

Information Theory

University of Amsterdam, **Master of Logic**, Fall 2016

Yfke Dulek and Christian Schaffner

Copyright © 2016 Yfke Dulek and Christian Schaffner

These lecture notes are licenced under the [Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International](#) license.

The source is available on [github](#) and we are grateful for any feedback including typos and suggestions for improvements! Feel free to send us [pull requests](#).

Parts of the text is based on, or taken verbatim, from ‘The Mathematical Theory of Information, and Applications’ (version 2.0) by Ronald Cramer and Serge Fehr [\[CF\]](#). The lay-out is based on the Legrand Orange Book (version 2.1.1), licenced under the [Creative Commons License BY-NC-SA 3.0](#).

Contents

1	Probability and Entropy	1
1.1	Preliminaries: Probability Theory	1
1.2	Some Important Distributions	4
1.3	Jensen's Inequality	5
1.4	Shannon Entropy	6
1.5	Conditional Entropy	8
1.6	Mutual Information	11
1.7	Relative entropy	12
1.8	Entropy Diagrams	13
1.9	Further Reading	13
2	Source Coding	15
2.1	Symbol Codes	15
2.2	Kraft's Inequality	17
2.3	Shannon's Source-Coding Theorem	20
2.4	Huffman Codes	21
2.5	Arithmetic Codes	23
2.6	Asymptotic Equipartition Property (AEP)	25
3	Perfectly Secure Encryption	31
3.1	Security	32
3.2	One-Time Pad	32
3.3	Lower Bound On Key Length	34

4	Stochastic Processes	35
4.1	Finite Markov Chains	35
4.2	Stochastic Processes	37
4.2.1	Types Of Stochastic Processes	38
4.3	Entropy rate	42
5	Noisy-Channel Coding	45
5.1	Channels	45
5.2	Codes	47
5.2.1	Repetition codes	48
5.2.2	The [7,4] Hamming code	49
5.2.3	Generalization: Linear Codes	51
6	Zero-Error Channel Coding	55
6.1	Confusability graphs	56
6.2	Multiple channel uses	57
6.3	Shannon capacity	60
7	The Noisy-Channel Coding Theorem	63
7.1	Channel capacity	63
7.2	Joint Asymptotic Equipartition Property	65
7.3	Shannon's Noisy-Channel Coding Theorem	69
7.4	Shannon's Noisy-Channel Coding Theorem: Converse	71

Chapter 1: Probability and Entropy

1.1 Preliminaries: Probability Theory

For this course, we will only be concerned with discrete probabilities. This section formalizes some notions you should already be familiar with: probability spaces, events and probability distributions.

Definition 1.1.1 — Probability space. A (discrete) probability space (Ω, \mathcal{F}, P) consists of a discrete, non-empty *sample space* Ω , an *event space* $\mathcal{F} \subseteq \mathcal{P}(\Omega)$ and a *probability measure* P which is a function $P : \Omega \rightarrow \mathbb{R}_{\geq 0}$ that satisfies

$$\sum_{\omega \in \Omega} P(\omega) = 1.$$

The event space \mathcal{F} is required to be non-empty and closed under intersection, union and complements. For convenience, we will most often assume that \mathcal{F} equals the powerset $\mathcal{P}(\Omega)$ of Ω , i.e. it contains all possible subsets of events, and therefore fulfils the required properties.

Definition 1.1.2 — Event. An event \mathcal{A} is an element of the event space $\mathcal{F} \subseteq \mathcal{P}(\Omega)$, i.e. a subset \mathcal{A} of the sample space Ω . Its probability is defined as

$$P[\mathcal{A}] := \sum_{\omega \in \mathcal{A}} P(\omega),$$

where by convention $P[\emptyset] = 0$.

As a notational convention, we write $P[\mathcal{A}, \mathcal{B}]$ for $P[\mathcal{A} \cap \mathcal{B}]$, and $P[\overline{\mathcal{A}}]$ for $P[\Omega \setminus \mathcal{A}]$.

Exercise 1.1.3 Prove the following identities (for arbitrary events $\mathcal{A}, \mathcal{B} \subseteq \Omega$):

$$P[\overline{\mathcal{A}}] = 1 - P[\mathcal{A}] \tag{1.1}$$

$$P[\mathcal{A} \cup \mathcal{B}] = P[\mathcal{A}] + P[\mathcal{B}] - P[\mathcal{A}, \mathcal{B}] \tag{1.2}$$

$$P[\mathcal{A}] = P[\mathcal{A}, \mathcal{B}] + P[\mathcal{A}, \overline{\mathcal{B}}]. \tag{1.3}$$

It is often useful to consider the probability of an event *given* that some other event happened:

Definition 1.1.4 — Conditional probability. For events \mathcal{A} and \mathcal{B} with $P[\mathcal{A}] > 0$, the conditional probability of \mathcal{B} given \mathcal{A} is defined as

$$P[\mathcal{B}|\mathcal{A}] := \frac{P[\mathcal{A}, \mathcal{B}]}{P[\mathcal{A}]}.$$

Example 1.1.5 — Fair die. We throw a six-sided fair die once, and consider the number that comes up. The sample space for this experiment is $\Omega = \{1, 2, 3, 4, 5, 6\}$, with event space $\mathcal{F} = \mathcal{P}(\Omega)$ and probability measure $P[i] = \frac{1}{|\Omega|} = \frac{1}{6}$ for all $i \in \Omega$ (this is a **uniform** probability measure). Consider the events $\mathcal{A} = \{2, 4, 6\}$ and $\mathcal{B} = \{3, 6\}$. Using the formulas in Definitions 1.1.2 and 1.1.4, we can compute the following probabilities:

$$P[\mathcal{A}] = \frac{1}{2} \quad (\text{the outcome is even})$$

$$P[\mathcal{B}] = \frac{1}{3} \quad (\text{the outcome is a multiple of 3})$$

$$P[\mathcal{A}, \mathcal{B}] = P[\{6\}] = \frac{1}{6} \quad (\text{the roll is even and a multiple of 3})$$

$$P[\mathcal{A}|\mathcal{B}] = \frac{1/6}{1/3} = \frac{1}{2} \quad (\text{the roll is even, given that it is a multiple of 3})$$

$$P[\mathcal{B}|\mathcal{A}] = \frac{1/6}{1/2} = \frac{1}{3} \quad (\text{the roll is a multiple of 3, given that it is even})$$

This example shows that in general, $P[\mathcal{A}|\mathcal{B}]$ is *not equal* to $P[\mathcal{B}|\mathcal{A}]$.

Definition 1.1.6 — Discrete Random Variable (RV). Let (Ω, \mathcal{F}, P) be a discrete probability space. A random variable X is a function $X : \Omega \rightarrow \mathcal{X}$ where \mathcal{X} is a set, and we may assume it to be discrete.

A *real* random variable is one whose image is contained in \mathbb{R} . A (The *image* and the *range* of a random variable X are given by the image and the range of X in the function-theoretic sense.) The image of a *binary* random variable is a set $\{x_0, x_1\}$ with only two elements.

Definition 1.1.7 — Probability distribution. Let X be a random variable. The probability distribution of X is the function $P_X : \mathcal{X} \rightarrow [0, 1]$ defined as

$$P_X(x) := P[X = x],$$

where $X = x$ denotes the event $\{\omega \in \Omega \mid X(\omega) = x\}$.

Alternatively, one can write $P_X(x) = P(X^{-1}(x))$ to express that the probability of x is precisely the P -measure of the pre-image of x under the random variable X .

Exercise 1.1.8 Verify that $(\mathcal{X}, \mathcal{P}(\mathcal{X}), P_X)$ is itself a probability space. ■

We say that P_X is a **uniform** distribution if the associated probability measure is uniform, i.e. $P_X(x) = \frac{1}{|\mathcal{X}|}$. The **support** of a random variable or a probability distribution is defined as $\text{supp}(P_X) := \{x \in \mathcal{X} \mid P_X(x) > 0\}$, the points of the range which have strictly positive probability. We often slightly abuse notation and write $\text{supp}(X)$ instead.

When given two or more random variables defined on the same probability space, we can consider the probability that each of the variables take on a certain value:

Definition 1.1.9 — Joint probability distribution. Let X and Y be two random variables defined on the same probability space, with respective ranges \mathcal{X} and \mathcal{Y} . The pair XY is a random variable with probability distribution $P_{XY} : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$ given by

$$P_{XY}(x, y) := P[X = x, Y = y].$$

This definition naturally extends to three and more random variables. Unless otherwise stated, a collection of random variables is assumed to be defined on the same (implicit) probability space, so that their joint distribution is always well-defined.

If $P_{XY} = P_X \cdot P_Y$, in the sense that $P_{XY}(x, y) = P_X(x)P_Y(y)$ for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$, then the random variables X and Y are said to be **independent**. If a set of variables X_1, \dots, X_n are all mutually independent and all have the same distribution (i.e. $P_{X_i} = P_{X_j}$ for all i, j), then they are **independent and identically distributed**, or **i.i.d.**

From a joint distribution, we can always find out the “original” (or **marginal**) distribution of one of the random variables (for example, X) by **marginalizing** out the variable that we want to discard (for example, Y):

$$P_X(x) = \sum_{y \in \mathcal{Y}} P_{XY}(x, y). \quad (1.4)$$

This marginalization process also works with more than two random variables.

Like events, probability distributions can also be conditioned on probabilistic events:

Definition 1.1.10 — Conditional probability distribution. If \mathcal{A} is an event with $P[\mathcal{A}] > 0$, then the conditional probability distribution of X given \mathcal{A} is given by

$$P_{X|\mathcal{A}}(x) = \frac{P[X = x, \mathcal{A}]}{P[\mathcal{A}]}.$$

If Y is another random variable and $P_Y(y) > 0$, then we write

$$P_{X|Y}(x|y) := P_{X|Y=y}(x) = \frac{P_{XY}(x, y)}{P_Y(y)}$$

for the conditional distribution of X , given $Y = y$.

Note that again, both $(\mathcal{X}, P_{X|\mathcal{A}})$ and $(\mathcal{X}, P_{X|Y=y})$ themselves form probability spaces. Note also that if X and Y are independent, then

$$P_{X|Y}(x|y) = \frac{P_{XY}(x, y)}{P_Y(y)} = \frac{P_X(x) \cdot P_Y(y)}{P_Y(y)} = P_X(x), \quad (1.5)$$

which aligns well with our intuition of independent variables: the distribution of X remains unchanged when Y is fixed to a specific value.

Example 1.1.11 — Fair die (continued). Consider again the throw of a six-sided fair die as in Example 1.1.5. Let the random variable X describe the number of (distinct) integer divisors for the outcome, that is

$$X(1) = 1 \quad X(2) = 2 \quad X(3) = 2 \quad X(4) = 3 \quad X(5) = 2 \quad X(6) = 4$$

X is a real random variable, with range $\mathcal{X} = \{1, 2, 3, 4\}$. The associated probability distribution

is

$$P_X(1) = P[\{1\}] = \frac{1}{6}, \quad P_X(2) = P[\{2, 3, 5\}] = \frac{1}{2}, \quad P_X(3) = P[\{4\}] = \frac{1}{6}, \quad P_X(4) = P[\{6\}] = \frac{1}{6}.$$

If we now condition on the event $\mathcal{A} = \{2, 4, 6\}$ (the outcome of the die being even), we get that

$$P_{X|\mathcal{A}}(1) = 0, \quad P_{X|\mathcal{A}}(2) = \frac{1}{3}, \quad P_{X|\mathcal{A}}(3) = \frac{1}{3}, \quad P_{X|\mathcal{A}}(4) = \frac{1}{3}$$

If X is a random variable and $f: \mathcal{X} \rightarrow \mathcal{Y}$ is a surjective function, then $f(X)$ is a random variable, defined by composing the map f with the map X . Its image is \mathcal{Y} . Clearly,

$$P_{f(X)}(y) = \sum_{x \in \mathcal{X}: f(x)=y} P_X(x). \quad (1.6)$$

For example, $1/P_X(X)$ denotes the real random variable obtained from another random variable X by composing with the map $1/P_X$ that assigns $1/P_X(x) \in \mathbb{R}$ to $x \in \mathcal{X}$.

Definition 1.1.12 — Expectation. The expectation of a *real* random variable X is defined as

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

Note that if X is not real, then we can still consider the expectation of some function $f: \mathcal{X} \rightarrow \mathbb{R}$, where

$$\mathbb{E}[f(X)] = \sum_{x \in \mathcal{X}} P_X(x) \cdot f(x). \quad (1.7)$$

Definition 1.1.13 — Variance. The variance of a *real* random variable X is defined as

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

The variation is a measure for the deviation of the mean. Hoeffding's inequality (here stated for binary random variables) states that for a list of i.i.d. random variables, the average of the random variables is close to the expectation, except with very small probability. We state it here without proof.

Theorem 1.1.14 — Hoeffding's inequality. Let X_1, \dots, X_n be independent and identically distributed binary random variables with $P_{X_i}(0) = 1 - \mu$ and $P_{X_i}(1) = \mu$, and thus $\mathbb{E}[X_i] = \mu$. Then, for any $\delta > 0$

$$P\left[\sum_i X_i > (\mu + \delta) \cdot n\right] \leq \exp(-2\delta^2 n).$$

1.2 Some Important Distributions

- The distribution of a biased coin with probability $P_X(1) = p$ to land heads, and a probability of $P_X(0) = 1 - p$ to land tails is called **Bernoulli(p) distribution**. Its entropy is given by the binary entropy $h(p)$. The expected value is $\mathbb{E}[X] = p$ and the variance is $\text{Var}[X] = p(1 - p)$.
- When n coins X_1, X_2, \dots, X_n are flipped independently and every X_i is Bernoulli(p) distributed, let $S = \sum_{i=1}^n X_i$ be their sum, i.e. the number of heads in n throws of a biased coin. Then, S has the **binomial(n, p) distribution**:

$$P_S(k) = \binom{n}{k} p^k (1 - p)^{n-k} \quad \text{where } k = 0, 1, 2, \dots, n. \quad (1.8)$$

From simple properties of the expected value and variance, one can show that $\mathbb{E}[S] = np$ and $\text{Var}[S] = np(1-p)$.

- The **geometric(p) distribution** of a random variable Y is defined as the number of times one has to flip a Bernoulli(p) coin before it lands heads:

$$P_Y(k) = (1-p)^{k-1}p \quad \text{where } k = 1, 2, 3, \dots \quad (1.9)$$

There is another variant of the geometric distribution used in the literature, where one excludes the final success event of landing heads in the counting:

$$P_Z(k) = (1-p)^k p \quad \text{where } k = 0, 1, 2, 3, \dots \quad (1.10)$$

While the expected values are slightly different, namely $\mathbb{E}[Y] = \frac{1}{p}$ and $\mathbb{E}[Z] = \frac{1-p}{p}$, their variances are the same $\text{Var}[Y] = \text{Var}[Z] = \frac{1-p}{p^2}$.

1.3 Jensen's Inequality

In the following, let \mathcal{D} be an interval in \mathbb{R} .

Definition 1.3.1 — Convex and concave functions. The function $f : \mathcal{D} \rightarrow \mathbb{R}$ is convex if for all $x_1, x_2 \in \mathcal{D}$ and all $\lambda \in [0, 1] \subset \mathbb{R}$:

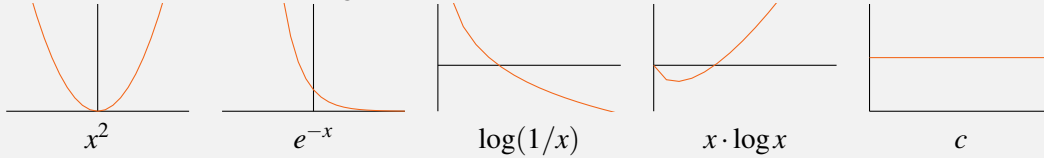
$$\lambda f(x_1) + (1-\lambda)f(x_2) \geq f(\lambda x_1 + (1-\lambda)x_2).$$

The function f is *strictly* convex if equality only holds when $\lambda \in \{0, 1\}$ or when $x_1 = x_2$.

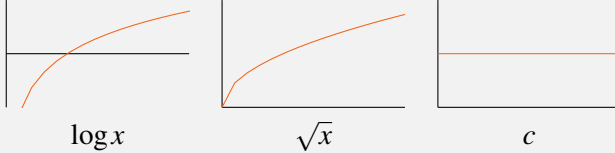
The function f is (strictly) concave if the function $-f$ is (strictly) convex.

Intuitively, a function is convex if any straight line drawn between two points $f(x_1)$ and $f(x_2)$ lies above the graph of f entirely. For a concave function, such a line must lie entirely beneath the graph.

Example 1.3.2 The following functions are convex (for $c \in \mathbb{R}$):



The following functions are concave (for $c \in \mathbb{R}$):



The following establishes a more formal method of proving the convexity of a function.

Proposition 1.3.3 Let $f : \mathcal{D} \rightarrow \mathbb{R}$. If \mathcal{D} is open, and for all $x \in \mathcal{D}$, the second order derivative $f''(x)$ exists and is non-negative (positive), then f is convex (strictly convex).

We omit the proof, which can be found in, for example, [CF] (Lemma 1).

Theorem 1.3.4 — Jensen's inequality. Let $f : \mathcal{D} \rightarrow \mathbb{R}$ be a convex function, and let $n \in \mathbb{N}$. Then for any $p_1, \dots, p_n \in \mathbb{R}_{\geq 0}$ such that $\sum_{i=1}^n p_i = 1$ and for any $x_1, \dots, x_n \in \mathcal{D}$ it holds that

$$\sum_{i=1}^n p_i f(x_i) \geq f\left(\sum_{i=1}^n p_i x_i\right).$$

If f is strictly convex and $p_1, \dots, p_n > 0$, then equality holds iff $x_1 = \dots = x_n$.

In particular, if X is a real random variable whose image \mathcal{X} is contained in \mathcal{D} , then

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X]),$$

and if f is strictly convex, equality holds iff there is a $c \in \mathcal{X}$ such that $X = c$ with probability 1.

Proof. The proof is by induction. The case $n = 1$ is trivial, and the case $n = 2$ is identical to the very definition of convexity. Suppose that we have already proved the claim up to $n - 1 \geq 2$. Assume, without loss of generality, that $p_n < 1$. Then:

$$\begin{aligned} \sum_{i=1}^n p_i f(x_i) &= p_n f(x_n) + \sum_{i=1}^{n-1} p_i f(x_i) \\ &= p_n f(x_n) + (1 - p_n) \sum_{i=1}^{n-1} \frac{p_i}{1 - p_n} f(x_i) \\ &\geq p_n f(x_n) + (1 - p_n) f\left(\sum_{i=1}^{n-1} \frac{p_i}{1 - p_n} x_i\right) && \text{(induction hypothesis)} \\ &\geq f\left(p_n x_n + (1 - p_n) \sum_{i=1}^{n-1} \frac{p_i}{1 - p_n} x_i\right) && \text{(definition of convexity)} \\ &= f\left(p_n x_n + \sum_{i=1}^{n-1} p_i x_i\right) \\ &= f\left(\sum_{i=1}^n p_i x_i\right). \end{aligned} \tag{1.11}$$

That proves the claim. As for the strictness claim, if x_1, \dots, x_n are not all identical, then either x_1, \dots, x_{n-1} are not all identical and the first inequality is strict by induction hypothesis, or $x_1 = \dots = x_{n-1} \neq x_n$ so that the second inequality is strict by the definition of convexity. ■

1.4 Shannon Entropy

In this section, we explore a measure for the amount of uncertainty of random variables. Consider some probabilistic event \mathcal{A} that occurs with probability $P[\mathcal{A}]$ for some probability measure P . The **surprisal value** $\log \frac{1}{P[\mathcal{A}]}$ indicates how surprised we should be when the event \mathcal{A} occurs: events with small probabilities yield high surprisal values, and vice versa. An event that occurs with certainty ($P_X(\mathcal{A}) = 1$) yields a surprisal value of 0. For a random variable X , we consider the *expected* surprisal value to be an indicator of how much uncertainty is contained in the variable, or how much information is gained by revealing the outcome. This expected surprisal value is more commonly known as the (Shannon) entropy¹ of a random variable:

¹Shannon once said: *My greatest concern was what to call it. I thought of calling it information, but the word was overly used, so I decided to call it uncertainty. When I discussed it with John von Neumann, he had a better idea. Von Neumann told me: "You should call it entropy, for two reasons. In the first place, your uncertainty function has been*

Definition 1.4.1 — Entropy. Let X be a random variable with image \mathcal{X} . The (Shannon) entropy $H(X)$ of X is defined as

$$H(X) := \mathbb{E} \left[\log \frac{1}{P_X(X)} \right] = \sum_{x \in \mathcal{X}} P_X(x) \cdot \log \frac{1}{P_X(x)} = - \sum_{x \in \mathcal{X}} P_X(x) \cdot \log P_X(x),$$

with the convention that the log function represents the *binary* logarithm \log_2 . As another convention, for $x \in \mathcal{X}$ with $P_X(x) = 0$, the corresponding argument in the summation is declared 0 (which is justified by taking a limit).

It is important to realize that the entropy of X is a function (solely) of the *distribution* P_X of X . However, it is customary to write $H(X)$ instead of the formally correct $H(P_X)$.

Exercise 1.4.2 Prove that $\lim_{p \rightarrow 0} p \log(p) = 0$. ■

Proposition 1.4.3 — Positivity. Let X be a random variable with image \mathcal{X} . Then

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

Equality on the left-hand side holds iff there exists $x \in \mathcal{X}$ with $P_X(x) = 1$ (and thus $P_X(x') = 0$ for all $x' \neq x$). Equality on the right-hand side holds iff $P_X(x) = 1/|\mathcal{X}|$ for all $x \in \mathcal{X}$.

Proof. The function $f : \mathbb{R}_{>0} \rightarrow \mathbb{R}$ defined by $y \mapsto \log y$ is strictly concave on $\mathbb{R}_{>0}$. Thus, by Jensen's inequality:

$$H(X) = \sum_{x \in \mathcal{X}} P_X(x) \cdot \log \frac{1}{P_X(x)} \leq \log \left(\sum_{x \in \mathcal{X}} 1 \right) = \log(|\mathcal{X}|). \quad (1.12)$$

Furthermore, since we may restrict the sum to all x with $P_X(x) > 0$, equality holds if and only if $\log(1/P_X(x)) = \log(1/P_X(x'))$, and thus $P_X(x) = P_X(x')$, for all $x, x' \in \mathcal{X}$.

Finally, for the characterization of the lower bound, it is obvious that $H(X) = 0$ if $P_X(x) = 1$ for some x , and, on the other hand, if $H(X) = 0$ then for any x with $P_X(x) > 0$ it must be that $\log(1/P_X(x)) = 0$ and hence $P_X(x) = 1$. ■

For a binary random variable X with image $\mathcal{X} = \{x_0, x_1\}$ and probabilities $P_X(x_0) = p$ and $P_X(x_1) = 1 - p$, we can write $H(X) = h(p)$, where h denotes the binary entropy function:

Definition 1.4.4 — Binary entropy function h . The binary entropy function is defined for $0 < q < 1$ as

$$h(q) := q \log \frac{1}{q} + (1 - q) \log \frac{1}{1 - q},$$

and is defined as $h(q) = 0$ for $q = 0$ or $q = 1$. The graph of h on the interval $[0, 1]$ is shown in Figure 1.1.

This binary entropy function is used, for example, to measure the entropy of a biased coin flip.

Example 1.4.5 Consider a random variable X with $\mathcal{X} = \{a, b, c\}$ and $P_X(a) = \frac{1}{2}$, $P_X(b) =$

used in statistical mechanics under that name, so it already has a name. In the second place, and more important, nobody knows what entropy really is, so in a debate you will always have the advantage."



Figure 1.1: The binary entropy function h as a function of the probability p .

$P_X(c) = \frac{1}{4}$. The entropy of X is

$$H(X) = \frac{1}{2} \log 2 + \frac{1}{4} \log 4 + \frac{1}{4} \log 4 = \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = \frac{3}{2}. \quad (1.13)$$

Another approach to computing the entropy of X by coming up with an appropriate underlying probability space (Ω, P) : we toss a fair coin twice, giving $\Omega = \{hh, ht, th, tt\}$ and $P(\omega) = \frac{1}{4}$ for all $\omega \in \Omega$. Then we define the function $X : \Omega \rightarrow \mathcal{X}$ as

$$X(hh) = X(ht) = a, \quad X(th) = b, \quad X(tt) = c.$$

This yields the correct distribution P_X . The following computation now leads to the entropy of X :

$$H(X) = h\left(\frac{1}{2}\right) + \frac{1}{2}h(0) + \frac{1}{2}h\left(\frac{1}{2}\right) = \frac{3}{2}. \quad (1.14)$$

The first coin toss determines whether the outcome is a (on heads h) or something else (on tails t). On heads, the second coin toss does not give any more information, whereas on tails, the second coin toss still decides between outcome b and outcome c . In general, the entropy of a random variable with probabilities p_1, \dots, p_n can be expressed as

$$\begin{aligned} H(p_1, \dots, p_k, p_{k+1}, \dots, p_n) &= h(p_1 + \dots + p_k) + \\ &\quad (p_1 + \dots + p_k)H\left(\frac{p_1}{p_1 + \dots + p_k} + \dots + \frac{p_k}{p_1 + \dots + p_k}\right) + \\ &\quad (p_{k+1} + \dots + p_n)H\left(\frac{p_{k+1}}{p_{k+1} + \dots + p_n} + \dots + \frac{p_n}{p_{k+1} + \dots + p_n}\right). \end{aligned} \quad (1.15)$$

1.5 Conditional Entropy

Let X be a random variable and \mathcal{A} an event. Applying Definition 1.4.1 to the conditional probability distribution $P_{X|\mathcal{A}}$ allows us to naturally define the entropy of X conditioned on the event \mathcal{A} , which

leads to the following notion:

Definition 1.5.1 — Conditional entropy. Let X and Y be random variables, with respective images \mathcal{X} and \mathcal{Y} . The conditional entropy $H(X|Y)$ of X given Y is defined as

$$H(X|Y) := \sum_{y \in \mathcal{Y}} P_Y(y) \cdot H(X|Y=y),$$

with the convention that the corresponding argument in the summation is 0 for $y \in \mathcal{Y}$ with $P_Y(y) = 0$, and where

$$H(X|\mathcal{A}) := \sum_{x \in \mathcal{X}} P_{X|\mathcal{A}}(x) \cdot \log \frac{1}{P_{X|\mathcal{A}}(x)}.$$

Note that conditional entropy $H(X|Y)$ is not the entropy of a probability distribution but an expectation: the average uncertainty about X when given Y . The following bound expresses that (on average!) additional information, i.e. knowing Y , can only *decrease* the uncertainty.

Proposition 1.5.2 Let X and Y be random variables with respective images \mathcal{X} and \mathcal{Y} . Then

$$0 \leq H(X|Y) \leq H(X)$$

Equality on the left-hand side holds iff X is determined by Y , i.e., for all $y \in \mathcal{Y}$, there is an $x \in \mathcal{X}$ such that $P_{X|Y}(x|y) = 1$. Equality on the right-hand side holds iff X and Y are independent.

Proof. The lower bound follows trivially from the definition and from Proposition 1.4.3, and so does the characterization of when $H(X|Y) = 0$. For the upper bound, note that

$$H(X|Y) = \sum_y P_Y(y) \sum_x P_{X|Y}(x|y) \log \frac{1}{P_{X|Y}(x|y)} = \sum_{x,y} P_{XY}(x,y) \log \frac{P_Y(y)}{P_{XY}(x,y)} \quad (1.16)$$

and

$$H(X) = \sum_x P_X(x) \log \frac{1}{P_X(x)} = \sum_{x,y} P_{XY}(x,y) \log \frac{1}{P_X(x)} \quad (1.17)$$

where the last equality is derived by marginalization. Note that in both expressions, we may restrict the sum to those pairs (x,y) with $P_{XY}(x,y) > 0$. Using Jensen's inequality, it follows that

$$\begin{aligned} H(X|Y) - H(X) &= \sum_{\substack{x,y \\ P_{XY}(x,y) > 0}} P_{XY}(x,y) \log \frac{P_X(x)P_Y(y)}{P_{XY}(x,y)} \\ &\leq \log \left(\sum_{\substack{x,y \\ P_{XY}(x,y) > 0}} P_X(x)P_Y(y) \right) \leq \log \left(\left(\sum_{x \in \mathcal{X}} P_X(x) \right) \left(\sum_{y \in \mathcal{Y}} P_Y(y) \right) \right) \\ &= \log 1 = 0. \end{aligned} \quad (1.18)$$

Note that in the second inequality, we replaced the summation over all (x,y) with $P_{XY}(x,y) > 0$ by the summation over all $(x,y) \in \mathcal{X} \times \mathcal{Y}$. Inequality then follows by the monotonicity of the logarithm function.

For the first inequality in (??), equality holds if and only if $P_{XY}(x,y) = P_X(x)P_Y(y)$ for all (x,y) with $P_{XY}(x,y) > 0$, and for the second inequality, equality holds if and only if $P_{XY}(x,y) = 0$ implies $P_X(x)P_Y(y) = 0$ for any $x \in \mathcal{X}$ and $y \in \mathcal{Y}$. It follows that $H(X|Y) = H(X)$ if and only if $P_{XY}(x,y) = P_X(x)P_Y(y)$ for all $(x,y) \in \mathcal{X} \times \mathcal{Y}$. ■

Proposition 1.5.3 — Chain Rule. Let X and Y be random variables. Then

$$H(XY) = H(X) + H(Y|X).$$

Proof. The chain rule is a simple matter of rewriting:

$$\begin{aligned} H(XY) &= - \sum_{x,y} P_{XY}(x,y) \log P_{XY}(x,y) \\ &= - \sum_{x,y} P_{XY}(x,y) \log (P_X(x) P_{Y|X}(y|x)) \\ &= - \sum_{x,y} P_{XY}(x,y) \log P_X(x) - \sum_{x,y} P_{XY}(x,y) \log P_{Y|X}(y|x) \\ &= - \sum_x P_X(x) \log P_X(x) - \sum_x P_X(x) \sum_y P_{Y|X}(y|x) \log P_{Y|X}(y|x) \\ &= H(X) + H(Y|X). \end{aligned} \tag{1.19}$$

This was to be shown. ■

The following inequality, also known as the ‘independence bound’, follows from the fact that $H(Y|X) \leq H(Y)$:

Corollary 1.5.4 — Subadditivity.

$$H(XY) \leq H(X) + H(Y).$$

Equality holds iff X and Y are independent.

Note that applying Definition 1.5.1 to the conditional distribution $P_{XY|\mathcal{A}}$ naturally defines $H(X|Y, \mathcal{A})$, the entropy of X given Y and conditioned on the event \mathcal{A} . Since the entropy is a function of the distribution of a random variable, the chain rule also holds when conditioning on an event \mathcal{A} . Furthermore, it holds that

$$H(X|YZ) = \sum_z P_Z(z) H(X|Y, Z=z), \tag{1.20}$$

which is straightforward to verify. With this observation, it is easy to see that the chain rule generalizes as follows.

Corollary 1.5.5 Let X , Y and Z be random variables. Then

$$H(XY|Z) = H(X|Z) + H(Y|XZ).$$

Inductively applying the (generalized) chain rule implies that for any sequence X_1, \dots, X_n of random variables:

$$H(X_1 \cdots X_n) = H(X_1) + H(X_2|X_1) + \cdots + H(X_n|X_{n-1} \cdots X_1). \tag{1.21}$$

Example 1.5.6 Consider the binary random variables X and Y , with joint distribution

$$P_{XY}(00) = \frac{1}{2}, \quad P_{XY}(01) = \frac{1}{4}, \quad P_{XY}(10) = 0, \quad P_{XY}(11) = \frac{1}{4}.$$

By marginalization, we find that $P_X(0) = \frac{3}{4}$ and $P_X(1) = \frac{1}{4}$, while $P_Y(0) = P_Y(1) = \frac{1}{2}$. This

allows us to make the following computations:

$$H(XY) = \frac{1}{2} \log 2 + \frac{1}{4} \log 4 + \frac{1}{4} \log 4 = \frac{3}{2} \quad (1.22)$$

$$H(X) = h\left(\frac{1}{4}\right) = h\left(\frac{3}{4}\right) \approx 0.81 \quad (1.23)$$

$$H(Y) = h\left(\frac{1}{2}\right) = 1 \quad (1.24)$$

$$H(X|Y) = H(XY) - H(Y) = \frac{1}{2} \quad (1.25)$$

$$H(Y|X) = H(XY) - H(X) \approx 0.69 \quad (1.26)$$

We also could have computed $H(X|Y)$ and $H(Y|X)$ directly through the definition of conditional entropy.

Note that for this specific distribution, $H(X|Y = 1) > H(X)$. It is important to remember that Proposition 1.5.2 only holds on average, not for specific values of Y . Note also that in this example, $H(X|Y) \neq H(Y|X)$.

1.6 Mutual Information

Definition 1.6.1 — Mutual information. Let X and Y be random variables. The mutual information $I(X;Y)$ of X and Y is defined as

$$I(X;Y) = H(X) - H(X|Y).$$

Thus, in a sense, mutual information reflects the reduction in uncertainty about X when given Y . Note the following properties of the mutual information:

$$I(X;Y) = H(X) + H(Y) - H(XY) \quad (\text{by chain rule}) \quad (1.27)$$

$$I(X;Y) = I(Y;X) \quad (\text{"symmetry"}) \quad (1.28)$$

$$I(X;Y) \geq 0 \quad (\text{by subadditivity}) \quad (1.29)$$

$$I(X;Y) = 0 \text{ iff } X \text{ and } Y \text{ are independent} \quad (1.30)$$

$$I(X;X) = H(X) \quad (\text{"self-information"}) \quad (1.31)$$

Applying Definition 1.6.1 to the conditional distribution $P_{XY|\mathcal{A}}$ naturally defines $I(X;Y|\mathcal{A})$, the mutual information of X and Y conditioned on the event \mathcal{A} .

Definition 1.6.2 — Conditional mutual information. Let X, Y, Z be random variables. Then the conditional mutual information of X and Y given Z is defined as

$$I(X;Y|Z) = \sum_z P_Z(z) I(X;Y|Z=z),$$

with the convention that the corresponding argument in the summation is 0 for z with $P_Z(z) = 0$.

Conditional mutual information has properties similar to the ones we saw above:

$$I(X;Y|Z) = I(Y;X|Z) \quad (1.32)$$

$$I(X;Y|Z) \geq 0 \quad (1.33)$$

$$I(X;Y|Z) = 0 \text{ iff } X \text{ and } Y \text{ are independent given } Z \quad (1.34)$$

Furthermore, the previous bounds $H(X) \geq 0$, $H(X|Y) \geq 0$, and $I(X;Y) \geq 0$, can all be seen as special cases of $I(X;Y|Z) \geq 0$. These bounds, and any bound they imply, are called **Shannon inequalities**.

It is important to realize that $I(X;Y|Z)$ may be larger or smaller than (or equal to) $I(X;Y)$. The following is easy to verify (and is sometimes used as definition of $I(X;Y|Z)$).

Proposition 1.6.3 Let X, Y, Z be random variables. Then

$$I(X;Y|Z) = H(X|Z) - H(X|YZ).$$

By this result, we obtain:

Corollary 1.6.4 — Chain rule for mutual information. Let W, X, Y and Z be random variables. Then

$$I(WX;Y|Z) = I(X;Y|Z) + I(W;Y|XZ).$$

Proof. The proof is a matter of writing out definitions and applying the generalized chain rule.

$$\begin{aligned} I(WX;Y|Z) &= H(WX|Z) - H(WX|YZ) \\ &= (H(X|Z) + H(W|XZ)) - (H(X|YZ) + H(W|XYZ)) \\ &= H(X|Z) - H(X|YZ) + H(W|XZ) - H(W|XYZ) \\ &= I(X;Y|Z) + I(W;Y|XZ). \end{aligned} \tag{1.35}$$

■

1.7 Relative entropy

A measure that is related to the mutual information is the relative entropy: it reflects how different two distributions are:

Definition 1.7.1 — Relative entropy. The relative entropy (or: **Kullback-Leibler divergence**) of two probability distributions P and Q over the same \mathcal{X} is defined by

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

where by convention, $\log \frac{p}{0} = \infty$ for all p .

Note that if $Q(x) = 0$ for some x with $P(x) > 0$, then $D(P||Q) = \infty$.

Exercise 1.7.2 Show that $I(X;Y) = D(P_{XY}||P_X \cdot P_Y)$. ■

This exercise, combined with the equality condition in Theorem 1.7.3 below, shows that the mutual information is a measure of ‘how independent’ the variables X and Y are: if $P_{XY} = P_X \cdot P_Y$, the variables are independent and their mutual information is zero.

Theorem 1.7.3 — Information inequality. For any two probability distributions P and Q defined on the same \mathcal{X} ,

$$D(P||Q) \geq 0.$$

Equality holds if and only if $P = Q$.

Proof. Left as an exercise. Hint: use Jensen's inequality. ■

Even though relative entropy is always nonnegative, it is not a proper distance measure, because it is not symmetric and does not satisfy the triangle inequality.

1.8 Entropy Diagrams

We finish this chapter by visually summing up the relations between entropy, joint entropy, conditional entropy, mutual information, and conditional mutual information. For two and three random variables, the relations between these different information-theoretic measures can be nicely represented by means of a Venn-diagram-like **entropy diagram**. The case of two random variables is illustrated in Figure 1.2 (left). From the diagram, one can for instance easily read off the relations $H(X|Y) \leq H(X)$, $I(X;Y) = H(X) + H(Y) - H(XY)$ etc. The case of three random variables is illustrated in Figure 1.2 (right). Also here, one can easily read off all the relations between the information-theoretic measures, like for instance $H(X|YZ) = H(X) - I(X;Z) - I(X;Y|Z)$, which is a relation that is otherwise not immediately obvious.



Figure 1.2: Entropy diagram for two (left) and three (right) random variables. The areas encompassed by the dotted lines represent $H(XY)$ and $H(XYZ)$, respectively.

One subtlety with the entropy diagram for three random variables is that the “area in the middle”, $R(X;Y;Z) = I(X;Y) - I(X;Y|Z)$, may be *negative*.

1.9 Further Reading

- Sections 2.1, 2.2, 3.1-3.3 of [CF]
- Sections 2.1, 2.2, 2.6 of [CT]
- For more background on probability theory, check for instance the [lecture script](#) of the Master of Logic course “Basic Probability:Theory” by Philip Schulz and Christian Schaffner.

Chapter 2: Source Coding

Suppose we sample x from a distribution P_X with image \mathcal{X} . In the context of data compression, P_X is typically called a **source** that emits value $x \in \mathcal{X}$ with probability $P_X(x)$. We want to compress (or encode) symbols x sampled from P_X in such a way that we can later decompress (or decode) it reliably, without losing any information about the value x .



A counting argument shows that it is possible to encode the elements of \mathcal{X} by bit strings of length n , where $n = \lceil \log(|\mathcal{X}|) \rceil$: we simply list all elements of \mathcal{X} , and use the (binary) index of x in the list as its encoding. Thus, to store or to transmit an element $x \in \mathcal{X}$, n bits of information always suffice. However, if not all $x \in \mathcal{X}$ are equally likely according to P_X , one should be able to exploit this to achieve codes with shorter *average* length. The idea is to use encodings of varying lengths, assigning shorter codewords to the elements in \mathcal{X} that have higher probabilities, and vice versa. The question we answer in this chapter is: how short can such a code be (on average over repeated samples x from P_X)?

We explore both **lossless** codes (where we want to recover the original data with certainty) and **lossy** codes (where with small probability, the data is lost).

2.1 Symbol Codes

We start by investigating codes that encode a source one symbol at a time. Later on, we will also see codes that group the source symbols together into blocks.

Definition 2.1.1 — Binary symbol code. Let P_X be the distribution of a random variable X (with image \mathcal{X}). A binary symbol code for P_X is an injective function $C : \mathcal{X} \rightarrow \{0, 1\}^*$.

The **extended** code $C^* : \mathcal{X}^* \rightarrow \{0, 1\}^*$ is defined by concatenation:

$$C^*(x_1, \dots, x_n) := C(x_1) \cdots C(x_n).$$

Here, $\mathcal{X}^* = \bigcup_{n \in \mathbb{N}} \mathcal{X}^n \cup \{\perp\}$, and \perp is the empty string.

We often refer to the set of codewords, $\mathcal{C} = \text{im}(C)$, as code and leave the actual encoding function C implicit.

The injectivity of C ensures that we can always uniquely decode $C(x)$. However, if one transmits a sequence $x_1, \dots, x_m \in \mathcal{X}$ (or stores them “sequentially”) by sending the concatenation $C(x_1, \dots, x_m)$, ambiguities may arise, namely in cases where it is possible to parse this long string in two consistent but different ways. Indeed, injectivity of the encoding function per se does not rule out that there exists a positive integer m' and elements $x'_1, \dots, x'_{m'} \in \mathcal{X}$ such that $C(x_1) \cdots C(x_m) = C(x'_1) \cdots C(x'_{m'})$. Of course, this problem can be circumvented by introducing a special separation symbol. However, such a symbol might not be available, and maybe even more importantly, even if an additional symbol *is* available, then one can often create a better code by using it as an ordinary code symbol (in addition to 0 and 1) rather than as a special separation symbol. This is why it is interesting to study the following class of symbol codes:

Definition 2.1.2 — Uniquely decodable code. A binary symbol code $C : \mathcal{X} \rightarrow \{0, 1\}^*$ is uniquely decodable if C^* is injective as well.

One convenient way to guarantee that a code is unique decodable is to require it to be prefix-free:

Definition 2.1.3 — Prefix-free code. A binary symbol code $C : \mathcal{X} \rightarrow \{0, 1\}^*$ is prefix-free (or: **instantaneous**) if for all $x, x' \in \mathcal{X}$ with $x \neq x'$, $C(x)$ is *not* a prefix of $C(x')$.

With a prefix-free encoding, the elements x_1, \dots, x_m can be uniquely recovered from $C(x_1) \cdots C(x_m)$, simply by reading the encoding from left to right one bit at a time: by prefix-freeness it will remain unambiguous as reading continues when the current word terminates and the next begins. This is a loose argument for the following:

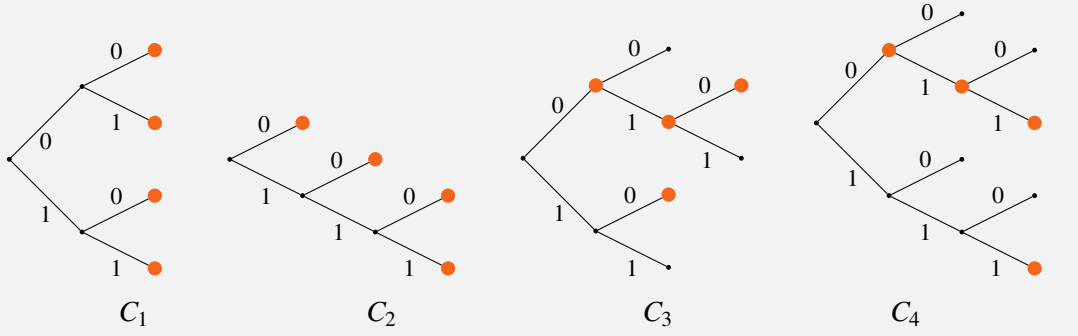
Proposition 2.1.4 If a code \mathcal{C} is prefix-free and $\mathcal{C} \neq \{\perp\}$ then \mathcal{C} is uniquely decodable.

The other direction does not hold: uniquely decodable codes need not be prefix-free. A prefix-free code is appealing from an efficiency point of view, as it allows to decode “on the fly”. For a general uniquely decodable code one may possibly have to inspect all bits in the entire string before being able to even recover the first word.

Example 2.1.5 The following are three codes for the source P_X , with $\mathcal{X} = \{a, b, c, d\}$:

x	$P_X(x)$	$C_1(x)$	$C_2(x)$	$C_3(x)$	$C_4(x)$
a	$1/2$	00	0	0	0
b	$1/4$	01	10	010	01
c	$1/8$	10	110	01	011
d	$1/8$	11	111	10	111

These codes can be visualised as binary trees, with marked codewords, as follows:



C_1 and C_2 are prefix-free, and therefore also uniquely decodable. C_3 is not uniquely decodable, as $C_3(ad) = C_3(b)$. C_4 is not prefix-free, but it is uniquely decodable, since it can be decoded from right to left (it is “postfix-free”). Note that the binary trees for the prefix-free codes C_1, C_2 only have codewords at the leaves. (The same holds for the postfix-free code C_4).

For efficiency reasons, we are often interested in the average (expected) length of a code C :

Definition 2.1.6 — Average length. Let $\ell(s)$ denote the length of a string $s \in \{0, 1\}^*$. The (average) length of a code C for a source P_X is defined as

$$\ell_C(P_X) := \mathbb{E}[\ell(C(X))] = \sum_{x \in \mathcal{X}} P_X(x) \ell(C(x)).$$

Example 2.1.7 For the codes from Example 2.1.5, we obtain the following average codeword lengths: $\ell_{C_1}(P_X) = 2$, $\ell_{C_2}(P_X) = \ell_{C_4}(P_X) = \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 = \frac{7}{4} = 1.75$ and $\ell_{C_3}(P_X) = \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 3 + \frac{1}{8} \cdot 2 + \frac{1}{8} \cdot 2 = \frac{7}{4} = 1.75$. We see that the codes C_2, C_3, C_4 have a smaller average codeword length, but C_2 and C_4 are preferred over C_3 because their unique decodability.

Notice that the individual codeword lengths of codes C_2 and C_4 correspond exactly to the surprisal values of P_X in bits, e.g. $\ell(C_2(b)) = \ell(C_4(b)) = 2 = -\log P_X(b)$. Therefore, the computations of the entropy $H(X)$ and of the average code length $\ell_{C_2}(P_X)$ are exactly the same, and we have that $H(X) = \ell_{C_2}(P_X) = \ell_{C_4}(P_X)$. We will see in Section 2.3 below that this property characterizes optimal codes.

Definition 2.1.8 — Minimal code length. The minimal code length of a source P_X is defined as

$$\ell_{\min}(P_X) := \min_{C \in \mathfrak{C}} \ell_C(P_X)$$

where \mathfrak{C} is some class of codes, for example the set of all prefix-free codes (resulting in $\ell_{\min}^{\text{p.f.}}$), or the set of all uniquely decodable codes (resulting in $\ell_{\min}^{\text{u.d.}}$).

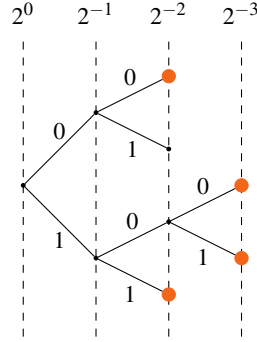
2.2 Kraft's Inequality

As argued, prefix-freeness is a nice feature, but it is also considerably more restrictive than mere unique decodability; thus, it is natural to ask: how much do we lose (in terms of the average codeword length) by requiring the encoding to be prefix-free rather than merely uniquely decodable? Surprisingly, the answer is: *nothing*. In this section, we will show that the length of an optimal prefix-free code and the length of an optimal uniquely decodable code coincide. In the next section, we will see that these lengths are essentially given by the Shannon entropy.

Theorem 2.2.1 — Kraft's inequality. There exists a prefix-free code with image $\mathcal{C} = \{c_1, \dots, c_m\}$ and codeword lengths $\ell_i := \ell(c_i)$, if and only if

$$\sum_{i=1}^m 2^{-\ell_i} \leq 1.$$

Proof. For the forward direction, suppose we have a prefix-free code \mathcal{C} with image $\mathcal{C} = \{c_1, \dots, c_m\}$ and codeword lengths $\ell_i := \ell(c_i)$. View this code as a tree, with codewords only on the leaves (but not necessarily all the leaves), and assign a weight of 2^{-d} to every node in the tree at depth d (including the leaves):



Note that the weight of each node is exactly the sum of the weight of its direct children, and thereby that the weight of the root is exactly the weight of all of the leaves. Since every codeword c_i resides on a leaf of depth ℓ_i (but not all leaves are necessarily occupied), the weight of the root is *at least* the sum of all the codeword weights:

$$\sum_{i=1}^m 2^{-\ell_i} \leq 2^0 = 1. \quad (2.1)$$

For the backward direction, we build a code $\mathcal{C} = \{c_1, \dots, c_m\}$ with $\ell(c_i) = \ell_i$ by selecting the appropriate leaves of a binary tree as codewords, assigning the most ‘expensive’ (i.e. those with small depth) first. We proceed by induction on the number of codewords, m :

For $m = 1$, the construction is clear: we can assign any string of length ℓ_1 to represent the single codeword.

For $m > 1$, assume without loss of generality that $\ell_1 \leq \dots \leq \ell_{m-1} \leq \ell_m$. We will first try to build a code with $(m-1)$ code words, of lengths $\ell_1, \dots, \ell_{m-2}, (\ell_{m-1} - 1)$. In order to be able to invoke the induction hypothesis, we do need to check that

$$\left(\sum_{i=1}^{m-2} 2^{-\ell_i} \right) + 2^{-(\ell_{m-1}-1)} \leq 1. \quad (2.2)$$

This can be seen by first noting that

$$1 > 1 - 2^{-\ell_m} \geq \sum_{i=1}^{m-1} 2^{-\ell_i} = \sum_{i=1}^m \frac{2^{\ell_{m-1}-\ell_i}}{2^{\ell_{m-1}}} = \frac{\sum_{i=1}^{m-1} 2^{\ell_{m-1}-\ell_i}}{2^{\ell_{m-1}}}. \quad (2.3)$$

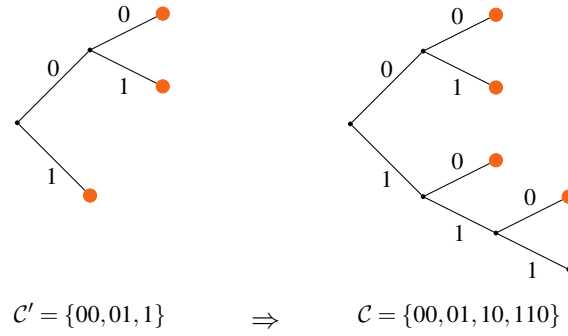
Thus, on the right-hand side, we end up with a fraction of the form $\frac{a}{b} < 1$, where a and b are both integers, with $a < b$. In this case it must follow that $a + 1 \leq b$, and hence $\frac{a+1}{b} \leq 1$. So we get

$$1 \geq \frac{\left(\sum_{i=1}^{m-1} 2^{\ell_{m-1}-\ell_i} \right) + 1}{2^{\ell_{m-1}}} = \left(\sum_{i=1}^{m-1} 2^{-\ell_i} \right) + 2^{-\ell_{m-1}}$$

$$\begin{aligned}
&= \left(\sum_{i=1}^{m-2} 2^{-\ell_i} \right) + 2 \cdot 2^{-\ell_{m-1}} \\
&= \left(\sum_{i=1}^{m-2} 2^{-\ell_i} \right) + 2^{-(\ell_{m-1}-1)},
\end{aligned} \tag{2.4}$$

as desired. So we invoke the induction hypothesis to create a prefix-free code $\mathcal{C}' = \{c'_1, \dots, c'_{m-1}\}$ with lengths $\ell_1, \dots, \ell_{m-2}, (\ell_{m-1} - 1)$. The new code \mathcal{C} is then constructed by setting $c_i = c'_i$ for all $i \leq m-2$, and furthermore setting $c_{m-1} = c'_{m-1}|0$ and $c_m = c'_{m-1}|10 \dots 0$, padding with enough zeroes to achieve $\ell(c_m) = \ell_m$. This new code is necessarily also prefix-free.

The following image illustrates the induction step for $\ell_1 = \ell_2 = \ell_3 = 2$ and $\ell_4 = 3$. The code is constructed from a prefix-free code \mathcal{C}' with code lengths 2, 2 and 1 by replacing the bottom codeword (of length 1) with two new codewords (of lengths 2 and 3).



■

A stronger version of Kraft's inequality holds as well, this time for uniquely decodable codes:

Theorem 2.2.2 — McMillan inequality. For a uniquely decodable code with image $\mathcal{C} = \{c_1, \dots, c_m\}$ and codeword lengths $\ell_i := \ell(c_i)$, it holds that

$$\sum_{i=1}^m 2^{-\ell_i} \leq 1.$$

Proof. Let \mathcal{C} be a uniquely decodable code as in the theorem statement. We can write

$$S := \sum_{c \in \mathcal{C}} \frac{1}{2^{\ell(c)}} = \sum_{\ell=L_{\min}}^{L_{\max}} \frac{n_{\ell}}{2^{\ell}} \tag{2.5}$$

where $L_{\min} = \min_{c \in \mathcal{C}} \ell(c)$, $L_{\max} = \max_{c \in \mathcal{C}} \ell(c)$, and $n_{\ell} = |\{c \in \mathcal{C} \mid \ell(c) = \ell\}|$. Furthermore, for any $k \in \mathbb{N}$, consider the k th power of S ,

$$S^k = \sum_{c_1, \dots, c_k \in \mathcal{C}^k} \frac{1}{2^{\ell(c_1) + \dots + \ell(c_k)}} = \sum_{\ell=kL_{\min}}^{kL_{\max}} \frac{n_{\ell}^{(k)}}{2^{\ell}} \tag{2.6}$$

where $n_{\ell}^{(k)}$ is defined as $n_{\ell}^{(k)} = |\{(c_1, \dots, c_k) \in \mathcal{C}^k \mid \sum_i \ell(c_i) = \ell(c_1 | \dots | c_k) = \ell\}|$. Note that

$$n_{\ell}^{(k)} = \sum_{x \in \{0,1\}^{\ell}} |\{(c_1, \dots, c_k) \in \mathcal{C}^k \mid c_1 | \dots | c_k = x\}| \leq \sum_{x \in \{0,1\}^{\ell}} 1 = 2^{\ell} \tag{2.7}$$

where the inequality follows from the unique decodability of \mathcal{C} . Thus, we can conclude that

$$S^k \leq (L_{\max} - L_{\min}) \cdot k \quad (2.8)$$

for all $k \in \mathbb{N}$, so S^k grows at most linearly in k , from which follows that $S \leq 1$ (for if not, S^k would grow exponentially in k). ■

Kraft's and McMillan's inequality together lead to the conclusion that the lengths of an optimal prefix-free code and an optimal uniquely decodable code coincide:

Corollary 2.2.3 Let P_X be a source. For every uniquely decodable code C , there exists a prefix-free code C' such that $\ell_C(P_X) = \ell_{C'}(P_X)$. Hence,

$$\ell_{\min}^{\text{p.f.}}(P_X) = \ell_{\min}^{\text{u.d.}}(P_X).$$

From now on, we will just write $\ell_{\min}(P_X)$ to denote either of these measures for average length. A code C for which $\ell_C(P_X) = \ell_{\min}(P_X)$ is called **optimal** for the source P_X .

2.3 Shannon's Source-Coding Theorem

We now know that prefix-free codes can achieve the same minimal code lengths for a source P_X as the more general class of uniquely decodable codes. How small is this minimal code length in general? In this section we explore the following relation between the minimal code length and the entropy of the source:

Theorem 2.3.1 — Shannon's source-coding theorem (for symbol codes). For any source P_X , we have the following bounds:

$$H(X) \leq \ell_{\min}(P_X) \leq H(X) + 1.$$

Proof. The proof relies on Kraft's inequality (Section 2.2). Let C be a code, and write ℓ_x for $\ell(C(x))$ as a notational convenience. For the lower bound, we have that

$$\begin{aligned} H(X) - \ell_C(P_X) &= - \sum_{x \in \mathcal{X}} P_X(x) \log(P_X(x)) - \sum_{x \in \mathcal{X}} P_X(x) \ell_x \\ &= \sum_{x \in \mathcal{X}} P_X(x) \left(-\log(P_X(x)) - \log(2^{\ell_x}) \right) \\ &= \sum_{x \in \mathcal{X}} P_X(x) \log \left(\frac{1}{P_X(x) \cdot 2^{\ell_x}} \right) \\ &\leq \log \left(\sum_{x \in \mathcal{X}} \frac{1}{2^{\ell_x}} \right) && \text{(by Jensen's inequality)} \\ &\leq \log(1) = 0 && \text{(by Kraft's inequality)} \end{aligned} \quad (2.9)$$

For the upper bound, let us denote by ℓ_x the surprisal value in bits rounded up to the next integer, i.e. for any $x \in \mathcal{X}$,

$$\ell_x := \left\lceil \log \frac{1}{P_X(x)} \right\rceil, \quad (2.10)$$

and note that

$$\sum_{x \in \mathcal{X}} 2^{-\ell_x} \leq \sum_{x \in \mathcal{X}} 2^{-\log \frac{1}{P_X(x)}} = \sum_{x \in \mathcal{X}} P_X(x) = 1. \quad (2.11)$$

Therefore, by Kraft's inequality, there exists a prefix-free code C such that $\ell(C(x)) = \ell_x$ for all $x \in \mathcal{X}$. This code satisfies

$$\begin{aligned}
 \ell_C(P_X) &= \sum_{x \in \mathcal{X}} P_X(x) \ell_x \\
 &\leq \sum_{x \in \mathcal{X}} P_X(x) \left(\log \frac{1}{P_X(x)} + 1 \right) \\
 &= - \sum_{x \in \mathcal{X}} P_X(x) \log P_X(x) + \sum_{x \in \mathcal{X}} P_X(x) \\
 &= H(X) + 1.
 \end{aligned} \tag{2.12}$$

We have thus constructed a code C with $\ell_C(P_X) \leq H(X) + 1$, so $\ell_{\min}(P_X) \leq H(X) + 1$. ■

2.4 Huffman Codes

Shannon's source-coding theorem shows us that in theory, the minimal code length for a source P_X is roughly $H(X)$. In this section we will investigate **Huffman codes**, which provide an explicit and neat construction for optimal prefix-free codes. A binary Huffman code for a source P_X is constructed by iteratively pairing the two symbols with the smallest probability together, building a binary tree on the way. This is best explained by example:

Example 2.4.1 — Binary Huffman code. Let the random variable X be given with $\mathcal{X} = \{a, b, c, d, e\}$ and $P_X(a) = P_X(b) = 0.25$, $P_X(c) = 0.2$, and $P_X(d) = P_X(e) = 0.15$. The following is a binary Huffman code for P_X :



We build up the tree from left to right, pairing the symbols (or groups of symbols) with smallest (combined) probabilities at every step. The codeword for every symbol is then determined by following the branches of the tree *from right to left* until the symbol is reached. Note that this way, the symbols with the smallest probabilities get assigned the longest codewords (paths).

The average codeword length for this code is

$$0.25 \cdot 2 + 0.25 \cdot 2 + 0.2 \cdot 2 + 0.15 \cdot 3 + 0.15 \cdot 3 = 2.3. \tag{2.13}$$

This is very close to the entropy $H(X) \approx 2.285$. The average codeword length lies between $H(X)$ and $H(X) + 1$.

The above was an example of how to construct *binary* Huffman codes. We can also generate Huffman codes for larger alphabets, resulting in ternary, quaternary, or, more generally, ***d*-ary Huffman codes**.

Example 2.4.2 — Ternary Huffman code. We build a Huffman code with the alphabet $\{0, 1, 2\}$ for the same distribution as in Example 2.4.1.

x	$P_X(x)$		code
a	0.25	0	0
b	0.25	1	1
c	0.2	0	20
d	0.15	1	21
e	0.15	2	22

Exercise 2.4.3 Use the above procedure to construct a ternary code for the source P_X with $\mathcal{X} = \{a, b, c, d, e, f\}$ and $P_X(a) = P_X(b) = 0.25$, $P_X(c) = 0.2$, $P_X(d) = P_X(e) = P_X(f) = 0.1$. Can you find another code with a smaller average codeword length? ■

We have to be careful, because with an alphabet size of greater than 2, the above procedure does not always give an optimal code! In fact, a d -ary code is only optimal if $|\mathcal{X}|$ is of the form $k(d-1) + 1$ for some $k \in \mathbb{N}$. This ensures that at every step, we can combine exactly d symbols to use the alphabet at full capacity. The ternary code in Example 2.4.2 is optimal because $|\mathcal{X}| = 5 = 2(3-1) + 1$, but the code you constructed in Exercise 2.4.3 is not. To remedy this, one can add one or more ‘dummy’ symbols to the source (each with probability zero) until an appropriate size of \mathcal{X} of the form $|\mathcal{X}| = k(d-1) + 1$ for some $k \in \mathbb{N}$ is reached. The codewords for those dummy symbols are discarded at the end.

It turns out that Huffman codes indeed have optimal code length (see [CT], Section 5.8).

Theorem 2.4.4 — Optimality of Huffman codes. Let P_X be a source, and let C^* be the associated Huffman code. For any other uniquely decodable code C' with the same alphabet,

$$\ell_{C^*}(P_X) \leq \ell_{C'}(P_X).$$

Example 2.4.5 — 20 questions. In the game of ‘20 questions’ (20Q) the goal is to identify an object from a set of objects using (at most) 20 yes/no questions. Please go and [play it online](#) right now if you have never heard of it before. Assume we know the probability distribution P_X over all possible objects, what is the most efficient sequence of questions to ask in order to determine the object? We can use Huffman coding to answer this question!

On the one hand, the Huffman code for P_X has optimal average code length and we could (in principle) ask the player about the first, second, third, etc. bit of the object in question (admittedly, these questions would be rather boring...). On the other hand, we can see the sequence of (questions and) answers as the ‘code’ for an object, because every object has a unique sequence of yes/no answers. The number of questions asked is the length of the codeword. By Shannon’s source-coding theorem,

$$H(X) \leq \text{expected number of questions} \leq H(X) + 1,$$

and the Huffman code procedure can be used to determine the optimal sequence of questions.

As we have seen, the average codeword length of Huffman codes is theoretically optimal. However,

Huffman codes (and symbol codes in general) still have a number of disadvantages:

- When compressing, for example, an English text symbol-by-symbol, the probability distribution for each position may depend on the string of text that precedes it: for example, the letter *n* is a lot more likely than the letter *a* if it comes after the string *informatio*. Given this change of distribution, the Huffman code may not produce the shortest possible code. This can be resolved by recomputing the Huffman code after every symbol, but this results in a lot of overhead.
- The average codeword length is upper bounded by $H(X) + 1$. This additive cost of 1 bit is fine when $H(X)$ is very large, but can be a significant overhead when $H(X)$ is small itself.

2.5 Arithmetic Codes

In this section, we study a different kind of code that can handle context-dependent distributions, unlike the Huffman code.

Definition 2.5.1 — Standard binary representation. The standard binary representation of a real number $r \in [0, 1)$ is a (possibly infinite) string of bits $c_1c_2\cdots$ such that

$$r = \sum_i c_i \cdot 2^{-i},$$

where by convention, 0 is represented by the string 0.

Not all reals in $[0, 1)$ have a finite representation, but any interval $[a, b)$ with $0 \leq a < b \leq 1$ contains at least one number with a finite binary representation.

Example 2.5.2 The following table lists some numbers $r \in [0, 1)$ and their standard binary representation.

r	binary representation of r
$1/2$	1
$1/3$	01010101...
$1/4$	01
$3/4$	11
$13/16$	1101
$13/32$	01101

Note that 1101 is also the binary form of the natural number 13. Adding a 0 on the left divides the represented value by 2.

These binary representations of numbers in the interval $[0, 1)$ give rise to a very elegant code:

Definition 2.5.3 — Arithmetic code. Given a source P_X with $\mathcal{X} = \{x_1, \dots, x_m\}$, construct the arithmetic code as follows. Divide the interval $[0, 1)$ into disjoint subintervals $I_{x_j} = [a_j, a_{j+1})$, where a_1, \dots, a_{m+1} are defined such that $a_{j+1} - a_j = P_X(x_j)$, and $a_1 = 0, a_{m+1} = 1$.

The encoding $AC(x_j)$ of the element x_j is the (shortest possible) standard binary representation of some number in the interval I_{x_j} .

Example 2.5.4 Let X be a random variable with $\mathcal{X} = \{1, 2, 3\}$ and $P_X(1) = \frac{1}{2}$, $P_X(2) = P_X(3) = \frac{1}{4}$. The arithmetic code is constructed by first determining the intervals:



This image results in the arithmetic code \mathcal{C} with $\mathcal{C}(1) = 0$ (the representation of 0), $\mathcal{C}(2) = 1$ (the representation of $\frac{1}{2}$), $\mathcal{C}(3) = 11$ (the representation of $\frac{3}{4}$).

For P_X , the codewords happen to fall exactly *on* the boundaries of the intervals. This is not always the case, however. The same code would have resulted from this procedure if we started with the random variable Y with $\mathcal{Y} = \{1, 2, 3\}$ and $P_Y(1) = P_Y(2) = 0.3$ and $P_Y(3) = 0.4$:



Proposition 2.5.5 For any (X, P_X) , the arithmetic code has average length $\ell_{AC}(P_X) \leq H(X) + 1$.

Proof. Let $x \in \mathcal{X}$, and define $\ell_x := \lceil \log(1/P_X(x)) \rceil$ to be the rounded surprisal value of x . Then

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

Therefore, since the size of the interval I_x is $P_X(x)$, there must exist $0 \leq s_x < 2^{\ell_x}$ such that $s_x \cdot 2^{-\ell_x}$ lies in the interval I_x . This number $s_x \cdot 2^{-\ell_x}$ has a binary representation of length $\ell_x \leq -\log P_X(x) + 1$. Repeating this argument for every x , we obtain

$$\ell_{AC}(P_X) = \mathbb{E}[\ell(AC(X))] = \sum_x P_X(x) \ell(AC(x)) \leq \sum_x P_X(x) (-\log P_X(x) + 1) = H(X) + 1. \quad (2.15)$$

■

As we can see from Example 2.5.4, this construction for arithmetic codes does not necessarily yield prefix-free codes. However, at the expense of one extra bit of code (on average), the construction can be adapted into a prefix-free code. One example of this is the **Shannon-Fano-Elias code** (see [CT], Section 5.9 or [Wikipedia](#)), which provides a more sophisticated way of selecting a number within each interval than simply selecting the number with the shortest binary representation. This alternative selection procedure ensures prefix-freeness. Another option is to select *binary intervals* within each interval:

Definition 2.5.6 — Binary interval. A binary interval is an interval of the form

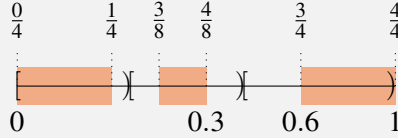
$$\left[\frac{s}{2^\ell}, \frac{s+1}{2^\ell} \right)$$

with $s, \ell \in \mathbb{N}$ and $0 \leq s < 2^\ell$. The **name** of the interval is the binary representation of s (as a natural number) padded with zeroes on the left to reach length ℓ .

Definition 2.5.7 — Arithmetic code (prefix-free version). The prefix-free arithmetic code is identical to Definition 2.5.3, except that the encoding $AC^{pf}(x_j)$ of the element x_j is now the name of the largest binary interval that fits entirely in I_{x_j} .

Similarly to Proposition 2.5.5, it can be shown that for any source P_X , $\ell_{AC^{pf}}(P_X) \leq H(X) + 2$. Note that we get prefix-freeness only at the expense of an extra bit on average.

Example 2.5.8 Let Y be the random variable as in Example 2.5.4, that is, $P_Y(1) = P_Y(2) = 0.3$ and $P_Y(3) = 0.4$. The prefix-free code for Y is constructed as follows:



This results in the codewords $\mathcal{C}(1) = 00$, $\mathcal{C}(2) = 011$, and $\mathcal{C}(3) = 11$.

The arithmetic code is slightly less efficient than the Huffman code in terms of average codeword length. A big advantage is the way it is able to adapt to changing distributions, such as when we are encoding a stream of English text. Suppose we are given the (not necessarily i.i.d.) random variables X_1, X_2, \dots, X_n , and we want to encode the source $P_{X_1 X_2 \dots X_n}$. We start by dividing the interval $[0, 1]$ into subintervals according to P_{X_1} . If, for example, the event $X_1 = b$ happens, we zoom into the interval corresponding to b , and subdivide *that* interval according to $P_{X_2|X_1}$, so that the sizes of these intervals add up to $P_{X_1}(b)$. The concept of arithmetic coding is exploited as an accessibility tool in the keyboard alternative **Dasher**, invented by the group of David MacKay at Cambridge University, UK.

2.6 Asymptotic Equipartition Property (AEP)

In this section, we consider the possibility of encoding blocks of symbols, rather than just one symbol at a time. We restrict our attention to sources that are *real* random variables. The following definition of converging random variables may remind you of a converging sequence of numbers. Recall that a sequence x_1, x_2, x_3, \dots of numbers converges to x if $\forall \varepsilon > 0 \exists n_0 \forall (n \geq n_0) : |x_n - x| < \varepsilon$. We denote this by writing $x_n \xrightarrow{n \rightarrow \infty} x$.

Definition 2.6.1 — Converging random variables. A sequence X_1, X_2, X_3, \dots of real random variables converges to a random variable X , if it satisfies one of the following definitions:

in probability (notation $X_n \xrightarrow{P} X$) if $\forall \varepsilon > 0, P[|X_n - X| > \varepsilon] \xrightarrow{n \rightarrow \infty} 0$

in mean square (notation $X_n \xrightarrow{m.s.} X$) if $\mathbb{E}[(X_n - X)^2] \xrightarrow{n \rightarrow \infty} 0$

almost surely (notation $X_n \xrightarrow{a.s.} X$) if $P[\lim_{n \rightarrow \infty} X_n = X] = 1$

where the definition of $X_n \xrightarrow{a.s.} X$ can be interpreted as $P[\{\omega \in \Omega \mid X_n(\omega) \xrightarrow{n \rightarrow \infty} X(\omega)\}] = 1$.

Example 2.6.2 Consider the following game: you flip a coin, and if it comes up heads you win one euro. You flip it again, and if it comes up head again, you win the double amount. You continue doing this, until a tails comes up: then you have to give back everything you won so far, and the game is over. To formalize this, let the sequence X_1, X_2, \dots denote the amount you won after 1, 2, etc coin flips. They are defined by $X_n = 2^n$ if you just flipped n heads in a row (probability $\frac{1}{2^n}$), and $X_n = 0$ otherwise. Furthermore, define X to be the constant-zero random variable with $P_X(0) = 1$.

- This sequence converges to X in probability: for every n and every $0 < \varepsilon \leq 1$, $P[|X_n - X| > \varepsilon] = P[X_n > \varepsilon] = \frac{1}{2^n}$, and from standard calculus we know that $\frac{1}{2^n} \xrightarrow{n \rightarrow \infty} 0$.
- This sequence does not converge to X in mean square, because $\mathbb{E}[(X_n - X)^2] = \mathbb{E}[(X_n)^2] = \frac{1}{2^n} \cdot 2^{2n} + (1 - \frac{1}{2^n}) \cdot 0 = 2^n$. Clearly, 2^n does not go to zero as n goes to infinity.
- This sequence almost surely converges to X : let Ω be the set of all possible (infinite) strings representing the coin toss outcomes (note that Ω itself is now also infinite). For any string ω that is *not* the all-heads string $hhhhh \dots$, we have that $X_n(\omega) \xrightarrow{n \rightarrow \infty} 0$, since there is an n_0 such that $X_{n_0}(\omega) = 0$ (namely if ω has a tails at the n_0 -th position), and remains zero for all $n > n_0$. Hence, $P[\{\omega \in \Omega \mid X_n(\omega) \xrightarrow{n \rightarrow \infty} X(\omega)\}] = P[\Omega \setminus \{hhhhh \dots\}] = 1$.

In general, the following implications hold (although their converses do not):

$$X_n \xrightarrow{m.s.} X \Rightarrow X_n \xrightarrow{P} X \quad (2.16)$$

$$X_n \xrightarrow{a.s.} X \Rightarrow X_n \xrightarrow{P} X \quad (2.17)$$

The following law states that if we sample several times from the same distribution, the average converges (in probability) to the expected value of the distribution.

Theorem 2.6.3 — Weak Law of Large Numbers. Let X_1, X_2, \dots be real i.i.d. random variables with mean $\mu = \mathbb{E}[X_i]$ and variance $\sigma^2 = \mathbb{E}[(X_i - \mu)^2] < \infty$. Define the random variables

$$S_n := \frac{1}{n} \sum_{i=1}^n X_i.$$

Then $S_n \xrightarrow{P} \mu$.

This important law has an entropy variant, which follows almost directly:

Theorem 2.6.4 — Asymptotic Equipartition Property (AEP). Let X_1, X_2, X_3, \dots be i.i.d. random variables with distribution P_X . Then

$$-\frac{1}{n} \log P_{X_1 \dots X_n}(X_1, \dots, X_n) \xrightarrow{P} H(X).$$

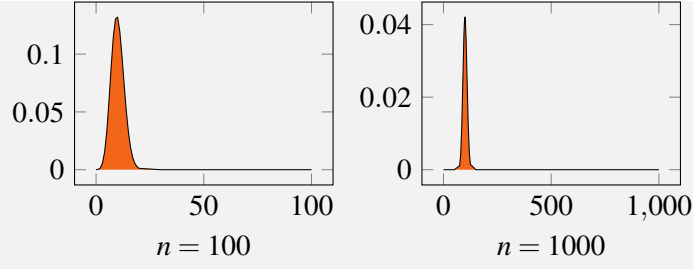
(Note that $P_{X_1 \dots X_n}(X_1, \dots, X_n)$ is itself a random variable, and $H(X)$ can be regarded as a constant random variable.)

Proof. Since the variables X_i are independent, so are the random variables $\log P_X(X_i)$. Then

$$\begin{aligned} -\frac{1}{n} \log P_{X_1 \dots X_n}(X_1, \dots, X_n) &= -\frac{1}{n} \sum_{i=1}^n \log P_X(x_i) \\ &\xrightarrow{P} -\mathbb{E}[\log P_X(X_i)] = H(X) \end{aligned} \quad (2.18)$$

by the weak law of large numbers. ■

Example 2.6.5 — Biased coin flip. Consider flipping a biased coin, with probability of heads being $p_1 = 0.1$ and probability of tails being $p_0 = 0.9$, and counting the number of heads that come up. The random variable X describing the outcome is distributed according to the binomial(n, p) distribution, where n is the number of coin flips. Below, the distribution is plotted for $n = 100$ and $n = 1000$:



We see that the variance of the sample mean, $\text{Var}[\frac{1}{n}X] = \frac{p(1-p)}{n}$, decreases as the number of samples increases. The weight of the distribution becomes centered around an increasingly narrow set of outcomes.

The above example leads us to defining the following subset of the image of a random variable:

Definition 2.6.6 — Typical set. The typical set $A_\varepsilon^{(n)}$ with respect to P_X is the set of strings $(x_1, \dots, x_n) \in \mathcal{X}^n$ such that

$$2^{-n(H(X)+\varepsilon)} \leq P_{X^n}(x_1, \dots, x_n) \leq 2^{-n(H(X)-\varepsilon)}.$$

The typical set is relatively small, but contains almost all of the probability mass. We start by establishing some general properties of typical sets.

Proposition 2.6.7 A typical set $A_\varepsilon^{(n)}$ satisfies the following:

1. For all $(x_1, \dots, x_n) \in A_\varepsilon^{(n)}$,

$$H(X) - \varepsilon \leq -\frac{1}{n} \log P_{X^n}(x_1, \dots, x_n) \leq H(X) + \varepsilon.$$

2. $P[A_\varepsilon^{(n)}] > 1 - \varepsilon$ (for large enough n).
3. $|A_\varepsilon^{(n)}| \leq 2^{n(H(X)+\varepsilon)}$.
4. $|A_\varepsilon^{(n)}| \geq (1 - \varepsilon)2^{n(H(X)-\varepsilon)}$ (for large enough n).

Proof.

1. This is immediate from the definition (take the logarithm and divide by $-n$, thereby reversing the inequalities).
2. This follows from the Asymptotic Equipartition Property: for all $\varepsilon > 0$, $P[|-\frac{1}{n} \log P_{X^n}(X_1, \dots, X_n) - H(X)| > \varepsilon] \xrightarrow{n \rightarrow \infty} 0$, that is,

$$\forall(\varepsilon > 0) \forall(\delta > 0) \exists n_0 \forall(n \geq n_0) P[|-\frac{1}{n} \log P_{X^n}(X_1, \dots, X_n) - H(X)| \leq \varepsilon] > 1 - \delta. \quad (2.19)$$

By choosing $\delta := \varepsilon$, the result follows from the first property.

3. First, observe that

$$1 = \sum_{\vec{x} \in \mathcal{X}^n} P_{X^n}(\vec{x}) \geq \sum_{\vec{x} \in A_\varepsilon^{(n)}} P_{X^n}(\vec{x}) \geq |A_\varepsilon^{(n)}| \cdot 2^{-n(H(X)+\varepsilon)}, \quad (2.20)$$

where the last inequality follows by the definition of typicality. The claim follows by multiplying both sides of the equation by $2^{n(H(X)+\varepsilon)}$.

4. By Property 2, we can choose an n large enough so that

$$1 - \varepsilon < P[A_\varepsilon^{(n)}] = \sum_{\vec{x} \in A_\varepsilon^{(n)}} P_{X^n}(\vec{x}) \leq |A_\varepsilon^{(n)}| \cdot 2^{-n(H(X)-\varepsilon)}, \quad (2.21)$$

where again, the last inequality follows by the definition of typicality. ■

Typical sets and their properties allow us to code a source P_X , in blocks of n symbols at a time, in either a lossy or a lossless way. For a lossy code, we notice that with overwhelming probability, a sequence of n iid samples from P_X is typical, so it suffices to assign binary labels of length (at most) $\lceil n(H(X) + \varepsilon) \rceil$ to the elements of $A_\varepsilon^{(n)}$, and assign some constant (dummy) codeword to all elements outside of the set. Decoding this dummy codeword will result in an error (data loss), but this error occurs with probability at most ε .

The above scheme can be extended to a lossless version by assigning longer labels to the elements outside of $A_\varepsilon^{(n)}$, for example binary labels of length $\lceil \log |\mathcal{X}|^n \rceil = \lceil n \log |\mathcal{X}| \rceil$. An extra ‘flag’ bit is needed to indicate whether the element is inside or outside the typical set. For large enough n , this code is quite efficient:

Theorem 2.6.8 Let X_1, \dots, X_n be i.i.d. real random variables with respect to the set \mathcal{X} , and distributed according to P_X . Let $\varepsilon > 0$. Then there exists a lossless code $\mathcal{X}^n \rightarrow \{0, 1\}^*$ such that, for sufficiently large n , $\mathbb{E}[\frac{1}{n} \ell(X^n)] \leq H(X) + \varepsilon$.

Proof. Consider the code described above: the code consist of a flag bit (indicating whether or not the element is inside the typical set), followed by either a short label (for elements in the typical set) or a longer one (for elements outside of it).

Let $\varepsilon' > 0$ (we will specify the value of ε' later). Let n be large enough such that $P[A_{\varepsilon'}^{(n)}] > 1 - \varepsilon'$ (see Proposition 2.6.7). Then

$$\begin{aligned}
 \mathbb{E}[\ell(X^n)] &= \sum_{\vec{x} \in \mathcal{X}^n} P_{X^n}(\vec{x}) \ell(\vec{x}) \\
 &= \sum_{\vec{x} \in A_{\varepsilon'}^{(n)}} P_{X^n}(\vec{x}) \ell(\vec{x}) + \sum_{\vec{x} \notin A_{\varepsilon'}^{(n)}} P_{X^n}(\vec{x}) \ell(\vec{x}) \\
 &\leq P[A_{\varepsilon'}^{(n)}] \cdot (\lceil n(H(X) + \varepsilon') \rceil + 1) + P[\overline{A_{\varepsilon'}^{(n)}}] \cdot (\lceil n \log |\mathcal{X}| \rceil + 1) \\
 &\leq P[A_{\varepsilon'}^{(n)}] \cdot (n(H(X) + \varepsilon') + 2) + P[\overline{A_{\varepsilon'}^{(n)}}] \cdot (n \log |\mathcal{X}| + 2) \\
 &\leq n(H(X) + \varepsilon') + \varepsilon' \cdot n \log |\mathcal{X}| + 2 \\
 &= n(H(X) + \varepsilon), \tag{2.22}
 \end{aligned}$$

where $\varepsilon = \varepsilon' + \varepsilon' \log |\mathcal{X}| + \frac{2}{n}$ (note that ε can be made arbitrarily small by choosing ε' and n wisely). The +1 in the first inequality is a consequence of the ‘flag’ bit. ■

For large enough blocks of symbols, typical sets thus allow the construction of an efficient code without the 1 bit of overhead that symbol codes may necessarily have. However, this efficiency is only guaranteed for ‘sufficiently large n ’, a rather theoretical condition that may not be achievable in practice.

We conclude this chapter by showing that the typical set is in a sense ‘optimal’, i.e. that picking a smaller set instead of the typical set does not allow for much shorter codewords on average in a lossy setting, not even if we allow rather large error probabilities by allowing about half of the elements to lie outside of the typical set.

Let us use the notation $B_\delta^{(n)}$ to denote the smallest subset of \mathcal{X}^n such that $P[B_\delta^{(n)}] > 1 - \delta$ (for some parameter $\delta > 0$). $B_\delta^{(n)}$ can be explicitly constructed by, for example, ordering \mathcal{X}^n in order of decreasing probability, and adding elements to $B_\delta^{(n)}$ until the probability threshold of $1 - \delta$ is reached. The following theorem states that even for large values of δ , we still need almost $nH(X)$ bits to denote an element from $B_\delta^{(n)}$.

Theorem 2.6.9 Let X_1, \dots, X_n be i.i.d. random variables distributed according to P_X . For any $\delta < \frac{1}{2}$, and any $\delta' > 0$, if $P[B_\delta^{(n)}] > 1 - \delta$, then

$$\frac{1}{n} \log |B_\delta^{(n)}| > H(X) - \delta',$$

for sufficiently large n .

Proof. Let $\delta, \varepsilon < \frac{1}{2}$, and consider some $B_\delta^{(n)}$ such that $P[B_\delta^{(n)}] > 1 - \delta$. We know that by Proposition 2.6.7, $P[A_\varepsilon^{(n)}] > 1 - \varepsilon$, for large enough n . Thus, by the union bound,

$$\begin{aligned} 1 - \varepsilon - \delta &< 1 - P[\overline{A_\varepsilon^{(n)}}] - P[\overline{B_\delta^{(n)}}] \\ &\leq 1 - P[\overline{A_\varepsilon^{(n)} \cup B_\delta^{(n)}}] \\ &= P[A_\varepsilon^{(n)} \cap B_\delta^{(n)}] \\ &= \sum_{\vec{x} \in A_\varepsilon^{(n)} \cap B_\delta^{(n)}} P_{X^n}(\vec{x}) \\ &\leq \sum_{\vec{x} \in A_\varepsilon^{(n)} \cap B_\delta^{(n)}} 2^{-n(H(X) - \varepsilon)} \\ &= |A_\varepsilon^{(n)} \cap B_\delta^{(n)}| \cdot 2^{-n(H(X) - \varepsilon)} \\ &\leq |B_\delta^{(n)}| \cdot 2^{-n(H(X) - \varepsilon)}. \end{aligned} \tag{2.23}$$

Rearranging this expression and taking the logarithm, we get

$$H(X) - \varepsilon + \frac{1}{n} \log(1 - \varepsilon - \delta) < \frac{1}{n} \log |B_\delta^{(n)}|. \tag{2.24}$$

If we now set $\delta' := \varepsilon - \frac{1}{n} \log(1 - \varepsilon - \delta)$, then

$$H(X) - \delta' < \frac{1}{n} \log |B_\delta^{(n)}|, \tag{2.25}$$

as desired. Observe that we can make the expression for δ' as small as desired by choosing a large enough n , even if δ is rather large. ■

Chapter 3: Perfectly Secure Encryption

Information theory is very useful when analyzing the security of perfectly secure encryption schemes. Consider a scenario where one party, Alice, wants to send a message m (sampled from some distribution P_M) to another party, Bob, over some public channel, for example the internet. Alice and Bob share a key k (sampled from another distribution P_K), which is a piece of information that is known only to them. Alice can use the key to encrypt her message (the **plaintext**), and Bob can use the same key to decrypt the **ciphertext** that Alice created, and read the message. The goal is to do this in such a way that if an eavesdropper (who usually goes by the name ‘Eve’) listens in on the channel and intercepts the encrypted message, she cannot derive any information about the message as long as she does not know the key k .



Let us formalize the above notion of encryption in the following definition:

Definition 3.0.1 — Encryption scheme. An encryption scheme for (the message) M consists of a key K and a ciphertext $C = \text{Enc}(M, K)$, such that

- $I(M; K) = 0$ (the key is independent of the message – this is a **setup assumption**), and
- $H(M|KC) = 0$ (given the key and the ciphertext, Bob can recover the original message)

Note that M , K and C are random variables.

Note that in order to satisfy the second requirement, the encryption function $\text{Enc}(\cdot, \cdot)$ needs to be injective: every message is mapped to a *unique* ciphertext.

3.1 Security

Definition 3.0.1 does not put any constraints on the amount of information that Eve can get from the ciphertext: we still need to explicitly require the scheme to be secure.

Definition 3.1.1 — Perfect security. An encryption scheme is perfectly secure if

$$I(M; C) = 0.$$

This is equivalent to saying $H(M|C) = H(M)$, or to saying that M and C are independent.

This type of security is also sometimes called perfect **information-theoretic security**, in order to stress that the ciphertext really does not contain *any* information about the plaintext message. Many commonly used encryption schemes do not provide this type of security. In **computationally secure** schemes, a lot of information about the message may be contained in the ciphertext, but it would take a ridiculous amount of resources (such as computation time or memory) to compute the information about the message from the ciphertext.

3.2 One-Time Pad

A classic example of a perfectly secure encryption scheme is the one-time pad.

Definition 3.2.1 — One-time pad (OTP). Let the message space \mathcal{M} be some additive group $(G, +)$. Define the random variable K to be uniformly distributed over the key space $\mathcal{K} = \mathcal{M}$, and define the ciphertext space to be $\mathcal{C} = \mathcal{M}$ as well. Define the encryption and decryption function as follows:

$$\text{Enc}(m, k) = m + k = c,$$

$$\text{Dec}(c, k) = c - k = m.$$

Here, $c - k$ stands for $c + (-k)$, where $-k$ is the additive inverse of k in the group $(G, +)$.

Example 3.2.2 — One-time pad for binary strings. The most common use of the one-time pad is for the group of binary strings under (bit-wise) addition modulo 2, i.e. $(\{0, 1\}^n, \oplus)$. In this group, every element is its own additive inverse, resulting in the encryption and decryption functions

$$\text{Enc}(m, k) = m \oplus k = c,$$

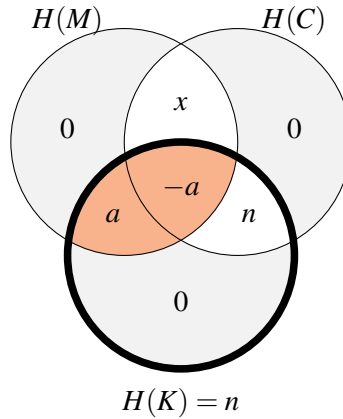
$$\text{Dec}(c, k) = c \oplus k = m.$$

For example, if $n = 4$, a possible message m is 0101, and a possible key k is 0110. The ciphertext c is $0101 \oplus 0110 = 0011$, and the decryption of c is again $0011 \oplus 0110 = 0101$, the original message m .

We can show that the one-time pad indeed satisfies the security definition 3.1.1 of the previous section.

Theorem 3.2.3 The one-time pad is perfectly secure.

Proof. Write $n = \log |G|$. We need to verify that $I(M;C) = 0$. We do so using a three-variable entropy diagram (see Section 1.8). We can already fill in the values $H(K) = n = \log |G|$ (because K is uniformly distributed), $H(M|CK) = H(C|MK) = H(K|MC) = 0$ (because each random variable is a function of the other two), and $I(M;K) = 0$ (this is our setup assumption).



Note that the area of $I(M;K) = I(M;K|C) + R(M;K;C)$ (shaded orange in the picture) as a whole is 0, but that does not mean that $I(M;K|C)$ and $R(M;K;C)$ themselves are zero, because $R(M;K;C)$ can be negative. We can conclude that there must be some (non-negative) real number $a \geq 0$ such that $I(M;K|C) = a$ and $R(M;K;C) = -a$. As the entropy of K has to be $H(K) = n$, we can furthermore conclude that $I(K;C|M) = n$.

From $I(M;C) \geq 0$ follows that $x \geq a$, and because $H(C) \leq n$, it follows that $x \leq a$ and hence, $x = a$ and $I(M;C) = 0$, as desired.



■

We have thus seen that the one-time pad provides perfect information-theoretic security. There is one enormous drawback to this encryption scheme though: the key needs to be as large as the message! To send a message of n bits, Alice needs to share n bits of key with Bob. It might be tempting for Alice to reuse the key k for several messages once she has shared it with Bob, but this is dangerous: Eve could, from two intercepted encryptions $(m_1 + k)$ and $(m_2 + k)$, recover

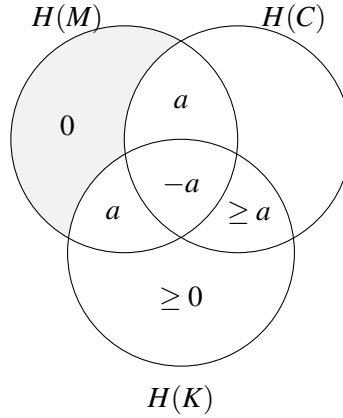
the difference of the two plaintext messages $m_1 + k - (m_2 + k) = m_1 - m_2$. Already the difference between two plaintext messages can reveal a lot of information about the individual messages, as illustrated in [this Cryptosmith blog post](#).

3.3 Lower Bound On Key Length

It turns out to be impossible to design an encryption scheme that provides both perfect security and short keys. So even though the one-time pad may seem inefficient, its key lengths are optimal for a perfectly secure scheme.

Theorem 3.3.1 — Shannon 1949 (Shannon49). For any perfectly secure encryption scheme, it holds that $H(K) \geq H(M)$.

Proof. Again, we turn to entropy diagrams. Write $a = I(M; C|K) \geq 0$. using the fact that $I(M; K) = 0$ (setup assumption) and $I(M; C) = 0$ (security), we can fill in the entropy diagram as follows:



Note that $I(K; C|M) \geq a$ follows from the fact that $I(C; K) \geq 0$ and $R(C; K; M) = -a$.

From the diagram, we observe that $H(M) = a$, and that $H(K) \geq a$. Hence, $H(K) \geq H(M)$. ■

Chapter 4: Stochastic Processes

So far in this course we have mostly seen independent random variables. In source coding, we designed our codes according to a single distribution P_X , and assumed that if we encoded a sequence of source symbols, the symbols in the sequence would be drawn independently according to P_X . In the real world, however, subsequent events are often dependent on each other. For example, in an English text, after observing a letter q , the next letter is much more likely to be u than it is to be r , even though in general the letter r is more prevalent in English text. The event of observing the letter q changes the probability distribution of the next letter.

4.1 Finite Markov Chains

We start by studying a restricted form of dependence between variables: Markov chains. Random variables forms a Markov chain if the distribution of each random variable depends only on the outcome of the random variable that directly precedes it.

Definition 4.1.1 — Markov chain (length 3). The random variables X, Y , and Z form a Markov chain (notation: $X \rightarrow Y \rightarrow Z$) if

$$P_{Z|XY} = P_{Z|Y}.$$

Exercise 4.1.2 Verify that for a Markov chain, $P_{XYZ} = P_X \cdot P_{Y|X} \cdot P_{Z|Y}$. ■

Proposition 4.1.3 For all random variables X, Y, Z , the following statements are equivalent:

- $X \rightarrow Y \rightarrow Z$,
- $Z \rightarrow Y \rightarrow X$,
- $P_{XZ|Y} = P_{X|Y} \cdot P_{Z|Y}$ (that is, $I(X; Z|Y) = 0$).

Because of the equivalence of items (a) and (b), we often write $X \leftrightarrow Y \leftrightarrow Z$ to denote a Markov chain.

Proof. We prove that (a) \Rightarrow (b), that (b) \Rightarrow (c), and that (c) \Rightarrow (a). All other directions follow from (combinations of) these three implications.

(a) \Rightarrow (b): Suppose $X \rightarrow Y \rightarrow Z$. Then

$$\begin{aligned} P_{X|YZ} &= \frac{P_{XYZ}}{P_{YZ}} = \frac{P_{XY} \cdot P_{Z|XY}}{P_{YZ}} = \frac{P_{XY} \cdot P_{Z|XY}}{P_Y \cdot P_{Z|Y}} \\ &= \frac{P_{XY} \cdot P_{Z|Y}}{P_Y \cdot P_{Z|Y}} = \frac{P_{XY}}{P_Y} = P_{X|Y}, \end{aligned} \quad (4.1)$$

where we go from the top to the bottom line using the assumption that $X \rightarrow Y \rightarrow Z$.

(b) \Rightarrow (c): Suppose $Z \rightarrow Y \rightarrow X$. Then

$$\begin{aligned} P_{XZ|Y} &= \frac{P_{XYZ}}{P_Y} = \frac{P_{X|YZ} \cdot P_{YZ}}{P_Y} \\ &= \frac{P_{X|Y} \cdot P_{YZ}}{P_Y} = P_{X|Y} \cdot \frac{P_{YZ}}{P_Y} = P_{X|Y} \cdot P_{Z|Y}, \end{aligned} \quad (4.2)$$

where we go from the top to the bottom line using the assumption that $Z \rightarrow Y \rightarrow X$.

(c) \Rightarrow (a): For the final implication, suppose that $P_{XZ|Y} = P_{X|Y} \cdot P_{Z|Y}$. Then

$$\begin{aligned} P_{Z|YX} &= \frac{P_{XYZ}}{P_{XY}} = \frac{P_{XZ|Y} \cdot P_Y}{P_{XY}} \\ &= \frac{P_{X|Y} \cdot P_{Z|Y} \cdot P_Y}{P_{XY}} = \frac{P_{XY} \cdot P_{Z|Y}}{P_{XY}} = P_{Z|Y}, \end{aligned} \quad (4.3)$$

and therefore $X \rightarrow Y \rightarrow Z$. ■

Example 4.1.4 — Whispering game. Alice, Bob and Charlie play a game: Alice comes up with a message, and whispers it into Bob's ear. Bob then proceeds to whisper the message into Charlie's ear, who says it out loud. Of course, Bob and Charlie may not hear the message correctly: there is some noise in the communication.

Let A be the random variable that contains the message that Alice picks. For simplicity, let's say she picks a uniformly random bit. Let B be the message that Bob heard: it is also a bit, but with probability 0.1, it is not the same as A . More formally,

$$\begin{aligned} P_{B|A}(0|0) &= P_{B|A}(1|1) = 0.9, \\ P_{B|A}(1|0) &= P_{B|A}(0|1) = 0.1. \end{aligned}$$

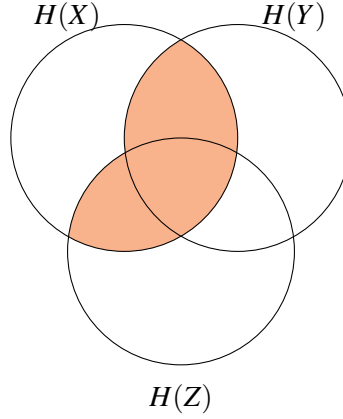
Finally, let C be the message that Charlie heard: again, let us suppose that it is unequal to the whispered message B with probability 0.1. For the conditional distribution of C , we see that $P_{C|AB}(0|00) = 0.9 = P_{C|B}(0|0)$, and similarly for all other possible values of A , B , and C . Thus, these random variables form a Markov chain: $A \rightarrow B \rightarrow C$.

Proposition 4.1.3 tells us that we can also look at this game in the converse direction: if we are curious about Alice's original message A , we only have to look at Bob's interpretation B of the message. Additionally knowing Charlie's interpretation C does not change the distribution of A .

Example 4.1.4 also exhibits an important property of Markov chains: you can only lose information down the line. Charlie's final message C does not contain any more information about Alice's original message A than what was already contained in Bob's message B . This observation is formalized in the following theorem:

Theorem 4.1.5 — Data-processing inequality. If $X \rightarrow Y \rightarrow Z$, then $I(X;Y) \geq I(X;Z)$. Equality holds if and only if $I(X;Y|Z) = 0$.

Proof. The following entropy diagram depicts the area $I(X;YZ)$:



From the diagram, we can see that

$$I(X;Z) + I(X;Y|Z) = I(X;YZ) = I(X;Y) + I(X;Z|Y). \quad (4.4)$$

Combining this with part (c) of Proposition 4.1.3, it follows that

$$I(X;Z) + I(X;Y|Z) = I(X;Y). \quad (4.5)$$

Since $I(X;Y|Z) \geq 0$, the result follows: $I(X;Z) \leq I(X;Y)$, with equality iff $I(X;Y|Z) = 0$. ■

Corollary 4.1.6 $I(X;Y) \geq I(X;g(Y))$ for any two random variables X and Y , and any function g on the range of Y .

We can extend the definition of Markov chains to more than three variables:

Definition 4.1.7 — Markov chain (length n). The random variables X_1, X_2, \dots, X_n form a Markov chain (notation: $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_n$) if for all $3 \leq i \leq n$,

$$P_{X_i|X_1 \dots X_{i-1}} = P_{X_i|X_{i-1}}.$$

Markov chains of length n exhibit similar properties to the properties we have seen for Markov chains of length 3. In particular, the reverse chain is also a Markov chain ($X_n \rightarrow \dots \rightarrow X_2 \rightarrow X_1$), and a more general form of the data-processing inequality holds in the sense that the further apart two variables are in the chain, the less correlated they are.

4.2 Stochastic Processes

Let us now consider more general sets of random variables than the three-variable Markov chains we saw in the previous section. In the rest of this chapter, we will consider infinite sequences of random variables, that can have varying degrees of independence. We start off with the most general definition of such an infinite process.

Definition 4.2.1 — Discrete-time stochastic process. A stochastic process is a sequence $\{X_i : \Omega \rightarrow \mathcal{X}\}$ of random variables indexed by $i \in \mathbb{N}$. The process is characterized by the

collection of probability distributions $P_{X_n|X_1 \dots X_{n-1}}$ for all $n \in \mathbb{N}$.

Equivalently, we can say that a stochastic process is characterized by the joint probability distributions $P_{X_1 \dots X_n}$ for all $n \in \mathbb{N}$.

Note that the random variables all have the same domain (the same sample space) and the same codomain. Sometimes this sample space can be infinite, as it is in the following example.

Example 4.2.2 — Repeatedly tossing a fair coin. Suppose you toss a fair coin an infinite amount of times, and at every step, you count the number of heads you have seen so far.

This experiment can be described as a stochastic process by letting each variable X_i denote the number of heads observed up until that toss. For example, a sequence of tosses THHTHT... results in the values $X_1 = 0, X_2 = 1, X_3 = 2, X_4 = 2, X_5 = 3, X_6 = 3, \dots$. The stochastic process is characterized by the probability distributions

$$P_{X_1}(0) = P_{X_1}(1) = \frac{1}{2} (P_{X_1}(x_1) = 0 \text{ for all other } x_1 \in \mathbb{N})$$

$$P_{X_{n+1}|X_1 \dots X_n}(x_{n+1}|x_1 \dots x_n) = \frac{1}{2} \text{ if } x_{n+1} = x_n \text{ or } x_{n+1} = x_n + 1, \text{ and } 0 \text{ otherwise.}$$

For this experiment, the sample space is the set of all possible outcomes for an infinite sequence of coin tosses,

$$\Omega = \{H, T\} \times \{H, T\} \times \{H, T\} \times \dots$$

In particular, $|\Omega|$ is uncountable. Properly defining a probability measure for Ω would require us to rethink many of our definitions that work only for finite sample spaces. That can be done (by using **σ -algebras** as event spaces), but that is beyond the scope of this course. For stochastic processes, we can always think of the sample space for X_n as a finite set (consisting of the first n samples only). In the current example, the probability distribution of X_n only depends on the first n coin tosses.

4.2.1 Types Of Stochastic Processes

Often, we will be interested in stochastic processes with specific properties, such as processes where all X_i are independent, or processes with a Markov-like property. We review several such properties in this section.

Definition 4.2.3 — Stationary process. A stochastic process is stationary if for all $n, k \in \mathbb{N}$,

$$P_{X_1 \dots X_n} = P_{X_{1+k} \dots X_{n+k}}.$$

Stationary processes are invariant under time shifts: when observing a subsequence of length n , it does not matter where in the process you look exactly.

Example 4.2.4 — i.i.d. process. Let X be a random variable. Consider a stochastic process $\{X_i\}$ where $P_{X_i} = P_X$ for all i . That is, the random variables in the sequence are all independent and identically distributed.

This process is stationary, which can be shown by induction on n . For $n = 1$, it is immediate from the definition of the process: $P_{X_1} = P_X = P_{X_{1+k}}$ for all $k \in \mathbb{N}$. For the induction step, suppose that for some $n \in \mathbb{N}$, for all k , $P_{X_1 \dots X_n} = P_{X_{1+k} \dots X_{n+k}}$. Then also for all k ,

$$P_{X_1 \dots X_{n+1}} = P_{X_1 \dots X_n} \cdot P_{X_{n+1}} \quad (\text{by independence})$$

$$\begin{aligned}
&= P_{X_1 \dots X_n} \cdot P_{X_{n+1+k}} && \text{(by identical distribution)} \\
&= P_{X_{1+k} \dots X_{n+k}} \cdot P_{X_{n+1+k}} && \text{(by induction hypothesis)} \\
&= P_{X_{1+k} \dots X_{n+1+k}} && \text{(by independence)}
\end{aligned}$$

Example 4.2.5 — Ten fair coins. Throw a fair coin 10 times: this gives the finite sample space $\{\text{H}, \text{T}\}^{10}$. Define the stochastic process $\{X_i\}$ by setting $X_i = 0$ if coin number $(i \bmod 10)$ is heads, and $X_i = 1$ otherwise.

This process is stationary. You can prove this yourself by observing that for all i , $X_i = X_{i+10} = X_{i+20} = \dots$, and using a similar induction argument as in Example 4.2.4.

Definition 4.2.6 — Markov process. A stochastic process is a Markov process (or: Markov chain) if for all $n \in \mathbb{N}$,

$$X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_n.$$

Exercise 4.2.7 Verify that for a Markov process, $P_{X_1 \dots X_n} = P_{X_1} \cdot P_{X_2|X_1} \cdot P_{X_3|X_2} \dots P_{X_n|X_{n-1}}$. ■

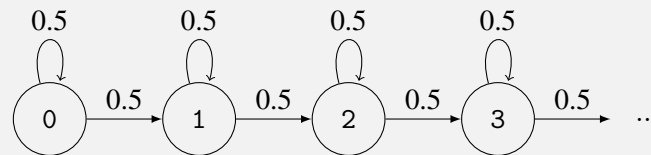
In a Markov process, the value of each step depends only on the previous value, similarly to the three-variable Markov chains from Section 4.1. The infinite sequence of coin tosses from Example 4.2.2 is an example of a Markov process: you can see this by inspecting the definition of $P_{X_{n+1}|X_1 \dots X_n}$, and noting that it is indeed independent of the values for X_1 to X_{n-1} . The process in this example even fulfills a stronger property: it is time invariant.

Definition 4.2.8 — Time-invariant Markov process. A Markov process is time invariant if for all $n \in \mathbb{N}$,

$$P_{X_{n+1}|X_n} = P_{X_2|X_1}.$$

Time-invariant Markov processes can be nicely visualized using state diagrams, where the **states** represent the values in \mathcal{X} , and the labels on the arrows represent the conditional probabilities. A time-dependent Markov process could also be visualized in this way, but the labels on the arrows would have to be functions of the time step i .

Example 4.2.9 — Repeatedly tossing a fair coin, continued. The time-invariant Markov process in Example 4.2.2 is represented by the following state diagram:

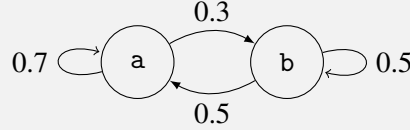


For example, the arrow from state 2 to state 3, represents the probability $P_{X_{n+1}|X_n}(3|2) = 0.5$.

In order to describe the stochastic process completely, we need to specify P_{X_1} , the **initial distribution**. In the current example we have $P_{X_1}(0) = P_{X_1}(1) = 0.5$: the process is equally likely to start out in state 0 (if the first toss is a tails) or state 1 (if the first toss is a heads).

Definition 4.2.10 — Finite-state Markov process. A time-invariant Markov process is finite-state if $|\mathcal{X}|$ is finite. [Yfke: Or should we say: if the set of elements \mathcal{X} that have non-zero probability is finite?]

Example 4.2.11 — A finite-state time-invariant Markov process. Consider the following state diagram, for a Markov process with initial distribution $P_{X_1}(a) = 1$ and $P_{X_1}(b) = 0$:



The process starts in state a with probability one. A possible run of the process would be aabaaabbabaa...

Note that in the state diagrams, the probabilities of the outgoing arrows add up to 1 for every state. This is necessary for a well-defined Markov process.

A time-invariant Markov process can alternatively be represented by a **transition matrix** P , where the entry P_{ij} represents the probability $P_{X_{n+1}|X_n}(j|i)$. For a finite-state process, the transition matrix is finite. In the transition matrix, the row entries have to sum up to 1.

Example 4.2.12 — A finite-state time-invariant Markov process, continued. The matrix representation of the process in Example 4.2.11 is given by

$$P = \begin{bmatrix} 0.7 & 0.3 \\ 0.5 & 0.5 \end{bmatrix}$$

where the state a is represented by the first row/column, and the state b by the second. Verify that the row entries indeed sum up to 1.

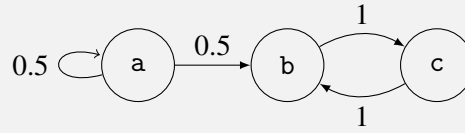
Definition 4.2.13 — Irreducible Markov process. A time-invariant Markov process is irreducible if every state is reachable from any other state in a finite number of steps.

The process in Example 4.2.11 is irreducible: the state a is reachable from b and vice versa. The process in Example 4.2.9 is not irreducible. For example, the state 0 is not reachable from the state 2.

Definition 4.2.14 — Aperiodic Markov process. A state in a time-invariant Markov process is aperiodic if the greatest common divisor of all path lengths from that state to itself is 1. The process itself is aperiodic if all states are aperiodic.

The processes in Example 4.2.9 and 4.2.11 are both aperiodic. In both examples, every state is reachable from itself with paths of any length, so the greatest common divisor of the path lengths is always 1. Below, we will present an example of a process that is periodic (i.e., not aperiodic).

Example 4.2.15 — Periodic Markov process. Consider the process that starts in state a with probability 1, and is further described by the following state diagram:



The state a is aperiodic, but b and c are not. The paths from state b to itself are all of even length, as are the paths from state c to itself. A typical run of this process might look like $aaaaabcbcbcbcb\cdots$; the period 2 is clearly visible in this run.

Suppose that we run the process of Example 4.2.11 for a very large number of steps, and wonder what the probability will be of observing an a at the next step. Given the initial distribution and the state diagram, we can compute the probability distribution for every X_i . In Figure 4.1, $P_{X_i}(a)$ is plotted for several values of i . The probability to observe an a seems to stabilize. This leads us to the following definition:

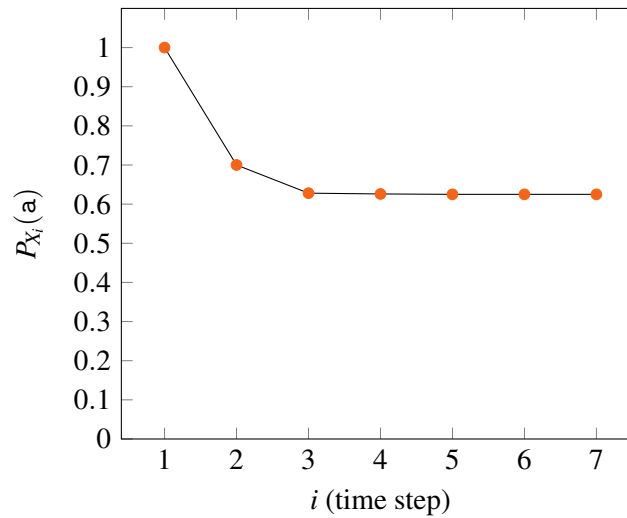


Figure 4.1: The probability that $X_i = a$ in the Markov process of Example 4.2.11, plotted for the first few values of i .

Definition 4.2.16 — Stationary distribution. A stationary distribution is a distribution P_{X_n} such that $P_{X_{n+1}} = P_{X_n}$.

If the initial distribution of a time-invariant Markov process is stationary, then the entire process is stationary in the sense of Definition 4.2.3.

Not every Markov process has a stationary distribution, and if it does, the distribution might not be unique. However, in some cases we can be sure that the stationary distribution exists and is unique.

Proposition 4.2.17 If a time-invariant Markov chain is finite-state, irreducible and aperiodic, then there exists a unique stationary distribution. Moreover, from any initial distribution, the distribution of X_n tends to the stationary distribution as $n \rightarrow \infty$.

Example 4.2.18 Consider again the process of Example 4.2.11. This process does have a stationary distribution, because it is finite-state, irreducible and aperiodic. Let μ_a denote the

probability of observing an a in the stationary distribution, and μ_b the probability of observing a b. Then the stationary distribution must satisfy the following set of equations:

$$\begin{aligned}\mu_a &= 0.7\mu_a + 0.5\mu_b, \\ \mu_b &= 0.3\mu_a + 0.5\mu_b, \\ \mu_a + \mu_b &= 1.\end{aligned}$$

Solving this set of equations gives the stationary distribution $\mu_a = 0.625$ and $\mu_b = 0.375$. This outcome matches our observations in Figure 4.1.

4.3 Entropy rate

Intuitively, some of the stochastic processes we have seen in the previous section are more predictable than others. The periodic Markov process from Example 4.2.15 is not so surprising anymore as soon as the first b is observed. In this section, we will study a measure for the unpredictability of a stochastic process: the entropy rate.

Definition 4.3.1 — Entropy rate. The entropy rate $H(\mathcal{X})$ of a stochastic process $\{X_i\}$ is

$$H(\mathcal{X}) := \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n),$$

if the limit exists, and undefined otherwise.

The entropy rate reflects the way in which the entropy of the sequence (observed so far) grows as n grows large.

Example 4.3.2 Consider a process $\{X_i\}$ where the X_i are i.i.d. fair coin tosses. Then

$$\begin{aligned}H(\mathcal{X}) &= \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1) \cdot H(X_1) \cdots H(X_n) \\ &= \lim_{n \rightarrow \infty} \frac{n}{n} H(X_1) \\ &= 1.\end{aligned}$$

So, every new coin toss increases the entropy of the entire observed sequence by 1.

We can also define an alternative measure of the unpredictability of a stochastic process, where we focus not on the amount of entropy in the entire sequence observed so far, but on the amount of entropy present in the current random variable, given the past sequence.

Definition 4.3.3 — Entropy rate given the past. The entropy rate given the past $H'(\mathcal{X})$ of a stochastic process $\{X_i\}$ is

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

if the limit exists, and undefined otherwise.

For all stationary processes, this alternative definition turns out to coincide with the original definition of entropy rate. In order to show this, we need an analytic statement about the convergence of sums.

Theorem 4.3.4 — Cesàro mean. If $\lim_{n \rightarrow \infty} a_n = a$ and $b_n = \frac{1}{n} \sum_{i=1}^n a_i$, then $\lim_{n \rightarrow \infty} b_n = a$.

Exercise 4.3.5 Prove Theorem 4.3.4. ■

Theorem 4.3.6 For a stationary process $\{X_i\}$, $H(\mathcal{X}) = H'(\mathcal{X})$ (and both exist).

Proof. We first show that $H(X_n | X_1, \dots, X_{n-1})$ is a non-increasing function of n :

$$\begin{aligned} H(X_n | X_1, \dots, X_{n-1}) &= H(X_{n+1} | X_2, \dots, X_n) && \text{(stationary)} \\ &\geq H(X_{n+1} | X_1, \dots, X_n). \end{aligned} \tag{4.6}$$

Combined with the fact that $H(X_n | X_1, \dots, X_{n-1})$ is lower bounded by 0, this implies that the limit $\lim_{n \rightarrow \infty} H(X_n | X_1, \dots, X_{n-1})$ must exist. It is $H'(\mathcal{X})$.

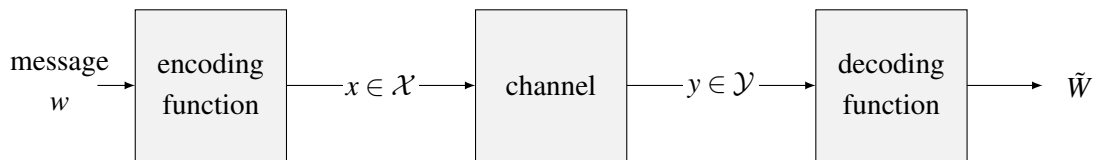
It remains to show that $H(\mathcal{X}) = H'(\mathcal{X})$:

$$\begin{aligned} H(\mathcal{X}) &= \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n H(X_i | X_1, \dots, X_{i-1}) \\ &= H'(\mathcal{X}). \end{aligned} \tag{4.7}$$

The final equality follows from Theorem 4.3.4, the Cesàro mean. ■

Chapter 5: Noisy-Channel Coding

In Chapter 2, we saw how to encode information from a source such that it can be stored or sent over some channel, and decoded at a later point in time. We assumed that the channel used to send the encoded information was perfect, meaning that no information got altered or lost while being sent over the channel. In this chapter, we consider a different setting, where the channel possibly contains some **noise**, that may convert the channel input x to some potentially different value y :



The goal is to design encoding and decoding functions that can resist this noise, so that the recovered message \tilde{W} is as close as possible to the original message w . The question is how short (efficient) such codes can be while still providing resistance to noise.

5.1 Channels

In this section, the notion of ‘channel’ will be made more precise, and different types of channels are considered.

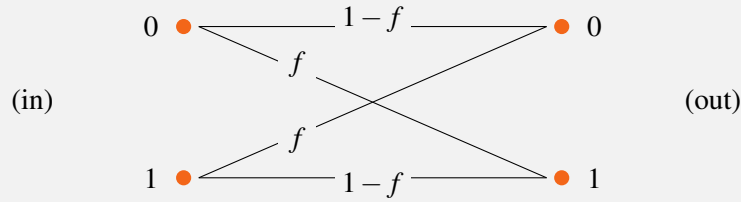
Definition 5.1.1 — Discrete channel. A (discrete) channel is a tuple $(\mathcal{X}, P_{Y|X}, \mathcal{Y})$ such that \mathcal{X} and \mathcal{Y} are finite sets, and $P_{Y|X} : \mathcal{X} \rightarrow [0, 1]$ is a probability distribution. \mathcal{X} represents the set of possible inputs, \mathcal{Y} the set of possible outputs, and $P_{Y|X}(y|x)$ is the probability of receiving output y given input x .

Note that if we give a distribution P_X for the set \mathcal{X} , this immediately determines a distribution P_Y for \mathcal{Y} by marginalization.

Example 5.1.2 — Binary symmetric channel (BSC). Define the binary symmetric channel with parameter $f \in [0, 1/2]$ by $\mathcal{X} = \mathcal{Y} = \{0, 1\}$ and

$$\begin{aligned} P_{Y|X}(0|0) &= P_{Y|X}(1|1) = 1 - f, \\ P_{Y|X}(0|1) &= P_{Y|X}(1|0) = f. \end{aligned}$$

With probability f , the input is flipped, and with probability $1 - f$, it remains unaffected. This channel can be represented visually as



Exercise 5.1.3 Verify that, up to relabeling of the outputs in \mathcal{Y} , a BSC with parameter $f \in [1/2, 1]$ is equivalent to a BSC with parameter $1 - f$. ■

We are interested in **memoryless** channels where the probability distribution of the output depends only on the current input. If the channel is used repeatedly, the channel distribution does not change depending on previous inputs and outputs. If we use a discrete memoryless channel n times, this can be regarded as the channel $(\mathcal{X}^n, P_{Y^n|X^n}, \mathcal{Y}^n)$ where

$$P_{Y^n|X^n}(\vec{y}|\vec{x}) = \prod_{i=1}^n P_{Y|X}(y_i|x_i), \quad (5.1)$$

Example 5.1.4 If we send 00110 over the binary symmetric channel from Example 5.1.2, we receive the output 01111 with probability

$$(1 - f) \cdot f \cdot (1 - f) \cdot (1 - f) \cdot f = (1 - f)^3 f^2.$$

In general, if a (memoryless) binary symmetric channel with noise parameter f is used n times, the probability of observing k bit flips is

$$\binom{n}{k} f^k (1 - f)^{n-k}.$$

Example 5.1.5 A binary symmetric channel (see Example 5.1.2) with parameter 0.1 is used to send 10^4 bits. The number of bit flips that result from the channel use are distributed according to the binomial($10^4, 0.1$) distribution: on expectation, $10^4 \cdot 0.1 = 1000$ bits are flipped, and the variance is $10^4 \cdot 0.1 \cdot 0.9 = 900$. Hence the **standard deviation** is $\sqrt{900} = 30$. With probability 99.7%, the number of bits that are flipped stays within three standard deviations of the mean, i.e. between 910 and 1090 bits are flipped with very high probability.

Definition 5.1.6 — Deterministic channel. A channel is deterministic if $H(Y|X) = 0$. In other words,

$$\forall x \in \mathcal{X} \exists y \in \mathcal{Y} : P_{Y|X}(y|x) = 1.$$

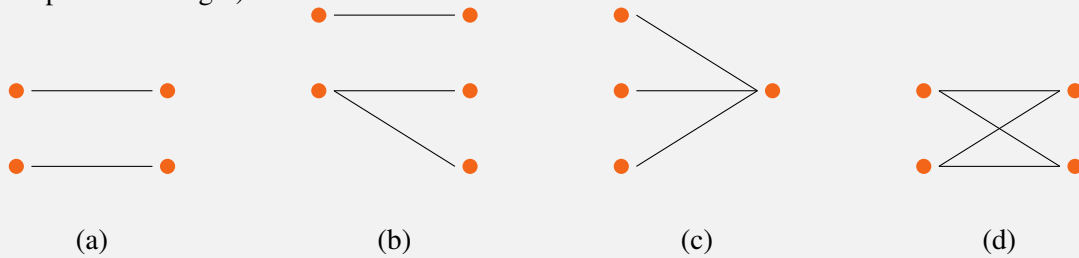
Definition 5.1.7 — Lossless channel. A channel is lossless (or **ideal**) if $H(X|Y) = 0$. In other words,

$$\forall y \in \mathcal{Y} \exists! x \in \mathcal{X} : P_{Y|X}(y|x) > 0.$$

(the notation $\exists! x$ means that there exists *exactly* one such x .)

In a deterministic channel, the output is completely determined by the input, whereas in a lossless channel, the input is completely determined by the output. A **noiseless** channel is a channel that is both deterministic and lossless.

Example 5.1.8 Consider the following channels. The connecting lines represent non-zero probabilities, but probabilities and inputs/outputs are not further specified (inputs are on the left, outputs on the right):



Channel (a) is both deterministic and lossless (and is therefore called a noiseless channel). Channel (b) is lossless but not deterministic, channel (c) is deterministic but not lossless, and channel (d) (the noisy binary channel) is neither.

5.2 Codes

In order to get as much information through a channel as possible, we can encode messages before sending them through the channel.

Definition 5.2.1 — Code. Let $M, n \in \mathbb{N}$. An (M, n) -code for the channel $(\mathcal{X}, P_{Y|X}, \mathcal{Y})$ consists of

- An index set $[M] = \{1, \dots, M\}$ representing the set of possible messages.
- A (possibly probabilistic) encoding function $\text{enc} : [M] \rightarrow \mathcal{X}^n$. This encoding function should be injective. n represents the number of channel uses we need to send a single message.
- A deterministic decoding function $\text{dec} : \mathcal{Y}^n \rightarrow [M]$. The set of all codewords, $\{\text{enc}(1), \dots, \text{enc}(M)\}$ is called the **codebook**.

The number of bits of information that are transmitted per channel use is captured by the following notion:

Definition 5.2.2 — Rate. The rate of an (M, n) -code is defined as

$$R := \frac{\log M}{n}.$$

Given a code for a specific channel, we can study the probability that an error occurs while

transmitting a message. For a fixed message m , this probability is represented by the following quantity.

Definition 5.2.3 — Probability of error. Given an (M, n) code for a channel $(\mathcal{X}, P_{Y|X}, \mathcal{Y})$, the probability of error λ_m is the probability that the decoded output is not equal to the input message m . More formally,

$$\lambda_w := P[\text{dec}(Y^n) \neq m \mid X^n = \text{enc}(w)].$$

Given this quantity, the **maximal probability of error** is defined as

$$\lambda^{(n)} := \max_{w \in [M]} \lambda_w.$$

Similarly, the **average probability of error** is defined as

$$p_e^{(n)} := \frac{1}{M} \sum_{w=1}^M \lambda_w.$$

5.2.1 Repetition codes

A simple and intuitive code is the n -bit **repetition code** R_n : a single message bit is encoded by simply repeating the bit n times. Decoding is done by ‘majority vote’, that is, $\text{dec}(y) = \text{MAJ}(y_1, \dots, y_n)$, which is 1 if and only if (strictly) more than half of the bits in y are 1s. In order to avoid “ties” in the decoding, repetition codes usually require that n is odd.

Example 5.2.4 — 3-bit repetition code. Consider the 3-bit repetition code R_3 . It is a $(2, 3)$ code with codebook $\{000, 111\}$. The rate of R_3 is $1/3$. The probability of error for the message $w = 0$, sent over a BSC with bit flip probability f , is

$$\begin{aligned} \lambda_0 &= P[\text{dec}(Y^n) \neq 0 \mid X^n = 000] \\ &= P[Y^n = 011 \cup Y^n = 101 \cup Y^n = 110 \cup Y^n = 111 \mid X^n = 000] \\ &= 3f^2(1-f) + f^3. \end{aligned}$$

A similar calculation shows that $\lambda_1 = \lambda_0$. Hence, the maximal and average probability of error are equal to λ_0 as well.

As a concrete example, if $f = 0.1$, the 3-bit repetition code has an error probability of approximately 0.03. Hence, the 3-bit repetition code provides an error probability that is about three times lower (3% instead of 10%), at the expense of a rate that is a factor 3 worse than simply sending the messages through the channel without encoding.

In general, the n -bit repetition code is a $(2, n)$ code with the relatively low rate of $1/n$ and an average/maximal probability of error of

$$\sum_{k=(n+1)/2}^n \binom{n}{k} f^k (1-f)^{n-k}$$

when used on a binary symmetric channel with bit flip probability f .

Example 5.2.5 — Application of the repetition code. A company producing hard drives wants to make sure that the probability that a customer experiences a storage failure is at most 0.1%. The company has 1000 customers, who each want to store 10 GB of data every day for

10 years without experiencing a failure. That means that the probability of a single stored bit being stored incorrectly has to be at most

$$\frac{1}{10 \cdot 8 \cdot 10^9} \cdot \frac{1}{10 \cdot 365} \cdot \frac{1}{1000} \cdot 0.001 \geq 10^{-20}.$$

The company tries to achieve that error probability using a repetition code R_n and some imperfect drives with bit flip probability $f = 0.1$. How big should n be to achieve $p_e^{(n)} \leq 10^{-20}$?

$$\begin{aligned} p_e^{(n)} &= \lambda_0 = P[MAJ(Y^n) \neq 0 \mid X^n = 0^n] \\ &= \sum_{k=(n+1)/2}^n \binom{n}{k} 0.1^k 0.9^{n-k} \\ &\approx \binom{n}{n/2} (0.1 \cdot 0.9)^{n/2} \\ &\approx 2^{n \cdot h(1/2)} \cdot 0.09^{n/2} \\ &= 2^n \cdot (\sqrt{0.09})^n \\ &= 0.6^n. \end{aligned}$$

In order for $0.6^n \leq 10^{-20}$ to be true, it has to be that $n \approx 90$.

5.2.2 The [7, 4] Hamming code

A more sophisticated error-correcting code is the [7, 4] **Hamming code**. It is a $(2^4, 7)$ code, meaning that it encodes a 4-bit message (there are 2^4 such messages) into 7 bits. The encoding function is defined as

$$\text{enc}(m_1 m_2 m_3 m_4) = m_1 m_2 m_3 m_4 t_5 t_6 t_7,$$

where the **parity bits**

$$t_5 = m_1 \oplus m_2 \oplus m_3,$$

$$t_6 = m_2 \oplus m_3 \oplus m_4,$$

$$t_7 = m_1 \oplus m_3 \oplus m_4$$

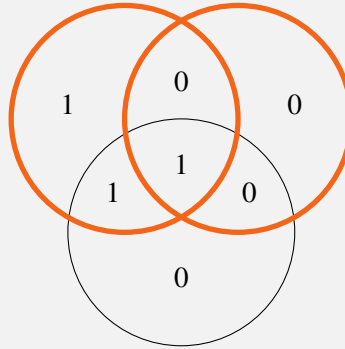
are appended at the right. Note that the choice of these parity bits differs throughout the literature.

Decoding is done by making sure that all parity bits check out, and if not, making the smallest possible number of bit flips such that they do. It is important to note that the parity bits themselves may have been flipped during the transmission through the channel. A visual way to perform this parity check is by using the following diagram:

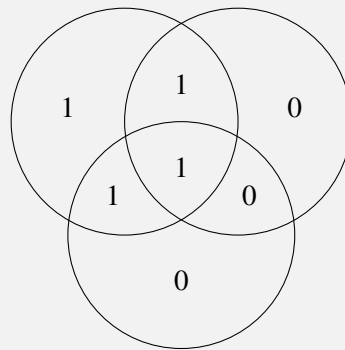


For all of the three circles, the parity bit should equal the parity of the three message bits in that circle.

Example 5.2.6 The decoding of the string 1010100 is performed as follows. First, we fill in the bits into the diagram:



We see that in two of the circles, the parity bit is incorrect. This is called the **error syndrome** (for a more formal definition of the error syndrome, see Section 5.2.3). By flipping only the bit m_2 (which is in the intersection of the top two circles, but not the bottom one), we can fix all the parity bits:



Hence, $\text{dec}(1010100) = 1110$.

The $[7,4]$ Hamming code can correctly decode if the codeword is corrupted in (at most) one place. However, if the codeword is corrupted in two or more places, the Hamming code will not decode correctly. For example, if $\text{enc}(0001) = 0001011$ is corrupted in two places to become 0011010, it will decode to $\text{dec}(0011010) = 0111$. The number of bit flips a code can correct for depends on the minimal distance between the words in the codebook:

Definition 5.2.7 — Hamming distance. The Hamming distance between two n -bit strings x and y is defined as

$$d(x, y) := \sum_{i=1}^n |x_i - y_i|.$$

One can equivalently define $d(x, y) = |x - y|$ where $|z|$ denotes the **Hamming weight** of a binary string: the number of ones in that string.

Definition 5.2.8 — Minimal distance. Given a code with codebook C , the minimal distance

of that code is defined as

$$d_{\min} := \min_{\substack{x, y \in C \\ x \neq y}} d(x, y).$$

By checking all pairs of codewords of the $[7, 4]$ Hamming code, one can verify that its minimal distance is 3 (for this reason, it is often called a $[7, 4, 3]$ code). Hence, if two bits in a codeword are flipped, it might be closer to some other codeword in terms of number of bit flips. By flipping a single bit, the channel output is (incorrectly) decoded into the message that corresponds to that other codeword.

Exercise 5.2.9 Compute the rate of the $[7, 4]$ Hamming code. ■

5.2.3 Generalization: Linear Codes

The repetition code and the $[7, 4]$ Hamming code both belong to the more general class of linear (error-correcting) codes, which we study in this section.

Definition 5.2.10 — Linear code. A code C is linear if any linear combination of codewords is also a codeword.

For the definition of linearity to make sense, addition and multiplication by constants needs to be defined on \mathcal{X}^n (formally, \mathcal{X}^n needs to be a **vector space**). In this case, C is linear if it is a **linear subspace** of \mathcal{X}^n . In the following, we will assume that $\mathcal{X} = \{0, 1\}$: in that case, we are talking about **binary codes** and addition is simply bitwise addition modulo 2 (which is the exclusive OR function). Note that for binary codes, addition and subtraction is the same operation, as $-1 = +1$ modulo 2.

Exercise 5.2.11 Check that the repetition code R_n and the $[7, 4]$ Hamming code are indeed linear. ■

For linear codes, the notation “ $[n, k, d]$ (linear) code” is often used to denote the fact that n bits are used to encode k message bits, with a minimal distance d . Be aware that in the notation of Definition 5.2.1, this would be a $(2^k, m)$ code. We use different types of brackets to differentiate between the two notations.

Definition 5.2.12 — Generator matrix. Given a $[n, k, d]$ linear code, the generator matrix G^T is an $n \times k$ matrix such that the columns c_1, \dots, c_k form a basis for C . The generator matrix can always be transformed into **systematic form**:

$$G^T = \begin{pmatrix} \mathbb{I}_k \\ P \end{pmatrix},$$

where P is some $(n - k) \times k$ matrix representing the parity bits of the code.

The generator matrix is used in the encoding function as follows:

$$\text{enc}(m) = G^T \cdot m.$$

The codebook C is the set $\{G^T \cdot m \mid m \in \mathcal{X}^k\}$.¹

¹The reason for the transposition in G^T is that historically, coding theorists preferred to use row vectors and matrix multiplication from the right instead of column vectors and multiplication from the left, which is more standard in other areas. Notice that for row vectors $c = m \cdot G$, we equivalently have column vectors $c^T = (m \cdot G)^T = G^T \cdot m^T$.

Example 5.2.13 — Generator matrix of the $[7,4]$ Hamming code. The following 7×4 matrix generates the $[7,4]$ Hamming code:

$$G^T = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ \hline 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \end{pmatrix}.$$

The generator matrix is given in systematic form. To encode the message 1010, we compute

$$G^T \begin{pmatrix} 1 \\ 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}.$$

The decoding operation can be described by defining the parity-check matrix:

Definition 5.2.14 — Parity-check matrix. Let C be a code with generator matrix

$$G^T = \begin{pmatrix} \mathbb{I}_k \\ P \end{pmatrix}.$$

Then the parity check matrix for C is given by

$$H := (-P \mid \mathbb{I}_{n-k}).$$

Note that if $\mathcal{X}^n = \{0, 1\}^n$, then P is a matrix with binary entries, and hence $-P = P$.

Definition 5.2.15 — Syndrome. Let C be a code with parity-check matrix H . The syndrome of a received codeword $c \in \mathcal{Y}^n$ is Hc .

For all $c \in C$, the syndrome is the all-zero vector, since $c = G^T m$ for some m , and

$$HG^T m = (-P \mid \mathbb{I}_{n-k}) \begin{pmatrix} \mathbb{I}_k \\ P \end{pmatrix} m = 0^{n-k}, \quad (5.2)$$

where we make use of **block matrix multiplication**.

The syndrome gives a lot of valuable information about where an error occurred. Suppose a codeword c is sent over a channel, and a single bit of c is flipped, at the i th position. The output of the channel is thus $c' = c + e_i$ (where e_i is the i th unit vector, which is 1 at position i and 0 elsewhere). The syndrome for this received output is

$$Hc' = H(c + e_i) = Hc + He_i = He_i, \quad (5.3)$$

which is the i th column of H . So by comparing the syndrome to the parity-check matrix H , we can find out which bit was most likely flipped.

Example 5.2.16 — Parity-check matrix of the $[7,4]$ Hamming code. The following 3×7 matrix is the parity-check matrix for the $[7,4]$ Hamming code:

$$H = \begin{pmatrix} 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 \end{pmatrix}.$$

When checking, for example, the codeword 1010100, a syndrome of 110 is observed, which means that the second bit (m_2) should be flipped to correct it: the decoding is 1110. Notice that all columns of H are different, and that there are $2^3 = 8$ possible syndromes we might observe. This aligns well with the fact that there are three parity bits that can all either be correct or incorrect.

Example 5.2.17 — The 5-bit repetition code. The generating matrix of R_5 , the 5-bit repetition code, is

$$G^T = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}.$$

The parity check matrix is

$$H = \begin{pmatrix} 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 \end{pmatrix}.$$

If the codeword 10100 is received, the syndrome is 1011. This is exactly the vector $e_1 + e_3$, the sum of the first and third column. Hence, we should flip the first and third bit to correct the error. Note that 1011 is also equal to $e_2 + e_4 + e_5$, so another option to arrive at a correct codeword is to flip the second, fourth and fifth bit. It is less likely, however, that three bit flips occurred instead of two.

In general, a code with minimal distance d can *detect* up to $d - 1$ bit flip errors, because in $d - 1$ bit flip ‘steps’, one can never arrive at another valid codeword. The code can accurately *correct* up to $\frac{d-1}{2}$ bit flip errors. This is because if all codewords are guaranteed to be at least d bit flips apart, the set of strings that result from $\frac{d-1}{2}$ bit flips on a codeword c_1 never overlaps with the set of strings that result from that number of bit flips on another codeword c_2 .



If every codeword needs an area around it that it does not share with any other codeword, this has consequences for the total number of codewords.

Proposition 5.2.18 — Hamming bound. If C is a binary code of block length n and minimum distance 3, then $|C| \leq \frac{2^n}{n+1}$.

Proof. For each $c \in C$, define the neighborhood of c to be $N(c) := \{y \in \{0, 1\}^n \mid d(x, y) \leq 1\}$. Every such neighborhood contains exactly $n + 1$ elements. Since $d = 3$, $N(c) \cap N(c') = \emptyset$ whenever $c \neq c'$. Hence,

$$2^n \geq \left| \bigcup_{c \in C} N(c) \right| = \sum_{c \in C} |N(c)| = |C| \cdot (n + 1), \quad (5.4)$$

and the result follows. ■

The $[7, 4]$ Hamming code is optimal in the sense that it achieves this Hamming bound: it is a code with block length 7 and minimal distance 3, so an upper bound to the codebook size is $\frac{2^7}{7+1} = 2^4$. The Hamming code achieves exactly this codebook size. There is an **entire family of Hamming codes** that are optimal in this sense. They form $[2^r - 1, 2^r - 1 - r, 3]$ codes for $r \in \mathbb{N}$. The columns of their parity-check matrix consist of all non-zero r -bit strings.

We conclude this section with two ways to determine the minimal distance of a code. Apart from the trivial way (determining the entire codebook and comparing all the codeword pairs), it turns out that it is already enough to consider just the Hamming weights of the (nonzero) codewords:

Proposition 5.2.19 For a code C , the minimal distance is equal to the minimal weight of the nonzero codewords.

Proof. The following derivation proves the claim:

$$d_{\min} = \min_{\substack{x, y \in C \\ x \neq y}} d(x, y) = \min_{\substack{x, y \in C \\ x \neq y}} \sum_{i=1}^n |x_i - y_i| = \min_{\substack{x, y \in C \\ x \neq y}} d(x - y, 0) = \min_{\substack{z \in C \\ z \neq 0}} d(z, 0) = \min_{\substack{z \in C \\ z \neq 0}} |z| \quad (5.5)$$

■

An equivalent way to determine the minimal distance of a code is possible if the parity check matrix is known.

Proposition 5.2.20 For a linear code C with parity check matrix H , the minimal distance d_{\min} equals the minimum number of columns of H that are linearly dependent.

Proof. Left as an exercise. ■

Chapter 6: Zero-Error Channel Coding

In this chapter, we consider the problem of using a noisy channel to transmit a message perfectly, i.e. with vanishing maximal probability of error. For some channels, for example the non-trivial binary symmetric channel with $f \notin \{0, 1\}$, it is not possible to send multiple different messages over the channel in this way. For other channels, an interesting question is: how many messages (or how much information) can be sent over this channel in an error-free way?

Example 6.0.1 Look again at channel (b) from Example 5.1.8, now with some specific input/output sets and probabilities:



We can send two messages, m_1 and m_2 , over the channel by defining $\text{enc}(m_1) = a$ and $\text{enc}(m_2) = b$. The decoding is defined as $\text{dec}(0) = m_1$, and $\text{dec}(1) = \text{dec}(2) = m_2$.

Example 6.0.2 — Noisy typewriter. The noisy typewriter channel is defined as follows:



There is a way to send two messages m_1 and m_2 error-free over this channel by defining $\text{enc}(m_1) = a$ and $\text{enc}(m_2) = b$. The decoding function is defined as $\text{dec}(a) = \text{dec}(b) = a$, and $\text{dec}(c) = \text{dec}(d) = c$ (note that the definition of $\text{dec}(e)$ is irrelevant, as this output symbol will never be observed).

Is there a way to encode three different messages in an error-free way? Upon inspection, we see that any encoding function $\text{enc}(\cdot)$ on three messages will map at least two messages to channel inputs that are **confusable** (i.e. are possibly mapped to the same channel output).

In general, it is not easy to tell directly from the channel how many messages can be perfectly transmitted. We invoke some graph theory to help us with the analysis.

6.1 Confusability graphs

Definition 6.1.1 — Graph. A (undirected simple) graph G consists of a set $V(G)$ of **vertices** and a set $V(E)$ of **edges**. The edges are unordered pairs of vertices (each edge connects two different vertices of the graph).

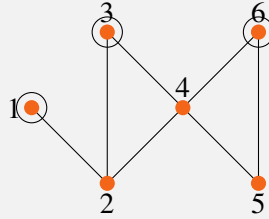
Definition 6.1.2 — Independence number. The independence number $\alpha(G)$ of a graph G is the size of the largest **independent set** of G , where an independent set of G is a set $S \subseteq V(G)$ such that

$$\forall x, x' \in S : (x, x') \notin E(G).$$

That is, an independent set S in G is a set of vertices such that there is no edge between any of the vertices.

Finding the independence number of a graph is an NP-hard problem, meaning there is no known efficient method for finding the independence number of an arbitrary graph.

Example 6.1.3 Consider the following graph with $V(G) = \{1, 2, 3, 4, 5, 6\}$ and $E(G) = \{\{1, 2\}, \{2, 3\}, \{2, 4\}, \{3, 4\}, \{4, 5\}, \{4, 6\}, \{5, 6\}\}$:



A maximal independent set $\{1, 3, 6\}$ is marked in the graph. Note that this set is not unique: $\{1, 3, 5\}$ is also an independent set of size 3. As there is no independent set of size 4, the independence number of this graph is 3.

The graph-theoretic notions just introduced can be used to analyze the number of channel inputs that can safely be used simultaneously in an error-free code.

Definition 6.1.4 — Confusability graph. Let $(\mathcal{X}, P_{Y|X}, \mathcal{Y})$ be a channel. The confusability graph G for the channel consists of

$$V(G) := \mathcal{X},$$

the set of input symbols of the channel, and

$$E(G) := \{\{x, x'\} \subset \mathcal{X} \mid x \neq x' \text{ and } \exists y \in \mathcal{Y} \text{ s.t. } P_{Y|X}(y|x) \cdot P_{Y|X}(y|x') > 0\}$$

is the set of input pairs that are confusable (because they reach a shared output symbol $y \in \mathcal{Y}$).

Example 6.1.5 — Confusability graph of the noisy typewriter. Consider again the noisy typewriter from Example 6.0.2. The confusability graph for that channel is



This graph is also known as C_5 , the circle of size 5. Its independence number is $\alpha(C_5) = 2$.

In the above example, the independence number of the confusability graph is exactly the number of messages that can be sent over the channel perfectly. This is no coincidence:

Proposition 6.1.6 Given a channel with confusability graph G , the maximal message set $[M]$ that can be communicated perfectly in a single channel use is of size $\alpha(G)$.

Proof. Let $(x, x') \in E(G)$, i.e. x and x' are confusable. They cannot both be used to send different messages, for suppose there are messages $m \neq m'$ such that $\text{enc}(m) = x$ and $\text{enc}(m') = x'$, then by definition of the confusability graph, there is a $y \in \mathcal{Y}$ such that x and x' are both mapped to y with nonzero probability. In order to correctly decode in all cases, it must therefore be that $\text{dec}(y) = m$ and $\text{dec}(y) = m'$, contradicting the assumption that $m \neq m'$. Therefore, the number of messages that can be sent over the channel cannot exceed the independence number $\alpha(G)$.

For the other direction, it is easy to find an encoding and decoding function for $\alpha(G)$ different messages. Let $\{x_1, x_2, \dots, x_{\alpha(G)}\}$ be a largest independent set of G . Define $\text{enc}(m_i) = x_i$ for all $i \in [\alpha(G)]$. Then for all $y \in \mathcal{Y}$, by definition of the confusability graph and the independent set, there is exactly one i such that $P_{Y|X}(y|x_i) > 0$. Define $\text{dec}(y) = m_i$. This code can send $\alpha(G)$ different messages over the channel without error. ■

6.2 Multiple channel uses

In the previous section we have found that the independence number of the confusability graph is a conclusive upper bound for the number of messages that can be sent over a channel *in a single use*. But it turns out that if the same channel is allowed to be used more than once, it is sometimes possible to achieve a higher rate, still error-free!

Example 6.2.1 — Two uses of the noisy typewriter. Let us reconsider the noisy typewriter from Examples 6.0.2 and 6.1.5, but this time we are allowed to use the channel twice. We can of

course use the following encoding function to encode $2 \times 2 = 4$ different messages:

$$\begin{aligned}\text{enc}(m_1) &= \text{aa} \\ \text{enc}(m_2) &= \text{ac} \\ \text{enc}(m_3) &= \text{ca} \\ \text{enc}(m_4) &= \text{cc}\end{aligned}$$

This encoding function is defined by simply concatenating the encoding function of Example 6.0.2. It has a rate of $\log(4)/2 = 1$, just like the single-use noisy typewriter channel code. There is, however, a code that is able to encode and correctly decode *five* different messages! It uses codewords aa, bc, ce, db, ed as codewords for messages m_1 through m_5 . To see that the decoding is error-free, observe that the channel maps the codewords to the following sets of possible outputs:

$$\begin{aligned}\text{aa} &\mapsto \{11, 22, 12, 21\} \\ \text{bc} &\mapsto \{23, 24, 33, 34\} \\ \text{ce} &\mapsto \{35, 45, 41, 31\} \\ \text{db} &\mapsto \{42, 52, 43, 53\} \\ \text{ed} &\mapsto \{54, 55, 14, 15\}\end{aligned}$$

All of these sets are non-overlapping, so none of the inputs are confusable. The rate of this code is $\log(5)/2 \approx 1.161$, which is strictly better than the single-use channel!

The above example shows that there are optimal codes that make better use of the channel if the channel is used more than once. Is there a structural way to find these codes? Let us expand our graph theoretic tools.

Definition 6.2.2 — Strong graph product. Let G, H be two graphs. We define the strong graph product $G \boxtimes H$ as follows. The set of vertices is

$$V(G \boxtimes H) := V(G) \times V(H).$$

The set of edges is

$$E(G \boxtimes H) := \{ \{(x, y), (x', y')\} \mid (x, y) \neq (x', y') \text{ and } (x = x' \text{ or } \{x, x'\} \in E(G)) \\ \text{and } (y = y' \text{ or } \{y, y'\} \in E(H)) \},$$

i.e. there is an edge between (x, y) and (x', y') if and only if the vertices of G are confusable (or equal) *and* the vertices of H are confusable (or equal).

For our application, we are often interested in the strong graph product of a graph G with itself, possibly many times. Therefore it is useful to work out the definition of $G^{\boxtimes n}$, based on Definition 6.2.2:

$$\begin{aligned}V(G^{\boxtimes n}) &= V(G) \times \cdots \times V(G) \\ E(G^{\boxtimes n}) &= \{ \{(x_1, \dots, x_n), (v_1, \dots, v_n)\} \mid (x_1, \dots, x_n) \neq (v_1, \dots, v_n) \text{ and } \forall i : x_i = v_i \vee \{x_i, v_i\} \in E(G) \}\end{aligned}$$

Example 6.2.3 — Binary symmetric channel. The confusability graph of the non-trivial binary symmetric channel (BSC) from Example 5.1.2 is C_2 :



This graph is also known as K_2 , the **complete graph** of size 2. $G \boxtimes G$ is the complete graph on 4 vertices K_4 :



In this case, using the channel twice does not help: the independence number of both graphs is 1, and the rate of any error-free code is therefore zero.

Exercise 6.2.4 Show that for all $n, m \in \mathbb{N}$, $\alpha(G \boxtimes H) \geq \alpha(G)\alpha(H)$. ■

Example 6.2.5 Consider the following graph H :



The strong graph product of H with the graph G from Example 6.2.3 is



The independence number of this graph is 2. The strong product of H with itself is



The independence number of this graph is 4. As the graphs get bigger, the independence number is increasingly hard to compute.

Exercise 6.2.6 Design a channel that has confusability graph H from Example 6.2.5. ■

6.3 Shannon capacity

In the previous section it became apparent that for some channels, two or more channel uses can provide a better rate than a single use of those channels. The maximum rate that can be achieved in this way is captured by the notion of Shannon capacity of the confusability graph of the channel.

Definition 6.3.1 — Shannon capacity. The Shannon capacity of a graph is

$$c(G) := \sup_{n \in \mathbb{N}} \frac{\log \alpha(G^{\boxtimes n})}{n}.$$

The Shannon capacity represents the maximum number of bits per channel use that can be perfectly communicated, over a channel with confusability graph G . Intuitively, allowing more channel uses can only increase the rate. This is captured by the following lemma.

Proposition 6.3.2

$$c(G) = \lim_{n \rightarrow \infty} \frac{\log \alpha(G^{\boxtimes n})}{n}.$$

Proof. First, note that $\alpha(G^{\boxtimes(k+\ell)}) \geq \alpha(G^{\boxtimes k})\alpha(G^{\boxtimes \ell})$ by Exercise 6.2.4. It then follows that

$$\log \alpha(G^{\boxtimes(k+\ell)}) \geq \log \alpha(G^{\boxtimes k}) + \log \alpha(G^{\boxtimes \ell}),$$

and so the sequence $(\log \alpha(G^{\boxtimes n}))_{n \in \mathbb{N}}$ is superadditive. Then by **Fekete's lemma**,

$$\lim_{n \rightarrow \infty} \frac{\log \alpha(G^{\boxtimes n})}{n}$$

exists and is equal to the supremum. ■

Example 6.3.3 We calculate the Shannon capacity of the graph H from Example 6.2.5.

$$\begin{aligned} c(H) &= \lim_{n \rightarrow \infty} \frac{\log \alpha(H^{\boxtimes n})}{n} \\ &= \lim_{n \rightarrow \infty} \frac{\log 2^n}{n} = 1. \end{aligned}$$

The second equality holds because the graph $H^{\boxtimes n}$ has a maximal independent set consisting of the ‘outermost’ corners of the graph (convince yourself of this for $n = 1, 2, 3$). There are 2^n such corners.

A channel with confusability graph H does not benefit from multiple channel uses: we can still only send a single bit of information over the channel on average.

The exact computational complexity of $c(G)$ is unknown. It is not even known to be a decidable problem. From Example 6.2.1, we know that $c(C_5)$ is at least $\frac{\log 5}{2}$, a fact that has been known

since Shannon showed it in 1956, but how can we tell whether it is bigger than that number? This question remained open until Lovasz showed in 1979 that $c(C_5) = \frac{\log 5}{2}$. For the relatively small circle of size 7, C_7 , the exact Shannon capacity remains unknown.

From the fact that $c(C_5) = \frac{\log 5}{2}$, we know that the 5-letter noisy typewriter does not benefit from more than two channel uses. In fact, all graphs for which the Shannon capacity is known attain that capacity with one, two or an infinite number of channel uses. It is not known if there is a deeper reason for this observed pattern.

In conclusion, the problem of zero-error channel coding studied in this chapter provides a nice connection between information theory and graph theory, a core topic of theoretical computer science and of great combinatorial interest. Zero-error information theory is an active area of research, in particular with respect to quantum variants of these problems.

Chapter 7: The Noisy-Channel Coding Theorem

In Chapter 6 we discovered that for some noisy channels, zero-error communication is very hard, or even impossible. For example, if Alice and Bob have to communicate over a binary symmetric channel (BSC, see Example 5.1.2) that has non-zero bit-flip probability, they cannot hope to do any zero-error communication, because the Shannon capacity of the BSC's confusability graph is zero.

In this chapter, Alice and Bob try to achieve a more modest (and realistic) goal. Consider again the setting where they communicate over a noisy, memoryless channel (as described in Section 5.1), but now they are content with a small probability of error. That is, Bob has to have high confidence that the message he decodes is correct, but he is allowed to be wrong with some small probability. Can Alice and Bob do more in this setting? For example, can they communicate information reliably over a noisy BSC?

7.1 Channel capacity

Central to our study of communication over noisy channels will be the concept of channel capacity. It reflects the maximum amount of information that could *in principle* be communicated with a single use of a channel.

Definition 7.1.1 — Channel capacity. The channel capacity C of a discrete, memoryless channel $(\mathcal{X}, P_{Y|X}, \mathcal{Y})$ is given by

$$C := \max_{P_X} I(X; Y).$$

The output distribution P_Y is fully determined by the input distribution P_X and the channel $P_{Y|X}$. Maximizing over P_Y rather than P_X would yield the same result for C . [Chris: what do you mean? How is P_{XY} defined given $P_{Y|X}$ and P_Y ? Bayes': $P_{X|Y} = P_{Y|X} \frac{P_X}{P_Y}$, we always seem to miss P_X .]

The following exercise asks you to check a basic property of channel capacity: the input alphabet size limits the amount of information the sender can communicate with a single use of the channel.

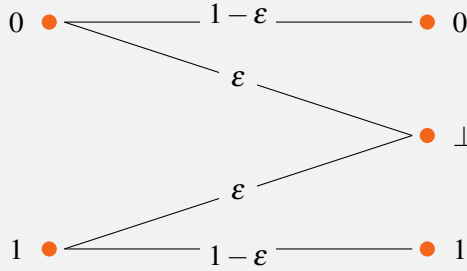
Exercise 7.1.2 Prove that $0 \leq C \leq \log |\mathcal{X}|$. ■

Example 7.1.3 — Capacity of a BSC. Recall the binary symmetric channel (BSC) from Example 5.1.2 with parameter $f \in [0, 1/2]$. The capacity of that channel is

$$\begin{aligned}
 \max_{P_X} I(X; Y) &= \max_{P_X} (H(Y) - H(Y|X)) \\
 &= \max_{P_X} \left(H(Y) - \sum_{x \in \mathcal{X}} P_X(x) \cdot H(Y|X=x) \right) \\
 &= \max_{P_X} \left(H(Y) - \sum_{x \in \mathcal{X}} P_X(x) \cdot h(f) \right) \\
 &= \max_{P_X} (H(Y) - h(f)) \\
 &= 1 - h(f).
 \end{aligned}$$

The last step follows because $H(Y)$ is maximized if Y is uniform, which is achievable by choosing X to be uniform.

Example 7.1.4 — Capacity of a BEC. Consider the binary erasure channel (BEC) with $\mathcal{X} = \{0, 1\}$ and $\mathcal{Y} = \{0, 1, \perp\}$, where \perp is the **erasure symbol**, and $\varepsilon \in [0, 1]$ is the **erasure probability**:



The channel capacity of the BEC is a function of ε , and computed as follows. Write p for $P_X(0)$.

$$\begin{aligned}
 \max_{P_X} I(X; Y) &= \max_{P_X} (H(X) - H(X|Y)) \\
 &= \max_p \left(h(p) - \sum_{y \in \mathcal{Y}} P_Y(y) \cdot H(X|Y=y) \right) \\
 &= \max_p (h(p) - P_Y(\perp) \cdot h(p)) \\
 &= \max_p (h(p)(1 - \varepsilon)) \\
 &= 1 - \varepsilon.
 \end{aligned}$$

Again, the last step follows because $H(X) = h(p)$ is maximized if X is uniform, hence $p = \frac{1}{2}$.

As stated, the channel capacity reflects the maximum amount of information that could theoretically be sent over a noisy channel per use of that channel. The question remains whether this capacity is **achievable** by an actual code in the following sense:

Definition 7.1.5 — Achievable rate. For a given channel, a rate R is achievable if there exists a sequence of $(2^{nR}, n)$ codes (for $n = 1, 2, 3, \dots$) such that $\lambda^{(n)} \xrightarrow{n \rightarrow \infty} 0$. Here, $\lambda^{(n)}$ is the maximum error probability as defined in Definition 5.2.3.

Note that any $(2^{nR}, n)$ code has rate R , since $\frac{1}{n} \cdot \log 2^{nR} = R$. Thus, a rate R is achievable if there exists a sequence of rate R codes such that increasing the block size reduces the maximum error probability asymptotically to zero.

The rest of this chapter will be concerned with proving a theorem due to Shannon, which states that any rate R that is strictly below the capacity C is achievable (Theorem 7.3.1), and conversely, that any rate strictly above is not achievable (Theorem 7.4.1). These theorems together are known as **Shannon's noisy-channel coding theorem**.

Very roughly, the first theorem follows from the asymptotic equipartition property (AEP) discussed in Chapter 2.6. Increasing n (the size of the input blocks) results in output blocks that are more likely to be typical. But can we be sure that no two typical input sequences result in the same typical output sequence? Such a collision would make it impossible to decode correctly. To show that this does not form a problem, we need stronger mathematical tools.

7.2 Joint Asymptotic Equipartition Property

In order to decode a channel output Y^n correctly, it is important that there is only one input codeword that is likely to have resulted in that specific output. The following concept will help us with formalizing this idea:

Definition 7.2.1 — Joint typicality. For a joint distribution P_{XY} , the jointly typical set $A_\varepsilon^{(n)}$ is defined as

$$A_\varepsilon^{(n)} := \{(x^n, y^n) \in \mathcal{X}^n \times \mathcal{Y}^n \mid \begin{aligned} &\left| -\frac{1}{n} \log P_{X^n}(x^n) - H(X) \right| < \varepsilon, \\ &\left| -\frac{1}{n} \log P_{Y^n}(y^n) - H(Y) \right| < \varepsilon, \\ &\left| -\frac{1}{n} \log P_{X^n Y^n}(x^n, y^n) - H(XY) \right| < \varepsilon \}. \end{aligned}$$

Here, $P_{X^n Y^n} := \prod_{i=1}^n P_{XY}(x_i, y_i)$.

A jointly typical set is a set where both elements of the pair are typical individually under the marginal distributions (see Definition 2.6.6), and their combination is typical under the joint distribution. To get some intuition for what a typical set looks like, we give two examples: one where X and Y are very correlated, and one where they are independent.

Example 7.2.2 Consider the almost fully correlated random variables X and Y with the following distribution:

	$X = 0$	$X = 1$
$Y = 0$	0.45	0.05
$Y = 1$	0.05	0.45

Note that the marginal distributions of X and Y are uniform, and so for all n and all $x^n \in \{0, 1\}^n$,

$$-\frac{1}{n} \log P_{X^n}(x^n) - H(X) = -\frac{1}{n} \log \frac{1}{2^n} - \log 2 = 0,$$

and a similar statement holds for Y . Hence, the first and second inequalities in Definition 7.2.1 are always satisfied for any $\varepsilon > 0$. What about the third inequality? The table below visualizes the probabilities of every pair for $n = 3$ and $\varepsilon = 0.35$: every square represents a pair of elements from \mathcal{X}^3 and \mathcal{Y}^3 . Bigger squares represent higher probabilities. The typical pairs are marked in orange.

		Y							
	000	001	010	100	011	101	110	111	X
	000	■	■	■	■	■	■	■	
	001	■	■	■	■	■	■	■	
	010	■	■	■	■	■	■	■	
	100	■	■	■	■	■	■	■	
	011	■	■	■	■	■	■	■	
	101	■	■	■	■	■	■	■	
	110	■	■	■	■	■	■	■	
	111	■	■	■	■	■	■	■	

As expected, pairs (x^3, y^3) with a higher overlap in bits have a higher probability of occurring. The jointly typical set has a very clear structure: only those pairs with $x^3 = y^3$ are included.

Increasing ε would result in a bigger typical set. For example, for $\varepsilon = 0.8$, all pairs that differ at at most one position (e.g., $(000, 001)$) are part of the jointly typical set. Increasing n would have the same effect: for example, for $n = 10$ and $\varepsilon = 0.35$, all pairs of sequences that differ at 0, 1, or 2 positions are part of the jointly typical set.

For big n and small ε we see an interesting pattern appear. Take, for example, $n = 10,000$ and $\varepsilon = 0.01$. A pair of sequences differing in $0 \leq m \leq 10,000$ positions occurs with probability $0.45^{10,000-m} \cdot 0.05^m$ under P_{XY} . The third inequality in Definition 7.2.1,

$$\left| -\frac{1}{10,000} \log(0.45^{10,000-m} \cdot 0.05^m) - H(XY) \right|,$$

is satisfied for $969 \leq m \leq 1031$, or around 10 percent of the positions. This might match your expectation, since the probability of sampling a non-matching pair of bits $((0,1)$ or $(1,0))$ from P_{XY} is 10 percent. Thus, these values of m result in very *typical* pairs of sequences. Note in particular that pairs of sequences that differ in less than 969 positions, which have the highest probability of occurring, are not included in the [Chris: jointly?] typical set! [Chris: hard to follow... is it worth trying to make a (huge) picture?] However, the number of such pairs is so small compared to the number of pairs which differ in 969 to 1031 positions, that their total probability is not relevant. As we will see shortly, the elements in the [Chris: jointly?] typical set together constitute almost all of the total probability.

Example 7.2.3 Consider the following distribution on X and Y :

	X = 0	X = 1
Y = 0	1/8	1/8
Y = 1	3/8	3/8

Note that X is uniform, and that X and Y are independent.

As in Example 7.2.2, since X is uniform, the first inequality in Definition 7.2.1 is always

satisfied for any $\varepsilon > 0$.

Furthermore, since X and Y are independent, the third inequality in Definition 7.2.1 is satisfied whenever the second inequality is, because

$$\begin{aligned} -\frac{1}{n} \log P_{X^n Y^n}(x^n, y^n) - H(XY) &= -\frac{1}{n} \log P_{X^n}(x^n) - \frac{1}{n} \log P_{Y^n}(y^n) - H(X) - H(Y) \\ &= -\frac{1}{n} \log P_{Y^n}(y^n) - H(Y). \end{aligned}$$

Hence, for this distribution, any pair where the second element is typical for Y is also jointly typical.

More concretely, consider $n = 3$ and $\varepsilon = 0.35$. For these values, the sequences 011, 101, and 110 are typical for Y . This results in the following jointly typical set (again, bigger squares represent bigger probabilities, and orange squares indicate elements of the typical set):

		Y							
		000	001	010	100	011	101	110	111
X	000	•	•	•	•	•	•	•	•
	001	■	■	■	■	■	■	■	■
	010	■	■	■	■	■	■	■	■
	100	■	■	■	■	■	■	■	■
	011	■	■	■	■	■	■	■	■
	101	■	■	■	■	■	■	■	■
	110	■	■	■	■	■	■	■	■
	111	■	■	■	■	■	■	■	■

Note that $P_Y(111) \approx 0.422$, which is higher than $P_Y(011) \approx 0.141$. However, 111 is not typical for Y while 011 is. 111's probability deviates too far from the 'benchmark' set by $H(Y) \approx 0.811$. Together, the sequences in the typical set have a probability of about 0.422 of occurring. As n increases, this number will grow, and the total probability of the elements that fall outside of the typical set (such as the all-1 string) will diminish.

In the above two examples, we have gathered some intuition about jointly typical sets. A (jointly) typical set is usually relatively small compared to the total sample space, but for large n , much of the probability lies on its elements. Comparing the above two examples, we also see the effect of mutual information between X and Y : if X and Y are highly correlated, the typical set has a diagonal structure, while if they are independent, the typical set is a rectangle. In the latter case, this means that independently selecting a typical X sample and a typical Y sample will certainly yield a jointly typical element (see Example 7.2.3). In the first case, there is a significant probability that sampling a typical X and a typical Y independently yields a pair that is not jointly typical (see Example 7.2.2).

These observations lead us to a joint version of Theorem 2.6.4, describing some properties of jointly typical sets:

Theorem 7.2.4 — Joint Asymptotic Equipartition Property (Joint AEP). Let $X^n Y^n$ be i.i.d. with respect to P_{XY} . Then:

1. $P_{X^n Y^n}(A_\varepsilon^{(n)}) \xrightarrow{n \rightarrow \infty} 1$.
2. $|A_\varepsilon^{(n)}| \leq 2^{n(H(XY) + \varepsilon)}$

3. For random variables \tilde{X} and \tilde{Y} , if $\tilde{X}^n \tilde{Y}^n$ is i.i.d. with respect to $P_{\tilde{X}} \cdot P_{\tilde{Y}}$, then

$$P[(\tilde{X}^n \tilde{Y}^n) \in A_\varepsilon^{(n)}] \leq 2^{-n(I(\tilde{X}; \tilde{Y}) - 3\varepsilon)}.$$

4. For random variables \tilde{X} and \tilde{Y} , if $\tilde{X}^n \tilde{Y}^n$ is i.i.d. with respect to $P_{\tilde{X}} \cdot P_{\tilde{Y}}$, then for large enough n ,

$$P[(\tilde{X}^n \tilde{Y}^n) \in A_\varepsilon^{(n)}] \geq (1 - \varepsilon) 2^{-n(I(\tilde{X}; \tilde{Y}) + 3\varepsilon)}.$$

Proof. We show the statements separately, for an arbitrary $X^n Y^n$ that is i.i.d. with respect to P_{XY} .

1. The first statement follows from the weak law of large numbers. Similarly to the proof of Theorem 2.6.4, we have that

$$\begin{aligned} -\frac{1}{n} \log P_{X^n}(X^n) &= -\frac{1}{n} \sum_{i=1}^n \log P_X(X_i) \\ &\xrightarrow{P} -\mathbb{E}[\log P_X(X)] = H(X). \end{aligned} \quad (7.1)$$

Writing out the definition of convergence in probability, we have that for all $\varepsilon > 0$, there exists an $n_1 \in \mathbb{N}$ such that for all $n > n_1$,

$$\begin{aligned} P\left[\left|-\frac{1}{n} \log P_{X^n}(X^n) - H(X)\right| \geq \varepsilon\right] &\leq P\left[\left|-\frac{1}{n} \log P_{X^n}(X^n) - H(X)\right| \geq \frac{\varepsilon}{3}\right] \\ &< \frac{\varepsilon}{3}. \end{aligned} \quad (7.2)$$

In the above, we can divide ε by 3 because the second inequality holds for *all* $\varepsilon > 0$, so in particular also for $\frac{\varepsilon}{3}$.

We can derive statements about Y^n and $X^n Y^n$ that are similar to (7.2). Doing so, we find $n_2, n_3 \in \mathbb{N}$ for Y^n and $X^n Y^n$, respectively. For all $\varepsilon > 0$, we can now choose $n_0 := \max\{n_1, n_2, n_3\}$, and find that for all $n > n_0$,

$$\begin{aligned} P[(X^n, Y^n) \in A_\varepsilon^{(n)}] &= 1 - P[(X^n, Y^n) \notin A_\varepsilon^{(n)}] \\ &\geq 1 - \left(P\left[\left|-\frac{1}{n} \log P_{X^n}(X^n) - H(X)\right| \geq \varepsilon\right] + \right. \\ &\quad \left. P\left[\left|-\frac{1}{n} \log P_{Y^n}(Y^n) - H(Y)\right| \geq \varepsilon\right] + \right. \\ &\quad \left. P\left[\left|-\frac{1}{n} \log P_{X^n Y^n}(X^n Y^n) - H(XY)\right| \geq \varepsilon\right] \right) \\ &\leq 1 - \left(\frac{\varepsilon}{3} + \frac{\varepsilon}{3} + \frac{\varepsilon}{3} \right) = 1 - \varepsilon. \end{aligned} \quad (7.3)$$

The first inequality is due to the union bound, and the second follows from Equation (7.2) and its analogues. By definition of convergence, the first statement is proven.

2. Observe that by rewriting the last line in the definition of $A_\varepsilon^{(n)}$ in Definition 7.2.1, one can derive that for every element $(x^n, y^n) \in A_\varepsilon^{(n)}$,

$$P_{X^n Y^n}(x^n, y^n) > 2^{-n(H(XY) + \varepsilon)}. \quad (7.4)$$

The second statement follows by rearranging the following inequality:

$$1 = \sum_{(x^n, y^n) \in \mathcal{X}^n \times \mathcal{Y}^n} P_{X^n Y^n}(x^n, y^n) \geq \sum_{(x^n, y^n) \in A_\varepsilon^{(n)}} P_{X^n Y^n}(x^n, y^n) \geq |A_\varepsilon^{(n)}| \cdot 2^{-n(H(XY) + \varepsilon)}. \quad (7.5)$$

3. First, we write out the probability of sampling an element in the typical set, again by rearranging the inequalities in Definition 7.2.1.

$$\begin{aligned}
 P[(\tilde{X}^n \tilde{Y}^n) \in A_\varepsilon^{(n)}] &= \sum_{(x^n, y^n) \in A_\varepsilon^{(n)}} P_{\tilde{X}^n \tilde{Y}^n}(x^n, y^n) \\
 &= \sum_{(x^n, y^n) \in A_\varepsilon^{(n)}} P_{X^n}(x^n) \cdot P_{Y^n}(y^n) \\
 &\leq |A_\varepsilon^{(n)}| \cdot 2^{-n(H(X)-\varepsilon)} \cdot 2^{-n(H(Y)-\varepsilon)}.
 \end{aligned} \tag{7.6}$$

Combining this with part (2) of this theorem, we get

$$\begin{aligned}
 P[(\tilde{X}^n \tilde{Y}^n) \in A_\varepsilon^{(n)}] &\leq 2^{n(H(XY)+\varepsilon)} \cdot 2^{-n(H(X)-\varepsilon)} \cdot 2^{-n(H(Y)-\varepsilon)} \\
 &= 2^{-n(-H(XY)+H(X)+H(Y)-3\varepsilon)} \\
 &= 2^{-n(I(X;Y)-3\varepsilon)}.
 \end{aligned} \tag{7.7}$$

This completes the proof of the third part.

4. By part (1), we have that $P[A_\varepsilon^{(n)}] \geq 1 - \varepsilon$ for large enough n . Thus,

$$\begin{aligned}
 1 - \varepsilon &\leq P[A_\varepsilon^{(n)}] \\
 &= \sum_{(x^n, y^n) \in A_\varepsilon^{(n)}} P_{X^n Y^n}(x^n, y^n) \\
 &\leq |A_\varepsilon^{(n)}| 2^{-n(H(XY)-\varepsilon)}.
 \end{aligned} \tag{7.8}$$

Rearranging this inequality and using the same type of derivation as in the proof of (3), it follows that

$$\begin{aligned}
 P[(\tilde{X}^n \tilde{Y}^n) \in A_\varepsilon^{(n)}] &= \sum_{(x^n, y^n) \in A_\varepsilon^{(n)}} P_{X^n}(x^n) \cdot P_{Y^n}(y^n) \\
 &\geq (1 - \varepsilon) \cdot 2^{n(H(XY)-\varepsilon)} \cdot 2^{-n(H(X)+\varepsilon)} \cdot 2^{-n(H(Y)+\varepsilon)} \\
 &= (1 - \varepsilon) \cdot 2^{-n(I(X;Y)+3\varepsilon)}.
 \end{aligned} \tag{7.9}$$

■

7.3 Shannon's Noisy-Channel Coding Theorem

With the tools from the previous section, we are ready to prove the forward direction of Shannon's noisy-channel coding theorem, which states that any rate strictly below the channel capacity is achievable:

Theorem 7.3.1 — Shannon's noisy-channel coding theorem (forward direction). For a discrete memoryless channel with capacity C , any rate $R < C$ is achievable.

Proof. Given a channel $(\mathcal{X}, P_{Y|X}, \mathcal{Y})$ with capacity $C = \max P_X I(X;Y)$, let $R < C$ and $\varepsilon > 0$. We will first show that for big enough n , a *randomly constructed* code with rate R has a low error probability. We will then argue that a low average error probability implies the existence of some *specific* code with low error probability.

Fix an input distribution P_X that maximizes $I(X;Y)$. For any n , construct a $(2^{nR}, n)$ -code \mathcal{C} [Chris: what if 2^{nR} is not an integer?] by choosing a codebook at random according to P_X . That

is, for every message $w \in [1, 2^{n \cdot R}]$, sample n times from the distribution P_X , creating a codeword $\mathcal{C}(w) = (\mathcal{C}_1(w), \mathcal{C}_2(w), \dots, \mathcal{C}_n(w))$ by concatenating the n independent samples $\mathcal{C}_i(w)$.

Since the channel is memoryless, if w is sent over the channel using \mathcal{C} , the output distribution Y^n is given by:

$$P_{Y^n|X^n}(y^n|\mathcal{C}(w)) = \prod_{i=1}^n P_{Y|X}(y_i | \mathcal{C}_i(w)). \quad (7.10)$$

What is the probability that the decoded message \hat{w} is incorrect, i.e., not equal to w ? This depends on the decoding method used by the receiver. The optimal decoding procedure is **maximum-likelihood decoding**, where the input message that is most likely with respect to $P_{X|Y}$ is selected as the decoding \hat{w} . However, it is hard to analyze the error probability for this decoding method. Instead, we will assume that the receiver applies **jointly typical decoding**, which has a slightly higher probability of decoding to the wrong message, but still small enough for our analysis. Jointly typical decoding works as follows: upon receiving an output y^n , the receiver looks for a *unique* message \hat{w} such that the pair $(\mathcal{C}(\hat{w}), y^n)$ is jointly typical. If there exists no such message, or if it is not unique, the receiver declares a failure by decoding to $\hat{w} = 0$ (which is always wrong because $w \in [1, 2^{n \cdot R}]$).

With this decoding procedure in mind, we will analyze the average error probability $P[\text{error}]$, where the average is taken over both the randomly constructed code \mathcal{C} and the uniformly randomly selected message w . Defining $\lambda_w(\mathcal{C}) := P[\hat{w} \neq w | \mathcal{C}(w) \text{ was sent over the channel}]$ to be the probability that a message w (encoded using \mathcal{C}) is decoded incorrectly, we get:

$$\begin{aligned} P[\text{error}] &= \sum_{\mathcal{C}} P[\mathcal{C}] \cdot \left(\sum_{w=1}^{2^{n \cdot R}} \frac{1}{2^{n \cdot R}} \cdot \lambda_w(\mathcal{C}) \right) \\ &= \frac{1}{2^{n \cdot R}} \sum_{w=1}^{2^{n \cdot R}} \sum_{\mathcal{C}} P[\mathcal{C}] \cdot \lambda_w(\mathcal{C}). \end{aligned} \quad (7.11)$$

Since we average over all randomly constructed codes \mathcal{C} , and the codewords for all messages are sampled independently, the value $\sum_{\mathcal{C}} P[\mathcal{C}] \cdot \lambda_w(\mathcal{C})$ does not depend on the particular message w . Hence if we set, for example, $w_0 = 1$, then for all $w \in [0, 2^{n \cdot R}]$,

$$\sum_{\mathcal{C}} P[\mathcal{C}] \lambda_w(\mathcal{C}) = \sum_{\mathcal{C}} P[\mathcal{C}] \lambda_{w_0}(\mathcal{C}). \quad (7.12)$$

This simplifies the calculation of $P[\text{error}]$ significantly:

$$\begin{aligned} P[\text{error}] &= \frac{1}{2^{n \cdot R}} \sum_{w=1}^{2^{n \cdot R}} \sum_{\mathcal{C}} P[\mathcal{C}] \cdot \lambda_{w_0}(\mathcal{C}) \\ &= \sum_{\mathcal{C}} P[\mathcal{C}] \cdot \lambda_{w_0}(\mathcal{C}). \end{aligned} \quad (7.13)$$

That is, the average probability of error is the probability (over the selection of the code \mathcal{C} , and over the randomness in the channel) that the message w_0 is decoded incorrectly. There are two possible reasons for an error in the decoding:

1. The output of the channel is not jointly typical with $\mathcal{C}(w_0)$. By the first item of the joint AEP (Theorem 7.2.4), this probability approaches zero as n goes to infinity. Hence, for big enough n , the probability of an error for this reason is smaller than ϵ .
2. There is some $w' \neq w_0$ such that the output of the channel is (also) jointly typical with $\mathcal{C}(w')$. Since \mathcal{C} is a random code (and so $\mathcal{C}(w')$ is independent from the channel output y^n), by the third item of the joint AEP the probability that this occurs is at most

$$\sum_{w' \neq w_0} 2^{-n(I(X;Y)-3\epsilon)} = (2^{n \cdot R} - 1) 2^{-n(I(X;Y)-3\epsilon)}. \quad (7.14)$$

We can thus bound the average probability of error, using the union bound and the bounds in the above analysis, by

$$\begin{aligned}
 P[\text{error}] &\leq \varepsilon + (2^{nR} - 1)2^{-n(I(X;Y) - 3\varepsilon)} \\
 &\leq \varepsilon + 2^{nR}2^{-n(I(X;Y) - 3\varepsilon)} \\
 &= \varepsilon + 2^{-n(I(X;Y) - R - 3\varepsilon)}
 \end{aligned} \tag{7.15}$$

As long as $R < I(X;Y)$, one can choose n large enough so that $P[\text{error}] \leq 2\varepsilon$.

This analysis upper bounds the (expected) average error probability for a random code \mathcal{C} . However, if this expected probability is low, there must be some specific code \mathcal{C}^* that also has low average error probability.

Finally, in \mathcal{C}^* , we aim to bound the *maximal* error probability, i.e., the probability of error for the worst message. We can do so by noting that at least half of the messages w has error probability $\lambda_w(\mathcal{C}^*) \leq 4\varepsilon$: if not, then the total error probability of these messages would already exceed $2^{nR} \cdot 2\varepsilon$, contradicting the upper bound of 2ε to the average error probability. Thus, we can construct a better code by discarding the worst half of the codewords, and using the remaining 2^{nR-1} codewords to construct a new code, with rate

$$\frac{\log(2^{nR-1})}{n} = \frac{n \cdot R - 1}{n} = R - \frac{1}{n} \tag{7.16}$$

and maximal probability of error $\lambda^{(n)} \leq 4\varepsilon$. ■

7.4 Shannon's Noisy-Channel Coding Theorem: Converse

In the previous section, we showed that any rate strictly below the channel capacity is achievable. Here, we show that one cannot do better: rates strictly above the channel capacity are not achievable. Specifically, codes with such rates suffer from non-negligible error probabilities.

Theorem 7.4.1 — Shannon's noisy-channel coding theorem (converse). On a discrete memoryless channel with capacity C , any code with rate $R > C$ has average probability of error $p_e^{(n)} \geq 1 - \frac{C}{R} - \frac{1}{nR}$.

Proof. For a code with rate $R > C$, let W be uniformly distributed over all possible messages, let X^n describe the encoding of the message (and the input to the channel), let Y^n describe the output of the channel, and \hat{W} the decoding of that output. The average probability of error, $p_e^{(n)}$, is equal to $P[W \neq \hat{W}]$, the probability that the original message differs from the decoded message. Note that $W \rightarrow X^n \rightarrow Y^n \rightarrow \hat{W}$ forms a Markov chain.

As a first step, we show that the mutual information between the message W and the channel output Y^n is upper bounded by $n \cdot C$, that is, there is a limit to the amount of information that can get through the channel. To see this, first observe that

$$\begin{aligned}
 H(Y^n W) &= H(Y^n W) + H(X^n | Y^n W) && \text{(since } W \rightarrow X^n) \\
 &= H(X^n Y^n W) && \text{(chain rule)} \\
 &= H(X^{n-1} Y^{n-1} W) + H(Y_n | X^n Y^{n-1} W) + H(X_n | X^{n-1} Y^{n-1} W) && \text{(chain rule)} \\
 &= H(X^{n-1} Y^{n-1} W) + H(Y_n | X^n Y^{n-1} W) && \text{(since } W \rightarrow X^n) \\
 &= H(X^{n-1} Y^{n-1} W) + H(Y_n | X_n) && \text{(memoryless)} \\
 &= \dots && \text{(repeat)}
 \end{aligned}$$

$$= H(W) + \sum_{i=1}^n H(Y_i | X_i). \quad (7.17)$$

Therefore,

$$\begin{aligned}
 I(W; Y^n) &= H(W) + H(Y^n) - H(Y^n W) && \text{(entropy diagram)} \\
 &= H(Y^n) - \sum_{i=1}^n H(Y_i | X_i) && \text{(by (7.17))} \\
 &\leq \sum_{i=1}^n H(Y_i) - \sum_{i=1}^n H(Y_i | X_i) \\
 &= \sum_{i=1}^n I(X_i; Y_i) \\
 &\leq n \cdot C.
 \end{aligned} \quad (7.18)$$

Now that we have established that $I(W; Y^n)$ is upper bounded by $n \cdot C$, we can show that the code with rate R induces a considerable error probability:

$$\begin{aligned}
 R &= \frac{\log |\mathcal{W}|}{n} \\
 &= \frac{1}{n} H(W) \\
 &= \frac{1}{n} (H(W | Y^n) + I(W; Y^n)) \\
 &\leq \frac{1}{n} (H(W | Y^n) + n \cdot C) \\
 &\leq \frac{1}{n} (1 + P[W \neq \hat{W}] \cdot n \log |\mathcal{W}| + n \cdot C) \\
 &= \frac{1}{n} + P[W \neq \hat{W}] \cdot R + C,
 \end{aligned} \quad (7.19)$$

where the second inequality is an application of Fano's inequality [Chris: reference to Chapter ? to be inserted]. Dividing both sides by R and rearranging, we get the desired inequality:

$$p_e^{(n)} = P[W \neq \hat{W}] \geq 1 - \frac{C}{R} - \frac{1}{nR}. \quad (7.20)$$

■

This theorem shows that if Alice and Bob try to communicate using a code with a rate $R > C$, their probability of error will be bounded away from zero by a constant factor of $1 - \frac{C}{R}$ (for big n , the last term in the inequality becomes insignificant). This error probability worsens for a bigger difference between R and C .