37)$$

$$= h(\epsilon_y) - \epsilon_y \sum_{x \neq y} \frac{P_{X|Y=y}(x)}{\epsilon_y} \log \frac{P_{X|Y=y}(x)}{\epsilon_y} \quad (2.38)$$

$$\leq h(\epsilon_y) + \epsilon_y \log(|X| - 1) \quad (2.39)$$

In the last step we simply used that the entropy is upper bounded by the logarithm of the support as shown in Chapter 1.

It remains to take an average of the above bound. Using (once again) concavity of the entropy, we find

$$\sum_y P_Y(y) H(X|Y = y) \leq \sum_y P_Y(y) h(\epsilon_y) + \epsilon_y \log |X| \quad (2.40)$$

$$\leq h\left(\sum_y P_Y(y) \epsilon_y\right) + \epsilon \log |X| \quad (2.41)$$

$$= h(\epsilon) + \epsilon \log |X|. \quad (2.42)$$

Finally, we can bound $h(\epsilon) \leq \log 2 = 1$ and $|X| - 1 \leq |X|$ to make the bound a bit simpler but still sufficient for most purposes. \square

Proof of converse of Theorem 2.18. Consider now any sequence of $(n, 2^{\lfloor nR \rfloor})$ -codes with encoders e_n and decoders d_n that satisfy $\epsilon_n \rightarrow 0$ as $n \rightarrow \infty$, where

$$\epsilon_n := P(\hat{X}^n \neq X^n) \quad (2.43)$$

with $\hat{X}^n = d_n(e_n(X^n))$. Fano's inequality applied to the estimation of X^n yields

$$H(X^n | \hat{X}^n) \leq \epsilon_n n \log |\mathcal{X}| + 1 \quad (2.44)$$

Since $H(X^n | M) \leq H(X^n | \hat{X}^n)$ by the data-processing inequality for the conditional entropy (see Corollary 1.25) applied to the decoder d_n , this can be relaxed to

$$H(X^n | M) \leq \epsilon_n n \log |\mathcal{X}| + 1 \quad (2.45)$$

Furthermore, using the dimension bound for $|M| = 2^{\lfloor nR \rfloor} \leq 2^{nR}$, and the definition of mutual information, we find

$$nR \geq H(M) \quad (2.46)$$

$$= I(X^n : M) + H(M | X^n) \quad (2.47)$$

$$= I(X^n : M) \quad (2.48)$$

$$= nH(X) - H(X^n | M) \quad (2.49)$$

We can now apply Eq. (2.45) to get

$$R \geq H(X) - \epsilon_n \log |\mathcal{X}| - \frac{1}{n}. \quad (2.50)$$

Since this inequality must hold for large n and $\epsilon_n \rightarrow 0$ as $n \rightarrow \infty$, we thus must have that

$$R \geq H(X). \quad (2.51)$$

Moreover, since this holds for any sequence of codes, we conclude that $R^*(\mathbf{X}) \geq H(X)$. \square

2.3.3 Proof of achievability and typical sets

The main idea in the proof of achievability is to only encode sequences of source outputs that are “typical” in a sense we will make precise below. The sets of typical sequences are chosen in such a way that two crucial properties hold:

- There are not too many such sequences so that we can encode them efficiently.
- We can safely ignore all the sequences that are not typical since they are guaranteed to only occur with very low probability.

Let us now make this more formal.

Definition 2.20. *The typical set of a DMS \mathbf{X} is defined as*

$$A_\epsilon^{(n)}(\mathbf{X}) := \left\{ x^n \in \mathcal{X}^n : \left| \frac{1}{n} \log \frac{1}{P_{X^n}(x^n)} - H(X) \right| \leq \epsilon \right\} \quad (2.52)$$

where

$$P_{X^n}(x^n) = P(X^n = x^n) = \prod_{i=1}^n P_X(x_i), \quad \forall x^n \in \mathcal{X}^n. \quad (2.53)$$

We call the elements of $A_\epsilon^{(n)}(\mathbf{X})$ *typical sequences* for the source \mathbf{X} . In words, this means that the typical sequences are those whose average log-likelihood (per symbol) is very close to the entropy of X , the i.i.d. output of the source. Note that

$$H(X) = \frac{1}{n} H(X_1 X_2 \dots X_n) \quad (2.54)$$

due to the additivity of entropy for product distributions, and thus we can alternatively interpret $H(X)$ as the entropy the source creates per symbol. (The latter interpretation is especially useful when we want to generalise the concept of typical sets beyond memoryless sources.)

The properties of the typical set mentioned above can now be formalised as follows:

Proposition 2.21 (Asymptotic equipartition property). *Let $\epsilon > 0$. The sequence of typical sets $A_\epsilon^{(n)}(\mathbf{X})$ for $n \in \mathbb{N}$ has the following properties:*

1. $H(X) - \epsilon \leq \frac{1}{n} \log \frac{1}{P_{X^n}(x^n)} \leq H(X) + \epsilon$ for all sequences $x^n \in A_\epsilon^{(n)}(\mathbf{X})$ and $n \in \mathbb{N}$.
2. $\lim_{n \rightarrow \infty} P[X^n \in A_\epsilon^{(n)}(\mathbf{X})] = 1$.
3. For all $n \in \mathbb{N}$, the size of the set satisfies

$$|A_\epsilon^{(n)}(\mathbf{X})| \leq 2^{n(H(X) + \epsilon)}. \quad (2.55)$$

The name *asymptotic equipartition property* alludes to the the first (and defining) property of the typical set, which ensures that all sequences in the set are approximately equally likely. More precisely, the definition implies that we have

$$\left| \frac{P_{X^n}(x^n)}{P_{X^n}(\tilde{x}^n)} \right| \leq \exp(2n\epsilon) \quad (2.56)$$

for all typical sequences $x^n, \tilde{x}^n \in A_\epsilon^{(n)}(\mathbf{X})$.

Proof. The first property is immediate from the definition (and holds for all n).

The second property follows from the weak law of large numbers. To see this, we consider the random variables X_i produced by the source and the new random variables $Z_i = \log \frac{1}{P_X(X_i)} - H(X)$.

Question 2.22. Verify that Z_i has zero mean and that the sequence (Z_1, Z_2, \dots, Z_n) is i.i.d..

We would now like to express the probability $P[X^n \in A_\epsilon^{(n)}(\mathbf{X})]$ in terms of the random variables Z_i we just introduced. We find the following sequence of equalities:

$$P[X^n \in A_\epsilon^{(n)}(\mathbf{X})] = P\left[\left|\frac{1}{n} \log \frac{1}{P_{X^n}(X^n)} - H(X)\right| \leq \epsilon\right] \quad (2.57)$$

$$= P\left[\left|\frac{1}{n} \log \frac{1}{\prod_{i=1}^n P_X(X_i)} - H(X)\right| \leq \epsilon\right] \quad (2.58)$$

$$= P\left[\left|\frac{1}{n} \sum_{i=1}^n \log \frac{1}{P_X(X_i)} - H(X)\right| \leq \epsilon\right] \quad (2.59)$$

$$= P\left[\left|\frac{1}{n} \sum_{i=1}^n Z_i\right| \leq \epsilon\right] \quad (2.60)$$

$$= 1 - P\left[\left|\frac{1}{n} \sum_{i=1}^n Z_i\right| > \epsilon\right] \quad (2.61)$$

Now since Z_i are i.i.d. and zero mean we can apply the weak law of large numbers (Proposition 0.22), which ensures that

$$\lim_{n \rightarrow \infty} P\left[\left|\frac{1}{n} \sum_{i=1}^n Z_i\right| > \epsilon\right] = 0, \quad (2.62)$$

and so, in particular, we have $\lim_{n \rightarrow \infty} P[X^n \in A_\epsilon^{(n)}(\mathbf{X})] = 1$.

The third property follows by a basic counting argument. Since every sequence in $A_\epsilon^{(n)}(\mathbf{X})$ has probability at least $2^{-n(H(X)+\epsilon)}$ by definition, there can be at most $2^{n(H(X)+\epsilon)}$ such sequences in the typical set as otherwise the total probability of all sequences would exceed 1. \square

Proof of achievability of Theorem 2.18. We fix $\epsilon > 0$ and construct encoders and decoders for each blocklength n . The main idea is to do a faithful encoding of all the sequences x^n in the typical set and essentially ignore and accept errors for sequences that are not typical.

To do this we index the elements of $A_\epsilon^{(n)}(X)$ in some order (say lexicographic). This simply means that to each sequence $x^n \in A_\epsilon^{(n)}(X)$, we assign a unique index $m(x^n) \in \{0, 1\}^{L_n}$. By the upper bound on the size of the typical set, we know that the codeword length can be bounded as $L_n \leq \lceil n(H(X) + \epsilon) \rceil \leq n(H(X) + \epsilon) + 1$. This assignment is known to both the encoder and the decoder.

We now design the encoder and decoder for blocklength n .

Encoder e_n : If the realised sequence is typical, i.e. $X^n \in A_\epsilon^{(n)}(X)$, then output the index $M = m(X^n)$. Otherwise set $M = 0^{L_n}$. In other words, the precise working of the

encoder is

$$e_n(x^n) = \begin{cases} m(x^n) & x^n \in A_\epsilon^{(n)}(X) \\ 0^{L_n} & x^n \notin A_\epsilon^{(n)}(X) \end{cases} \quad (2.63)$$

Decoder: Given m , output x^n such that $m = m(x^n)$. (This choice is unique.) In other words, look in the table for the sequence that corresponds to the index m .

We can now compute the probability of error for this encoder and decoder. Suppose first that the realised sequence is typical. We will never make an error because the sequence that is output coincides with the emitted sequence of the DMS, by construction of the encoder and decoder. Thus, we can only make an error if the emitted source sequence is atypical. Hence,

$$\epsilon_n = P(\hat{X}^n \neq X^n) \leq P(X^n \notin A_\epsilon^{(n)}(X)). \quad (2.64)$$

and in particular $\lim_{n \rightarrow \infty} \epsilon_n = 0$. This implies that the sequence of codeword lengths L_n is achievable. Moreover, since $\frac{L_n}{n} \leq H(X) + \epsilon + \frac{1}{n} \leq H(X) + 2\epsilon$ for large enough n , we can conclude that the rate $H(X) + 2\epsilon$ is achievable. Since $\epsilon > 0$ is arbitrarily small, we have

$$R^*(\mathbf{X}) = \inf\{R : R \text{ is achievable}\} \leq H(X). \quad (2.65)$$

□

2.3.4 Strong converse via typical sets

We can also argue with typical sets to get a stronger statement for our converse bound.

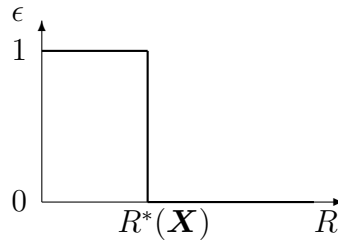


Figure 2.4: Illustration of the strong converse property. For any rate sequence of $(n, 2^{nR})$ codes with rate $R < R^*(\mathbf{X}) = H(X)$, the asymptotic error $\lim_{n \rightarrow \infty} P(\hat{X}^n \neq X^n)$ necessarily converges to 1.

In fact, the above converse proof is quite conservative. Even if we allow that

$$\limsup_{n \rightarrow \infty} P(\hat{X}^n \neq X^n) \leq \epsilon \quad (2.66)$$

for any $0 \leq \epsilon < 1$, it turns out that the ϵ -optimal rate must be no smaller than $H(X)$. We state this formally as follows:

Theorem 2.23. *For any sequence of $(n, 2^{nR})$ -codes it holds that if $R < H(X)$, then necessarily $P(\hat{X}^n \neq X^n) \rightarrow 1$ as $n \rightarrow \infty$.*

This theorem removes all hope to devise a more efficient source coding scheme that can beat the compression rate $H(X)$ by allowing some small error.

One way to prove this statement is to expand our characterisation of the typical set.

Proposition 2.24 (Asymptotic equipartition property, continued). *Let $\epsilon, \mu > 0$. Then there exists an N_0 such that for all $n \geq N_0$, the following holds:*

1. $P[X^n \in A_\epsilon^{(n)}(\mathbf{X})] \geq 1 - \mu.$

2. *The size of the set satisfies*

$$|A_\epsilon^{(n)}(X)| \geq (1 - \mu)2^{n(H(X) - \epsilon)}. \quad (2.67)$$

Proof. The first statement is a simple consequence of the second property in Proposition 2.21, i.e. since $\lim_{n \rightarrow \infty} P[X^n \in A_\epsilon^{(n)}(\mathbf{X})] = 1$ there must exist an N_0 with the desired property by definition of the limit. The second statement now again follows by a counting argument. Since we know that every sequence $x^n \in A_\epsilon^{(n)}(X)$ satisfies $P_{X^n}(x^n) \leq 2^{-n(H(X) - \epsilon)}$ we will need at least $(1 - \mu)2^{n(H(X) - \epsilon)}$ such elements to reach a total probability of $1 - \mu$ as stipulated by the first property above. \square

This allows us to lay out the idea for a proof of Theorem 2.23. Let us assume that $R < H(X)$, and define $\epsilon = \frac{1}{2}(H(X) - R)$. Then using the bound in Eq. (2.67), we find

$$2^{nR} = 2^{n(H(X) - 2\epsilon)} \leq \frac{2^{-n\epsilon}}{1 - \mu} |A_\epsilon^{(n)}(X)|. \quad (2.68)$$

For sufficiently large n this implies that $|M| < |A_\epsilon^{(n)}(X)|$, and in fact $|M|$ is smaller by a factor that grows exponentially in n . This implies that we can only faithfully represent a smaller and smaller fraction of all typical source sequences in M . Moreover, since all typical sequences are almost equiprobable this induces an error approaching 1 exponentially fast.

We finally give an alternative and more formal proof of Theorem 2.23 that uses the structure of encoder and decoder more explicitly.

Proof of Theorem 2.23. We assume $R < H(X)$ and try to give a lower bound on the probability of error.

Fix an arbitrary encoder $e_n : \mathcal{X}^n \rightarrow \{1, 2, \dots, 2^{nR}\}$. This encoder induces a partition of the space \mathcal{X}^n into disjoint subsets $\mathcal{D}_m \subset \mathcal{X}^n, m \in \{1, 2, \dots, 2^{nR}\}$ defined as $\mathcal{D}_m = \{x^n : f(x^n) = m\}$. The best decoder (the one that minimizes the error probability) is the one that uses the following rule given message m :

$$d_n(m) = \operatorname{argmax}_{x^n \in \mathcal{D}_m} P_{X^n}(x^n). \quad (2.69)$$

(This is known as the maximum likelihood decoder as the decoder maximizes the likelihood of the data.) We also denote the resulting random variable by $\hat{X}^n = d_n(e_n(X^n))$. Fix now $\epsilon > 0$ small enough such that $R < H(X) - 2\epsilon$ and also fix $\mu > 0$.

Question 2.25. *Argue why such an $\epsilon > 0$ always exists.*

The error probability can be bounded as

$$\epsilon_n = 1 - P(\hat{X}^n = X^n) \quad (2.70)$$

$$= 1 - \sum_{m=1}^{2^{nR}} P_{X^n}(d_n(m)) \quad (2.71)$$

$$= 1 - \sum_{m: d_n(m) \in \mathcal{A}_\epsilon^{(n)}(X)} P_{X^n}(d_n(m)) - \sum_{m: d_n(m) \notin \mathcal{A}_\epsilon^{(n)}(X)} P_{X^n}(d_n(m)) \quad (2.72)$$

$$\stackrel{(a)}{\geq} 1 - \sum_{m: d_n(m) \in \mathcal{A}_\epsilon^{(n)}(X)} P_{X^n}(d_n(m)) - \mu \quad (2.73)$$

$$\stackrel{(b)}{\geq} 1 - \sum_{m: d_n(m) \in \mathcal{A}_\epsilon^{(n)}(X)} 2^{-n(H(X)-\epsilon)} - \mu \quad (2.74)$$

$$\geq 1 - \sum_{m=1}^{2^{nR}} 2^{-n(H(X)-\epsilon)} - \mu \quad (2.75)$$

$$\stackrel{(c)}{=} 1 - 2^{-n(H(X)-R-\epsilon)} - \mu \quad (2.76)$$

where (a) follows because the probability of the atypical set is smaller than μ for sufficiently large n according to Proposition 2.24, (b) follows since the probability of sequences in the typical set is at most $2^{-n(H(X)-\epsilon)}$ and (c) because the number of messages is 2^{nR} . Hence, since $R < H(X) - 2\epsilon$, we find

$$\liminf_{n \rightarrow \infty} \epsilon_n \geq \liminf_{n \rightarrow \infty} 1 - \mu - 2^{-n\epsilon} = 1 - \mu. \quad (2.77)$$

Since μ can be arbitrarily small, we in fact have $\lim_{n \rightarrow \infty} P_e^{(n)} = 1$, concluding the proof. \square

Chapter 3

Cryptography: randomness extraction

[Week 6]

Intended learning outcomes:

- You can compute and use guessing probability and min-entropy.
- You can construct a randomness extractor using a family of hash functions.
- You understand that deterministic functions cannot increase entropy.

3.1 Problem setup

One of the most prominent concepts in cryptography is randomness, and it lies at the core of information-theoretic security. To understand, for example, whether a given bit string is random, we do not want to look at a particular instance of the string (although that is interesting as well and leads ultimately to the notion of algorithmic randomness) but instead want to check that the process that created the bit string selected it at random. Similarly and maybe even more evidently, the concept of a *secret* bit string cannot be defined unless we look at the process by which a random variable is produced. If the random process is such that the bit string is independent of any side information held by an eavesdropper, then secrecy (relative to that eavesdropper) can be claimed.

In the following we say that a random variable Z on an alphabet \mathcal{Z} is close to uniformly random if its pmf is close to a uniform pmf in total variation distance, i.e. if

$$\delta_{\text{tvd}}(P_Z, U_Z) := \frac{1}{2} \|P_Z - U_Z\|_1 = \frac{1}{2} \sum_{z \in \mathcal{Z}} |P_Z(z) - U_Z(z)| \quad (3.1)$$

is small, where U_Z denotes the uniform distribution on \mathcal{Z} . In the next chapter we learn more about the total variation distance and its use in statistics.

Question 3.1. *Verify that the total variation distance is a metric: is it symmetric, is it positive and zero only if the two distributions are equal, and does it satisfy the triangle inequality?*

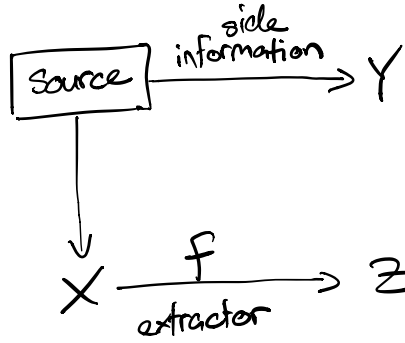


Figure 3.1: The setup of randomness extraction. A source produces random variables X and Y . An extractor f is used to create a new random variable Z that is (close to) uniform and independent of Y .

The total variation distance satisfies the data-processing inequality, i.e. for any channel $W_{Y|X}$ and any two pmfs P_X and Q_X , we have

$$\delta_{\text{td}}(P_X, Q_X) \geq \delta_{\text{td}}(P_Y, Q_Y), \quad (3.2)$$

where $P_Y(y) = \sum_x P_X(x)W_{Y|X}(y|x)$ and $Q_Y(y) = \sum_x Q_X(x)W_{Y|X}(y|x)$. This can be understood as saying that after we apply a channel $W_{Y|X}$, that is, introduce some noise, the output distributions are generally closer than the input distributions. So in a sense the two distributions have become more difficult to distinguish after applying the channel. We will verify this property in the homework.

We can now extend the definition of uniformity to the case where some side information Y on Z is available, and we want to make sure that the randomness is in fact not only uniform but also independent of the side information. This leads us to the following more general definition.

Definition 3.2. Let P_{ZY} be a joint pmf for two random variables Z on \mathcal{Z} and Y on \mathcal{Y} . For any $\epsilon \in (0, 1)$, we say that Z is ϵ -uniformly random and independent of Y if

$$\delta_{\text{td}}(P_{ZY}, U_Z \times P_Y) \leq \epsilon. \quad (3.3)$$

We will now consider the task of randomness extraction, namely the task of creating approximately uniform and independent random variables from a random source X that is generally neither uniform nor independent of Y . In cryptography the i.i.d. assumption (as appears for example in memoryless sources) is often not very natural since we often cannot guarantee that a random source is exactly memoryless and thus we want to ensure that our extraction scheme works even if we do not make any assumptions on the structure of the source. See Figure ?? for a schematic.

This is generally difficult: one thing we can immediately notice is that if one output of the source is very likely, for example if it appears with probability 0.5, then we can produce

exactly one bit of perfect randomness from this source (the new uniform random variable would be the indicator function for this event, which takes the value 0 and 1 with probability 0.5 each.), and this is in fact the best we can hope for. The maximal probability over any output of the source thus appears prominently in the analysis of randomness extraction, even in the approximate case, and we will introduce it formally in the next section in terms of guessing probability and min-entropy.

Let us now formally define a randomness extractor for a fixed source, which takes X and produces a bit string Z that is uniformly random and independent of Y .

Definition 3.3. An (ϵ, ℓ) -extractor for a source X with side information Y governed by a pmf P_{XY} is a function $f : \mathcal{X} \rightarrow \{0, 1\}^\ell$ such that

$$\delta_{\text{td}}(P_{ZY}, U_Z \times P_Y) \leq \epsilon \quad \text{where} \quad Z = f(X) \quad (3.4)$$

and thus P_{ZY} is the distribution induced by f , i.e. $P_{ZY}(z, y) = \sum_{x: f(x)=z} P_{XY}(x, y)$.

Question 3.4. Why should we not allow random functions/channels as extractors here?

We may now ask for the maximum length ℓ of such a approximately uniform and independent string of bits. For this purpose we define

$$\ell_\epsilon^*(X|Y)_P := \max \{ \ell \in \mathbb{N} : \exists \text{ an } (\epsilon, \ell)\text{-extractor for } P_{XY} \}. \quad (3.5)$$

We will now find bounds on this quantity from above and below that hold for arbitrary distributions P_{XY} . These bounds will be in terms of the smooth min-entropy of the source, which we will introduce in the next section. In the homework we will also consider the special case where these sources are memoryless.

3.2 Guessing probability and min-entropy

In this section we consider a source with side information given by a joint pmf P_{XY} on two random variables X on \mathcal{X} and Y on \mathcal{Y} . We characterise our source using the concepts of guessing probability and min-entropy. The guessing probability of X given Y is the probability that an observer with access to Y can correctly guess the value of X . It is not difficult to find the optimal strategy for this task: given a sample $y \in \mathcal{Y}$, the observer will simply choose its guess as

$$\hat{x} = \operatorname{argmax}_{x \in \mathcal{X}} P_{X|Y}(x|y). \quad (3.6)$$

The average probability of guessing the correct value of X is thus given by the guessing probability as defined in the following.

Definition 3.5. Let P_{XY} be a joint pmf as above. The guessing probability of X conditioned on Y is defined as

$$p_{\text{guess}}(X|Y)_P := \sum_{y \in \mathcal{Y}} P_Y(y) \max_{x \in \mathcal{X}} P_{X|Y}(x|y). \quad (3.7)$$

Moreover, the conditional min-entropy of X conditioned on Y is defined as

$$H_{\min}(X|Y)_P := -\log p_{\text{guess}}(X|Y)_P. \quad (3.8)$$

The min-entropy belongs to a class of Rényi entropies that have found widespread use in information theory, and we will explore that connection in the homework. For now let us just point out that it is always smaller than the Shannon entropy.

Lemma 3.6. *For any joint pmf P_{XY} , we have $H_{\min}(X|Y) \leq H(X|Y)$.*

Proof. To see this, we use Jensen's inequality on the convex function $t \mapsto -\log t$ to find

$$H_{\min}(X|Y) = -\log \left(\sum_{y \in \mathcal{Y}} P_Y(y) \max_{x \in \mathcal{X}} P_{X|Y}(x|y) \right) \quad (3.9)$$

$$\leq \sum_{y \in \mathcal{Y}} P_Y(y) \min_{x \in \mathcal{X}} \left(-\log P_{X|Y}(x|y) \right) \quad (3.10)$$

$$\leq \sum_{y \in \mathcal{Y}} P_Y(y) \sum_{x \in \mathcal{X}} P_{X|Y}(x|y) \left(-\log P_{X|Y}(x|y) \right) = H(X|Y), \quad (3.11)$$

where for the second inequality we used the fact that the minimum over x is smaller than the expectation over x under the distribution $P_{X|Y}(x|y)$. \square

We will state our results using a variation of the min-entropy, the *smooth min-entropy*, which is the maximum of the min-entropy over a set of distributions that are close to P_{XY} in total variation distance.

Definition 3.7. *Let P_{XY} a joint pmf and $\epsilon \in [0, 1)$. We define the smooth min-entropy of X conditioned on Y as*

$$H_{\min}^{\epsilon}(X|Y)_P := \max \left\{ H_{\min}(X|Y)_{\tilde{P}} : \delta_{\text{td}}(\tilde{P}_{XY}, P_{XY}) \leq \epsilon \right\}. \quad (3.12)$$

3.3 Achievability via two-universal hash functions

There are several ways to construct extractors, including using the property of typical sets that all its elements are approximately equally likely. Here we follow a different approach (which is quite standard in cryptography) and use a random construction based on hash functions. In particular, we consider a family of two-universal hash functions $\{f_s\}_{s \in \mathcal{S}}$ where $f_s : \mathcal{X} \rightarrow \{0, 1\}^{\ell}$. They are parametrised by a seed s and have the property that

$$\sum_{s \in \mathcal{S}} \frac{1}{|\mathcal{S}|} 1\{f_s(x) = f_s(x')\} = \frac{1}{2^{\ell}} \quad \forall x \neq x'. \quad (3.13)$$

Such families of hash functions exist if we choose \mathcal{S} large enough.

Question 3.8. *Can you come up with such a family?*

Let us know apply a function f_S from a two-universal family of hash functions for $S \in \mathcal{S}$ chosen uniformly at random to get an output $Z = f_S(X)$.

Theorem 3.9. *Let $\{f_s\}_s$ be a two-universal family of hash functions. Using the notation introduced above, we have*

$$\sum_{s \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \delta_{\text{td}}(P_{ZY}^s, U_Z \times P_Y) \leq \frac{1}{2} \sqrt{2^{\log \ell - H_{\min}(X|Y)_P}}, \quad (3.14)$$

where $P_{ZY}^s(z, y) = \sum_{x \in \mathcal{X}: f_s(x)=z} P_{XY}(x, y)$.

Proof. Without loss of generality we can assume that the marginal P_Y has full support as otherwise we can just remove unused symbols. Using the Cauchy-Schwarz inequality in Lemma 0.24, we can bound

$$2\delta_{\text{td}}(P_{ZY}^s, U_Z \times P_Y) = \|P_{ZY}^s - U_Z \times P_Y\|_1 \quad (3.15)$$

$$= \left\| \left(1_Z \times P_Y^{\frac{1}{2}}\right) \cdot \left(1_Z \times P_Y^{-\frac{1}{2}}\right) \cdot (P_{ZY}^s - U_Z \times P_Y) \right\|_1 \quad (3.16)$$

$$\leq \left\| 1_Z \times P_Y^{\frac{1}{2}} \right\|_2 \left\| \left(1_Z \times P_Y^{-\frac{1}{2}}\right) \cdot (P_{ZY}^s - U_Z \times P_Y) \right\|_2 \quad (3.17)$$

$$= \sqrt{\sum_{z,y} P_Y(y)} \sqrt{\sum_{z,y} P_Y^{-1}(y) (P_{ZY}^s(z, y) - U_Z(z) P_Y(y))^2} \quad (3.18)$$

$$= \sqrt{2^\ell} \sqrt{\sum_{z,y} P_Y^{-1}(y) P_{ZY}^s(z, y)^2 - 2 \cdot 2^{-\ell} P_{ZY}^s(z, y) + 2^{-2\ell} P_Y(y)} \quad (3.19)$$

$$= \sqrt{\sum_{z,y} 2^\ell P_Y^{-1}(y) P_{ZY}^s(z, y)^2 - 1}, \quad (3.20)$$

where we introduced the vector 1_Z for which $1_Z(z) = 1$ for all $z \in \mathcal{Z}$, and recall that $U_Z(z) = 2^{-\ell}$ is the uniform distribution on \mathcal{Z} .

Using Jensen's inequality, the expectation over the seed S of the above quantity can then be bounded as

$$\sum_{s \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \delta_{\text{td}}(P_{ZY}^s, U_Z \times P_Y) \leq \frac{1}{2} \sqrt{2^\ell \sum_{s \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \sum_{z,y} P_Y^{-1}(y) P_{ZY}^s(z, y)^2 - 1} \quad (3.21)$$

It remains to analyse the term

$$\sum_{s \in \mathcal{S}} \frac{2^\ell}{|\mathcal{S}|} \sum_{z,y} P_Y^{-1}(y) P_{ZY}^s(z, y)^2 \quad (3.22)$$

$$= \sum_{s \in \mathcal{S}} \frac{2^\ell}{|\mathcal{S}|} \sum_y P_Y(y) \sum_{x,x'} 1\{f_s(x) = f_s(x')\} P_{X|Y}(x'|y) P_{X|Y}(x|y) \quad (3.23)$$

$$= \sum_y P_Y(y) \sum_{x \neq x'} P_{X|Y}(x'|y) P_{X|Y}(x|y) + 2^\ell \sum_{x,y} P_Y(y) P_{X|Y}(x|y)^2 \quad (3.24)$$

$$\leq 1 + 2^\ell \sum_y P_Y(y) \max_x P_{X|Y}(x|y) \quad (3.25)$$

$$= 1 + 2^\ell p_{\text{guess}}(X|Y)_P. \quad (3.26)$$

Hence, plugging this into Eq. (3.20), we arrive at the desired bound. \square

We then arrive at the following result.

Theorem 3.10. *Consider $\epsilon \in (0, 1)$ and a source with pmf P_{XY} . As long as*

$$\ell \leq H_{\min}^{\frac{\epsilon}{4}}(X|Y)_P - 2 \log \frac{1}{\epsilon}, \quad (3.27)$$

there exists an (ϵ, ℓ) -extractor for P_{XY} . Or, in other words,

$$\ell_\epsilon^*(X|Y)_P \geq H_{\min}^{\frac{\epsilon}{4}}(X|Y)_P - 2 \log \frac{1}{\epsilon} - 1. \quad (3.28)$$

Proof. Let \tilde{P}_{XY} denote the distribution that achieves the maximum for the smooth min-entropy, i.e. $H_{\min}^{\frac{\epsilon}{4}}(X|Y)_P = H_{\min}(X|Y)_{\tilde{P}}$. Theorem 3.9 applied for the source \tilde{P}_{XY} with the above choice of ℓ yields

$$\sum_{s \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \delta_{\text{tvd}}(\tilde{P}_{ZY}^s, U_Z \times \tilde{P}_Y) \leq \frac{\epsilon}{2}. \quad (3.29)$$

Hence, there is at least one seed value s for which this bound holds, and it remains to show that f_s constitutes an (ϵ, ℓ) -extractor. However, $\delta_{\text{tvd}}(\tilde{P}_{XY}, P_{XY}) \leq \frac{\epsilon}{4}$ implies $\delta_{\text{tvd}}(\tilde{P}_Y, P_Y) \leq \frac{\epsilon}{4}$ and $\delta_{\text{tvd}}(\tilde{P}_{ZY}^s, P_{ZY}^s) \leq \frac{\epsilon}{4}$ by the data-processing inequality. And hence, using the triangle inequality twice we have

$$\delta_{\text{tvd}}(P_{ZY}^s, U_Z \times P_Y) \quad (3.30)$$

$$\leq \delta_{\text{tvd}}(\tilde{P}_{ZY}^s, P_{ZY}^s) + \delta_{\text{tvd}}(\tilde{P}_{ZY}^s, U_Z \times \tilde{P}_Y) + \delta_{\text{tvd}}(U_Z \times \tilde{P}_Y, U_Z \times P_Y) \quad (3.31)$$

$$\leq \epsilon. \quad (3.32)$$

\square

3.4 Converse via an entropy inequality

The converse relies on a generalization of the following lemma, which states that applying a function to a random variable cannot increase the uncertainty about it.

Lemma 3.11. *Let $f : \mathcal{X} \rightarrow \mathcal{Z}$ be a function. Then $H(X) \geq H(f(X))$ and $H_{\min}(X) \geq H_{\min}(f(X))$.*

Proof. Let $Z = f(X)$. The joint distribution $P_{XZ}(x, z) = P_X(x) 1\{f(x) = z\}$ satisfies

$$H(XZ) = \sum_{x,z} P_{XZ}(x, z) \log \frac{1}{P_{XZ}(x, z)} = \sum_x P_X(x) \log \frac{1}{P_X(x)} = H(X). \quad (3.33)$$

Hence, we can conclude that $H(X) = H(XZ) = H(Z) + H(X|Z) \geq H(Z)$.

The proof for the min-entropy cannot rely on the chain rule but by inspecting the definition of the respective guessing probabilities,

$$p_{\text{guess}}(X) = \max_x P_X(x) \quad \text{and} \quad p_{\text{guess}}(Z) = \max_z P_Z(z) = \max_z \sum_{x:f(x)=z} P_X(x), \quad (3.34)$$

we see that the second term is always at least as large as the first one, i.e. $p_{\text{guess}}(Z) \geq p_{\text{guess}}(X)$. This coincides with our intuition that the input of a function is at least as hard to guess as its output, since once you guessed the input you can get the output by just applying the function. The relation for the min-entropy then follows. \square

It is really important that in the statement we only allow for deterministic functions, as otherwise the equality in Eq. (3.33) would not hold.

Question 3.12. *Give an example where the inequality is violated by a probabilistic function.*

For our argument we need something similar to the above lemma, but for min-smooth entropy and with side information. Here again we can intuitively argue that it is at least as difficult to guess the input of a function as it is to guess the output (with equality if the function is injective). Formally, we can show the following:

Lemma 3.13. *Let $\epsilon \in [0, 1]$ and $f : \mathcal{X} \rightarrow \mathcal{Z}$ be a function. Then, $H_{\min}^\epsilon(X|Y) \geq H_{\min}^\epsilon(f(X)|Y)$.*

In the proof we will make the assumption that f is surjective. This can be avoided, but since it is not really restrictive we made it here to allow for a streamlined presentation.

Proof. The function f can be interpreted as a channel, $W_{Z|X}(z|x) = \delta_{z,f(x)}$. We can define an inverse channel $\widetilde{W}_{X|ZY}$ of that recovers the distribution P_{XY} by Bayes' rule:

$$\widetilde{W}_{X|ZY}(x|z, y) = \frac{P_{XZ|Y}(x, z|y)}{P_{Z|Y}(z|y)} = \frac{\delta_{z,f(x)} P_{X|Y}(x|y)}{\sum_{x':f(x')=z} P_{X|Y}(x'|y)} \quad (3.35)$$

Question 3.14. *Can you see what goes wrong here if the function is not surjective?*

Since this channel only maps z to values of x with $f(x) = z$ it is in fact a proper right-inverse of $W_{Z|X}$ in the following sense. For any pdf Q_{ZY} on the output we define \tilde{Q}_{ZY} as the distribution resulting from first applying $\tilde{W}_{X|YZ}$ and then $W_{Z|X}$ to Q_{ZY} . We then find that for all z, y ,

$$\tilde{Q}_{ZY}(z, y) = \sum_{x'} W_{Z|X}(z|x') \sum_{z'} \tilde{W}_{X|YZ}(x'|z', y) Q_{ZY}(z', y) \quad (3.36)$$

$$= \frac{\sum_{x', z'} \delta_{z, f(x')} \delta_{z', f(x')} P_{X|Y}(x'|y) Q_{ZY}(z', y)}{\sum_{x': f(x')=z} P_{X|Y}(x'|y)} = Q_{ZY}(z, y). \quad (3.37)$$

Now let us assume that the distribution Q_{ZY} is optimal for the smooth min-entropy $H_{\min}^\epsilon(Z|Y)_P$, i.e. $H_{\min}^\epsilon(Z|Y)_P = H_{\min}(Z|Y)_Q$. We can then construct

$$Q_{XY}(x, y) = \sum_{z'} \tilde{W}_{X|YZ}(x|z', y) Q_{ZY}(z', y). \quad (3.38)$$

Note now that due to Eq. (3.37) the pdf $Q_{ZY}(z, y)$ is recovered by applying the function f on the register X . By the data-processing inequality for the total variation distance we have $\delta_{\text{tvd}}(Q_{XY}, P_{XY}) \leq \delta_{\text{tvd}}(Q_{ZY}, P_{ZY}) \leq \epsilon$. Hence,

$$H_{\min}^\epsilon(X|Y)_P \geq H_{\min}(X|Y)_Q = -\log p_{\text{guess}}(X|Y)_Q. \quad (3.39)$$

Now we simply use Lemma 3.11 to show that

$$p_{\text{guess}}(X|Y)_Q = \sum_y Q_Y(y) p_{\text{guess}}(X)_{Q^y} \quad (3.40)$$

$$\leq \sum_y Q_Y(y) p_{\text{guess}}(Z)_{Q^y} = p_{\text{guess}}(Z|Y)_Q, \quad (3.41)$$

where $Q_X^y(x) = Q_{X|Y}(x|y)$ and $Q_Z^y(z) = Q_{Z|Y}(z|y)$, respectively. Combining this with Eq. (3.39) yields the desired bound:

$$H_{\min}^\epsilon(X|Y)_P \geq H_{\min}(Z|Y)_Q = H_{\min}^\epsilon(Z|Y)_P. \quad (3.42)$$

□

Now we are ready to provide an upper bound on the amount of randomness that can be extracted from a source. It matches the lower bound that we derived using two-universal hash functions, and thus we know that this construction was essentially optimal.

Theorem 3.15. *Consider $\epsilon \in (0, 1)$ and a source with pmf P_{XY} . Then, any (ϵ, ℓ) -extractor for P_{XY} must satisfy*

$$\ell \leq H_{\min}^\epsilon(X|Y)_P \quad (3.43)$$

Or, in other words, we have $\ell_\epsilon^(X|Y)_P \leq H_{\min}^\epsilon(X|Y)_P$.*

Proof. Let us assume there exists a function f that constitutes an (ϵ, ℓ) -extractor. We then necessarily have

$$\delta_{\text{td}}(P_{ZY}, U_Z \times P_Y) \leq \epsilon \quad (3.44)$$

for $Z = f(X)$, and thus

$$H_{\min}^\epsilon(Z|Y)_P \geq H_{\min}(Z|Y)_{U \times P} = H_{\min}(Z)_U = \ell, \quad (3.45)$$

where we simply evaluated the min-entropy for the distribution $U_Z \times P_Y$, which is ϵ -close to the distribution P_{ZY} . Combining this with Lemma 3.13 yields the bound $H_{\min}^\epsilon(X|Y)_P \geq \ell$, and since this holds for any (ϵ, ℓ) -extractor we have shown the desired statement. \square

Chapter 4

Information theory in statistics: hypothesis testing

[Week 7]

Intended learning outcomes:

- You can compute the minimal error probability in binary hypothesis testing with known priors, and understand its relationship with the total variation distance.
- You can compute the Chernoff exponent.
- You understand the setup of asymmetric binary hypothesis testing and Stein's lemma.

4.1 Binary hypothesis testing

We consider binary hypothesis testing where we try to distinguish between two models of a random process. The random process produces a sequence of random variables $\mathbf{X} = (X_1, X_2, \dots)$ that are independently drawn from some (unknown) probability distribution $Q \in \mathcal{P}(\mathcal{X})$, where we take \mathcal{X} to be any discrete set. Consider the hypothesis test

$$\begin{aligned} H_0 : Q &= P_0 \\ H_1 : Q &= P_1, \end{aligned} \tag{4.1}$$

where $P_0, P_1 \in \mathcal{P}(\mathcal{X})$ are two candidate probability distributions (or models) of the random process. Our goal is to deduce, from the observation of the random sequence \mathbf{X} , which of the two hypothesis is correct. H_0 is usually called the *null-hypothesis* and H_1 the *alternate hypothesis*.

A (deterministic) *test* for the sequence $X^{(n)} = (X_1, X_2, \dots, X_n)$ is a region $\mathcal{A}_n \subset \mathcal{X}^n$. We say that the alternate hypothesis is *accepted* for this test if the observed sequence satisfies

$(x_1, x_2, \dots, x_n) \in \mathcal{A}_n$, and it is *rejected* otherwise. If the alternate hypothesis is rejected the null-hypothesis is maintained. We can then define two kinds of errors:

$$\alpha_n(\mathcal{A}_n) := P_0^n(\mathcal{A}_n) \quad (4.2)$$

$$\beta_n(\mathcal{A}_n) := 1 - P_1^n(\mathcal{A}_n) \quad (4.3)$$

The *error of the first kind* or *type-1 error*, $\alpha_n(\mathcal{A}_n)$, considers the acceptance of the alternate hypothesis even if the null-hypothesis is true. The *error of the second kind* or *type-2 error*, $\beta_n(\mathcal{A}_n)$, considers the rejection of the alternate hypothesis even though it is true.

Question 4.1. *Can you formulate this problem in the general framework of probability theory as covered in Chapter 0? What is a test in this framework?*

Ideally we would like to devise a test such that both of these errors are small, and get smaller as n increases. We could, for example, try to compute the optimal average error (assuming a *uniform prior* on the two distributions):

$$\epsilon_{\text{sym},n}^* := \frac{1}{2} \min_{\mathcal{A}_n \subset \mathcal{X}^n} \left(\alpha_n(\mathcal{A}_n) + \beta_n(\mathcal{A}_n) \right). \quad (4.4)$$

Here uniform prior means that the probability we assign to the two hypotheses prior to observing the random sequence is equal, and thus $\epsilon_{\text{sym},n}^*$ is indeed the probability of making a wrong decision. However, as their names indicates, often these two hypotheses are not treated on the same footing. Indeed, the question can be easily generalised to the case when the prior over the two hypothesis is not uniform.

Example 4.2. *Assume the alternate hypothesis is that a patient is suffering from COVID-19, and the null-hypothesis is that this is not the case. The error of the first kind is then a false positive and the error of the second kind is a false negative. If we devise a test distinguishing these two hypothesis we are probably more tolerant of false positives than false negatives.*

For such cases it is natural to look at an asymmetric setting, and define, for all $\epsilon \in (0, 1)$,

$$\beta_n^*(\epsilon) := \min\{\beta_n(\mathcal{A}_n) : \alpha_n(\mathcal{A}_n) \leq \epsilon\}, \quad (4.5)$$

where \mathcal{A}_n runs through all subsets of \mathcal{X}^n . Asymmetric hypothesis testing also allows us to deal with the situation when we do not know the prior probabilities. In that case the sum (or probabilistic mixture) of the two errors does not make sense and we need to look at the errors independently. We can, however, still ask the question how these two errors trade off against each other. This is done by analysing $\beta_n^*(\epsilon)$, and in particular by looking at its asymptotics for large n .

4.2 Symmetric hypothesis testing

The first natural question is to evaluate $\epsilon_{\text{sym},n}^*$ for $n = 1$. With that in hand, we will then give a bound on the asymptotic error when $n \rightarrow \infty$.

4.2.1 Total variation distance

We will see that $\epsilon_{\text{sym},n}^*$ can be expressed in terms of the *total variation distance* (tvd) between two pmfs P_0 and P_1 , which is given by

$$\delta_{\text{tvd}}(P_0, P_1) := \frac{1}{2} \sum_{x \in \mathcal{X}} |P_0(x) - P_1(x)|. \quad (4.6)$$

The total variation distance vanishes if and only if $P_0 = P_1$ and it reaches its maximum 1 when P_0 and P_1 are orthogonal, that is, when for every $x \in \mathcal{X}$ either $P_0(x) = 0$ or $P_1(x) = 0$. We can alternatively express the TVD using the following variational formulae, which motivate its name.

Lemma 4.3. *For any two pmfs P_0 and P_1 , following relations hold:*

$$\delta_{\text{tvd}}(P_0, P_1) = \max_{\mathcal{A} \subseteq \mathcal{X}} \left(\sum_{x \in \mathcal{A}} P_0(x) - P_1(x) \right) \quad (4.7)$$

$$= \max_{\mathcal{A} \subseteq \mathcal{X}} \left(\sum_{x \in \mathcal{A}} P_1(x) - P_0(x) \right). \quad (4.8)$$

Proof. To see this equivalence, first note that

$$\sum_{x \in \mathcal{X}} P_0(x) - P_1(x) = 0 \quad (4.9)$$

by normalisation, and thus, for any set $\mathcal{A} \subseteq \mathcal{X}$, we have

$$\sum_{x \in \mathcal{A}} P_0(x) - P_1(x) = \sum_{x \in \mathcal{A}^c} P_1(x) - P_0(x). \quad (4.10)$$

Specifically, for the set $\mathcal{A} = \{x \in \mathcal{X} : P_0(x) > P_1(x)\}$ that is optimal for the maximisation in Eq. (4.7), we find

$$\sum_{x \in \mathcal{A}} |P_0(x) - P_1(x)| = \sum_{x \in \mathcal{A}^c} |P_1(x) - P_0(x)| \quad (4.11)$$

and thus

$$\max_{\mathcal{A} \subseteq \mathcal{X}} \left(\sum_{x \in \mathcal{A}} P_0(x) - P_1(x) \right) = \sum_{x \in \mathcal{A}} |P_0(x) - P_1(x)| = \frac{1}{2} \sum_{x \in \mathcal{X}} |P_0(x) - P_1(x)|. \quad (4.12)$$

□

The total variational distance is closely related to the Schatten 1-norm, which is defined for any vectors v_0 and v_1 that do not necessary need to be normalised. It is defined as

$$\|v_0 - v_1\|_1 := \sum_{x \in \mathcal{X}} |v_0(x) - v_1(x)|. \quad (4.13)$$

Hence, in particular, $\delta_{\text{tvd}}(P_0, P_1) = \frac{1}{2} \|P_0 - P_1\|_1$. We can now state the following result for binary hypothesis testing with general priors.

Proposition 4.4. *Let H_0 (with probability P_0) and H_1 (with probability P_1) have prior probabilities p and $1 - p$, respectively. The minimal probability of error for the binary hypothesis test with $n = 1$, denoted $\epsilon_{p,1}^*$, is given by*

$$\epsilon_{p,1}^* = \frac{1}{2} \left(1 - \|pP_0 - (1-p)P_1\|_1 \right). \quad (4.14)$$

In particular, if $p = 1 - p = \frac{1}{2}$, we have $\epsilon_{\text{sym},1}^* = \frac{1}{2} (1 - \delta_{\text{tvd}}(P_0, P_1))$. This gives a clear operational interpretation for the total variation distance, which is a widely used distance measure in statistics. On the one hand, when $P_0 = P_1$ the total variation distance vanishes and the best thing we can do is a random guess. On the other hand, when P_0 and P_1 are orthogonal, then we can distinguish them perfectly and the error vanishes.

Proof. First we observe the following relations:

$$\epsilon_{p,1}^* = \min_{\mathcal{A} \subset \mathcal{X}} (pP_0(\mathcal{A}) + (1-p)P_1(\mathcal{A}^c)) \quad (4.15)$$

$$= p - \max_{\mathcal{A} \subset \mathcal{X}} (pP_0(\mathcal{A}^c) - (1-p)P_1(\mathcal{A}^c)) \quad (4.16)$$

$$= 1 - p - \max_{\mathcal{A} \subset \mathcal{X}} ((1-p)P_1(\mathcal{A}) - pP_0(\mathcal{A})) \quad (4.17)$$

Combining the two relations we get

$$2\epsilon_{p,1}^* = 1 - \max_{\mathcal{A} \subset \mathcal{X}} ((1-p)P_1(\mathcal{A}) - pP_0(\mathcal{A})) - \max_{\mathcal{A} \subset \mathcal{X}} (pP_0(\mathcal{A}^c) - (1-p)P_1(\mathcal{A}^c)) \quad (4.18)$$

At this point we can determine which sets achieve the maximum in these two optimisation. Clearly both are achieved by the set $\mathcal{A} = \{x \in \mathcal{X} : (1-p)P_1(x) \geq pP_0(x)\}$. The above expression then simplifies to

$$2\epsilon_{p,1}^* = 1 - \sum_{x \in \mathcal{A}} (1-p)P_1(x) - pP_0(x) - \sum_{x \in \mathcal{A}^c} pP_0(x) - (1-p)P_1(x) \quad (4.19)$$

$$= 1 - \sum_{x \in \mathcal{A}} |pP_0(x) - (1-p)P_1(x)| - \sum_{x \in \mathcal{A}^c} |pP_0(x) - (1-p)P_1(x)| \quad (4.20)$$

$$= 1 - \sum_{x \in \mathcal{X}} |pP_0(x) - (1-p)P_1(x)| \quad (4.21)$$

$$= 1 - \|pP_0 - (1-p)P_1\|_1. \quad (4.22)$$

Dividing both sides by 2 then yields the desired result. \square

4.2.2 Chernoff exponent

When we look at n i.i.d. copies of the sample, distributed according to P_0^n or P_1^n , respectively, we make the at first sight surprising observation that these two distributions get closer and closer to orthogonal as $n \rightarrow \infty$ (unless $P_0 = P_1$, of course). Or, in other words, the total variation distance between P_0^n and P_1^n converges to 1 as $n \rightarrow \infty$. Or, in other words, the symmetric error $\epsilon_{\text{sym},n}^*$ converges to zero. We are also interested how fast this convergence to zero is. We will show the following bound:

Proposition 4.5. *For any two pmfs P_0 and P_1 , we have*

$$\lim_{n \rightarrow \infty} -\frac{1}{n} \log \epsilon_{\text{sym},n}^* \geq C(P_0, P_1), \quad (4.23)$$

where we introduced the Chernoff distance or Chernoff exponent,

$$C(P_0, P_1) := -\min_{0 \leq \lambda \leq 1} \log \sum_{x \in \mathcal{X}} P_0(x)^\lambda P_1(x)^{1-\lambda}. \quad (4.24)$$

Note that this actually corresponds to an asymptotic upper bound on the probability of error, so it says that there exists a sequence of tests for which the error drops as $2^{-nC(P_0, P_1)}$. It turns out (but we will not show this here) that this is optimal, i.e., that equality holds in Eq. (4.23).

Proof. To show the bound on the error probability, we argue that

$$2\epsilon_{\text{sym},n}^* = \min_{\mathcal{A}_n \subset \mathcal{X}^n} P_0^n(\mathcal{A}_n) + P_1^n(\mathcal{A}_n^c) \quad (4.25)$$

$$= \sum_{x^n \in \mathcal{X}^n} \min\{P_0^n(x^n), P_1^n(x^n)\} \quad (4.26)$$

$$\leq \sum_{x^n \in \mathcal{X}^n} P_0^n(x^n)^\lambda P_1^n(x^n)^{1-\lambda} \quad (4.27)$$

$$= \sum_{x_1 \in \mathcal{X}} \dots \sum_{x_n \in \mathcal{X}} P_0(x_1)^\lambda \dots P_0(x_n)^\lambda P_1(x_1)^{1-\lambda} \dots P_1(x_n)^{1-\lambda} \quad (4.28)$$

$$= \left(\sum_{x \in \mathcal{X}} P_0(x)^\lambda P_1(x)^{1-\lambda} \right)^n \quad (4.29)$$

We now take the logarithm on both sides and divide through n to get

$$\frac{1}{n} \log \epsilon_{\text{sym},n}^* \leq \log \sum_{x \in \mathcal{X}} P_0(x)^\lambda P_1(x)^{1-\lambda} - \frac{1}{n} \quad (4.30)$$

The last term vanishes in the limit $n \rightarrow \infty$, and thus the bound in (4.23) follows by maximising the right-hand side over all choices of $\lambda \in [0, 1]$. \square

Question 4.6. *What changes when we do the same analysis for $\epsilon_{p,n}^*$?*

4.3 Asymmetric hypothesis testing and the information spectrum method

For simplicity we assume that $D(P_0 \| P_1) < \infty$ in the following, as otherwise by definition of the relative entropy there are some $x \in \mathcal{X}$ with $P_0(x) > 0$ but $P_1(x) = 0$, and, as we will see in the homework, it is possible to come up with tests that have $\beta_n^*(\epsilon) = 0$ for large enough n .

Under this assumption, our goal is to show that regardless of the constant upper bound ϵ on the type-I error, the type-II error behaves as

$$\beta_n^*(\epsilon) \approx 2^{-nD(P_0\|P_1)}, \quad (4.31)$$

where the approximation is up to factors that are sub-exponential in n . This means that the optimal exponential rate at which the type-II error approaches zero is determined by the relative entropy (in first order), thus giving the relative entropy $D(P_0\|P_1)$ a clear operational interpretation in statistics. Let us restate this as a theorem:

Theorem 4.7 (Chernoff-Stein Lemma). *For every $\epsilon \in (0, 1)$,*

$$\lim_{n \rightarrow \infty} -\frac{1}{n} \log \beta_n^*(\epsilon) = D(P_0\|P_1) \quad (4.32)$$

For the proof we will use the information spectrum method (see [2] for more information on that technique). Consider a random variable X that takes values on \mathcal{X} and two pmfs $P_0, P_1 \in \mathcal{P}(\mathcal{X})$ as above. Recall that the log-likelihood ratio for the two pmfs is the random variable

$$Z = \log \frac{P_0(X)}{P_1(X)}, \quad (4.33)$$

where X (and thus Z) is distributed according to P_0 . The log-likelihood ratio is an important random variable in the analysis of many different information processing tasks.

Question 4.8. *Verify that the expectation value of Z under P_0 is $D(P_0\|P_1)$.*

We now introduce the following expression:

$$D_s^\epsilon(P_0\|P_1) := \sup \{R \in \mathbb{R} : P_0[Z \leq R] \leq \epsilon\} \quad (4.34)$$

This quantity looks complicated at first sight, but it simply evaluates exactly where (the value R) we need to cut off the pmf for the *log-likelihood ratio*, Z , so that the probability that $Z \leq R$ is at most ϵ . One could also simply see it as an inverse of the cumulative distribution function of Z . (This result can also be interpreted as a consequence of the Neyman-Pearson lemma, which states that all tests optimising the two types of errors are threshold tests for the log-likelihood ratio.)

Lemma 4.9. *Let $n \in \mathbb{N}$, $\epsilon \in (0, 1)$ and $\delta \in (0, 1 - \epsilon)$. The following two inequalities hold:*

$$D_s^\epsilon(P_0^n\|P_1^n) \leq -\log \beta_n^*(\epsilon) \leq D_s^{\epsilon+\delta}(P_0^n\|P_1^n) + \log \frac{1}{\delta}. \quad (4.35)$$

For $n = 1$ this gives bounds on asymmetric hypothesis testing for any two distributions P_0 and P_1 , without using the i.i.d. structure. If one plugs in n -fold i.i.d. distributions instead this recovers the result for general n . This is an example of a *one-shot bound*, a generic bound on an information-theoretic quantity that can then be easily statistically analysed by taking advantage of an i.i.d. or similar structure.

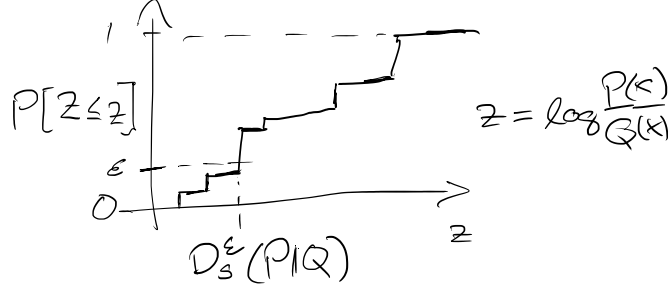


Figure 4.1: Example of information spectrum. The plot shows the cumulative distribution of $Z = \log \frac{P(X)}{Q(X)}$ and the value of $D_s^\epsilon(P\|Q)$.

Proof. To get the lower bound, we construct a test $\mathcal{A}_{R,n} := \{x^n \in \mathcal{X}^n : P_0^n(x^n) \leq 2^R P_1^n(x^n)\}$, where $R \in \mathbb{R}$ still needs to be chosen. These tests, which simply check if the log-likelihood ratio exceeds R , are called Neyman-Pearson tests and are known to be the most powerful tests. We will actually show this optimality here implicitly, as our converse bound will match the achievability we get using this test.

Let us choose $R = D_s^\epsilon(P_0^n\|P_1^n) - \mu$ for some $\mu > 0$ that can be chosen arbitrarily small. The reason we need this small slack $\mu > 0$ is simply that by definition of the supremum in (4.34) this ensures that we have $\alpha_n(\mathcal{A}_{R,n}) = P_0^n(\mathcal{A}_{R,n}) \leq \epsilon$ for any $\mu > 0$, while the same might not necessarily be true at $\mu = 0$. (Recall that the supremum can be taken at the boundary of an open set.) Moreover, we have

$$\beta_n(\mathcal{A}_{R,n}) = P_1^n(\mathcal{A}_{R,n}^c) \quad (4.36)$$

$$= \sum_{x^n \in \mathcal{X}^n} P_1^n(x^n) 1\{P_0^n(x^n) > 2^R P_1^n(x^n)\} \quad (4.37)$$

$$\leq 2^{-R} \sum_{x^n \in \mathcal{X}^n} P_0^n(x^n) 1\{P_0^n(x^n) > 2^R P_1^n(x^n)\} \quad (4.38)$$

$$\leq 2^{-R}. \quad (4.39)$$

This directly implies that $\beta_n^*(\epsilon) \leq 2^{-R}$, or, equivalently,

$$-\log \beta_n^*(\epsilon) \geq D_s^\epsilon(P_0^n\|P_1^n) - \mu. \quad (4.40)$$

Since this holds for all $\mu > 0$ we get the desired inequality.

To get the upper bound, let \mathcal{A}_n be the optimal test for $\beta_n^*(\epsilon)$, i.e. we have $\alpha_n(\mathcal{A}_n) \leq \epsilon$ and $\beta_n(\mathcal{A}_n) = \beta_n^*(\epsilon)$. Recall also the definition of the log-likelihood ratio, $Z = \log \frac{P_0^n(X^n)}{P_1^n(X^n)}$. Using these properties we can establish the following sequence of inequalities:

$$1 - P_0^n(Z \leq R) = P_0^n\left(\log \frac{P_0^n(X^n)}{P_1^n(X^n)} > R\right) \quad (4.41)$$

$$= \sum_{x^n \in \mathcal{X}^n} P_0^n(x^n) 1\{P_0^n(x^n) > 2^R P_1^n(x^n)\} \quad (4.42)$$

$$\geq \sum_{x^n \in \mathcal{X}^n} (P_0^n(x^n) - 2^R P_1^n(x^n)) 1\{P_0^n(x^n) > 2^R P_1^n(x^n)\} \quad (4.43)$$

$$\geq \sum_{x^n \in \mathcal{X}^n} (P_0^n(x^n) - 2^R P_1^n(x^n)) 1\{x^n \in \mathcal{A}_n^c\} \quad (4.44)$$

$$= P_0^n(\mathcal{A}_n^c) - 2^R P_1^n(\mathcal{A}_n^c) \quad (4.45)$$

$$= 1 - \alpha_n(\mathcal{A}_n) - 2^R \beta_n(\mathcal{A}_n) \quad (4.46)$$

$$\geq 1 - \epsilon - 2^R \beta_n^*(\epsilon). \quad (4.47)$$

The critical step is to get from Eq. (4.43) to Eq. (4.44). To verify this, note that the test $1\{P_0^n(x^n) > 2^R P_1^n(x^n)\}$ is actually the one that maximises the sum since it cuts out all negative contributions. Any other test, including \mathcal{A}_n^c , can thus only reduce the sum.

Now, if we choose $R = \log \delta - \log \beta_n^*(\epsilon)$, the above implies that

$$P_0^n(Z \leq R) \leq \epsilon + \delta \quad (4.48)$$

and thus we have $D_s^{\epsilon+\delta}(P_0^n \| P_1^n) \geq R \geq -\log \beta_n^*(\epsilon) + \log \delta$, which is the desired upper bound. \square

In the homework you have shown that the following asymptotic limit holds.¹

Lemma 4.10. *Let $P_0, P_1 \in \mathcal{P}(X)$ such that $D(P_0 \| P_1) < \infty$ and $\epsilon \in (0, 1)$. Then,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} D_s^\epsilon(P_0^n \| P_1^n) = D(P_0 \| P_1). \quad (4.49)$$

As an aside, we can evaluate the quantity on the left, $\frac{1}{n} D_s^\epsilon(P_0^n \| P_1^n)$, even to higher orders in n using the central limit theorem. While we will not need this here, analysing such higher order terms has been a fruitful area of research recently as it allows us to make more precise statements about optimal errors for smaller n , and thus for practical settings where we are far from the asymptotic setting of very large n .

Proof of Theorem 4.7. The proof of the theorem is evident once we combine Lemma 4.9 and Lemma 4.10. Namely, from Lemma 4.9 we get

$$\frac{1}{n} D_s^\epsilon(P_0^n \| P_1^n) \leq -\frac{1}{n} \log \beta_n(\epsilon) \leq \frac{1}{n} D_s^{\epsilon+\delta}(P_0^n \| P_1^n) + \frac{1}{n} \log \frac{1}{\delta} \quad (4.50)$$

and in the limit $n \rightarrow \infty$ both the lower and upper bound converge to the relative entropy by Lemma 4.10. \square

¹It is essentially a direct consequence of the law of large numbers applied for the random variable $Z = \sum_{i=1}^n \log P_0(X_i) - \log P_1(X_i)$, where X_i are i.i.d. distributed according to the law P_0 .

Chapter 5

Error correcting codes

[Week 8]

Intended learning outcomes:

- You will be familiar with the concept of error correcting codes and can compute their rate and minimal distance.
- You can construct a linear code from its generator matrix or parity check matrix, and understand the finite field arithmetic of F_q for $q \in \{2, 3, 4, 5\}$.
- You understand the basic concepts behind Reed-Solomon codes and can construct .

5.1 Definitions and bounds on codebook size

Error correcting codes are a very rich topic and are studied by communication engineers, computer scientists and mathematicians alike. In computer science, for example, they are used in complexity theory, cryptography, and the study of pseudo-randomness. We can only touch the very surface of this theory here. We will first discuss some general properties of codes and particularly linear codes, and then move on to describe one widely used class of codes in more detail, the Reed-Solomon family of codes.

In the following we will consider codewords that are strings of a fixed length, on some alphabet Σ . The following notions are useful.

Definition 5.1. *The Hamming weight of a string Σ^n is defined as $|\{i : x_i \neq 0\}|$, i.e., the number of nonzero elements of x^n . The Hamming distance between two strings $x^n, y^n \in \Sigma^n$ is defined as*

$$\delta(x^n, y^n) = |\{i : x_i \neq y_i\}|, \quad (5.1)$$

i.e., the number of locations where the strings differ.

We will now introduce the notion of an *error correction code*, or simply code for the remainder of this chapter.

Definition 5.2. An error correction code C of length n over a finite alphabet Σ is a subset of Σ^n . The elements of C are called codewords, and C is sometimes also called the codebook. We will use the following properties and definitions:

- An error correction code is a binary code if $\Sigma = \{0, 1\}$. (We will mostly consider binary codes in the following.)
- A binary code is a linear code if C is a subspace of $\{0, 1\}^n$. This means that for any two codewords $c, c' \in C$, the bitwise XOR of c and c' , denoted $c \oplus c'$, is an element of C as well. In particular, the all zero vector is in C .
- The size of the codebook is denoted by $|C|$.
- The rate of the code is defined as

$$R(C) = \frac{\log |C|}{n \log |\Sigma|} \quad (5.2)$$

- The minimal distance of a code C , denoted $d(C)$, is defined as

$$d(C) = \min_{\substack{c, c' \in C \\ c \neq c'}} \delta(c, c') \quad (5.3)$$

Question 5.3. Consider a binary code $C \subseteq \{0, 1\}^n$ where each codeword is constructed by adding a parity bit to a bit string of length $n - 1$. Is this a linear code? What can you say about its minimum distance?

The following relationships between minimal distance of a binary code and its use for error correction are rather immediate. Consider a binary code with minimum distance $2t + 1$. Such a code can be used to

- Detect up to $2t$ bit flip errors.
- Correct up to t bit flip errors.
- Correct up to $2t$ erasures.

In the erasure model the decoder is informed which bits of the codeword are faulty.

The following bounds establish a relationship between these parameters, limiting the size of the code in terms of the other parameters.

Lemma 5.4 (Hamming bound). *Let C be a binary code with block length n and distance d . Then,*

$$|C| \leq \frac{2^n}{\sum_{i=0}^{\lfloor \frac{d-1}{2} \rfloor} \binom{n}{i}} \quad (5.4)$$

In particular, when $d = 3$ we have $|C| \leq \frac{2^n}{n+1}$.

Proof. For every codeword c define its neighbourhood $N(c, r)$ as all the string that differ from c in at most r locations. Setting $r = \lfloor \frac{d-1}{2} \rfloor$, we note that $N(c, r) \cap N(c', r) = \emptyset$ for any two distinct codewords c and c' . Moreover, we have

$$|N(c, r)| = \sum_{i=0}^r \binom{n}{i}. \quad (5.5)$$

Hence, we can write

$$2^n \geq \left| \bigcup_{c \in C} N(c, r) \right| = \sum_{c \in C} |N(c, r)| = |C| \sum_{i=0}^r \binom{n}{i}. \quad (5.6)$$

Solving this for $|C|$ yields the desired inequality. \square

Note that for this bound to hold with equality the space must be exactly filled out by these neighbourhood balls. We all codes for which this is true *perfect codes*. The following bound applies to all codes, not only binary codes.

Lemma 5.5 (Singleton bound). *Let C be a code with block length n and distance d on an alphabet with $|\Sigma| = q$. Then, we must have*

$$|C| \leq q^{n-d+1}. \quad (5.7)$$

Proof. First observe that there are q^n possible codewords. Let C be an arbitrary code of minimum distance d . Clearly, all codewords $c \in C$ are distinct. Moreover, if we puncture the code by deleting the first $d-1$ letters of each codeword, then all resulting codewords must still be pairwise different. The newly obtained codewords each have length $n - (d-1) = n-d+1$, and thus, there can be at most q^{n-d+1} of them. \square

5.2 Linear codes

We will often be interested in binary codes, but it is important to note that these ideas can all be extended to the case where Σ is any finite field. Linear codes are defined on a field F_q , and codewords of length n are vectors in F_q^n . Since linear codes form subspaces we can express every codeword as a linear combination of a basis of codewords. We denote by k the *dimension* of the subspace, or the minimal number of codewords needed to for a basis. A linear code with a k -dimensional subspace of an n -dimensional space is referred to as a $[n, k]_q$ -code. Furthermore, if it has minimum distance d , we call it an $[n, k, d]_q$ -code. We usually drop the subscript q when it is clear from context, e.g. when we are discussing binary codes.

Definition 5.6. *Let C be an $[n, k]$ -code. A matrix $G \in F_q^{n \times k}$ is said to be a generator matrix for C if its k columns span C .*

Using the generator matrix we can encode any binary string $x \in F_q^k$ into a codeword $c \in \Sigma^n$ by the matrix multiplication $c = Gx$. Note that a linear code admits different generator matrices, corresponding to the different choices of basis for the code as a vector space. This corresponds to different encodings of the messages into codewords, with the same fixed set of codewords.

Example 5.7. Consider the binary repetition code for $n = 3$ comprised of the codewords 000 and 111. The generator matrix for this code is $G = (1, 1, 1)^T$.

For linear codes we can give a bound on the codebook size—the Singleton bound simply evaluates to $k \leq n - d + 1$ or

$$d \leq n - k + 1. \quad (5.8)$$

There are two generic ways to characterise a subspace:

- By specifying a basis of the subspace, as we have done above using the generator matrix.
- By specifying a basis of the orthogonal subspace.

For linear codes that orthogonal subspace is spanned by vectors that are orthogonal to the linear subspace spanned by the codewords. Those vectors can be interpreted as parity checks.

Definition 5.8. Let C be an $[n, k]$ -code. A matrix $H \in F_q^{(n-k) \times n}$ is said to be a parity check matrix for C if $Hc = 0$ for every $c \in C$.

Example 5.9. The Hamming code is a binary $[7, 4, 3]$ -code given by codewords of the form

$$x_1, \quad x_2, \quad x_3, \quad x_4, \quad x_2 \oplus x_3 \oplus x_4, \quad x_1 \oplus x_3 \oplus x_4, \quad x_1 \oplus x_2 \oplus x_4. \quad (5.9)$$

A possible generator matrix for this code is given by

$$G = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \end{pmatrix} \quad (5.10)$$

A possible parity check matrix is given by

$$H = \begin{pmatrix} 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{pmatrix} \quad (5.11)$$

The Hamming code is a perfect code since $2^k = \frac{2^n}{n+1}$ for $k = 4$ and $n = 7$.

Definition 5.10. The dual of a binary $[n, k]$ -code C , the $[n, n - k]$ -code C^\perp , is the space spanned by all codewords $c' \in F_q^n$ such that

$$\sum_{i=1}^n c_i c'_i = 0 \quad (5.12)$$

for all $c \in C$.

From the definition we can see that $G^\perp = H^T$ and $H^\perp = G^T$. In particular, the dual of a dual code is the code itself.

5.3 Reed-Solomon codes

Reed-Solomon were first used to do error correction for the Voyager program and became really widespread in their use to protect against errors on compact discs. They are still used in two-dimensional bar codes like QR codes.

The Reed-Solomon code is actually a family of codes, where every code is characterised by three parameters: an alphabet size q , a block length n , and a message length k , with $k < n \leq q$. In this code a message $m = (m_0, m_1, \dots, m_{k-1}) \in F_q^k$ is first mapped to a polynomial $p_m(x)$ with $x \in F_q$ of degree $k - 1$ given by

$$p_m(x) = \sum_{i=0}^{k-1} m_i x^i. \quad (5.13)$$

The codeword for m is then obtained by evaluating p_m at n different points $x_i \in F_q$ for $i \in [n]$, i.e.

$$C(m) = (p_m(x_1), \dots, p_m(x_n)). \quad (5.14)$$

This constitutes a linear code with the generator matrix $G \in F_q^{k \times n}$ given by

$$G = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ \vdots & \vdots & & \vdots \\ x_1^{k-1} & x_2^{k-1} & \dots & x_n^{k-1} \end{pmatrix}^T \quad (5.15)$$

The basic idea is that a polynomial of order $k - 1$ is uniquely specified if we know its value at k points or more.

Lemma 5.11. The above Reed-Solomon code is a $[n, k, n - k + 1]_q$ -code.

Proof. The code is an $[n, k]_q$ -code by construction. The only property we need to show here is that the minimal distance of the code is given by $d = n - k + 1$.

On the one hand, we can generally write $\delta(c_1, c_2) = \delta(c_1 - c_2, 0)$ for any two codewords $c_1, c_2 \in C$, where $c_1 - c_2$ is also a codeword since C is linear. Hence, the minimal distance of the codebook is simply given by the minimal number of nonzero elements in any codeword that is not the all zero codeword. But since for any $m \neq 0$, the polynomial $p_m(x)$ is nontrivial and of order $k - 1$, we know that it has at most $k - 1$ roots, so we must have $d \geq n - (k - 1) = n - k + 1$.

On the otherhand, we observe that by the Singleton bound we must have $d \leq n - k + 1$, therefore concluding the proof. \square

If we choose $q = 2^n$ we can interpret the Reed-Solomon code as a binary code. For $q = 2^2$, $n = 4$ and $k = 2$ we get the following mappings, where each symbol $\{0, 1, 2, 3\}$ can be interpreted as a binary sequence $\{00, 01, 10, 11\}$. The codewords are constructed by evaluating the polynomial at the points $\{0, 1, 2, 3\}$, using the multiplication and addition rules for F_4 discussed above. For example, we get

$$0011 = 03 \rightarrow 3 + 0x \rightarrow 3333 = 11111111 \quad (5.16)$$

$$1001 = 21 \rightarrow 1 + 2x \rightarrow 1320 = 01111000 \quad (5.17)$$

$$1010 = 22 \rightarrow 2 + 2x \rightarrow 2013 = 10000111. \quad (5.18)$$

The above should be read as “initial bit string” = “written in F_4 by interpreting it as binary representation” \rightarrow “corresponding polynomial of degree 1” \rightarrow “polynomial evaluated at $x = \{0, 1, 2, 3\}$ ” = “encoded bit string using binary representation”.

5.4 Low density parity check (LDPC) codes

LDPC codes are linear codes with the property that the parity-check matrix H is sparse, i.e., the individual parities that need to be checked involve only a few of the message bits. The sparsity of H allows for a relatively efficient iterative decoding heuristic called *belief propagation*. Essentially this algorithm updates a probabilistic guess (in the spirit of Bayesian updates) of the codeword bits by punishing assignments that do not satisfy the checksums. We do not have time to go into this in detail but it turns out that this heuristic works extremely well in practice, even though we cannot generally prove its convergence. On the other hand, finding the message with the maximum likelihood, as an optimal decoder would do, is a problem that we do not know how to solve in time polynomial in the block length n . These efficiency considerations are important since this implies that LDPC codes can be used for large block lengths, and in particular can be used to approach the capacity of a communication channel. We will cover the capacity of channels in the next lecture.

5.4.1 Decoding with belief propagation

TBD

Chapter 6

Noisy channel coding

[Week 9–12]

Intended learning outcomes:

- You can compute the channel mutual information, suitably simplifying the calculation if the channel exhibits symmetry.
- You understand the formal setup of the noisy channel coding problem, and are familiar with the binary symmetric channel (BSC), binary erasure channel (BEC) and additive white Gaussian noise (AWGN) channel.
- You know the difference between asymptotic and one-shot bounds and can derive the former from the latter.
- You can determine the type of a sequence and compute the empirical distribution, and construct random codes.
- You can determine if a DMS can be transmitted through a DMC using the source-channel separation theorem.

Book reference: Chapter 7 in Cover & Thomas [1].

6.1 Channel mutual information

To quote Shannon from his pivotal paper “A Mathematical theory of communication”:

The fundamental problem of communication is that of reproducing at one point either exactly or approximately a message selected at another point.

The basic setup of the communication problem consists of a source that generates digital information which is to be reliably communicated to a destination through a channel, preferably in the most efficient manner possible. The destination could be spatially or temporally separated.

In this chapter we will first learn how to transmit a source that produces messages with uniform probability from some set of messages and then argue that the optimal strategy to transmit an arbitrary source is to first compress it (which makes it approximately uniform) and then send it over the channel. The latter is called the source-channel separation theorem since it allows to treat channel coding and source coding independently as two separate tasks, without loss of efficiency — at least in the asymptotic limit of large block lengths.

Before we state the main theorem we want to explore the following quantity:

Definition 6.1. *We fix alphabets of input symbols, \mathcal{X} , and output symbols, \mathcal{Y} . Let W be a channel, a stochastic map represented as a conditional probability distribution $W_{Y|X}(y|x)$. The channel mutual information of W is defined as*

$$I(W) := \max_{P_X \in \mathcal{P}(\mathcal{X})} I(X : Y), \quad \text{where } P_{XY}(x, y) = P_X(x)W_{Y|X}(y|x) \quad (6.1)$$

is the joint distribution of channel input and output.

This is the maximal mutual information between channel input and output. The quantity is often called “channel capacity” in the literature, and we will see that it in fact corresponds to the maximal rate at which information can be transmitted over the channel in the next section. However, we prefer to keep a semantic difference between information quantities, like the channel mutual information, and operational quantities, like the channel capacity. Only through the study of information theory do we actually establish their equivalence, and usually only in special cases, e.g. for memoryless channels in this case.

The optimisation is well-behaved since the underlying function is concave in P_X , as the following lemma shows.

Lemma 6.2. *For a fixed stochastic map W , the mutual information between channel input and output, $I(X : Y)$, is concave in the marginal pmf P_X .*

Proof. We have $I(X : Y) = H(Y) - H(Y|X)$, which we may write as

$$I(X : Y) = H(P_Y) - \sum_x P_X(x)H(Y|X = x), \quad (6.2)$$

where $P_Y(y) = \sum_x P_X(x)P_{Y|X}(y|x)$. By concavity of the entropy function we now see that the first term is concave in P_X . The second term is linear and thus concave in P_X as well. \square

One very important consequence of this concavity is that we can simplify the optimisation for symmetric channels. In strongest symmetry we consider here is one where every permutation of input symbols can be “undone” by doing a respective permutation of the output symbols of the channel.

Proposition 6.3. *Consider a channel W such that for every permutation π on \mathcal{X} there exists a permutation π' on \mathcal{Y} such that $W(y|x) = W(\pi'(y)|\pi(x))$ for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$. Then, the uniform pmf on X is achieving the channel mutual information.*

Proof. Let us first simplify the notation a bit by making the alphabet more concrete. Without loss of generality we can pick $\mathcal{X} = \{0, 1, \dots, d-1\}$ for some $d \in \mathbb{N}$. Consider the permutations $\pi_k : x \rightarrow x + k \pmod d$; by assumption of the proposition there then exists a permutations π'_k on \mathcal{Y} such that $W(y|x) = W(\pi'_k(y)|\pi_k(x))$.

Take fix any pmf P_X . Our first claim is that the pmfs $P_X^k(x) = P_X(\pi_k(x))$ achieve the same mutual information as P_X for any $k \in \mathcal{X}$. To verify this, we write $I(X : Y)_{P^k} = H(P_Y^k) - \sum_x P_X^k(x) H(P_{Y|X=x})$, thus we can show equivalence for each of the two terms separately. We have

$$P_Y^k(y) = \sum_{x \in \mathcal{X}} P_X^k(x) W(y|x) \quad (6.3)$$

$$= \sum_{x \in \mathcal{X}} P_X(\pi_k(x)) W(\pi'_k(y)|\pi_k(x)) \quad (6.4)$$

$$= \sum_{x \in \mathcal{X}} P_X(x) W(\pi'_k(y)|x) = P_Y(\pi'_k(y)), \quad (6.5)$$

and, thus, the probabilities in P_Y^k and P_Y are the same (although in different order), and thus we have $H(P_Y) = H(P_Y^k)$. Moreover, the condition on $W(y|x)$ directly implies that $H(P_{Y|X=x}) = H(P_{Y|X=x'})$ for any $x, x' \in \mathcal{X}$, and we have thus established the identity

$$I(X : Y)_P = I(X : Y)_{P_k} \quad \forall k \in \mathcal{X}. \quad (6.6)$$

Next we use this identity to write

$$I(X : Y)_P = \sum_{k \in \mathcal{X}} \frac{1}{d} I(X : Y)_{P_k}.$$

This can be interpreted as the expected value of $I(X : Y)$, where the pmf is chosen uniformly at random from amongst the permuted pmfs P_X^k . However, since the mutual information is concave in the pmf, we have that

$$\sum_{k \in \mathcal{X}} \frac{1}{d} I(X : Y)_{P_k} \leq I(X : Y)_Q,$$

where

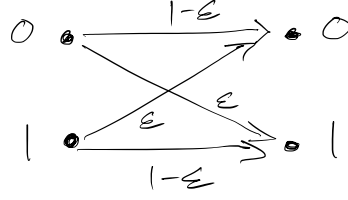
$$Q_X(x) = \sum_{k \in \mathcal{X}} \frac{1}{d} P_X^k(x) = \frac{1}{d} \sum_{k \in \mathcal{X}} P_X^k(\pi_k(x)) = \frac{1}{d}$$

is the expected pmf of X , with the expectation taken over the above-mentioned uniform distribution on the pmfs.

This implies that the symmetrized distribution Q_X always leads to a mutual information that is at least as high as what can be achieved by P_X . We can thus conclude that the distribution that maximises the channel mutual information is of the form Q_X . \square

We will now consider two very prominent examples of communication channels and compute their channel mutual information.

1. The *binary symmetric channel* (BSC) takes a binary input to a binary output. The bit is flipped with a certain probability, here denoted ϵ , and otherwise left intact:



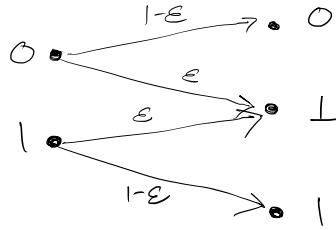
The conditional probability distribution of the channel is given by

$$W_{\text{BSC}(\epsilon)}(y|x) = (1 - \epsilon)1\{x = y\} + \epsilon 1\{x \neq y\}. \quad (6.7)$$

The channel mutual information for the BSC is easy to evaluate, even without invoking Proposition 6.3—which does clearly apply here. Let us simply note that $H(Y|X = x) = h(\epsilon)$, the binary entropy evaluated for ϵ , and this is independent of $x \in \{0, 1\}$. Hence the mutual information is given by $I(X : Y) = H(Y) - h(\epsilon)$, which is maximised when Y is uniformly distributed. This is achieved when X is uniformly distributed itself. Hence, the channel mutual information is given by

$$I(W_{\text{BSC}(\epsilon)}) = 1 - h(\epsilon). \quad (6.8)$$

2. The *binary erasure channel* (BEC) takes a binary input to a ternary output, $\{0, 1, \perp\}$. The output \perp has probability ϵ on either input, and otherwise the input symbol remains unaffected. Essentially this is a channel that flags errors:



The conditional probability distribution of the channel is given by

$$W_{\text{BEC}(\epsilon)}(y|x) = (1 - \epsilon)1\{x = y\} + \epsilon 1\{y = \perp\}. \quad (6.9)$$

By Proposition 6.3 we can again argue that the maximising input distribution is the uniform distribution. And we get the output distribution

$$P_Y(y) = \begin{cases} \frac{1}{2}(1 - \epsilon) & \text{if } y \in \{0, 1\} \\ \epsilon & \text{if } y = \perp \end{cases} \quad (6.10)$$

We again have $H(Y|X = x) = h(\epsilon)$ independent of x and can then compute

$$I(W_{\text{BEC}(\epsilon)}) = H(Y) - h(\epsilon) \quad (6.11)$$

$$= -2 \cdot \frac{1}{2}(1 - \epsilon) \log \frac{1}{2}(1 - \epsilon) - \epsilon \log \epsilon - h(\epsilon) \quad (6.12)$$

$$= (1 - \epsilon) + h(\epsilon) - h(\epsilon) \quad (6.13)$$

$$= 1 - \epsilon \quad (6.14)$$

For the general case the problem is a bit more difficult, but concavity ensures that if we find a local maximum for the mutual information then that maximum is in fact global. Based on this, there are algorithms that can compute the channel mutual information efficiently for any stochastic map.

The following expression for the channel mutual information is useful to know.

Proposition 6.4. *For any stochastic map W , we have*

$$I(W) = \min_{Q \in \mathcal{P}(\mathcal{Y})} \max_{x \in \mathcal{X}} D(W_{Y|X}(\cdot|x) \| Q_Y). \quad (6.15)$$

Proof. We first write, using the definition of the channel mutual information,

$$I(W) = \max_{P_X \in \mathcal{P}(\mathcal{X})} D(P_{XY} \| P_X \times P_Y) = \max_{P_X \in \mathcal{P}(\mathcal{X})} \min_{Q_Y \in \mathcal{P}(\mathcal{Y})} D(P_{XY} \| P_X \times Q_Y), \quad (6.16)$$

where the second equality comes from the fact that $D(P_{XY} \| P_X \times Q_Y) = D(P_{XY} \| P_X \times P_Y) + D(P_Y \| Q_Y)$ and the minimum is thus achieved for $Q_Y = P_Y$. If we further rewrite

$$D(P_{XY} \| P_X \times Q_Y) = \sum_{x \in \mathcal{X}} P_X(x) D(W_{Y|X}(\cdot|x) \| Q_Y) \quad (6.17)$$

we realise that this quantity is linear in P_X and convex in Q_Y . The idea then is to use Sion's minimax theorem [7], which states that the minimum and maximum in the above expressions can be interchanged. Hence, we get

$$I(W) = \min_{Q_Y \in \mathcal{P}(\mathcal{Y})} \max_{P_X \in \mathcal{P}(\mathcal{X})} \sum_{x \in \mathcal{X}} P_X(x) D(W_{Y|X}(\cdot|x) \| Q_Y). \quad (6.18)$$

Finally note that the maximum in the above expression is taken for a P_X that is concentrated on a single point. This yields the expression in (6.15). \square

6.2 The channel coding theorem

Let us now move on to a more operational description of the channel coding problem. As we have seen a noisy channel can be described by a conditional probability distribution $W_{Y|X}$. If such a channel can be used multiple times, without any memory effects, we speak of a *discrete memoryless channel* (DMC). We will not consider more complicated channels that change over time or have memory effects here, and thus our definition is restricted to the discrete memoryless case.

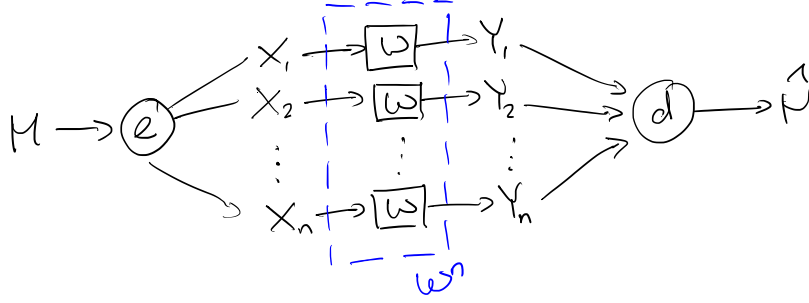


Figure 6.1: **The channel coding setup.** The figure depicts the setup for block length n . The encoder, e , takes a message M and encodes it into n channel input symbols, X_1, X_2, \dots, X_n . The channels act independently on these input symbols. The channel outputs Y_1, Y_2, \dots, Y_n are decoded using the decoder, d , to an estimate \hat{M} of M .

Definition 6.5. A discrete memoryless channel \mathbf{W} is fully characterised by a stochastic map $W = W_{Y|X}$. For any $n \in \mathbb{N}$, the stochastic map W^n takes a sequence of input symbols $x^n \in \mathcal{X}^n$ to a sequence of output symbols $y^n \in \mathcal{Y}^n$ such that

$$P[Y^n = y^n | X^n = x^n] = W^n(y^n | x^n) = \prod_{i=1}^n W_{Y|X}(y_i | x_i). \quad (6.19)$$

Definition 6.6. An $(\epsilon, |M|, n)$ -channel code for a channel $W_{Y|X}$ is comprised of an encoder function $e : [|M|] \rightarrow \mathcal{X}^n$ and a decoder function $d : \mathcal{Y}^n \rightarrow [|M|]$, so that the Markov chain $M \rightarrow X^n \rightarrow Y^n \rightarrow \hat{M}$ where $X = e(M)$ and $\hat{M} = d(Y)$ satisfies

$$P[M \neq \hat{M}] \leq \epsilon \quad (6.20)$$

when M follows the uniform distribution. Here n is called the block length, $|M|$ is the number of possible messages, and ϵ is the allowed probability of error.

This allows us to define the concept of achievable rates and capacity of a DMC.

Definition 6.7. We say that a rate R is achievable for a DMC \mathbf{W} if there exists a sequence of $(\epsilon_n, \lceil 2^{nR} \rceil, n)$ codes for all $n \in \mathbb{N}$ such that

$$\lim_{n \rightarrow \infty} \epsilon_n = 0 \quad (6.21)$$

The capacity of \mathbf{W} , denoted $C(\mathbf{W})$, is the supremum over all achievable rates R .

The main theorem of this section now relates the channel capacity of a DMC with the maximal mutual information of the underlying stochastic map.

Theorem 6.8. For any DMC \mathbf{W} with stochastic map W , we have

$$C(\mathbf{W}) = I(W). \quad (6.22)$$

We will prove this theorem in several steps. First we will derive an upper bound on the cardinality of the message set that holds even for a single use of the channel, the so-called meta-converse. From this we will then prove the converse (upper bound on the rate) and finally show how this rate can be achieved.

6.2.1 The meta-converse

The task in noisy channel coding is to transmit a message reliably over a DMC. We will for now assume that the message is uniformly distributed over some set of messages, but this assumption will be relaxed in the next section. Same as with source coding, we will eventually consider this problem in an asymptotic scenario where the the number of times the channel can be used, n , is taken to infinity. However, some results can conveniently be stated in a *one-shot setting*, without such a limit in mind, and we will do this here. Let us first define the notion of a code in the one-shot setting.



The condition on the distribution of M implies that ϵ enforces an *average error* criterion. We could alternatively also require that the probability of error is small for any distribution of M , which would enforce a *maximum error* criterion.

Our first result is a bound on the cardinality of the message set. This result is called the *meta-converse* as it can be used to derive various different fundamental limits (or converse bounds).¹

Theorem 6.9. *For any $(\epsilon, |M|, 1)$ -channel code for a stochastic map $W_{Y|X}$, we must have*

$$|M| \leq \min_{Q_Y \in \mathcal{P}(Y)} \max_{P_X \in \mathcal{P}(X)} \frac{1}{\beta_\epsilon^*(P_{XY} \| P_X \times Q_Y)}, \quad (6.23)$$

where $P_{XY}(x, y) = P_X(x)W_{Y|X}(y|x)$ is the joint distribution of channel input and output and $\beta_\epsilon^*(P_{XY} \| P_X \times Q_Y)$ is the minimal error of the second kind, as defined as in Eq. (4.5), for the hypothesis testing problem where H_0 is P_{XY} and H_1 is $P_X \times Q_Y$.

The way to think about this hypothesis test is the following. The null hypothesis is that P_{XY} are in fact the channel input and output of our channel W , and the alternative hypothesis is that the output Q_Y has been produced independently of the input P_X , i.e. that it is the output of a channel that is completely useless for information transmission.

We provide the proof here for the special case where the encoder and decoder are deterministic and the encoder is furthermore injective, i.e. the function e uniquely maps messages

¹This result was only established quite recently, by Polyanskiy-Poor-Verdú [5]. So even though information theory (and in particular channel coding) is by now very established, some progress can still be made when it comes to simplifying mathematical proofs and presenting them in a unified way.

to channel inputs. These assumptions make the proof a bit simpler but are not really restrictive. It is easy to verify that for any code using randomness there is a deterministic one performing at least equally well, and non-injective codes give up on distinguishing certain messages from the beginning, which can only give an advantage in extreme regimes where we tolerate large error and want to transmit more messages than fit into the channel alphabet.

Proof of Theorem 6.9. Since we are considering an (ϵ, M) -channel code we must have $P[M \neq \hat{M}] \leq \epsilon$. We now consider the hypothesis testing problem at hand, where H_0 is P_{XY} and H_1 is $P_X \times Q_Y$ for some arbitrary distribution $Q_Y \in \mathcal{P}(\mathcal{Y})$. For this problem we take the test

$$\mathcal{A} = \{(x, y) \in \mathcal{X} \times \mathcal{Y} : x \neq e(d(y))\} \quad (6.24)$$

We can then compute the errors of the first and second kind for this test

$$\alpha(\mathcal{A}) = P_{XY}(\mathcal{A}) = P[X \neq e(d(Y))] = P[e(M) \neq e(\hat{M})] \leq P[M \neq \hat{M}] \leq \epsilon \quad (6.25)$$

and, furthermore,

$$\beta(\mathcal{A}) = P_X \times Q_Y(\mathcal{A}^c) \quad (6.26)$$

$$= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} P_X(x) Q_Y(y) 1\{x = e(d(y))\} \quad (6.27)$$

$$= \sum_{m \in [M]} \frac{1}{|M|} \sum_{y \in \mathcal{Y}} Q_Y(y) 1\{e(m) = e(d(y))\} \quad (6.28)$$

$$= \frac{1}{|M|} \sum_{y \in \mathcal{Y}} Q_Y(y) \sum_{m \in [M]} 1\{e(m) = e(d(y))\} \quad (6.29)$$

$$\leq \frac{1}{|M|} \sum_{y \in \mathcal{Y}} Q_Y(y) = \frac{1}{|M|}. \quad (6.30)$$

Hence, we can deduce that $\beta_\epsilon^*(P_{XY} \| P_X \times Q_Y) \leq \frac{1}{|M|}$. Since we do not know the distribution P_X — it depends on the specific encoding map used — we prefer to loosen the bound to

$$|M| \leq \max_{P_X \in \mathcal{P}(\mathcal{X})} \frac{1}{\beta_\epsilon^*(P_{XY} \| P_X \times Q_Y)}, \quad (6.31)$$

and since this holds for all $Q_Y \in \mathcal{P}(\mathcal{Y})$ the desired inequality holds. \square

6.2.2 Proof of converse and types

[This proof itself will not be part of the final exam; however the concept of types and empirical distribution should be known however.]

In this subsection we will show that $C(\mathbf{W}) \leq I(W)$. In fact, we will show something much stronger. We will show that for any sequence of $(\epsilon, \lceil 2^{nR} \rceil, n)$ -channel codes for a DMC \mathbf{W} with $\epsilon \in (0, 1)$ fixed, we must have $R \leq I(W)$. Hence, even if we allow for a nonzero

error asymptotically, the maximal rate is still bounded by the channel mutual information. This is what is called a strong converse for channel coding.²

For this we will need generic bounds that hold for all codes, like the one established in the meta-converse. However, now we want to do this for n channel uses instead of just a single channel use. But note that we can always interpret the setup for block length n as one single super-channel W^n that takes a full string X^n as input and outputs a string Y^n . For this super-channel W^n we can apply the meta-converse. Since $|M| = \lceil 2^{Rn} \rceil$ for a rate R code, the meta-converse implies that

$$R \leq \min_{Q_{Y^n} \in \mathcal{P}(\mathcal{Y}^n)} \max_{P_{X^n} \in \mathcal{P}(\mathcal{X}^n)} \frac{1}{n} \log \frac{1}{\beta_\epsilon^*(P_{X^n Y^n} \| P_{X^n} \times Q_{Y^n})}, \quad (6.32)$$

where

$$P_{X^n Y^n}(x^n, y^n) = P_{X^n}(x^n) \prod_{i=1}^n W_{Y|X}(y_i | x_i) \quad (6.33)$$

is the joint distribution of channel inputs and outputs when the stochastic map W^n is applied to an arbitrary input distributions P_{X^n} . This upper bound holds for any particular choice of Q_{Y^n} , so in particular we may choose $Q_{Y^n}(y^n) = Q_Y^n$ to be i.i.d..

We can now reuse the tools we introduced to analyse hypothesis testing in Chapter 4. In particular, Lemma 4.9 allows us to bound β_ϵ^* in terms of the information spectrum. If we apply the lemma to the above situation, we get

$$R \leq \min_{Q_Y \in \mathcal{P}(\mathcal{Y})} \max_{P_{X^n} \in \mathcal{P}(\mathcal{X}^n)} \frac{1}{n} D_s^\mu(P_{X^n Y^n}^n \| P_{X^n} \times Q_Y^n) + \log \frac{1}{\delta}, \quad (6.34)$$

where $\delta \in (0, 1 - \epsilon)$ and we set $\mu = \epsilon + \delta$. Recall that the information spectrum is given by

$$D_s^\mu(P_{X^n Y^n}^n \| P_{X^n} \times Q_Y^n) \quad (6.35)$$

$$= \sup \left\{ R \in \mathbb{R} : P_{X^n Y^n} \left[\log \frac{P_{X^n Y^n}(X^n, Y^n)}{P_{X^n}(X^n) Q_Y^n(Y^n)} \leq R \right] \leq \mu \right\} \quad (6.36)$$

$$= \sup \left\{ R \in \mathbb{R} : P_{X^n Y^n} \left[\sum_{i=1}^n \log \frac{W_{Y|X}(Y_i | X_i)}{Q_Y(Y_i)} \leq R \right] \leq \mu \right\} \quad (6.37)$$

$$= \sup \left\{ R \in \mathbb{R} : \sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) W_{Y^n | X^n = x^n} \left[\sum_{i=1}^n \log \frac{W_{Y|X}(Y_i | x_i)}{Q_Y(Y_i)} \leq R \right] \leq \mu \right\} \quad (6.38)$$

$$\leq \max_{x^n \in \mathcal{X}^n} \sup \left\{ R \in \mathbb{R} : W_{Y^n | X^n = x^n} \left[\sum_{i=1}^n \log \frac{W_{Y|X}(Y_i | x_i)}{Q_Y(Y_i)} \leq R \right] \leq \mu \right\}. \quad (6.39)$$

The last inequality holds since if $W_{Y^n | X^n = x^n}[\dots]$ does not exceed μ when we take the expectation over P_{X^n} , then there must (at least) exist one sequence x^n for which the expression

²The proof strategy we follow here is from [8].

does not exceed μ . We can thus relax the bound by just requiring the inequality to hold for this particular sequence x^n , yielding an upper bound. Since we don't know how this sequence looks like, we further relax the bound by just maximising over all sequences x^n .

Let us now analyse the expression

$$W_{Y^n|X^n=x^n} \left[\sum_{i=1}^n \log \frac{W_{Y|X}(Y_i|x_i)}{Q_Y(Y_i)} \leq R \right]. \quad (6.40)$$

First, we note that the expression does not actually depend on the full sequence x^n , it only depends on how many times each element of \mathcal{X} appears in the sequence x^n . This is called the *type* of x^n . We will also introduce the *empirical distribution* of a sequence x^n , which also only depends on the number of times each symbol appears, and is defined as

$$P_X^{x^n}(x) = \frac{1}{n} |\{i \in [n] : x_i = x\}|. \quad (6.41)$$

The random variable in (6.40) is a sum of independent (but not identical!) random variables. Let us compute its expectation and variance, which are

$$\sum_{i=1}^n \mathbb{E} \left[\log \frac{W_{Y|X}(Y_i|x_i)}{Q_Y(Y_i)} \right] = \sum_{i=1}^n D(W(\cdot|x_i)||Q) = n \sum_{x \in \mathcal{X}} P_X^{x^n}(x) D(W(\cdot|x)||Q) \quad (6.42)$$

$$\sum_{i=1}^n \text{Var} \left[\log \frac{W_{Y|X}(Y_i|x_i)}{Q_Y(Y_i)} \right] \leq n \underbrace{\max_{x \in \mathcal{X}} \text{Var} \left[\log \frac{W_{Y|X}(Y_i|x)}{Q_Y(Y_i)} \right]}_{=: \sigma^2}, \quad (6.43)$$

where σ^2 is some constant. So if we set $R = n(\sum_{x \in \mathcal{X}} P_X^{x^n}(x) D(W(\cdot|x)||Q) + \nu)$ for some small $\nu > 0$, Chebyshev's inequality yields

$$W_{Y^n|X^n=x^n} \left[\sum_{i=1}^n \log \frac{W_{Y|X}(Y_i|x_i)}{Q_Y(Y_i)} \geq R \right] \leq \frac{\sigma^2}{\nu^2 n}, \quad (6.44)$$

or, equivalently,

$$W_{Y^n|X^n=x^n} \left[\sum_{i=1}^n \log \frac{W_{Y|X}(Y_i|x_i)}{Q_Y(Y_i)} \leq R \right] \geq 1 - \frac{\sigma^2}{\nu^2 n}, \quad (6.45)$$

If we choose n large enough, then $1 - \frac{\sigma^2}{\nu^2 n}$ is always larger than μ . Hence, we can deduce that

$$\sup \left\{ R \in \mathbb{R} : W_{Y^n|X^n=x^n} \left[\sum_{i=1}^n \log \frac{W_{Y|X}(Y_i|x_i)}{Q_Y(Y_i)} \leq R \right] \leq \mu \right\} \quad (6.46)$$

$$\leq n \left(\sum_{x \in \mathcal{X}} P_X^{x^n}(x) D(W(\cdot|x)||Q) + \nu \right) \quad (6.47)$$

$$\leq n \max_{P_X \in \mathcal{P}(\mathcal{X})} \sum_{x \in \mathcal{X}} P_X(x) D(W(\cdot|x)||Q) + n\nu \quad (6.48)$$

$$\leq n \max_{x \in \mathcal{X}} D(W(\cdot|x)||Q) + n\nu. \quad (6.49)$$

This expression now no longer depends on x^n . Plugging it into Eq. (6.34), we find

$$R \leq \min_{Q_Y \in \mathcal{P}(\mathcal{Y})} \max_{x \in \mathcal{X}} D(W(\cdot|x)||Q) + \nu + \frac{1}{n} \log \frac{1}{\delta} \quad (6.50)$$

Since $\frac{1}{n} \log \frac{1}{\delta} \leq \nu$ for sufficiently large n and furthermore ν can be chosen arbitrarily small, we can conclude that the inequality

$$R \leq \min_{Q_Y \in \mathcal{P}(\mathcal{Y})} \max_{x \in \mathcal{X}} D(W(\cdot|x)||Q) \quad (6.51)$$

must hold. But according to Proposition 6.4, this minimax expression is exactly the channel mutual information, concluding our proof.

6.2.3 Proof of achievability and random codes

We will need the following technical lemma³.

Lemma 6.10. *Let $t \geq 0$ and $s \in [0, 1]$. Then,*

$$1 - \frac{s}{s+t} \leq 2(1-s) + 4t. \quad (6.52)$$

Proof. We note that since $0 \leq \left(1 - \left(\frac{1}{\sqrt{s+t}} - 1\right)\right)^2$, we have

$$2\left(\frac{1}{\sqrt{s+t}} - 1\right) \leq 1 + \left(\frac{1}{\sqrt{s+t}} - 1\right)^2. \quad (6.53)$$

We may thus write

$$1 - \frac{s}{s+t} = \frac{t}{s+t} \quad (6.54)$$

$$= t + 2t\left(\frac{1}{\sqrt{s+t}} - 1\right) + t\left(\frac{1}{\sqrt{s+t}} - 1\right)^2 \quad (6.55)$$

$$\leq 2t + 2t\left(\frac{1}{\sqrt{s+t}} - 1\right)^2 \quad (6.56)$$

$$\leq 2t + 2(s+t)\left(\frac{1}{s+t} + 1 - \frac{2}{\sqrt{s+t}}\right). \quad (6.57)$$

Now it remains to bound $\sqrt{s+t} \geq \sqrt{s} \geq s$ since $s \leq 1$. Thus,

$$(s+t)\left(\frac{1}{s+t} + 1 - \frac{2}{\sqrt{s+t}}\right) = 1 + s + t - 2\sqrt{s+t} \leq 1 + s + t - 2s = 1 + t - s \quad (6.58)$$

Plugging this into (6.57) yields the desired inequality. \square

³This is taken from [3].

We again first analyse the channel coding problem in the one-shot setting where the channel is only used once.

Theorem 6.11. *There exists an $(6\epsilon, |M|, 1)$ -channel code for a stochastic map $W_{Y|X}$ as long as the code parameters satisfy*

$$|M| \leq \epsilon \cdot \frac{1}{\beta_\epsilon^*(P_{XY} \| P_X \times P_Y)} \quad (6.59)$$

for some pmf $P_X \in \mathcal{P}(\mathcal{X})$.

Proof. We now construct a code for a single use of the channel. First, we fix any distribution $P_X \in \mathcal{P}(\mathcal{X})$. From this we generate $|M|$ codewords independently by picking them from the distribution P_X , i.e. the output of the decoder, $E(m)$, is itself a random variable following the distributions P_X for each message m . The decoder is constructed as follows. Consider the binary hypothesis testing problem between $H_0 : P_{XY}$ and $H_1 : P_X \times P_Y$. The optimal test $\mathcal{A} \subset \mathcal{X} \times \mathcal{Y}$ satisfies

$$P_{XY}(\mathcal{A}) \leq \epsilon \quad \text{and} \quad (P_X \times P_Y)(\mathcal{A}^c) = \beta_\epsilon^*(P_{XY} \| P_X \times P_Y). \quad (6.60)$$

From this we construct the sets $\mathcal{A}_x = \{y \in \mathcal{Y} : (x, y) \notin \mathcal{A}\}$ for all $x \in \mathcal{X}$. For a fixed encoder $E = e$, the decoder is probabilistic. Given a channel output y it assigns $\hat{M} = m$ with probability

$$P[\hat{M} = m | Y = y] = \frac{1\{y \in \mathcal{A}_{e(m)}\}}{\sum_{m'} 1\{y \in \mathcal{A}_{e(m')}\}} = \frac{1\{y \in \mathcal{A}_{e(m)}\}}{1\{y \in \mathcal{A}_m\} + \sum_{m' \neq m} 1\{y \in \mathcal{A}_{e(m')}\}}. \quad (6.61)$$

Let us now analyse the probability of error for this code, first for a fixed set of codewords (or fixed encoder, e) and fixed message m .

$$P[M \neq \hat{M} | M = m, E = e] = 1 - \sum_{y \in \mathcal{Y}} W(y|e(m)) P[\hat{M} = m | Y = y] \quad (6.62)$$

$$= \sum_{y \in \mathcal{Y}} W(y|e(m)) \left(1 - \frac{1\{y \in \mathcal{A}_{e(m)}\}}{1\{y \in \mathcal{A}_m\} + \sum_{m' \neq m} 1\{y \in \mathcal{A}_{e(m')}\}} \right) \quad (6.63)$$

We can now use Lemma 6.10 to bound this as

$$P[M \neq \hat{M} | M = m, E = e] \quad (6.64)$$

$$\leq \sum_{y \in \mathcal{Y}} W(y|e(m)) \left(2 \cdot 1\{y \notin \mathcal{A}_{e(m)}\} + 4 \cdot \sum_{m' \neq m} 1\{y \in \mathcal{A}_{e(m')}\} \right). \quad (6.65)$$

We may now take the average over all encoders e , so that $e(m)$ and $e(m')$ are independent and follow the distribution P_X . This gives the following bound

$$P[M \neq \hat{M} | M = m] \quad (6.66)$$

$$\leq \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} P_X(x) W(y|x) \left(2 \cdot 1\{y \notin \mathcal{A}_x\} + 4 \underbrace{(|M| - 1)}_{\leq |M|} \sum_{x' \in \mathcal{X}} P_X(x') 1\{y \in \mathcal{A}_{x'}\} \right) \quad (6.67)$$

and we note that the bound no longer depends on the choice of m , i.e. Eq. (6.67) is in fact an upper bound on $P[M \neq \hat{M}]$. Let us now investigate the two summands in (6.67) individually. We first observe that

$$\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} P_X(x) W(y|x) 1\{y \notin \mathcal{A}_x\} = P_{XY} [1\{(x, y) \in \mathcal{A}\}] = P_{XY} [A] \leq \epsilon \quad (6.68)$$

by definition of the sets \mathcal{A}_x and \mathcal{A} . We can also evaluate

$$\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} P_X(x) W(y|x) \sum_{x' \in \mathcal{X}} P_X(x') 1\{y \in \mathcal{A}_{x'}\} = \sum_{y \in \mathcal{Y}} P_Y(y) \sum_{x' \in \mathcal{X}} P_X(x') 1\{(x', y) \in \mathcal{A}\} \quad (6.69)$$

$$= (P_X \times P_Y)[\mathcal{A}^c] \quad (6.70)$$

$$= \beta_\epsilon^*(P_{XY} \| P_X \times P_Y). \quad (6.71)$$

Summarising this, we find that

$$P[M \neq \hat{M}] \leq 2\epsilon + 4|M|\beta_\epsilon^*(P_{XY} \| P_X \times P_Y) \quad (6.72)$$

So, in particular, as long as we choose $|M| \leq \epsilon \cdot \beta_\epsilon^*(P_{XY} \| P_X \times P_Y)^{-1}$, we achieve $P[M \neq \hat{M}] \leq 6\epsilon$, as required. \square

Based on this, we will now show that any rate $R < I(W)$ is achievable by again considering the one-shot results applied to the super-channel W^n . Theorem 6.11 stipulates that there exists a code with $\lceil 2^{nR} \rceil$ codewords and error 6ϵ as long as

$$\lceil 2^{nR} \rceil \leq \epsilon \cdot \frac{1}{\beta_\epsilon^*(P_{X^n Y^n} \| P_{X^n} \times P_{Y^n})} \quad (6.73)$$

for some input distribution $P_{X^n} \in \mathcal{P}(\mathcal{X}^n)$. Or, using that $\log(x-1) \geq \log x - 1$ for $x > 2$,

$$R \leq \frac{1}{n} \left(\log \frac{1}{\beta_\epsilon^*(P_{X^n Y^n} \| P_{X^n} \times P_{Y^n})} + \log \epsilon - 1 \right). \quad (6.74)$$

If we further require that P_{X^n} is i.i.d. we get a more stringent requirement on R , but it makes our analysis simpler. We can then use that, for every $P_X \in \mathcal{P}(\mathcal{X})$, we have by the Chernoff-Stein's Lemma (cf. Theorem 4.7) that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \frac{1}{\beta_\epsilon^*(P_{X^n Y^n} \| P_X^n \times P_Y^n)} = D(P_{XY} \| P_X \times P_Y) = I(X : Y) \quad (6.75)$$

And thus, using the definition of the channel mutual information, we find that as long as $R \leq I(W) - \mu$ for some $\mu > 0$, and for any error $\epsilon > 0$, there exists code with 2^{nR} messages for sufficiently large n . This implies that $I(W) - \mu$ is an achievable rate, and since this is true for any $\mu > 0$ we have shown the desired bound, $C(\mathbf{W}) \geq I(W)$.

6.2.4 Maximum probability of error

Consider a code with $|M|$ codewords. So far we have used the *average probability of error* as a metric for our codes, namely we required that

$$P[\hat{M} = M] = \sum_{m \in [|M|]} \frac{1}{|M|} P[\hat{M} = m | M = m] \quad (6.76)$$

vanishes asymptotically. Sometimes we would however like to impose an even stricter condition, namely that the *maximum probability of error*, given by

$$\max_{m \in [|M|]} P[\hat{M} = m | M = m], \quad (6.77)$$

vanishes asymptotically. Because the condition is stricter our converse bounds still hold even with this new definition of error; however, the random codes we constructed so far do not necessarily lead to a small maximum probability of error.

The following lemma allows us to construct codes that overcome this.

Lemma 6.12. *Given an $(\epsilon, |M|, 1)$ -average error channel code, there exists also a $(2\epsilon, \frac{|M|}{2}, 1)$ -maximum error channel code.*

The proof uses expurgation of bad codewords.

Proof. By definition of the $(\epsilon, |M|, 1)$ -average error channel code, we have

$$\sum_{m \in [|M|]} \frac{1}{|M|} P[\hat{M} = m | M = m] \leq \epsilon \quad (6.78)$$

Hence, there must be a subset $M_{\text{good}} \subseteq [|M|]$ of size at least $\frac{|M|}{2}$ with

$$P[\hat{M} = m | M = m] \leq 2\epsilon \quad \forall m \in M_{\text{good}} \quad (6.79)$$

as otherwise the inequality in Eq. (6.78) cannot hold. The codewords in M_{good} constitute an $(2\epsilon, \frac{|M|}{2}, 1)$ -maximum error channel code. \square

6.3 Source-channel separation theorem

We have until now covered the case where a message that is uniformly chosen from a set needs to be transmitted through the noisy channel. Does anything change when instead we want to transmit a general source? The setting is the same as with channel coding, except that now for each block length n we want to transmit a memoryless source given by i.i.d. $Z^n = (Z_1, Z_2, \dots, Z_n)$. A code for block length n is given by an encoder $e_n : \mathcal{Z}^n \rightarrow \mathcal{X}^n$ and a decoder $d_n : \mathcal{Y}^n \rightarrow \mathcal{Z}^n$ and our goal is to find a sequence of such codes that satisfy

$$\lim_{n \rightarrow \infty} P[\hat{Z}^n = Z^n] \rightarrow 0 \quad (6.80)$$

Here we want to show the following theorem

Theorem 6.13. *Given a DMS \mathbf{Z} and DMC \mathbf{W} , there exists a sequence of codes satisfying with asymptotically vanishing error if $H(Z) < I(W)$. Moreover, if $H(Z) > I(W)$ such a sequence of codes cannot exist.*

When $H(Z) < I(W)$ we can simply compress the source at a rate $R = H(Z) + \epsilon$ and then transmit it over the channel at the same rate $R = I(W) - \epsilon$, where we choose $\epsilon = \frac{1}{2}(I(W) - H(Z))$. That is, we first apply the encoder for source compression, transmit the compressed source through the channel using a maximum probability of error channel code, and finally decompress the source at the receiver. The error of such a scheme is simply the sum of the individual errors of the source compression code and the channel code, both of which vanish asymptotically by the source and channel coding theorems, respectively.

The second statement of this theorem, which is conceptually more interesting, shows that such a separate treatment of compression and channel coding is in fact optimal (at least when we only look at the first order asymptotics). We will only give a formal proof of the second statement.

Proof. Assume $H(Z) - I(W) = \nu > 0$. If there is a sequence of codes with asymptotically vanishing error then for every $\epsilon > 0$ there must be a block length n such that $P[\hat{Z}^n \neq Z^n] \leq \epsilon$. For such a code, by Fano's inequality, we have

$$H(Z^n) - I(Z^n : Y^n) = H(Z^n | Y^n) \leq H(Z^n | \hat{Z}^n) \leq 1 + \epsilon n \log |Z| \quad (6.81)$$

We can now evaluate $H(Z^n) = nH(Z)$ since the source is i.i.d., and furthermore

$$I(Z^n : Y^n) \leq I(W^n) = nI(W). \quad (6.82)$$

To verify the first inequality we simply note that $Z^n \rightarrow X^n \rightarrow Y^n$ form a Markov chain, which implies that $I(Z^n : Y^n) \leq I(X^n : Y^n)$. Maximising $I(X^n : Y^n)$ over all input distributions then yields the inequality. To verify the equality $I(W^n) = nI(W)$ for $n = 2$ we simply note that for two channel inputs X_1 and X_2 following any distribution $P_{X_1 X_2}$ and two channel outputs Y_1 and Y_2 produced by a DMC applied to X_1 and X_2 , respectively, we have

$$I(X_1, X_2 : Y_1, Y_2) = I(X_1 : Y_1, Y_2) + I(X_2 : Y_1, Y_2 | X_1) \quad (6.83)$$

$$= I(X_1 : Y_1) + I(X_2 : Y_2 | X_1) \quad (6.84)$$

$$\leq 2I(W). \quad (6.85)$$

Hence this equality in particular holds for the the distribution $P_{X_1 X_2}$ achieving $I(W^2)$. This argument can then be chained to show $I(W^n) \leq nI(W)$ and the other direction and thus equality clearly holds since we can take any i.i.d. distribution when optimising $I(W^n)$.

Finally, combining Eqs. (6.81) and (6.82) yields

$$\epsilon \log |Z| \geq H(Z) - I(W) - \frac{1}{n} = \nu - \frac{1}{n} \quad (6.86)$$

but since for large enough n the term on the right-hand side is strictly positive, ϵ is bounded away from zero, leading to a contradiction. \square

6.4 Gaussian channels

We have not discussed continuous variables in any detail and indeed all our proofs so far have assumed that the random variables take values in a finite alphabet. We will now however explore one very important channel that is continuous, the additive white Gaussian noise (AWGN) channel. This channel takes an input $X \in \mathbb{R}$ and outputs

$$Y = X + Z, \quad (6.87)$$

where Z follows a Gaussian distribution with mean 0 and standard deviation σ , and is independent of X . The channel behaviour can thus be characterised by the conditional pdf

$$w_{Y|X}(y|x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-x)^2}{2\sigma^2}}, \quad (6.88)$$

which is the Gaussian pdf with mean x and standard deviation σ_N . We can now ask the usual question about this channel—at what rate can we transmit information over it? It turns out that without further restrictions the answer to this that we can transmit as much information as we want, even through a single use of the channel. We simply map the messages to a lattice of values $x_m \in \mathbb{R}$ that are sufficiently separated so that even after the noise is added $w(y|x_m)$ and $w(y|x_{m'})$ only have small overlap for distinct messages m and m' . If the grid distance is chosen to be 6σ , for example, we will get a decoding error that is lower than 0.5% by such a construction. And we can get an arbitrarily small error by spreading the lattice further.

Question 6.14. *Formally construct such an encoder and decoder and compute the probability of error.*

In practical applications an AWGN for example arises when we encode information in an electromagnetic field, and X_i and Y_i for each channel use $i \in [n]$ are then simply amplitudes of the field. On the other hand, the energy stored in the field grows with the square of the amplitude and needs to be invested by the sender of the electromagnetic pulse. It is natural to restrict how much energy per channel-use, or power, is available at the source.⁴ Formally, this is done by requiring that every codeword $\mathbf{x} = (x_1, x_2, \dots, x_n)$ satisfies

$$\frac{1}{n} \sum_{i=1}^n x_i^2 \leq P \quad (6.89)$$

This communication channel models many practical channels, including wireless and satellite links. The noise may be due to a variety of microscopic reasons; however, the central limit ensures that collectively they resembles an additive noise with a Gaussian distribution.

We will now analyse the channel capacity of the AWGN channel under the above constraint. To do this, we however need to first introduce the notion of differential entropy.

⁴The nomenclature makes sense since power is energy per time unit, and channel uses are temporally separated in this context.

6.4.1 Differential entropy and mutual information

Definition 6.15. Let X be a real-valued continuous random variable with support on S and pdf p_X . The differential entropy of X is defined as

$$h(X) = - \int_S p_X(x) \log p_X(x) dx. \quad (6.90)$$

It is worth noting that this integral does not always exist and might in fact be infinite in many cases.

Question 6.16. For maths enthusiasts: Construct an example with a valid pdf for which the integral diverges and one for which it becomes negative.

A class of distributions for which it is relatively well-behaved is the uniform distribution, where $p_X(x) = \frac{1}{a}$ in an interval $[0, a]$ and zero elsewhere. In this case it is easy to verify that we have $h(X) = \log a$. An other interesting case is the case of Gaussian distribution with

$$p_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}. \quad (6.91)$$

In this case we can evaluate the differential entropy as follows

$$h(X) = - \int p_X(x) \log p_X(x) dx \quad (6.92)$$

$$= \int p_X(x) \left(\frac{(x-\mu)^2}{2\sigma^2} \log e + \log \sqrt{2\pi\sigma^2} \right) dx \quad (6.93)$$

$$= \mathbb{E}[(x-\mu)^2] \frac{\log e}{2\sigma^2} + \frac{1}{2} \log (2\pi\sigma^2) \quad (6.94)$$

$$= \frac{1}{2} \log (2e\pi\sigma^2). \quad (6.95)$$

But note that since the differential entropy can get negative it is hard to give it operational meaning.

One thing we can immediately observe is that the differential entropy is independent of the mean of X . This is true more generally: $h(X) = h(X + a)$ for any constant a , which can be verified by a simple change of variable. Note, however, that $h(cX) = h(X) + \log |c|$, so the entropy is not invariant under rescaling. This can be seen as an analogue of the invariance of the entropy of discrete random variables under relabelings.

We can define conditional entropy and mutual information analogously to the discrete case.

Definition 6.17. Let X and Y be real-valued continuous random variables with joint pdf p_{XY} . The conditional differential entropy of X given Y is defined as

$$h(X|Y) = - \int_S p_{XY}(x, y) \log p_{X|Y}(x|y) dx dy, \quad (6.96)$$

where $p_{X|Y}(x|y) = \frac{p_{XY}(x,y)}{p_Y(y)}$ is the conditional pdf and $p_Y(y) = \int p_{XY}(x,y)dx$ the marginal pdf on Y . Moreover, the mutual information between X and Y is defined as

$$I(X : Y) = \int p_{XY}(x, y) \log \frac{p_{XY}(x, y)}{p_X(x)p_Y(y)} dx dy. \quad (6.97)$$

As expected, we can again decompose $I(X : Y) = h(X) - h(X|Y)$, for example.

The mutual information for continuous variables is very naturally linked to mutual information for discrete variables. To see this, consider the discrete random variables X^Δ that takes values $x_r = r\Delta$ for $r \in \mathbb{Z}$ with probability

$$P_{X^\Delta}(x_r) = \int_{x_r-\Delta/2}^{x_r+\Delta/2} p_X(x) dx. \quad (6.98)$$

This is simply a discretised version of X , where everything in an interval of length Δ is coarse-grained into a single discrete value. For sufficiently small Δ continuity of $p_X(x)$ implies that $P_{X^\Delta}(x_r) \rightarrow \Delta \cdot p_X(x_r)$, and thus we find, under some mild regularity assumptions,

$$I(X : Y) = \int p_{XY}(x, y) \log \frac{p_{XY}(x, y)}{p_X(x)p_Y(y)} dx dy \quad (6.99)$$

$$= \lim_{\Delta \rightarrow 0} \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} \Delta^2 p_{XY}(x_i, y_j) \log \frac{p_{XY}(x_i, y_j)}{p_X(x_i)p_Y(y_j)} \quad (6.100)$$

$$= \lim_{\Delta \rightarrow 0} \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} P_{X^\Delta Y^\Delta}(x_i, y_j) \log \frac{P_{X^\Delta Y^\Delta}(x_i, y_j)}{P_{X^\Delta}(x_i)P_{Y^\Delta}(y_j)} \quad (6.101)$$

$$= \lim_{\Delta \rightarrow 0} I(X^\Delta : Y^\Delta), \quad (6.102)$$

where the first equality is simply by definition of a Riemann integral.

Question 6.18. *Argue that this implies that the mutual information is always non-negative even for continuous variables.*

A corresponding result does not hold for differential entropy, instead a different normalisation is required. This is evident simply from the fact that fine-graining a random variable will generally strictly increase its entropy, and thus the entropy would always diverge to infinity in the above limit.

Finally, we can define the relative entropy between two pdfs p_X and q_X as

$$D(p_X \| q_X) = \int p_X(x) \log \frac{p_X(x)}{q_X(x)} dx. \quad (6.103)$$

The same argument we used in Chapter 1, based on Jensen's inequality, reveals that

$$D(p_X \| q_X) \geq 0 \quad (6.104)$$

for all pairs of pdfs.

6.4.2 Channel coding theorem for the AWGN channel

Theorem 6.19. *The capacity of the AWGN channel \mathbf{W} with variance σ^2 and power constraint P is given by*

$$C(\mathbf{W}) = \frac{1}{2} \log \left(1 + \frac{P}{\sigma^2} \right). \quad (6.105)$$

The expression $\frac{P}{\sigma^2}$ is called the *signal-to-noise ratio* (SNR).

This expression does not remind us of the usual channel coding theorem, but this is only because it is already simplified for the channel at hand. Let us thus first show the following identity.

Lemma 6.20. *For an AWGN channel W with variance σ^2 and power constraint P , we have*

$$\max_{\substack{p_X \in \mathcal{P}(\mathbb{R}) \\ \mathbb{E}[X^2] \leq P}} I(X : Y) = \frac{1}{2} \log \left(1 + \frac{P}{\sigma^2} \right) \quad (6.106)$$

where $p_{XY}(x, y) = p_X(x)w_{Y|X}(y|x)$ as usual.

Due to the constraint on the codewords the optimisation is now not over all input distributions but only such distributions that satisfy the required bound on the expectation of X^2 .

Proof. We may rewrite $I(X : Y) = h(Y) - h(Y|X)$ where

$$h(Y|X) = h(X + Z|X) = h(Z) = \frac{1}{2} \log(2\pi e\sigma^2) \quad (6.107)$$

can be simplified immediately. We now make the following observations. We have

$$\mathbb{E}[Y^2] = \mathbb{E}[X^2 + 2XZ + Z^2] = \mathbb{E}[X^2] + 2\mathbb{E}[X]\mathbb{E}[Z] + \mathbb{E}[Z^2] \leq P + \sigma^2 \quad (6.108)$$

since X and Z are independent and $\mathbb{E}[Z] = 0$. So we have a bound on the variance of Y —does this allow us to conclude anything about its differential entropy?

We now argue that $h(Y)_p$ cannot exceed the entropy of a Gaussian ϕ_Y with the same standard deviation as p_Y . To see this, we write

$$0 \leq D(p_Y \parallel \phi_Y) = -h(Y)_p + \int p_Y(y) \log \frac{1}{\phi_Y(y)} dy \quad (6.109)$$

but since $\log \phi_Y(y)$ is quadratic in y , we may replace $p_Y(y)$ with $\phi_Y(y)$ in the latter integral. This yields $h(Y)_p \leq h(Y)_\phi$. Thus, we can in particular write

$$h(y)_p \leq \frac{1}{2} \log(2\pi e(P + \sigma^2)). \quad (6.110)$$

Hence, we can conclude that

$$I(X : Y) \leq \frac{1}{2} \log(2\pi e(P + \sigma^2)) - \frac{1}{2} \log(2\pi e\sigma^2) = \frac{1}{2} \log \left(1 + \frac{P}{\sigma^2} \right) \quad (6.111)$$

Finally, we can see that equality can be achieved by choosing p_X Gaussian with standard deviation P (and zero mean) as the input distribution. \square

It now remains to show that the channel capacity equals the power-restricted channel mutual information, i.e.,

$$C(\mathbf{W}) = \max_{\substack{p_X \in \mathcal{P}(\mathbb{R}) \\ \mathbb{E}[X^2] \leq P}} I(X : Y) \quad (6.112)$$

For the converse, we can actually largely build on the proof we already have — we will thus only sketch the argument here. We will ignore all technical aspects that come from the fact that we are now dealing with pdfs instead of pmfs and focus on the main idea. In the meta-converse, in the last step we introduced a maximisation over all channel input distributions: if we have restrictions on which codewords are allowed, we can also restrict the distribution there. Thus, when we apply the meta-converse for n channels, we now get

$$R \leq \min_{Q_{Y^n} \in \mathcal{P}(\mathcal{Y}^n)} \max_{P_{X^n}} \frac{1}{n} \log \frac{1}{\beta_\epsilon^*(P_{X^n Y^n} \| P_{X^n} \times Q_{Y^n})}, \quad (6.113)$$

where we optimise over pdfs P_{X^n} has support only on codewords x^n that satisfy the constraint $\frac{1}{n} \sum_{i=1}^n x_i^2 \leq P$. This will help us in Eq. (6.39), where we can now restrict the optimisation over sequences x^n with the above property as well. It remains now only to note that the empirical distributions corresponding to these sequences satisfy

$$\mathbb{E}[X^2] = \sum_x P_X^{x^k}(x) x^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 \leq P. \quad (6.114)$$

The converse can thus be generalised along these lines to the continuous case, although significant care has to be taken when we go from discrete to continuous variables.

A very rough sketch can also be drawn up for achievability. The critical part here is that we choose codewords X^n at random using the i.i.d. law $P_X(x)^n$ and for some distribution with $\mathbb{E}[X^2] \leq P - \epsilon$, their power consumptions satisfies

$$\frac{1}{n} \sum_{i=1}^n X_i^2 \rightarrow \mathbb{E}[X^2] \leq P - \epsilon, \quad (6.115)$$

as $n \rightarrow \infty$ by the weak law of large numbers, and, thus

$$P \left[\frac{1}{n} \sum_{i=1}^n X_i^2 \leq P \right] \rightarrow 1. \quad (6.116)$$

Thus, with high probability the codewords constructed in our random coding scheme satisfy the power constraint. In case an invalid codeword is chosen by the random process we simply discard it and count this as an error.

Chapter 7

Learning theory: Multiarmed stochastic bandits

[Week 13]

Intended learning outcomes:

- You understand how multiarmed stochastic bandits can be used to investigate exploration vs. exploitation tradeoffs.
- You can lower bound the minimax regret by construction adversarial distributions.

7.1 Problem setup and objective

Multiarmed stochastic bandits are an example of an unsupervised learning problem, where decisions have to be made under uncertainty. Multiarmed stochastic bandits are used to investigate tradeoffs between exploration and exploitation. Exploration here means that we want to learn properties of the various arms (namely their expected rewards) by observing samples so as to find the arm with the highest expected reward. Exploitation means that we mostly want to play the arms which we think have the highest expected rewards. But clearly at the start we do not know yet which arms have higher rewards, so some exploration is necessary before exploitation can occur.

Before we continue let us first formally define the problem. A *multiarmed stochastic bandit* is given by a set \mathcal{A} that is called the action set (which we here assume to be finite). Furthermore, the bandit is in an *environment*, a collection of pdfs $\nu = \{p_a : a \in \mathcal{A}\}$. The environment $\nu \in \mathcal{P}$ is unknown but taken from some set of potential environments \mathcal{P} . At each round $t \in [n]$ a learner chooses an action $a_t \in \mathcal{A}$ and receives a reward $X_t \in \mathbb{R}$ according to the (a priori unknown) pdf p_{a_t} .

- A *policy* π is a set of conditional probability distributions $\pi_t(A_t|X_1, A_1, \dots, X_{t-1}, A_{t-1})$ for $t \in [n]$ that determines the action the player takes at round t .

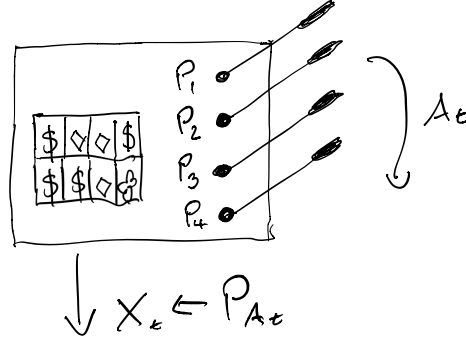


Figure 7.1: A bandit with 4 arms, each inducing a distribution P_i on the reward. Each round $t \in [n]$ an action A_t is chosen and the reward X_t is governed by the distribution P_{A_t} .

- The joint pdf of X_1, \dots, X_n and $A_1 \dots A_n$ is then given by

$$p(x_1, a_1, \dots, x_n, a_n) = \prod_{t=1}^n \pi_t(a_t | x_1, a_1, \dots, x_{t-1}, a_{t-1}) p_{a_t}(x_t) \quad (7.1)$$

- The *expected reward* of an action $a \in \mathcal{A}$ is defined as $\mu_a = \int p_{a_t}(x) dx$.
- The *maximal expected reward* is defined as $\mu^* = \max_{a \in \mathcal{A}} \mu_a$.
- The *expected regret* for a policy π on a bandit ν is defined as

$$R_n(\pi, \nu) := n\mu^* - \mathbb{E} \left[\sum_{t=1}^n X_t \right] \quad (7.2)$$

- The *worst-case regret* of a policy π over the environments \mathcal{P} is defined as

$$R_n^*(\pi, \mathcal{P}) := \sup_{\nu \in \mathcal{P}} R_n(\pi, \nu) \quad (7.3)$$

- Finally, the *minimax regret* for \mathcal{P} is defined as

$$R_n^*(\mathcal{P}) := \inf_{\pi} R_n^*(\pi, \mathcal{P}) = \inf_{\pi} \sup_{\nu \in \mathcal{P}} R_n(\pi, \nu). \quad (7.4)$$

Question 7.1. Show that we have $R_n(\pi, \nu) \geq 0$ for all π . Which policy π achieves the minimal regret for a fixed (and known) ν ?

Our objective is to find how the minimax regret scales with the number of samples n . As we have seen with the other problems we discussed in this module, there are two directions we have to approach this from. On the one hand, we might want to come up with good policies

that achieve a small worst-case regret. This is in some sense analogous to the “achievability” problem in source or channel coding, where we also need exhibit a code with the desired properties. The regret for any policy obviously gives an upper bound on the minimax regret. On the other hand, we also want to find lower bounds on the minimax regret that hold for all policies. That is analogous to the “converse” direction in channel or source coding. If the upper and lower bounds match then we know that our policy is optimal.

During this lecture we will only be able to derive a lower bound and you will have to believe me that this lower bound is in fact tight. Finding a good policy and analysing it is a bit outside the realm of traditional information theory and more similar to the task of algorithm design and analysis in computer science. If you are interested in this topic, please have a look at the recent book by Lattimore and Szepesvári [4].

7.2 A lower-bound on minimax regret

We will need a couple of lemmas for our proof, so we will discuss them first before we present and proof the main result.

In this section we will show the following theorem:

Theorem 7.2. *Let $k \geq 1$, $n \geq k - 1$ and let \mathcal{P} be a class of environments with k arms that allows for rewards with Gaussian pdfs $\mathcal{N}(\cdot; \nu, 1)$ for all $\nu \in [0, 1]$. Then,*

$$R_n^*(\mathcal{P}) \geq \frac{1}{27} \sqrt{(k-1)n}. \quad (7.5)$$

This means that, independent of the policy chosen, the worst-case regret scales at least as the square root of the number of trials and the number of arms, or $R_n^*(\pi, \mathcal{P}) \geq \frac{1}{27} \sqrt{(k-1)n}$ for any policy π . To show this, given any π , we will have to construct at least one environment ν that exhibits a regret that allows for this lower bound. So, what we will need to show is

$$\forall \pi, \exists \nu : R_n(\pi, \nu) \geq \frac{1}{27} \sqrt{(k-1)n}. \quad (7.6)$$

The proof of this statement follows in the next few sections.

7.2.1 Decomposing the regret

The first lemma allows us to decompose the regret:

Lemma 7.3. *Define $\Delta_a = \mu^* - \mu_a$ and let $T_a(n) = \sum_{t=1}^n 1\{A_t = a\}$ be the random variable counting the number of times the action $a \in \mathcal{A}$ is chosen. Then,*

$$R_n(\pi, \nu) = \sum_{a \in \mathcal{A}} \Delta_a \mathbb{E}[T_a(n)]. \quad (7.7)$$

Proof. Starting from the definition of the regret and the fact that $\sum_{a \in \mathcal{A}} \mathbb{E}[T_n(a)] = n$ and $\sum_{a \in \mathcal{A}} 1\{A_t = a\} = 1$, we have

$$R_n(\pi, \nu) := n\mu^* - \mathbb{E} \left[\sum_{t=1}^n X_t \right] \quad (7.8)$$

$$= \sum_{a \in \mathcal{A}} \mathbb{E}[T_n(a)]\mu^* - \mathbb{E} \left[\sum_{t=1}^n \sum_{a \in \mathcal{A}} 1\{A_t = a\} X_t \right] \quad (7.9)$$

$$= \sum_{a \in \mathcal{A}} \left(\mathbb{E}[T_n(a)]\mu^* - \mathbb{E} \left[\sum_{t=1}^n 1\{A_t = a\} X_t \right] \right). \quad (7.10)$$

Now we observe that in the latter expectation the action a is fixed and thus the random variables X_t are drawn independently from p_a . The expectation value thus factorises to

$$\mathbb{E} \left[\sum_{t=1}^n 1\{A_t = a\} X_t \right] = \mathbb{E} \left[\sum_{t=1}^n 1\{A_t = a\} \right] \mu_a = \mathbb{E}[T_n(a)]\mu_a, \quad (7.11)$$

which implies the desired result when plugged into (7.10) \square

7.2.2 Constructing worst-case environments

Without loss of generality we can take the action set to be $\mathcal{A} = [k]$. Let us introduce a vector of means $\mu \in [0, 1]^k$ and then the corresponding environment as

$$v_\mu = \{\mathcal{N}(\cdot; \mu_1, 1), \mathcal{N}(\cdot; \mu_2, 1), \dots, \mathcal{N}(\cdot; \mu_k, 1)\}. \quad (7.12)$$

To show the lower bound in Theorem 7.2, for every policy π , it is sufficient to find two vectors μ and μ' such that $\max\{R_n(\pi, v_\mu), R_n(\pi, v_{\mu'})\} \geq \frac{1}{27} \sqrt{(k-1)n}$, or, alternatively,

$$R_n(\pi, v_\mu) + R_n(\pi, v_{\mu'}) \geq \frac{2}{27} \sqrt{(k-1)n}. \quad (7.13)$$

Let us now construct such vectors for a fixed policy π . We choose $\mu = (\Delta, 0, \dots, 0)$ for some $\Delta \in (0, 1]$ still to be specified. Clearly the optimal policy for v_μ would always choose the first action. Now consider a run of the algorithm with policy π on environment v_μ . Then we define

$$i_* := \operatorname{argmin}_{i \in [k] \setminus \{1\}} \{\mathbb{E}[T_n(i)]\}, \quad (7.14)$$

the action that is chosen the least number of times (in expectation) by the policy π when run on the environment v_μ . Since $\sum_{i=1}^k \mathbb{E}[T_n(i)] = n$ we must have

$$\mathbb{E}[T_n(i_*)] \leq \frac{n}{k-1}. \quad (7.15)$$

by the Pigeonhole principle.¹

We can now define our alternative environment as

$$\mu' = (\Delta, 0, \dots, 0, \underbrace{2\Delta}_{\text{at position } i_*}, 0, \dots, 0). \quad (7.16)$$

Note that this is a worst-case scenario in the sense that if we run policy π and think that the environment is μ instead of μ' then our regret would be maximal as we play i_* the least often. Obviously our policy should be clever enough so that we at some point learn that the environment is μ' and adapt our action choices accordingly, but as we will see by choosing Δ small we can make it very difficult to distinguish the two cases.

According to Lemma 7.3 we can decompose and then bound the two regrets as

$$R_n(\pi, v_\mu) = \Delta \sum_{i \in [k] \setminus \{1\}} \mathbb{E}[T_i(n)] = \Delta (n - \mathbb{E}[T_1(n)]) , \quad (7.17)$$

$$R_n(\pi, v_{\mu'}) = \Delta \mathbb{E}'[T_1(n)] + 2\Delta \sum_{i \in [k] \setminus \{1, i_*\}} \mathbb{E}'[T_i(n)] \geq \Delta \mathbb{E}'[T_1(n)] . \quad (7.18)$$

Here we used \mathbb{E}' to denote the expectation under the distribution p' induced by π and $v_{\mu'}$. We can further bound

$$\mathbb{E}[T_1(n)] \leq \frac{n}{2} p \left[T_1(n) < \frac{n}{2} \right] + n p \left[T_1(n) \geq \frac{n}{2} \right] = \frac{n}{2} \left(1 + p \left[T_1(n) \geq \frac{n}{2} \right] \right) \quad (7.19)$$

and, using Markov's inequality, $\mathbb{E}'[T_1(n)] \geq \frac{n}{2} p' \left[T_1(n) \geq \frac{n}{2} \right]$. This yields

$$R_n(\pi, v_\mu) \geq \frac{n\Delta}{2} \left(1 - p \left[T_1(n) \geq \frac{n}{2} \right] \right) = \frac{n\Delta}{2} p \left[T_1(n) < \frac{n}{2} \right] \quad (7.20)$$

$$R_n(\pi, v_{\mu'}) \geq \frac{n\Delta}{2} p' \left[T_1(n) \geq \frac{n}{2} \right] . \quad (7.21)$$

And, thus,

$$R_n(\pi, v_\mu) + R_n(\pi, v_{\mu'}) \geq \frac{n\Delta}{2} \left(p \left[T_1(n) < \frac{n}{2} \right] + p' \left[T_1(n) \geq \frac{n}{2} \right] \right) . \quad (7.22)$$

7.2.3 Lower-bounding the regret

This is really where the information theory tools come in! The next lemma gives us a lower bound on the sum of the probabilities of two complementary events, evaluated on two (generally different) distributions. When these distributions are similar we expect that the sum of probabilities is close to 1.

¹Indeed, if on the contrary $\mathbb{E}[T_n(i_*)] > \frac{n}{k-1}$, then we get a contradiction since

$$n = \sum_{i=1}^k \mathbb{E}[T_n(i)] \geq \sum_{i=2}^k \mathbb{E}[T_n(i)] \geq (k-1) \mathbb{E}[T_n(i_*)] > (k-1) \frac{n}{k-1} = n .$$

Lemma 7.4 (Bretagnolle-Huber inequality). *Let p and q be two pdfs for the same random variable X taking values on \mathcal{X} . For any $A \subset \mathcal{X}$, we have*

$$p(A) + q(A^c) \geq \frac{1}{2} \exp(-D(p\|q)) , \quad (7.23)$$

where $D(p\|q) = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} dx$.

Proof. We first note that

$$p(A) + q(A^c) = \int_A p(x) dx + \int_{A^c} q(x) dx \geq \int_{\mathcal{X}} \min\{p(x), q(x)\} dx . \quad (7.24)$$

Furthermore, using the Cauchy-Schwartz inequality, we find

$$\left(\int_{\mathcal{X}} \sqrt{p(x)q(x)} dx \right)^2 = \left(\int_{\mathcal{X}} \sqrt{\min\{p(x), q(x)\} \max\{p(x), q(x)\}} dx \right)^2 \quad (7.25)$$

$$\leq \left(\int_{\mathcal{X}} \min\{p(x), q(x)\} dx \right) \left(\int_{\mathcal{X}} \max\{p(x), q(x)\} dx \right) \quad (7.26)$$

$$\leq 2 \int_{\mathcal{X}} \min\{p(x), q(x)\} dx . \quad (7.27)$$

Thus, combining this with Eq. (7.24), we have

$$p(A) + q(A^c) \geq \frac{1}{2} \exp \left(2 \log \int_{\mathcal{X}} p(x) \sqrt{\frac{q(x)}{p(x)}} dx \right) . \quad (7.28)$$

Finally, using Jensen's inequality for the logarithm, we arrive at

$$p(A) + q(A^c) \geq \frac{1}{2} \exp \left(2 \int_{\mathcal{X}} p(x) \log \sqrt{\frac{q(x)}{p(x)}} dx \right) \quad (7.29)$$

$$= \frac{1}{2} \exp \left(- \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} dx \right) \quad (7.30)$$

$$= \frac{1}{2} \exp(-D(p\|q)) . \quad (7.31)$$

□

We can now continue with our derivation of a lower bound. Applying this Lemma to our bound in Eq. (7.22), we find

$$R_n(\pi, v_\mu) + R_n(\pi, v_{\mu'}) \geq \frac{n\Delta}{4} \exp(-D(p\|p')) , \quad (7.32)$$

and it then remains to evaluate this relative entropy for our two distributions.

Lemma 7.5 (Divergence decomposition lemma). *Let π be any policy and p and p' be the distributions induce by the environments (p_1, p_2, \dots, p_k) and $(p'_1, p'_2, \dots, p'_k)$, respectively. Then,*

$$D(p\|p') = \sum_{i=1}^k \mathbb{E}[T_i(n)] D(p_i\|p'_i). \quad (7.33)$$

Proof. Recall that $p(x_1, a_1, \dots, x_n, a_n) = \prod_{t=1}^n \pi_t(a_t|x_1, a_1, \dots, x_{t-1}, a_{t-1})p_{a_t}(x_t)$. Then,

$$D(p\|p') = \mathbb{E} \left[\log \frac{p(X_1, A_1, \dots, X_n, A_n)}{p'(X_1, A_1, \dots, X_n, A_n)} \right] \quad (7.34)$$

$$= \mathbb{E} \left[\sum_{t=1}^n \log \frac{p_{A_t}(X_t)}{p'_{A_t}(X_t)} \right] \quad (7.35)$$

$$= \sum_{t=1}^n \sum_{i=1}^k P[A_t = i] \mathbb{E} \left[\log \frac{p_{A_t}(X_t)}{p'_{A_t}(X_t)} \middle| A_t = i \right] \quad (7.36)$$

$$= \sum_{i=1}^k \mathbb{E} \left[\sum_{t=1}^n \delta_{A_t, i} \right] \mathbb{E} \left[\log \frac{p_i(X_t)}{p'_i(X_t)} \middle| A_t = i \right] \quad (7.37)$$

$$= \sum_{i=1}^k \mathbb{E}[T_i(n)] D(p_i\|p'_i), \quad (7.38)$$

where we used the law of total expectation to get the third equality and in the penultimate step we used that X_t is drawn independently from p_i once $A_t = i$ is fixed. \square

Applying this to our particular situation we find that

$$D(p\|p') = \sum_{i=1}^k \mathbb{E}[T_i(n)] D(\mathcal{N}(\cdot; \mu_i, 1) \| \mathcal{N}(\cdot; \mu'_i, 1)) \quad (7.39)$$

$$= \mathbb{E}[T_{i^*}(n)] \frac{(2\Delta)^2}{2} \quad (7.40)$$

$$\leq \frac{2\Delta^2 n}{k-1}, \quad (7.41)$$

where we realised that μ_i and μ'_i only differ at $i = i^*$ and evaluated the relative entropy for two normal Gaussian distributions with different means. In the last step we used Eq. (7.15).

Question 7.6. *Verify that $D(\mathcal{N}(\cdot; \mu, 1) \| \mathcal{N}(\cdot; \mu', 1)) = \frac{(\mu - \mu')^2}{2}$.*

Substituting this into Eq. (7.32), we find

$$R_n(\pi, v_\mu) + R_n(\pi, v_{\mu'}) \geq \frac{n\Delta}{4} \exp \left(-\frac{2\Delta^2 n}{k-1} \right). \quad (7.42)$$

Choosing now $\Delta = \sqrt{\frac{k-1}{4n}}$ yields

$$R_n(\pi, v_\mu) + R_n(\pi, v_{\mu'}) \geq \sqrt{n(k-1)} \underbrace{\frac{1}{8} \exp\left(-\frac{1}{2}\right)}_{\geq \frac{2}{27}}, \quad (7.43)$$

which is what we set out to show.

Bibliography

- [1] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. Wiley, 1991.
- [2] T. S. Han. *Information-Spectrum Methods in Information Theory*. Applications of Mathematics. Springer, 2002.
- [3] Masahito Hayashi and Hiroshi Nagaoka. General Formulas for Capacity of Classical-Quantum Channels. *IEEE Transactions on Information Theory*, 49(7):1753–1768, jul 2003.
- [4] T. Lattimore and C. Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- [5] Yury Polyanskiy, H. Vincent Poor, and Sergio Verdú. Channel Coding Rate in the Finite Blocklength Regime. *IEEE Transactions on Information Theory*, 56(5):2307–2359, may 2010.
- [6] C. Shannon. A Mathematical Theory of Communication. *Bell System Technical Journal*, 27:379–423, 1948.
- [7] Maurice Sion. On General Minimax Theorems. *Pacific Journal of Mathematics*, 8:171–176, 1958.
- [8] Marco Tomamichel and Vincent Y. F. Tan. A Tight Upper Bound for the Third-Order Asymptotics for Most Discrete Memoryless Channels. *IEEE Transactions on Information Theory*, 59(11):7041–7051, nov 2013.