

University of Cambridge

Part II of the Mathematical Tripos

Coding and Cryptography

Lectured by Rachel Camina, Lent 2024–25

Notes by Avish Kumar

`ak2461@cam.ac.uk`

<https://ak1089.github.io/maths/notes>

Version 1.11

These notes are unofficial and may contain errors. While they are written and published with permission, they are not endorsed by the lecturer or University.

Contents

1	Noiseless Coding	3
1.1	The Coding Problem	3
1.2	Communication Channels	3
1.3	Strings and Alphabets	5
1.4	Mathematical Entropy	7
1.5	Optimal Noiseless Coding	9
1.6	Coding Sequences	12
2	Error Control Codes	14
2.1	Binary Codes	14
2.2	Bounds on Codes	18
2.3	Operational Channel Capacity	21
2.4	Informational Channel Capacity	25
2.5	Shannon's Noisy Coding Theorem	28
2.6	The Kelly Criterion	29
3	Linear Codes	31
3.1	Codes and Vector Spaces	31

1 Noiseless Coding

1.1 The Coding Problem

The general problem of coding is that of transmitting a message across a communication channel. For example, if we wish to send an email containing a message $m = \text{"Call me!"}$, we may encode this as a sequence of binary strings using the standard ASCII format.

Under this code, $f(\text{"C"}) = 1000011$. In fact, each character is mapped to seven binary digits, or *bits*. The entire message is the concatenation of these strings: "Call me!" becomes:

$$f^*(\text{"Call me!"}) = 100001111100001110110011011000100000110110111001010100001.$$

Definition 1.1 (Source, Encoder, Channel, Receiver, Decoder)

More generally, we have a *source*, often called Alice. She uses an *encoder* to convert plaintext messages into encoded messages. These encoded messages are sent through a *channel*: this channel may be *noisy*, and introduce errors into the code. The encoded message is received by a *receiver* Bob, who uses a *decoder* to convert it back into the original plaintext.

Given a source and a channel, described probabilistically, we want to design an encoder and decoder in order to transmit source information across the channel. We might want certain properties:

1. Economy: we would like to minimise the amount of unnecessary information sent: the code should not be too long, as it wastes time and money.
2. Reliability: the decoder should be able to successfully decipher the plaintext with very high probability, or mistakes should be detectable.
3. Privacy: we may want only someone with the decoder to be able to read the message.

Accomplishing this last desideratum is the aim of *cryptography* in particular, while *coding* deals with the first two. How might we achieve these?

Remark 1.2 (Economy and Reliability)

Morse code is economic in that it attempts to minimise message length. This is done by giving the shorter codes to letters which are used more frequently. For example, $E = \cdot$, while $Q = - - - \cdot - -$.

The ISBN system for numbering books is reliable. Each book has a unique ten-digit ISBN: the first nine digits encode information about the book (its publisher, ID, and region) while the last digit is a *check digit* chosen such that $10a_1 + 9a_2 + \dots + 2a_9 + a_{10} \equiv 0 \pmod{11}$.

This has robust error-detection capabilities. For example, a transposition of any two digits a_i and a_j will add $(j - i)(a_j - a_i) \not\equiv 0 \pmod{11}$ to the left hand side. This means a single error can be detected, since the congruence will be broken.

1.2 Communication Channels

Definition 1.3 (Channel)

A *communication channel* takes letters from an input alphabet $\Sigma_1 = \{a_1, \dots, a_r\}$ and emits letters from an output alphabet $\Sigma_2 = \{b_1, \dots, b_s\}$. It is determined by the probabilities

$$\mathbb{P}[y_1 \dots y_k \text{ emitted} \mid x_1 \dots x_k \text{ input}] \text{ where } x_i \in \Sigma_1^*, y_i \in \Sigma_2^*.$$

Note: The important feature of a channel is that it is not necessarily perfect! Much like in the real world, where we deal with problems like TV static or data corruption, the channels we will study are affected by noise.

Definition 1.4 (Discrete Memoryless Channel)

A *discrete memoryless channel* over a finite alphabet is a channel for which the probabilities

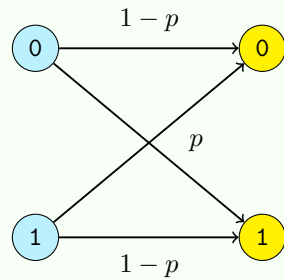
$$P_{ij} = \mathbb{P}[b_j \text{ received} \mid a_i \text{ sent}]$$

are the same every time the channel is used, independent of past and future channel use. This is the *memoryless property*, while the discrete nature is given by the alphabets.

We often identify the channel with its *channel matrix* P , which is the $r \times s$ matrix with entries p_{ij} equal to those probabilities. Note that the rows of P , but not necessarily its columns, sum to 1, and all entries are non-negative: we thus say that P is a *stochastic matrix*.

Example 1.5 (Binary Symmetric/Erasure Channel)

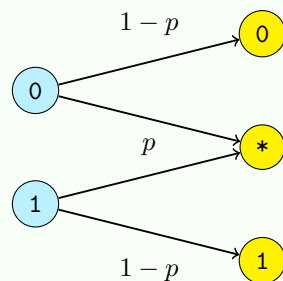
For example, a Binary Symmetric Channel with probability $0 \leq p \leq 1$ of error is a DMC over the binary alphabet $\Sigma_1 = \Sigma_2 = \{0, 1\}$. In particular, any bit sent has a probability p of being flipped by the channel due to noise. This can be seen in the below diagram.



$$P = \begin{pmatrix} 1-p & p \\ p & 1-p \end{pmatrix}$$

Usually, we assume $p < 0.5$. If $p > 0.5$, then we can just pre-flip the bits sent to reduce the error probability (since any bit is likely flipped back by the channel). If $p = 0.5$, then every single bit received is equally likely to be 0 and 1, independently of what was actually transmitted, so the channel is entirely useless (pure noise).

A Binary Erasure Channel is similar, taking $\Sigma_1 = \{0, 1\}$ but $\Sigma_2 = \{0, 1, *\}$, with the $*$ understood to be an *erasure*. Each bit transmitted has a probability $0 \leq p \leq 1$ of being erased, making it unreadable, and is transmitted correctly otherwise (never flipped), giving:



$$P = \begin{pmatrix} 1-p & p & 0 \\ 0 & p & 1-p \end{pmatrix}$$

where the columns correspond to 0, *, and then 1.

Definition 1.6 (Capacity)

The *capacity* of a channel is the highest rate at which information can be reliably transmitted over the channel. Here, the rate is measured as units of information per unit time: for a binary channel, this might be the number of decoded bits per transmitted bit. High reliability is achieved by an arbitrarily low probability of error.

1.3 Strings and Alphabets

We frequently work with alphabets, which are simply sets of elements called letters or characters. These are the building blocks of a language: the input alphabet of a code is the set of atoms which are encoded into something else.

Definition 1.7 (String, Concatenation, Length)

For an alphabet Σ , we define the set of Σ -strings to be $\Sigma^* = \bigcup_{n \geq 0} \Sigma^n$. These are usually written as concatenations rather than tuples, so that the set of binary alphabet strings is

$$\Sigma_{01}^* = \{\varepsilon, 0, 1, 00, 01, 10, 11, 000, 001, 010, 011, 100, 101, 110, 111, 0000, 0001, \dots\}$$

The *length* of a string is the number of letters contained. Here ε is the empty string with length 0. If $x = x_1x_2 \dots x_r$ and $y = y_1y_2 \dots y_s$ are Σ -strings, their *concatenation* is given by $xy = x_1 \dots x_r y_1 \dots y_s$.

For two alphabets Σ_1 and Σ_2 , a *code* is a function $f : \Sigma_1 \rightarrow \Sigma_2^*$. The strings $\{f(x) : x \in \Sigma_1\}$, or the image of f , are called codewords.

Example 1.8 (Polybius Square)

For example, the Polybius Square is a cipher developed by Ancient Greek polymath Polybius, who created a way to encode Greek as numbers.

The input alphabet Σ_1 was the 24 Greek letters α to ω , while the output alphabet Σ_2 was the set $\{1, 2, 3, 4, 5\}$. Each letter was mapped to precisely two digits from 1 to 5 for easy transmission (using every pair except 55). This made the codewords the set

$$\{1 \dots 5\}^2 = \{11, 12, 13, 14, 15, 21, 22, \dots, 52, 53, 54\}.$$

Note: For English-language codes, we do not necessarily have Σ_1 being $\{a \dots z\}$ the set of letters. The domain of the code function is the set of atoms of the code, which is often pairs of letters, or even more commonly entire words.

We apply a code by encoding $x_1x_2 \dots x_n \in \Sigma_1^*$ as $f(x_1)f(x_2) \dots f(x_n) \in \Sigma_2^*$. This extends f the code function from atoms to entire words in the input language, which we call $f^* : \Sigma_1^* \rightarrow \Sigma_2^*$.

However, not every function $f : \Sigma_1 \rightarrow \Sigma_2^*$ works as a code.

Definition 1.9 (Decipherable)

A code f is *decipherable* if f^* is injective, so that every string in Σ_2^* arises from at most one message. Without this condition, the output of encoding might have come from multiple possible inputs, and we would have no way of knowing which when decoding it.

Proposition 1.10 (Decipherability requires injectivity)

A decipherable code f requires f injective. However, this is not a sufficient condition.

Proof: Firstly, if f is not injective, then $f(x) = f(y)$ where $x \neq y \in \Sigma_1$. But then the encoding of x and y when treated as members of Σ_1^* is the same, violating injectivity of f^* .

However, this is not a sufficient condition. Suppose $\Sigma_1 = \{1, 2, 3, 4\}$ and $\Sigma_2 = \{0, 1\}$. Define

$$f(1) = 0 \quad f(2) = 1 \quad f(3) = 00 \quad f(4) = 01$$

so that f is injective, but $f^*(1112) = 0001 = f^*(34)$, so f^* is not. \square

How do we construct decipherable codes? There are a few basic properties of codes which guarantee decipherability (none of these are necessary, but are all sufficient) provided that f is injective.

1. A *block code* is a code where all codewords are of the same length. For example, the Polybius cipher had all codewords of length 2, and so it can be decoded by considering the output as a list of length-2 strings.
2. A *comma code* reserves one letter in Σ_2 to act as the comma, which appears at the end of every output of f and nowhere else. It thus delimits words in Σ_2^* , so we know where each letter in the original input to the code was mapped.
3. A *prefix-free* (or *instantaneous*) code is a code where no codeword is a prefix of any other codeword: for any $x, y \in \Sigma_1$, we have $f(x) \neq f(y)\alpha$ for any $\alpha \neq \varepsilon \in \Sigma_2^*$.

Note: In fact, block codes and comma codes are special cases of prefix-free codes.

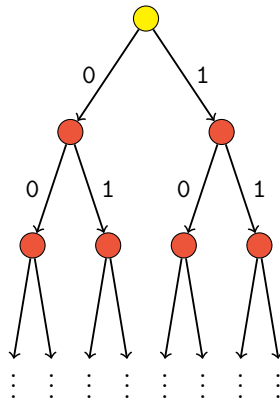
Theorem 1.11 (Kraft's Inequality)

Let $\Sigma_1 = \{x_1, \dots, x_m\}$ and $|\Sigma_2| = a$. Then a prefix-free code $f : \Sigma_1 \rightarrow \Sigma_2^*$ with word lengths s_1, \dots, s_m (where $|f(x_i)| = s_i$) exists if and only if

$$\sum_{i=1}^m a^{-s_i} \leq 1$$

Proof: (\Rightarrow) Consider an infinite tree where each node has a descendants corresponding to the letters of Σ_2 . Then each codeword corresponds to precisely one of these nodes, where the path to the node spells out the codeword along the branches taken.

Assuming f is prefix-free, no codeword is the ancestor of any other. View the tree as a network, with water being pumped in at the root, where each node divides the flow equally between each descendant. The total amount of water extracted at the codewords is therefore the sum of a^{-s_i} , which is at most the total amount of water pumped in, demonstrating the inequality.



(\Leftarrow) Conversely, we can construct a prefix-free code with word lengths $s_1 < s_2 < \dots < s_m$. Choose codewords sequentially, ensuring that any previous codewords are not prefixes. Suppose that the r^{th} codeword has no valid code available. Then constructing the tree above gives

$$\sum_{i=1}^{r-1} a^{-s_i} = 1 \implies \sum_{i=1}^m a^{-s_i} > 1$$

which contradicts our assumption. \square

Theorem 1.12 (McMillan / Karush)

Every decipherable code satisfies Kraft's inequality.

Proof: Let $f : \Sigma_1 \rightarrow \Sigma_2^*$ be a decipherable code with word lengths $s_1 \dots s_m$, where $s = \max s_i$. For any $r \in \mathbb{N}$, we must have

$$\left(\sum_{i=1}^m a^{-s_i} \right)^r = \left(\sum_{\ell=1}^{rs} b_\ell a^{-\ell} \right)$$

where b_ℓ is the number of ways of choosing r codewords with total length ℓ . Since f is decipherable, we know that $b_\ell \leq |\Sigma_2|^\ell = a^\ell$, since no string can be the encoding of more than one set of codewords. This means that we can write

$$\left(\sum_{i=1}^m a^{-s_i} \right)^r \leq \left(\sum_{\ell=1}^{rs} a^\ell a^{-\ell} \right) = rs \implies \left(\sum_{i=1}^m a^{-s_i} \right) \leq (rs)^{1/r}$$

But this is true for any r , and $(rs)^{1/r} \rightarrow 1$ as $r \rightarrow \infty$. Therefore the left hand side of the inequality is at most 1, which is exactly the statement of Kraft's inequality. \square

As a result, we mostly restrict our attention to prefix-free codes.

1.4 Mathematical Entropy

Entropy is a measure of “randomness” or “uncertainty”. Suppose a random variable X takes values $x_1 \dots x_n$ with probabilities $p_1 \dots p_n$, where we have $0 \leq p_i \leq 1$ for all i and $\sum p_i = 1$. Loosely, the *entropy* $H(X)$ is the expected number of yes/no questions required to determine the value of X .

This is not a formal definition yet: we consider some examples first.

Example 1.13 (Basic Entropy Examples)

Let's consider X taking values $x_1 \dots x_4$.

If $p_1 = p_2 = p_3 = p_4 = 1/4$, then asking precisely two yes/no questions can consistently determine the value of X . For example, the two questions could be “is $X \in \{x_1, x_2\}$?” and “is $X \in \{x_1, x_3\}$?”. Here, this means $H(X) = 2$ directly.

Now, suppose $(p_1, p_2, p_3, p_4) = (1/2, 1/4, 1/8, 1/8)$. Then we could ask the question “ $X = x_1$?” and finish with one question half the time. If the answer is no, we ask “ $X = x_2$?” and again be done half the time (so a quarter overall). Failing that, we ask “ $X = x_3$?” and know the value of X with certainty after three questions.

This gives $H(X) = 1 \times 1/2 + 2 \times 1/4 + 3 \times 1/8 + 3 \times 1/8 = 7/4 < 2$, and so the first random variable is “more random”. This aligns with our intuition.

This gives us enough to write down our formal definition.

Definition 1.14 (Entropy)

For a random variable X taking values $x_1 \dots x_n$ with probabilities $p_1 \dots p_n$, where we have $0 \leq p_i \leq 1$ for all i and $\sum p_i = 1$, the *entropy* $H(X)$ is defined to be

$$H(X) = H(p_1, \dots, p_n) = - \sum_{i=1}^n p_i \log p_i.$$

Note: In this course, we always consider the logarithm to be defined as \log_2 (the logarithm with base 2), rather than the natural logarithm with base e , due to our focus on binary.

Note: This definition breaks down if $p_i = 0$ for some i : in this case, we take $p_i \log p_i = 0$ as a convention, since we could have excluded p_i and x_i for an equivalent distribution.

Corollary: As $p_i \log p_i \leq 0$ for $0 \leq p_i \leq 1$, the entropy $H(X) \geq 0$.

Example 1.15 (Entropy of a Biased Coin)

Toss a biased coin which lands heads with probability p and tails with probability $1 - p$. Then

$$h(p) = H(p, 1 - p) = -p \log p - (1 - p) \log(1 - p)$$

Plotting this, we get an arch-shaped curve with $h(0) = h(1) = 0$, since the outcome is certain and thus there is no randomness. The graph is symmetric, which makes sense, and we get a peak at $p = 1/2$, which is the case for a fair coin (which is therefore maximal entropy).

We now prove a result which will come up frequently in the study of entropy.

Theorem 1.16 (Gibbs' Inequality)

Let $\mathbf{p} = (p_1 \dots p_n)$ and $\mathbf{q} = (q_1 \dots q_n)$ be probability distributions. Then

$$- \sum_{i=1}^n p_i \log p_i \leq - \sum_{i=1}^n p_i \log q_i.$$

with equality if and only if $\mathbf{p} = \mathbf{q}$.

Proof: Since $\log x = \ln(x)/\ln(2)$, we may prove the equality using \ln in place of \log , and dividing through both sides afterwards. Note that $\ln x \leq x - 1$ with equality if and only if $x = 1$.

Let $I = \{1 \leq i \leq n : p_i > 0\}$ be the set of nontrivial indices. Then

$$\ln(q_i/p_i) \leq q_i/p_i - 1 \quad \forall i \in I.$$

and therefore we have

$$\sum_{i \in I} p_i \ln(q_i/p_i) \leq \sum_{i \in I} q_i - \sum_{i \in I} p_i = \sum_{i \in I} q_i - 1 \leq 0$$

Rearranging this inequality yields

$$- \sum_{i \in I} p_i \ln p_i \leq - \sum_{i=1}^n p_i \ln p_i \leq - \sum_{i=1}^n p_i \ln q_i$$

as required. Equality is only possible if we had equality in the first line, with all $i \in I$ satisfying $\ln(q_i/p_i) = q_i/p_i - 1 \implies q_i/p_i = 1 \implies q_i = p_i$ as desired, and thus the proof holds. \square

Corollary: $H(p_1 \dots p_n) \leq \log n$ with equality if and only if $p_i = 1/n$ for all i .

1.5 Optimal Noiseless Coding

Recall that the coding problem considers alphabets Σ_1 and Σ_2 of sizes $m, a \geq 2$ respectively. When considering a channel, we model the source as a sequence of random variables X_1, X_2, \dots which take values in Σ_1 .

Definition 1.17 (Memoryless Source)

A *Bernoulli* (or *memoryless*) source is a sequence X_1, X_2, \dots of independently and identically distributed random variables.

Definition 1.18 (Expected Word Length, Optimal Code)

Consider a memoryless source. Let $\Sigma_1 = \{\mu_1 \dots \mu_m\}$ and define $p_i = \mathbb{P}[X_1 = \mu_i]$. The expected word length S of a code $f : \Sigma_1 \rightarrow \Sigma_2^*$ with word length $s_1 \dots s_m$ is therefore

$$\mathbb{E}[S] = \sum_{i=1}^m p_i s_i$$

The code f is then said to be *optimal* if it has the shortest possible expected word length among decipherable codes: that is, if it minimises $\mathbb{E}[S]$.

This brings us to one of the most important results in information theory.

Theorem 1.19 (Shannon's Noiseless Coding Theorem)

The expected word length of an optimal decipherable code $f : \Sigma_1 \rightarrow \Sigma_2^*$ satisfies

$$\frac{H(X)}{\log(a)} \leq \mathbb{E}[S] < \frac{H(X)}{\log(a)} + 1.$$

This theorem was proved in 1948 by Claude Shannon, the father of information theory. It is also known by several other names, like *Shannon's Source Coding Theorem for Symbol Codes*.

Proof: The lower bound is given by combining Gibbs' and Kraft's inequalities (1.16 and 1.11), taking $q_i = a^{-s_i}/c$, where $c = \sum a^{-s_i} \leq 1$ is such that $\sum q_i = 1$. Then

$$\begin{aligned} H(X) &= - \sum_{i=1}^m p_i \log p_i \\ &\leq - \sum_{i=1}^m p_i \log q_i \quad (\text{by Gibbs'}) \\ &= - \sum_{i=1}^m p_i \log(a^{-s_i}/c) \\ &= \log a \sum_{i=1}^m p_i s_i + \sum_{i=1}^m p_i \log c \\ &= \mathbb{E}[S] \times \log a + \log c \\ &\leq \mathbb{E}[S] \times \log a \end{aligned}$$

where the last line follows by $c \leq 1$ implying $\log c \leq 0$. Dividing through by $\log a$ yields the result. We achieve this lower bound only if $q_i = p_i$, that is if $p_i = a^{-s_i}$ for some integers s_i for all p .

In fact, this lower bound must hold for *all* decipherable codes, by McMillan / Karush (1.12).

For the upper bound, we take $s_i = \lceil -\log_a p_i \rceil$. We have $s_i < -\log_a p_i + 1 \implies a^{-s_i} \leq p_i$. This means we satisfy Kraft's inequality (1.11), since

$$\sum_{i=1}^m a^{-s_i} \leq \sum_{i=1}^m p_i = 1$$

and therefore there is a prefix-free code with word lengths $s_1 \dots s_m$. Also,

$$\mathbb{E}[S] = \sum_{i=1}^m p_i s_i < \sum_{i=1}^m p_i (-\log_a p_i + 1) = \frac{H(X)}{\log a} + 1,$$

so our upper bound holds. \square

Example 1.20 (Shannon-Fano Coding)

This is an example of a code which follows from the above proof. We set $s_i = \lceil -\log_a p_i \rceil$ and construct a prefix-free code with word lengths $s_1 \dots s_m$ sorted in ascending order, ensuring that previous codewords are not used as prefixes.

Suppose our source emits words $\mu_1 \dots \mu_5$ with probabilities 0.4, 0.2, 0.2, 0.1, and 0.1. Then we construct the binary Shannon-Fano code (with $a = 2$) by taking $s_i = \lceil -\log_2 p_i \rceil$, which is equal to 2, 3, 3, 4, and 4 respectively.

We then have a lot of freedom. At each stage, we may choose *anything* which does not contain a previous word as a prefix. For example, set $\mu_1 \mapsto 00$. Then we may choose $\mu_2 \mapsto$ anything of length 3 which does not begin 00, say 010. Similarly, $\mu_3 \mapsto 100$, then $\mu_4 \mapsto 1100$ and lastly $\mu_5 \mapsto 1110$, which is a prefix-free and thus decipherable code.

The expected word length $\mathbb{E}[S] = 2.8$. For comparison, the entropy $H(X) \approx 2.122$.

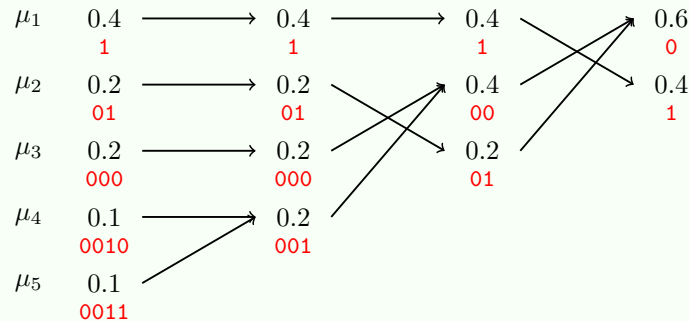
Note: The Shannon-Fano code is not always optimal. However, the next one is!

Example 1.21 (Huffman Coding)

We define Huffman Coding inductively. For simplicity, we take the binary case $a = 2$ again. Order the p_i in descending order $p_1 > \dots > p_m$. Now:

1. If $m = 2$, then assign $s_1 = 0$ and $s_2 = 1$.
2. If $m > 2$, then find the Huffman code f' for $m - 1$ words: $\mu_1 \dots \mu_{m-2}$ emitted with probabilities $p_1 \dots p_{m-2}$, and a new word ν with probability $p_{m-1} + p_m$. Then, assign words $\mu_1 \dots \mu_{m-2}$ the same codes, and set $f(\mu_{m-1}) = f'(\nu)0$ and $f(\mu_m) = f'(\nu)1$.

This gives a prefix-free code. When some of the p_i are equal, then we can choose how to order them, so Huffman codes are not necessarily unique. In our previous example, this gives:



Now, $\mathbb{E}[S] = 2.2$: notably better than the Shannon-Fano code and remarkably close to $H(X)$.

Let us prove an auxiliary lemma first to show the optimality of the Huffman scheme.

Proposition 1.22 (Sorting and Almost-Equality)

Suppose that $\mu_1 \dots \mu_m$ are emitted with probabilities $p_1 \dots p_m$. Let f be an optimal prefix-free code with word lengths $s_1 \dots s_m$. Then

1. If $p_i > p_j$, then $s_i \leq s_j$.
2. There are two codewords of maximal length which are equal up to the last letter.

Proof: (1) is obvious: simply swap the codewords $f(\mu_i)$ and $f(\mu_j)$. This strictly decreases the expected word length, contradicting the optimality of f .

(2) is less obvious. Suppose it is false. Then either there is only one codeword of maximal length, or any two codewords of maximal length differ before the last digit. In either case, remove the last letter of each codeword of maximal length. This maintains the prefix-free condition, and shortens the expected word length, again contradicting the optimality of f . \square

Note: Shannon-Fano satisfies the first of these properties, but not necessarily the second. In our example, $f(\mu_4) = 1100$ and $f(\mu_5) = 1110$ did not satisfy this condition.

Now, we may prove our target theorem.

Theorem 1.23 (Huffman Coding Optimal)

The Huffman coding scheme is optimal (1.18): for words $\mu_1 \dots \mu_m$ emitted with probabilities $p_1 \dots p_m$, it minimises the expected word length $\mathbb{E}[S]$.

Proof: We prove this for the binary $a = 2$ case by showing via induction on m that any Huffman code of size m is optimal. The $m = 2$ case is obvious: the codewords 0 and 1 are minimal.

Suppose $m > 2$. The inductive source emits $\mu_1 \dots \mu_{m-2}$ with probabilities $p_1 \dots p_{m-2}$, and a new word ν with probability $p_{m-1} + p_m$. The Huffman code f_{m-1} is optimal for this source. Now, we construct the Huffman code f_m of size m by extending f_{m-1} . The expected word length satisfies:

$$\mathbb{E}[S_m] = \mathbb{E}[S_{m+1}] + p_{m-1} + p_m$$

Let f'_m be an optimal code for X_m which is prefix-free. Proposition 1.22 then yields that the codewords associated with μ_{m-1} and μ_m are of maximal length and differ only in the last letter. Say these are $y0$ and $y1$ for some string $y \in \Sigma_2^*$. We define a code f'_{m-1} for X_{m-1} with

$$\begin{aligned} f'_{m-1}(\mu_i) &= f'_m(\mu_i) : 1 \leq i \leq m-2, \\ f'_{m-1}(\nu) &= y. \end{aligned}$$

Then f'_{m-1} is a prefix-free code and the expected-word length satisfies

$$\mathbb{E}[S'_m] = \mathbb{E}[S'_{m-1}] + p_{m-1} + p_m$$

By the induction hypothesis, f_{m-1} is optimal, so $\mathbb{E}[S_{m-1}] \leq \mathbb{E}[S'_{m-1}]$. Putting this all together, we obtain $\mathbb{E}[S_m] \leq \mathbb{E}[S'_m]$: that is, f_m has word length less than or equal to that of an optimal code for the source X_m . Therefore f_m must itself be optimal, as required. \square

Note: Not all optimal codes are Huffman, but (from the proof of the above) it can be seen that for any optimal sequence of word lengths $s_1 \dots s_m$ associated with $p_1 \dots p_m$, there is a Huffman code which results in these word lengths.

1.6 Coding Sequences

In Shannon's Noiseless Coding Theorem (1.19) we don't always attain the lower bound $H(X)/\log a$. However, by coding longer *sequences* we can make our code more efficient and closer to this bound.

Example 1.24 (Motivation for Coding Sequences)

Suppose we have a memoryless (1.17) source which emits μ_1 with probability $3/4$ and μ_2 otherwise. The optimal binary code assigns $\mu_1 \mapsto 0$ and $\mu_2 \mapsto 1$.

Consider strings of length 2. We have $\mathbb{E}[S^2] = 2$, since every two-letter input sequence codes to precisely two output letters. Can we beat this? Yes, if we code strings of length 2 directly.

Consider the 4 "letters" $\mu_1\mu_1$, $\mu_1\mu_2$, $\mu_2\mu_1$, and $\mu_2\mu_2$. These have probabilities $9/16$, $3/16$, $3/16$, and $1/16$ respectively. If we apply the Huffman algorithm, we obtain word lengths of 1, 2, 3, and 3, which maps to an expected word length of $\mathbb{E}[S^2] = 27/16 < 2$.

This saving came from mapping $\mu_1\mu_1$, which is a very common sequence, to just one bit. This would not have been possible without word combination! The idea is thus to split our sequences into high probability typical sequences and low probability atypical sequences.

If a coin with $\mathbb{P}[\text{Heads}] = p$ is tossed N times, we expect approximately Np heads and $(1-p)N$ tails. A particular sequence of precisely this many heads and tails has probability:

$$p^{pN}(1-p)^{(1-p)N} = 2^{N(p \log p + (1-p) \log(1-p))} = 2^{-NH(X)},$$

where X is the result of an individual coin toss. So with high probability, we will get a typical sequence, and its probability will be close to $2^{-NH(X)}$. Can we formalise this idea?

Definition 1.25 (Asymptotic Equipartition Property)

A source $X_1, X_2, X_3 \dots$ satisfies the *Asymptotic Equipartition Property* (AEP) with constant $H \geq 0$ if for all $\varepsilon > 0$ we have that $\exists N$ s.t. $(\forall n > N, \exists T_n \subseteq \Sigma^n)$ with:

1. $\mathbb{P}[(X_1 \dots X_n) \in T_n] > 1 - \varepsilon$
2. $2^{-n(H+\varepsilon)} \leq p(x_1, \dots, x_n) \leq 2^{-n(H-\varepsilon)}$ for all $(x_1, \dots, x_n) \in T_n$

This T_n is called a *typical set*. We then encode the high probability typical sequences carefully and encode the low probability atypical sequences arbitrarily.

Remark 1.26 (AEP Helpful)

For any given $\varepsilon > 0$ and sufficiently large n , we have $p(x_1, \dots, x_n) \geq 2^{-n(H+\varepsilon)}$ for all $\mathbf{x} \in T_n$. Summing over these \mathbf{x} , we obtain:

$$1 \geq \mathbb{P}[(X_1 \dots X_n) \in T_n] \geq 2^{-n(H+\varepsilon)} |T_n| \implies |T_n| \leq 2^{n(H+\varepsilon)}$$

We encode each of these sequences into some r -length string, so we require $a^r > |T_n|$. For atypical sequences, we encode by prefixing a string of length r (not already used) with a string of length n . Then

$$\mathbb{E}[S_n] \leq \frac{\lceil n(H+\varepsilon) \rceil}{\log a} + \delta n$$

This is close to the Shannon bound (1.19)! We get $\mathbb{E}[S^n]/n \leq H/\log a + \delta'$, where we can make the δ' small, yielding a compact encoding of n -strings.

Now, we consider a property of sources, related to our intuition about the AEP.

Definition 1.27 (Reliable Encodability, Information Rate)

A source X_1, X_2, \dots is *reliably encodable* at rate r if for each n there is $A_n \subseteq \Sigma^n$ with:

1. $\log |A_n| \times (1/n) \rightarrow r$ as $n \rightarrow \infty$.
2. $\mathbb{P}[(X_1, \dots, X_n) \in A_n] \rightarrow 1$ as $n \rightarrow \infty$

The *information rate* H of a source is the infimum of all rates at which it is reliably encodable. Then nH is roughly the number of bits required to encode (X_1, \dots, X_n) .

Theorem 1.28 (Shannon's First Coding Theorem, 1948)

If a source satisfies the asymptotic equipartition property (1.25) with constant H , then the source has information rate (1.27) equal to H .

Proof: Omitted. □

We now present an alternative definition of the asymptotic equipartition property (1.25).

Definition 1.29 (Asymptotic Equipartition Property)

A source X_1, X_2, \dots satisfies the AEP if for some $H \geq 0$, we have

$$-\frac{1}{n} \log p(x_1, \dots, x_n) \xrightarrow{\mathbb{P}} H \text{ as } n \rightarrow \infty$$

where the arrow refers to *convergence in probability*. This allows the source to take very different values for large n , but only on a set of small probability.

Remark 1.30 (Weak Law of Large Numbers)

Recall that the *weak law of large numbers* states that for any independently and identically distributed sequence of random variables X_1, X_2, \dots with finite expected value $\mathbb{E}[X_i] = \mu$:

$$\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{\mathbb{P}} \mu \text{ as } n \rightarrow \infty.$$

We can apply this to our toy model of a memoryless (1.17) source, since $p(X_i)$ are iid. random variables, and $p(x_1, \dots, x_n) = p(x_1) \times \dots \times p(x_n)$, yielding:

$$-\frac{1}{n} \log p(x_1, \dots, x_n) = -\frac{1}{n} \sum_{i=1}^n \log p(x_i) \xrightarrow{\mathbb{P}} \mathbb{E}[-\log p(X)] = H(X) \text{ as } n \rightarrow \infty.$$

Thus any memoryless source satisfies the AEP with constant $H = H(X)$.

Corollary: A memoryless source has information rate equal to its entropy $H(X)$.

2 Error Control Codes

2.1 Binary Codes

Recall our initial schematic of a code from 1.1. In the previous chapter, we considered the problem of sending coded messages when the channel was *noiseless*, that is when we had a perfect guarantee that the message we sent would be accurately received. Now, we consider the case when this does not hold, because our channel is *noisy*. Examples of such channels can be found in 1.5.

Definition 2.1 (Binary Code)

An $[n, m]$ *binary code* is a subset $C \subseteq \{0, 1\}^n$ of size m . We say that C has length n . The elements of C are called *codewords*.

Note: By this definition, since all elements of C have length n , C is a *block code*, where all m of the codewords are of equal length. As seen previously, all block codes are prefix-free.

We use an $[n, m]$ -code to send one of m possible messages through a binary symmetric channel (1.5) making n uses of the channel.

Definition 2.2 (Information Rate)

The *information rate* of an $[n, m]$ binary code C is defined as $\rho(C) = \log(m)/n$.

Corollary: Since $C \subseteq \{0, 1\}^n$ is of size m , $\rho(C) \leq 1$, with equality if and only if $C = \{0, 1\}^n$ (or equivalently if $m = 2^n$). Similarly, a code of size $m = 1$ has information rate 0.

The error rate depends on the *decoding rule*. We consider three possible rules:

1. The *ideal observer* decoding rule decodes $x \in \{0, 1\}^n$ as the codeword $c \in C$ which maximises the probability $\mathbb{P}[c \text{ sent} \mid x \text{ received}]$.
2. The *maximum likelihood* decoding rule decodes $x \in \{0, 1\}^n$ as the codeword $c \in C$ which maximises the probability $\mathbb{P}[x \text{ received} \mid c \text{ sent}]$.
3. The *minimum distance* decoding rule decodes $x \in \{0, 1\}^n$ as the codeword $c \in C$ which has the fewest digits changed: that is, minimising $\#\{1 \leq i \leq n : x_i \neq c_i\}$.

Note: For each of these, some convention is needed in case of a tie (when the codeword chosen is not unique). We could choose one at random, or arbitrarily yet consistently, or ask for the message to be sent again.

Proposition 2.3 (Decoder Agreement 1)

If all messages in C are equally likely to be sent, then the *ideal observer* decoder method and the *maximum likelihood* decoder method agree on how to decode any received message.

Proof: By Bayes' rule, we can calculate the probability:

$$\mathbb{P}[c \text{ sent} \mid x \text{ received}] = \frac{\mathbb{P}[c \text{ sent}, x \text{ received}]}{\mathbb{P}[x \text{ received}]} = \frac{\mathbb{P}[c \text{ sent}]}{\mathbb{P}[x \text{ received}]} \times \mathbb{P}[x \text{ received} \mid c \text{ sent}]$$

but having received any x , this last fraction is equal for all c , and so the two probabilities must be equal to each other. The methods thus assign equal “scores” to all codewords, so must agree. \square

Now, we use the *minimum distance* rule as the basis for a definition.

Definition 2.4 (Hamming Distance)

For $x, y \in \{0, 1\}^n$, the *Hamming distance* between x and y is the scoring rule used by the minimum distance observer $d(x, y) = \#\{1 \leq i \leq n : x_i \neq y_i\}$. Notice that this is a metric!

Proposition 2.5 (Decoder Agreement 2)

If $p < 1/2$, then the *maximum likelihood* decoder method and the *minimum distance* decoder method agree on how to decode any received message.

Proof: Suppose $d(x, c) = r$. Then we can calculate the probability explicitly as:

$$\mathbb{P}[x \text{ received} \mid c \text{ sent}] = p^r (1-p)^{n-r} = (1-p)^n \times \left(\frac{p}{1-p}\right)^{n-r}$$

When $p < 1/2$, this last fraction is less than 1. Therefore choosing c to maximise this probability is the same as choosing c to minimise $d(x, c)$. \square

Note: As mentioned in 1.5, we usually take $p < 1/2$ in general. If $p > 1/2$, then our bit is flipped most of the time, so it would make more sense to send the opposite bit, in which case our bit is now mostly correct (as if $p < 1/2$). If $p = 1/2$, we have an entirely useless channel which simply outputs a stream of random bits with no correlation to what we sent, which is uninteresting.

Example 2.6 (Encoding Codewords)

Suppose we have the codewords “000” and “111” which are sent with probabilities $\alpha = 0.9$ and $1 - \alpha = 0.1$ respectively. We use a BSC with error probability $p = 1/4$.

If we receive the string “110”, how should we decode it?

Clearly, the *minimum distance* decoder (and therefore the *maximum likelihood* decoder too, by Proposition 2.5) will decode the string as “111”.

However, the *ideal observer* decoder will calculate the odds ratio as:

$$\underbrace{9 : 1}_{\text{prior of 000 vs 111}} \times \underbrace{(3/64) : (9/64)}_{\text{odds of two flips vs one flip}} = \underbrace{3 : 1}_{\text{posterior odds ratio}}$$

giving a probability of 3/4 that “000” was sent, and thus choosing it as the most likely of the two codewords to have been sent, having received “110”.

Note: The *ideal observer* rule is also known as the *minimum error* rule. It seems much better, but it requires knowing the prior probabilities of each codeword being sent. From now on, we use the other two methods, since they are equivalent.

Definition 2.7 (Error Detecting/Correcting)

C is *d-error detecting* if changing at most d letters of a codeword never produces a different codeword. Equivalently, this is the *minimum separation distance* of C .

C is *e-error correcting* if the knowledge that the string received has at most e errors is sufficient to determine with certainty which codeword was sent.

We often consider the *repetition code* of length n . Here, the codewords we want to send are just the n -long strings of all 0s and all 1s, where we simply repeat a single bit we want to send. This is an $[n, 2]$ binary code. We can detect $n - 1$ errors, and correct anything less than $n/2$ errors, which is fairly good for a code! Unfortunately, the information rate (2.2) is only $1/n$.

Note: The information rate seems like a good definition! In this example we used n bits of channel space to send 1 bit of actual information, and had an information rate of $1/n$. In fact, this holds in general: the information rate can be thought of as the “bits per bit” of a code.

Example 2.8 (Simple Parity Check Code)

The *simple parity check code* of length n , also known as the *paper tape code*, is another common code example. Here, we view the first $n - 1$ bits as the actual information to communicate, and use the last bit as a free bit to enforce the rule that the total number of 1s in the codeword is even. That is:

$$C = \left\{ (x_1 \dots x_n) \in \{0, 1\}^n : \sum_{i=1}^n x_i \equiv 0 \pmod{2} \right\}$$

is the set of codewords, which is an $[n, 2^{n-1}]$ code. It is 1-error detecting, but it cannot correct any errors (0-error correcting). Its information rate is $1 - 1/n$, which is a lot better.

Note: Suppose we change our code C to use the same permutation to reorder each codeword. Then we get a code with the same information rate, error detection capabilities, and so forth. We say such a code is *permutationally equivalent*.

In the 1940s, Richard Hamming was working at Bell Labs on an old computer which used punch cards to store and run code. Since users of punch cards were prone to making errors, there were safety checks built in to the machines, so that they could detect malformed input, and loudly alert the operators with bright flashing lights and loud noises. Hamming was frustrated by this, and was said to have remarked “Damn it, if the machine can detect the error, why can’t it correct it?”

This experience influenced him to create the original *error-correcting Hamming code*.

Example 2.9 (Hamming’s Original 1950 Code)

Let $C \in \{0, 1\}^7$ be defined by the 7-tuples which satisfy the congruences:

$$c_1 + c_3 + c_5 + c_7 \equiv 0 \pmod{2}$$

$$c_2 + c_3 + c_6 + c_7 \equiv 0 \pmod{2}$$

$$c_4 + c_5 + c_6 + c_7 \equiv 0 \pmod{2}$$

Since there are three of these congruences, the size of C is $2^{7-3} = 16$. This means that C is a $[7, 16]$ code, and thus has information rate $4/7$.

Suppose we receive some $x \in \{0, 1\}^7$. Then we form the *syndrome* $z_x = (z_1, z_2, z_4)$, where:

$$z_1 = x_1 + x_3 + x_5 + x_7$$

$$z_2 = x_2 + x_3 + x_6 + x_7$$

$$z_4 = x_4 + x_5 + x_6 + x_7$$

with addition taken modulo 2. For any $x \in C$, by construction we have $z_x = (0, 0, 0)$.

If $d(x, c) = 1$ for some $c \in C$, then the place where they differ is given by $z_1 + 2z_2 + 4z_4$. This is because if $x = c + e_i$, where e_i is a vector with all 0s except for a 1 in the i^{th} place, then the syndrome of x is the syndrome of e_i , which is the binary expansion of i for all $1 \leq i \leq 7$.

This is because, for example, x_3 appears in the definitions of z_1 and z_2 only, since $3 = 110_2$ and so only bits 1 and 2 would be affected.

Thus our code C corrects any single error! However, it doesn’t correct two: for example, the string $1110000 \in C$ could be corrupted to 1000000 , which would be decoded as 0000000 .

Now recall that in the definition of the Hamming distance (2.4), we stated that this was a metric. We now prove this formally.

Proposition 2.10 (Hamming Metric)

The Hamming distance $d(x, y) = \# \{1 \leq i \leq n : x_i \neq y_i\}$ is a metric on $\{0, 1\}^n$.

Proof: Clearly $d(x, y) \geq 0$, as the count of a set. If $d(x, y) = 0$, then x and y differ in zero places, and so they must be the same: conversely, $d(x, x)$ is clearly 0. Also, the symmetry of the definition gives us the relation $d(x, y) = d(y, x)$. So we only need to show the triangle inequality, using:

$$\{1 \leq i \leq n : x_i \neq z_i\} \subseteq \{1 \leq i \leq n : x_i \neq y_i\} \cup \{1 \leq i \leq n : y_i \neq z_i\}$$

which yields $d(x, z) \leq d(x, y) + d(y, z)$. Equivalently, it is the sum metric on n copies of the discrete metric on $\{0, 1\}$, which is therefore a metric itself. \square

Definition 2.11 (Minimum Distance)

The *minimum distance* of a code $C \subseteq \{0, 1\}^n$ is the smallest Hamming distance between two distinct codewords. An $[n, m]$ code with minimum distance d is sometimes referred to as an $[n, m, d]$ code. For example, Hamming's original code is a $[7, 16, 3]$ code.

Note: We have $m \leq 2^n$, with equality if and only if $C = \{0, 1\}^n$. This is called the *trivial code*. From this definition, we can prove some bounds on how “good” a code can be (2.7).

Proposition 2.12 (Error Bounds)

Let C be a code with minimum distance $d = d(C)$. Then:

- (i) C can always detect up to $d - 1$ errors, but not necessarily d errors.
- (ii) C can always correct $\lfloor \frac{d-1}{2} \rfloor$ errors, but not necessarily more.

Proof: Suppose $x \in \{0, 1\}^n$ and $c \in C$ with $1 \leq d(x, c) \leq d - 1$. Then $x \notin C$, as otherwise the minimum distance would not be d . Therefore C can detect up to $d - 1$ errors. But if $c_1, c_2 \in C$ with $d(c_1, c_2) = d$, then c_1 can be corrupted to c_2 with just d errors, which the code would not be able to detect. So d errors cannot always be detected.

Now, let $e = \lfloor \frac{d-1}{2} \rfloor$, so $e \leq \frac{d-1}{2} \leq e + 1$, or equivalently $2e < d \leq e + 1$.

Then take $x \in \{0, 1\}^n$. If there is some $c_1 \in C$ with $d(x, c_1) \leq e$, we want to show $d(x, c_2) > e$ for all $c_2 \neq c_1$ in C . This is given directly by the triangle inequality:

$$d(x, c_2) \geq d(c_1, c_2) - d(x, c_1) \geq d - e > e.$$

Thus C is e -error correcting. However, take $c_1, c_2 \in C$ with $d(c_1, c_2) = d$. Let x differ from c_1 in precisely $e + 1$ places where c_1 and c_2 also differ. Then $d(x, c_1) = e + 1$, and we have

$$d(x, c_2) = d - (e + 1) \leq e + 1$$

so both c_1 and c_2 can be corrupted to x with $e + 1$ errors. Therefore C cannot correct $e + 1$ errors, and so our bounds are tight. \square

Corollary: The *repetition code* is an $[n, 2, n]$ code, so detects $n - 1$ errors and corrects $\lfloor \frac{n-1}{2} \rfloor$.

Corollary: The *paper tape code* is an $[n, 2^{n-1}, 2]$ code, so detects one error but corrects none.

Corollary: The *original Hamming code* is a $[7, 16, 3]$ code, as mentioned earlier.

Given an $[n, m, d]$ code, we might want to transform it: making the code “safer” by enabling it to detect or correct more errors, but at the cost of increasing the length of the codewords. Conversely, we may want to go the other way, sacrificing robustness for efficiency.

Definition 2.13 (Parity Extension, Punctured Code, Shortened Code)

Let C be an $[n, m, d]$ code. Then the *parity extension* of C is

$$\bar{C} = \left\{ (c_1, c_2, \dots, c_n, \sum_{i=1}^n c_i) : (c_1 \dots c_n) \in C \right\},$$

where the addition is taken modulo 2. That is, the new code is the old code, where each of the codewords has an extra bit added as a parity check. This makes \bar{C} an $[n+1, m, d']$ code, where $d \leq d' \leq d+1$, depending on the parity of d .

This code is longer but potentially more error-detecting. Conversely, the *punctured code* goes the other direction, and deletes the i^{th} letter from each codeword. This forms a code which is one bit shorter, but possibly combines two codewords, unless no two codewords differ only in this letter. A sufficient condition to ensure that this does not happen is to enforce $d \geq 2$.

Similarly, we define the *shortened code* for a fixed $a \in \{0, 1\}$ and $1 \leq i \leq n$. We take all the codewords in C , and remove the i^{th} letter, given that it is an a . For some choice of a , this will retain at least $\lceil m/2 \rceil$ codewords from the original C .

2.2 Bounds on Codes

Now, we try and find bounds on codes with certain nice or maximal properties. Recall that the Hamming distance (2.4) is a metric. As in any metric space, this allows us to define a *ball*.

Definition 2.14 (Hamming Ball)

Let $x \in \{0, 1\}^n$ with $r \geq 0$. The *closed Hamming Ball* with centre x and radius r is:

$$B_r(x) = B(x, r) = \{y \in \{0, 1\}^n : d(x, y) \leq r\}$$

The *volume* of the ball is the size of the set. This is given by:

$$V(n, r) = \sum_{i=0}^r \binom{n}{i}$$

which is independent of x .

This definition allows us to quantify precisely how error correcting (2.7) a code can be.

Proposition 2.15 (Hamming’s Bound)

If $C \subseteq \{0, 1\}^n$ is e -error correcting, then we must have

$$|C| \leq \frac{2^n}{V(n, e)}.$$

Proof: Since C is e -error correcting, the Hamming balls $B_e(c)$ are pairwise disjoint for each $c \in C$. (If not, then there would be $x \in B_e(c_1) \cap B_e(c_2)$: that is, a string obtained via at most e errors on two different words.) Then $|C| \times V(n, e) \leq 2^n$, which proves the bound. \square

Note: An $[n, m]$ code (2.1) which can correct e errors is called *perfect* if this bound is tight: that is, if we have $m = 2^n/V(n, e)$.

Corollary: If $2^n/V(n, e) \notin \mathbb{Z}$ then no perfect e -error correcting code of length n can exist.

Corollary: Hamming's original $[7, 16, 3]$ code (2.9) can correct $e = 1$ errors, so we can calculate $2^7 = 128$ and $V(n, e) = 1 + 7 = 8$, so as $16 = 128/8$ the code is perfect.

Corollary: Since a perfect e -error correcting code has balls which cover the entire set, *any* instance of $e + 1$ errors will always be in another ball, and therefore must be decoded incorrectly.

Definition 2.16 (Maximal Code Size)

We define the *maximal code size* with parameters n and d to be

$$A(n, d) = \max \{m \in \mathbb{N}_0 : \text{there exists some } [n, m, d] \text{ code}\}.$$

Corollary: $A(n, 1) = 2^n$: any distance is allowed, so we can have all codewords.

Corollary: $A(n, n) = 2$: every codeword must be distinct in every bit, so there can be only two.

Proposition 2.17 (Gilbert-Shannon-Varshamov Bound)

For any n and d , we have the *GSV lower bound* and the *Hamming upper bound*:

$$\frac{2^n}{V(n, d-1)} \leq A(n, d) \leq \frac{2^n}{V(n, \lfloor \frac{d-1}{2} \rfloor)}.$$

Proof: (GSV) Let $C \subseteq \{0, 1\}^n$ be a code of length n and minimum distance d of maximal size. Then there cannot be $x \in \{0, 1\}^n$ with $d(x, c) \geq d$ for all $c \in C$, otherwise we could take $C \cup \{x\}$ to be a larger code.

Thus the union of the $B(c, d-1)$ cover $\{0, 1\}^n$, and so we must have $2^n \leq |C| \times V(n, d-1)$, since the number of total strings of length n is at most the number covered by $|C|$ balls, each of size $V(n, d-1)$. This proves the bound. \square

Note: We omit the proof of the upper bound, known as *Hamming's bound*.

Example 2.18 (GSV Bound)

Take $n = 10$ and $d = 3$, so that $2^n = 1024$. Then we have:

$$\begin{aligned} V(n, 1) &= 1 + 10 = 11 \\ V(n, 2) &= 1 + 10 + 45 = 56 \\ \implies 1024/56 &\leq A(10, 3) \leq 1024/11. \end{aligned}$$

These bounds work out to around 18.3 and 93.1 respectively, so we have

$$19 \leq A(10, 3) \leq 93.$$

So these bounds are not very tight! In fact, $A(10, 3) = 72$, discovered in 1999.

Note: In general, calculating specific values of $A(n, d)$ is an open problem.

There also exist asymptotic versions of these bounds! Loosely, as n grows to infinity, we wish to find the maximal size of a code of length n which can correct a fraction δ of the errors.

Proposition 2.19 (Asymptotic GSV)

For $0 < \delta < 1/2$, we define the limiting code size to be

$$\alpha(\delta) = \limsup_{n \in \mathbb{N}} n^{-1} \log A(n, \delta n).$$

Now, recall from 1.15 the notation $h(\delta) = -\delta \log \delta - (1 - \delta) \log(1 - \delta)$. Then we must have the asymptotic bounds:

$$1 - h(\delta) \leq \alpha(\delta) \leq 1 - h(\delta/2).$$

In particular, we claim that:

- (i) $\log V(n, \lfloor n\delta \rfloor) \leq n \times h(\delta)$, and
- (ii) $\log A(n, \lfloor n\delta \rfloor) \geq n \times (1 - h(\delta))$.

Proof: (i) Since $0 < \delta < 1/2$, $\delta < 1 - \delta$. Now, notice that we have

$$\begin{aligned} 1 &= (\delta + (1 - \delta))^n = \sum_{i=0}^n \binom{n}{i} \delta^i (1 - \delta)^{n-i} \\ &\geq \sum_{i=0}^{\lfloor n\delta \rfloor} \binom{n}{i} \delta^i (1 - \delta)^{n-i} \\ &= (1 - \delta)^n \sum_{i=0}^{\lfloor n\delta \rfloor} \binom{n}{i} \left(\frac{\delta}{1 - \delta} \right)^i \\ &\geq (1 - \delta)^n \sum_{i=0}^{\lfloor n\delta \rfloor} \binom{n}{i} \left(\frac{\delta}{1 - \delta} \right)^{n\delta} \\ &= \delta^{n\delta} (1 - \delta)^{n(1-\delta)} \sum_{i=0}^{\lfloor n\delta \rfloor} \binom{n}{i} \\ &= \delta^{n\delta} (1 - \delta)^{n(1-\delta)} V(n, \lfloor n\delta \rfloor) \end{aligned}$$

Taking logs then gives us the inequality

$$\begin{aligned} 0 &\geq n\delta \log \delta + n(1 - \delta) \log(1 - \delta) + \log V(n, \lfloor n\delta \rfloor) \\ nh(\delta) &\geq \log V(n, \lfloor n\delta \rfloor). \end{aligned}$$

(ii) Now, the GSV bound gives us:

$$\begin{aligned} A(n, \lfloor n\delta \rfloor) &\geq \frac{2^n}{V(n, \lfloor n\delta \rfloor) - 1} \geq \frac{2^n}{V(n, \lfloor n\delta \rfloor)} \\ \implies \log A(n, \lfloor n\delta \rfloor) &\geq n - \log V(n, \lfloor n\delta \rfloor) \end{aligned}$$

From subtracting the first inequality from n , we get:

$$\begin{aligned} n - nh(\delta) &\leq n - \log V(n, \lfloor n\delta \rfloor) \\ \implies n(1 - h(\delta)) &\leq n - \log V(n, \lfloor n\delta \rfloor) \leq \log A(n, \lfloor n\delta \rfloor) \end{aligned}$$

This proves the desired inequality! □

Heuristically, we can interpret δ as the fraction of errors which our code can correct. $A(n, \lfloor n\delta \rfloor)$ is then the maximum size of a code with length n which is capable of this, and the limit supremum bounds the asymptotic behaviour of such a code.

2.3 Operational Channel Capacity

Now, we consider channels again, in order to describe properties of channels by considering the best possible codes which can be used to transmit messages across them. As usual, we will mainly be considering discrete memoryless channels (1.4).

Denote $|\Sigma| = q$: usually, we take $q = 2$ for the output alphabet (that is, a binary code). A code of length n is then a subset $C \subseteq \Sigma^n$. For each code, a decoding rule is chosen: here, we focus on the *minimum distance* rule (2.5).

Definition 2.20 (Operational Channel Capacity)

We define $\hat{e}(C) = \max_{c \in C} \mathbb{P}[\text{error} \mid c \text{ sent}]$ to be the *maximum error probability* of a code C .

A channel can *transmit reliably* at rate $0 \leq R \leq 1$ if there exists some infinite sequence of codes C_1, C_2, \dots with C_n a code of length n and size $\lfloor 2^{nR} \rfloor$, such that $\hat{e}(C_n) \rightarrow 0$ as $n \rightarrow \infty$.

The *operational capacity* of a channel is then the supremum over all rates R such that the channel can transmit reliably at rate R .

Note: Later, we will define the *informational* channel capacity (Definition 2.31). In fact, these two definitions will coincide exactly, which we prove (Theorem 2.32).

Note: The information rate (2.2) of C_n is $\log \lfloor 2^{nR} \rfloor / n$, which is bounded by and tends to R .

Proposition 2.21 (Error Rate Bound)

Let $\varepsilon > 0$. Consider a BSC (1.5) with error probability p being used to send n binary digits. Then we must have

$$\lim_{n \rightarrow \infty} \mathbb{P}[\text{number of errors} \geq n \times (p + \varepsilon)] = 0.$$

Proof: Let $\mu_1 \dots \mu_n$ be a sequence of independent identically distributed random variables, taking the values representing whether an error occurs in the i^{th} position:

$$\mu_i = \begin{cases} 1 & i^{\text{th}} \text{ digit mistransmitted} \\ 0 & \text{otherwise} \end{cases}$$

Then we have $\mathbb{P}[\mu_i = 1] = p$ for all i : in particular, $\mathbb{E}[\mu_i] = p$. The probability

$$\mathbb{P}[\text{number of errors} \geq n \times (p + \varepsilon)] = \mathbb{P}\left[\sum_{i=1}^n \mu_i \geq n(p + \varepsilon)\right] \leq \mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \mu_i\right| \geq p + \varepsilon\right]$$

But the right hand side tends to 0 as $n \rightarrow \infty$, by the weak law of large numbers. Therefore the left hand side must too. \square

This allows us to prove a nice result about operational channel capacity.

Proposition 2.22 (Nonzero Channel Capacity)

The operational channel capacity (2.20) of a binary symmetric channel (1.5) with an error probability of $p < 1/4$ is not zero.

Proof: We choose δ with $2p < \delta < 1/2$, and prove that there is reliable transmission (2.20) at a rate $1 - h(\delta) > 0$. Let C_n be the largest code of length n and minimum distance $\lfloor n\delta \rfloor$. Then:

$$|C_n| = A(n, \lfloor n\delta \rfloor) \geq 2^{n(1-h(\delta))} = 2^{nR}$$

due to Gilbert-Shannon-Varshamov (2.19).

Replacing C_n by a subcode gives us $|C_n| \leq \lfloor 2^{nR} \rfloor$ with a minimum distance still at least $\lfloor n\delta \rfloor$. Using minimum distance decoding, we see that the maximum error probability is at most:

$$\hat{e}(C_n) \leq \mathbb{P} \left[\text{the BSC makes at least } \left\lfloor \frac{\lfloor n\delta \rfloor - 1}{2} + 1 \right\rfloor \text{ errors} \right]$$

which is at most the probability it makes at least $(n\delta - 1)/2$ errors. As $p < 1/4$ we may choose $\varepsilon > 0$ such that $p + \varepsilon < \delta/2$. Then

$$\frac{n\delta - 1}{2} = n \left(\frac{\delta}{2} - \frac{1}{2n} \right) > n(p + \varepsilon)$$

for sufficiently large n , and this means $\hat{e}(C_n)$ is at most the probability of making $n(p + \varepsilon)$ errors. But by the previous proposition, this tends to 0 as $n \rightarrow \infty$, and so $\hat{e}(C_n)$ does too.

Therefore there is a sequence of codes such that C_n has length n and size $\lfloor 2^{n(1-h(\delta))} \rfloor$ such that the maximum error probability $\hat{e}(C_n)$ tends to 0.

Thus the channel transmits reliably at rate $R = 1 - h(\delta) > 0$. Therefore it has some operational channel capacity at least this large: in particular, it is not zero. \square

Now, we return to the idea of entropy, and expand our understanding from the basic case (1.14) of one variable to multiple variables.

Definition 2.23 (Joint Entropy)

Let X and Y be random variables taking values in Σ_1 and Σ_2 . The *joint entropy* of X and Y is then given by

$$H(X, Y) = - \sum_{x \in \Sigma_1} \sum_{y \in \Sigma_2} p_{xy} \log p_{xy} \text{ where } p_{xy} = \mathbb{P}[X = x, Y = y].$$

Proposition 2.24 (Joint Entropy Inequality)

The joint entropy is at most the sum of the individual entropies: we have

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

with equality if and only if X and Y are independent.

Proof: Let $\Sigma_1 = \{x_1 \dots x_m\}$ and $\Sigma_2 = \{y_1 \dots y_n\}$. Define $p_{ij} = \mathbb{P}[X = x_i, Y = y_j]$ like before, with $p_i = \mathbb{P}[X = x_i]$ and $q_j = \mathbb{P}[Y = y_j]$. Then by Gibbs' inequality (1.16) we have

$$- \sum_{i,j} p_{ij} \log p_{ij} \leq - \sum_{i,j} p_{ij} \log(p_i q_j) = - \sum_i \left(\sum_j p_{ij} \right) \log p_i - \sum_j \left(\sum_i p_{ij} \right) \log q_j$$

Notice that the sum across the j of p_{ij} is just p_i , and the sum across the i is q_j . Thus

$$H(X, Y) = - \sum_{i,j} p_{ij} \log p_{ij} \leq - \sum_i p_i \log p_i - \sum_j q_j \log q_j = H(X) + H(Y)$$

with equality if and only if the two probability distributions coincide, that is $p_{ij} = p_i q_j$ for all i and j . But this is the same as the random variables being independent, proving the result. \square

Sometimes, knowing the value of a random variable gives you information about another random variable. In fact, we can quantify precisely how much information!

Definition 2.25 (Conditional Entropy)

The *conditional entropy* of a random variable X given the event $\{Y = y\}$ is given by:

$$H(X | Y = y) = - \sum_{x \in \Sigma_1} \mathbb{P}[X = x | Y = y] \log \mathbb{P}[X = x | Y = y]$$

The conditional entropy of a random variable X given another random variable Y is:

$$H(X | Y) = \sum_{y \in \Sigma_2} \mathbb{P}[Y = y] \times H(X | Y = y)$$

which is the “expected” conditional entropy, given some value of Y .

Proposition 2.26 (Conditional Entropy Equality)

The joint entropy $H(X, Y)$ is equal to $H(X | Y) + H(Y)$.

Proof: Use Bayes’ rule to rewrite the conditional entropy as:

$$\begin{aligned} H(X | Y) &= - \sum_{y \in \Sigma_2} \sum_{x \in \Sigma_1} \mathbb{P}[X = x | Y = y] \times \mathbb{P}[Y = y] \times \log \mathbb{P}[X = x | Y = y] \\ &= - \sum_{y \in \Sigma_2} \sum_{x \in \Sigma_1} \mathbb{P}[X = x, Y = y] \times \log \frac{\mathbb{P}[X = x, Y = y]}{\mathbb{P}[Y = y]} \\ &= - \sum_{y \in \Sigma_2} \sum_{x \in \Sigma_1} p_{xy} \times \log p_{xy} + \sum_{y \in \Sigma_2} \left(\sum_{x \in \Sigma_1} p_{xy} \right) \times \log q_y \\ &= H(X, Y) - H(Y) \end{aligned}$$

which proves the result. □

Corollary: $H(X | Y) \leq H(X)$ with equality if and only if X and Y are independent.

Example 2.27 (Joint and Conditional Entropy)

Suppose we throw a fair six-sided die. Define X to be the value shown, and define

$$Y = \begin{cases} 0 & X \text{ even} \\ 1 & X \text{ odd.} \end{cases}$$

Then $H(X, Y) = H(X)$, since Y is fully determined by X , and this is $\log_2 6$. The entropy of Y is simply $H(Y) = \log_2 2 = 1$.

Then the conditional entropies are:

1. $H(X | Y) = H(X, Y) - H(Y) = \log 6 - 1 = \log 3$.
2. $H(Y | X) = H(X, Y) - H(X) = \log 6 - \log 6 = 0$.

Both of these make sense! Given Y , there are 3 possible equally likely values for X . However, given X , the value of Y is totally determined, so there is no “remaining” randomness.

Note: X and Y having zero covariance is a necessary but not sufficient condition for independence. However, if $H(X | Y) = H(X)$, then independence really is always attained!

Note: In the definition of conditional and joint entropy, we did not use the actual values of X and Y , so we may replace random variables X and Y with vectors $\mathbf{x} = (x_1 \dots x_r)$ and $\mathbf{y} = (y_1 \dots y_s)$.

Proposition 2.28 (Double Conditional Entropy)

The conditional entropy $H(X | Y)$ is at most $H(X | Y, Z) + H(Z)$.

Proof: We expand $H(X, Y, Z)$ in two different ways.

1. $H(X, Y, Z) = H(Z | X, Y) + H(X | Y) + H(Y)$.
2. $H(X, Y, Z) = H(X | Y, Z) + H(Z | Y) + H(Y)$.

Since the entropy $H(Z | X, Y) \geq 0$, we must have

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

which completes the proof. \square

Theorem 2.29 (Fano's Inequality)

Suppose X and Y are random variables taking values in Σ , with $|\Sigma| = m$. Let $p = \mathbb{P}[X \neq Y]$. Then we must have

$$H(X | Y) \leq H(p) + p \log(m - 1).$$

Proof: Let Z be the indicator variable for $X \neq Y$, so that $\mathbb{E}[Z] = \mathbb{P}[Z = 0] = 1 - \mathbb{P}[Z = 1] = p$. Then by 2.28, we must have:

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

Here, $H(Z) = h(p)$. Now, we can take the first term on the right hand side and write:

$$\begin{aligned} H(X | Y = y, Z = 0) &= 0 \\ H(X | Y = y, Z = 1) &\leq \log(m - 1) \end{aligned}$$

The first line is because $Z = 0$ means $X = Y = y$ with certainty: there is no entropy. The second line is bounded by $\log(m - 1)$ because there are $m - 1$ choices for X remaining, so the maximum possible entropy is if they are all equally likely. Then

$$\begin{aligned} H(X | Y, Z) &= \sum_{y,z} \mathbb{P}[Y = y, Z = z] \times H(X | Y = y, Z = z) \\ &\leq \sum_y \mathbb{P}[Y = y, Z = 1] \times \log(m - 1) \\ &= \mathbb{P}[Z = 1] \times \log(m - 1) \end{aligned}$$

which proves the inequality, since $\mathbb{P}[Z = 1] = p$. \square

Note: We often interpret X and Y to be the input and output of a channel respectively. Then p is the probability that the transmission is incorrect.

Definition 2.30 (Mutual Information)

For X and Y random variables, the *mutual information* is given by

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

This is the amount of information about X conveyed by Y .

Corollary: We can also write this as $I(X, Y) = H(X) + H(Y) - H(X, Y) \geq 0$ by 2.26 and 2.24, so this definition is symmetric with $I(X, Y) = 0$ if and only if X and Y are independent.

2.4 Informational Channel Capacity

We now consider a different definition of the channel capacity, which also measures how good a channel is at transmitting. In the previous section, we introduced the *operational* channel capacity (Definition 2.20), which loosely defines the limiting behaviour of a channel: as the length of codes grows, the information rate tends to R . Now, we consider the *informational* channel capacity.

Definition 2.31 (Informational Channel Capacity)

Once again, consider a discrete memoryless channel (1.4). Let X take values in an alphabet Σ_1 of size m , with probabilities $p_1 \dots p_m$, and let Y be the random variable representing the channel's output when the input is X .

The *informational channel capacity* is $\max_X \{I(X, Y)\}$, where this maximum is taken over all possible random variables X as defined above.

Note: Since this capacity is a maximum and not a property of any particular input random variable, it depends only on the *channel matrix*.

Note: We are maximising over all probabilities $\mathbf{p} \in \{(p_1 \dots p_m) : p_i \geq 0, \sum p_i = 1\}$, which is a compact set, since it is closed and bounded in \mathbb{R}^m . As the function $\mathbf{p} \mapsto I(X, Y)$ is continuous, the maximum is therefore attained by the Extreme Value Theorem.

Theorem 2.32 (Shannon's Noisy Coding Theorem)

The operational channel capacity (2.20) and informational channel capacity (2.31) are in fact the same for all discrete memoryless channels.

Note: This is Shannon's second coding theorem, with the first being the *noiseless* version (1.19). We prove some cases in §2.5, first computing the capacity of certain channels assuming this result.

Example 2.33 (BSC Channel Capacity)

Suppose we have a BSC (1.5) with error probability $0 \leq p < 1/2$. Then the input X can be defined by $\mathbb{P}[X = 0] = 1 - \alpha$ and $\mathbb{P}[X = 1] = \alpha$. The output Y is then:

$$\begin{aligned}\mathbb{P}[Y = 0] &= (1 - p)(1 - \alpha) + p\alpha \\ \mathbb{P}[Y = 1] &= p(1 - \alpha) + p(1 - \alpha)\end{aligned}$$

since the probability of mistransmission is p . Recall that $h(\delta) = -\delta \log \delta - (1 - \delta) \log(1 - \delta)$. Then we have to maximise the mutual information over α , which we can calculate to be:

$$\begin{aligned}\text{capacity } C &= \max_{0 \leq \alpha \leq 1} I(X, Y) \\ &= \max_{0 \leq \alpha \leq 1} (H(Y) - H(Y | X)) \\ &= \max_{0 \leq \alpha \leq 1} (h(p(1 - \alpha) + p(1 - \alpha)) - h(p)) \\ &= 1 - h(p), \text{ attained when } \alpha = 1/2 \\ &= 1 + p \log p + (1 - p) \log(1 - p)\end{aligned}$$

In fact, recalling Proposition 2.22, we already had the bound $C \geq 1 - h(\delta)$ for all $2p < \delta < 1/2$. This was useful for $p < 1/4$: now, we have “the same bound” but with the error probability doubled! This definition also works for $p \geq 1/2$, even though we ignore these cases.

Example 2.34 (BEC Channel Capacity)

Now, we consider a BEC (Binary Erasure Channel, also from 1.5) with erasure probability p . The input is again parameterised by α in the same way, but now:

$$\begin{aligned}\mathbb{P}[Y = 0] &= (1 - p)(1 - \alpha) \\ \mathbb{P}[Y = 1] &= (1 - p)\alpha \\ \mathbb{P}[Y = *] &= p\end{aligned}$$

with $*$ being the erasure character. Now, if $Y = 0$ or $Y = 1$, we know X with certainty, since the bit is never fully “flipped”, only erased. Thus $H(X | Y = 0) = H(X | Y = 1) = 0$, and:

$$H(X | Y = *) = - \sum_x \mathbb{P}[X = x | Y = *] \log \mathbb{P}[X = x | Y = *]$$

By Bayes’ rule, we can see that:

$$\mathbb{P}[X = 0 | Y = *] = \frac{\mathbb{P}[X = 0, Y = *]}{\mathbb{P}[Y = *]} = \frac{(1 - \alpha)p}{p} = 1 - \alpha$$

and similarly $\mathbb{P}[X = 1 | Y = *] = \alpha$. This is fairly obvious: erasure is symmetric, so you gain no information over the prior. Therefore $H(X | Y = *) = h(\alpha)$, so $H(X | Y) = ph(\alpha)$. So:

$$\begin{aligned}\text{capacity } C &= \max_{0 \leq \alpha \leq 1} I(X, Y) \\ &= \max_{0 \leq \alpha \leq 1} (H(Y) - H(Y | X)) \\ &= \max_{0 \leq \alpha \leq 1} (h(\alpha) + ph(\alpha)) \\ &= (1 - p) \max_{0 \leq \alpha \leq 1} h(\alpha) \\ &= 1 - p, \text{ again attained when } \alpha = 1/2\end{aligned}$$

Thus $1 - p$ is the capacity of the channel.

Corollary: A BSC with error probability p has capacity $1 - h(p)$, and a BEC with erasure probability q has capacity $1 - q$. Thus it is equally bad to “lose” a proportion $h(p)$ bits as it is to flip a proportion p of bits. Since $h(p) > 2p$ for $0 < p < 1/2$, we can say that flipping a bit is in fact over *twice* as bad as losing it!

Note: This makes sense: with erasures, we at least know where our errors are coming from.

Definition 2.35 (Channel Extension)

We model using a channel n times as the n^{th} extension. That is, we replace the input and output alphabets by $\Sigma'_1 = \Sigma_1^n$ and $\Sigma'_2 = \Sigma_2^n$. Then the channel probabilities are:

$$\mathbb{P}[y_1 \dots y_n \text{ received} | x_1 \dots x_n \text{ sent}] = \prod_{i=1}^n \mathbb{P}[y_i \text{ received} | x_i \text{ sent}]$$

by memorylessness of the channel yielding independence.

Note: We interpret this as sending a block of n characters. The independence of the X_i states that in fact every letter is independent of all other letters. In real life, this is usually not true!

Remark 2.36 (Entropy of the English Language)

The 26 letters of the English language are obviously neither equiprobable nor independent. What is the *actual* information rate of English, assuming we consider only the 26 letters and excluding other characters?

Of course, the maximum entropy would be $\log_2(26)$, if all probabilities were the same: this is around 4.70. But of course, this isn't true. Samuel Morse (who invented Morse code) wanted to assign probabilities to letters to make his code shorter. He estimated these by counting the letters in sets of *printer's type*, which was a set of metal blocks used for traditional ink pressing in the 1800s.

Each letter was provided in different quantities for printing, with the quantities intended to approximate their use in printing. Treating these as probabilities, this distribution implies an entropy of around 4.22 (90% of the maximum entropy). Modern estimates of frequency from a much larger corpus of text gives a similar estimate of 4.14 bits.

However, the letters are also not independent! Claude Shannon was the first to estimate the true entropy, in a 1950 paper entitled *Prediction and Entropy of Printed English*. He found an entropy of around 1 bit per letter, so a “redundancy” of 75% (equivalently, an information rate of around 0.25). In fact, even using only the previous eight letters, the entropy is only 2.3 bits. This estimate is also fairly accurate compared to more modern ones!

Proposition 2.37 (Scalar Capacity)

If a DMC has informational channel capacity (2.31) C , then the n^{th} extension of the channel has information capacity nC .

Proof: Take the random variable input $\mathbf{X} = (X_1 \dots X_n)$ which produces as output the random variable $\mathbf{Y} = (Y_1 \dots Y_n)$. Then consider the entropy:

$$H(\mathbf{Y} | \mathbf{X}) = \sum_{\mathbf{x} \in \Sigma_1^n} \mathbb{P}[\mathbf{X} = \mathbf{x}] \times H(\mathbf{Y} | \mathbf{X} = \mathbf{x})$$

Since the channel is memoryless, each Y_i is independent of everything except the corresponding X_i . Therefore we can write the entropy as the sum:

$$H(\mathbf{Y} | \mathbf{X} = \mathbf{x}) = \sum_{i=1}^n H(Y_i | \mathbf{X} = \mathbf{x}) = \sum_{i=1}^n H(Y_i | X_i = x_i).$$

Therefore we can write the overall conditional entropy as

$$H(\mathbf{Y} | \mathbf{X}) = \sum_{i=1}^n \sum_{\mu \in \Sigma} H(Y_i | X_i = \mu) \times \mathbb{P}[X_i = \mu] = \sum_{i=1}^n H(Y_i | X_i)$$

So $H(\mathbf{Y} | \mathbf{X})$ is the sum of $H(Y_i | X_i)$. We know that we can bound the entropy $H(\mathbf{Y})$ from above by the sum of the n entropies $H(Y_i)$ for $1 \leq i \leq n$, which means that:

$$\begin{aligned} I(\mathbf{X}, \mathbf{Y}) &= H(\mathbf{Y}) - H(\mathbf{Y} | \mathbf{X}) \\ &\leq \sum_{i=1}^n H(Y_i) - \sum_{i=1}^n H(Y_i | X_i) \\ &\leq \sum_{i=1}^n I(X_i, Y_i) \end{aligned}$$

But this is at most nC , attained when $H(\mathbf{Y})$ is equal to the sum of the $H(Y_i)$. So when the Y_i are all independent, we have a channel capacity of nC , as required. \square

2.5 Shannon's Noisy Coding Theorem

We now consider Shannon's Noisy Coding Theorem (Theorem 2.32) in more detail, and prove it. This theorem states that the *operational* and *informational* definitions of channel capacity, as given in 2.20 and 2.31, in fact coincide.

At first, this result is surprising, as the operational channel capacity is defined in terms of reliable transmission, while the informational channel capacity is given in terms of the seemingly unrelated constructs of entropy and mutual information.

Proposition 2.38 (DMC Direction 1)

For a discrete memoryless channel (1.4), the operational channel capacity is no greater than the informational channel capacity.

Proof: Let C be the informational capacity, and suppose by way of contradiction we can transmit reliably at some rate $R > C$. Take the sequence of codes C_1, C_2, \dots with each C_n of length n and size $\lfloor 2^{nR} \rfloor$ and maximum error probability $\hat{e}(C_n) \rightarrow 0$ as $n \rightarrow \infty$.

Consider the definition of $\hat{e}(C_n)$, compared to the simple error probability:

$$e(C_n) = \frac{1}{|C_n|} \sum_{c \in C_n} \mathbb{P}[\text{error} \mid c \text{ sent}] \leq \hat{e}(C_n).$$

Take X to be the random variable input of the channel, distributed uniformly over C_n . Let Y be the random variable output when X is transmitted and decoded. Then $e(C_n) = \mathbb{P}[X \neq Y] = p_n$.

Now, since X is the uniform distribution, we have $H(X) = \log |C_n|$. For sufficiently large n , this is at least $nR - 1$, since $|C_n| = \lfloor 2^{nR} \rfloor$. Also, $H(X \mid Y) \leq h(p_n) + p_n \log(|C_n| - 1)$, by Fano's inequality (Theorem 2.29). Thus the mutual information is at most:

$$nC \geq I(X, Y) = H(X) - H(X \mid Y) \geq (nR - 1) - (1 + p_n nR)$$

since $\log(|C_n| - 1) \geq \log 2^{nR} = nR$, and $h(p_n) \leq 1$, where the fact that the capacity is nC follows from Proposition 2.37. But then rearranging yields

$$p_n \geq \frac{n(R - C) - 2}{nR} = 1 - (C/R) - (2/nR) \rightarrow 1 - (C/R).$$

Thus p_n tends to $1 - (C/R) > 0$ as $n \rightarrow \infty$, since we assumed that $C < R$. But we established that $p_n \leq \hat{e}(C_n)$, so then $\hat{e}(C_n)$ cannot tend to 0, contradicting reliable transmission! \square

Proposition 2.39 (BSC Error Probability)

For a binary symmetric channel (1.5) with error probability p , take any $R < 1 - h(p)$. Then there is some sequence of codes C_1, C_2, \dots with C_n of length n and size $\lfloor 2^{nR} \rfloor$ such that the average error probability $e(C_n) \rightarrow 0$ as $n \rightarrow \infty$.

Proof: The idea of this proof is to use a *random code*. Without loss of generality, assume $p < 1/2$. Then there is some $\varepsilon > 0$ with $R < 1 - H(p + \varepsilon)$. We use minimum distance decoding, making an arbitrary choice in case of a tie.

Let $m = \lfloor 2^{nR} \rfloor$, and pick an $[n, m]$ code C_n at random. That is, we pick each of the possible codes $C_n \subseteq \{0, 1\}^n$ at random with equal probability 2^n choose m .

Now, choose $1 \leq i \leq m$ at random, each with probability $1/m$. We send c_i through the channel, and get output Y . It suffices to show that the probability $\mathbb{P}[Y \text{ not decoded as } c_i] \rightarrow 0$ as $n \rightarrow \infty$.

Let $r = \lfloor n(p + \varepsilon) \rfloor$. Then we can split the incorrect decoding probability into two cases:

$$\mathbb{P}[Y \text{ not decoded as } c_i] = \underbrace{\mathbb{P}[c_i \notin B_r(Y)]}_{\text{too many errors}} + \underbrace{\mathbb{P}[B_r(Y) \cap C_n \not\supseteq \{c_i\}]}_{\text{some other codeword}}.$$

The first case can be written as $\mathbb{P}[d(c_i, Y) > r] = \mathbb{P}[\text{channel makes more than } r \text{ errors}]$. But this tends to 0 as $n \rightarrow \infty$, by Proposition 2.21.

Now, consider the second case. For any $j \neq i$, the randomness of the code yields

$$\mathbb{P}[c_j \in B(Y, r) \mid c_i \in B(Y, r)] = \frac{V(n, r) - 1}{2^n - 1} \leq \frac{V(n, r)}{2^n}.$$

Summing this expression over the $m - 1 \leq 2^{nR}$ other codewords and using Proposition 2.19 yields:

$$\mathbb{P}[B_r(Y) \cap C_n \not\supseteq \{c_i\}] \leq \frac{(m - 1)V(n, r)}{2^n} \leq \frac{2^{nR}V(n, r)}{2^n} \leq 2^{nR} \times 2^{nH(p+\varepsilon)} \times 2^{-n} = 2^{n(R - (1 - H(p+\varepsilon)))}$$

which tends to 0 as $n \rightarrow \infty$, since $R < 1 - H(p + \varepsilon)$ by assumption! \square

Note: This is *not* the condition for reliable transmission! For that, we require the maximum error probability $\hat{e}(C_n)$ to tend to 0, but here we have only bounded the average error probability. To salvage this proof, we simply throw out the worst half of the codewords!

Proposition 2.40 (BSC Direction 2)

For a binary symmetric channel, the operational channel capacity is at least the informational channel capacity. In particular, if the error probability is p , then let $R < 1 - h(p)$. There is then a sequence of codes with $\hat{e}(C_n) \rightarrow 0$ as $n \rightarrow \infty$.

Proof: Choose R' strictly between R and $C = 1 - h(p)$. Use the previous proposition to construct a sequence of codes C'_n which have average error probability $e(C'_n)$ tending to 0, with the size of each code being $\lfloor 2^{nR'} \rfloor$.

Then, sort the codewords in C'_n by their error probability $\mathbb{P}[\text{mistransmitted} \mid c \text{ sent}]$ and throw out the worse half. This gives a code C_n with $\hat{e}(C_n) \leq 2e(C'_n)$. Therefore $\hat{e}(C_n)$ tends to 0 as $n \rightarrow \infty$, and also $2^{nR'-1} = 2^{n(R'-1/n)} > 2^{nR}$ for sufficiently large n .

We can replace C_n by a subcode of size $\lfloor 2^{nR} \rfloor$ for sufficiently large n , and any code at all for n before this point, to obtain a sequence of codes of the right size with maximum error probability tending to 0 as required. Therefore C transmits reliably at rate R .

But this is true for all $R < 1 - h(p) = C$, and so the supremum of these rates is C . Therefore the operational channel capacity is at least this supremum, and so at least the informational channel capacity, exactly as required. \square

Note: Since these proofs used random codes, they were entirely non-constructive. In practice, we build redundancy into our codes to transmit at or below a desired error probability.

2.6 The Kelly Criterion

In 1956, John Larry Kelly Jr. was working at Bell Labs, and published *A New Interpretation of Information Rate*, a paper which applied the lessons of noisy coding and transmission rates to something very different: gambling.

The game proceeds as follows. Every day at noon, you may make a bet for any amount $\$K$ of your choice (provided you have the capital), and give this money to your friend. Your friend keeps this money and tosses a biased coin which lands on heads with probability p and tails otherwise. If heads, you receive $\$K \times u$ in return.

The question is: what is the optimal strategy? Clearly, this depends on the probability p with which you win the game, and also the proportional payout u .

1. Clearly, if $pu < 1$, then the expected value of this game is $Kpu - K < 0$, and so the game is negative in expectation. Therefore you should not take the bet.
2. If $pu = 1$, then the game is a martingale: the expected value is zero, so the game is fair. Any sort of loss aversion (which people tend to have) leads to a recommendation of not playing.
3. What if $pu > 1$? Well, the expected value is positive, but simply betting all your money isn't necessarily a good strategy. In fact, with probability 1, you will go broke eventually, and in fact will go broke with probability $1 - (1 - p)^n$ after n days.

Now, we can write down a recurrence. Suppose our fortune after n days is Z_n , where $Z_0 = 1$ is our initial wealth (starting capital), and we bet a proportion w of our wealth daily. Then we have:

$$Z_{n+1} = Z_n \times Y_{n+1} \text{ where } Y_{n+1} = \begin{cases} uw + (1 - w) & \text{if the } n + 1^{\text{st}} \text{ toss is a head} \\ (1 - w) & \text{if the } n + 1^{\text{st}} \text{ toss is a tail} \end{cases}$$

Now, we apply the weak law of large numbers, as in 1.30, noticing that $Z_n = Y_1 \times Y_2 \times \cdots \times Y_n$, and taking the sequence of independent and identically distributed random variables to be $\log Y_i$.

$$\mathbb{P} \left[\left| \frac{1}{n} \log Z_n - \mathbb{E}[\log Y_1] \right| > \varepsilon \right] \rightarrow 0.$$

So to maximise Z_n (and hence $\frac{1}{n} \log Z_n$) in the long run, we maximise

$$\begin{aligned} f(w) &= \mathbb{E}[\log Y_1] = p \log(uw + 1 - w) + (1 - p) \log(1 - w) \\ f'(w) &= \frac{(pu - 1) - (u - 1)w}{((u - 1)w + 1)(1 - w)} \end{aligned}$$

If $u < 1$ this is negative: we don't bet, and in fact should take the other side of the bet if we can! Now assume $u \geq 1$. If $pu \leq 1$, then $f(w)$ is decreasing for $w \geq 0$, so the same applies.

Now, if $pu > 1$, we take a maximum at:

$$w_0 = \frac{pu - 1}{u - 1}$$

which is therefore the proportion of our wealth we should bet! When $u = 2$, which corresponds to “even odds”, we therefore bet money if and only if $p > 1/2$: that is, if the game is biased in our favour. This again matches our heuristic.

Kelly showed how to interpret this using information theory. In his model, the gambler receives information about the game (in his example, a horse race) over a noisy channel. As in Shannon's Noisy Coding Theorem (Theorem 2.5), information can be transmitted close to the channel capacity (Definition 2.31) with negligible risk of error, in the long run.

This means if the game lasts for a sufficiently long time, the gambler can increase their fortune at arbitrarily close to this optimal rate with very high probability!

3 Linear Codes

So far, we have considered binary codes as being arbitrary subsets $C \subseteq \{0, 1\}$. Now, instead we insist on some extra structure.

Definition 3.1 (Field Of Two Elements)

We define $\mathbb{F}_2 = \{0, 1\}^n$ to be a *field* over two elements: 0 and 1. Addition and multiplication are possible in this field modulo 2. We have $0 + 0 = 1 + 1 = 0$, and $0 + 1 = 1 + 0 = 1$. We also have $0 \times 0 = 0 \times 1 = 1 \times 0 = 0$, and $1 \times 1 = 1$.

3.1 Codes and Vector Spaces

We can consider vector spaces over the field \mathbb{F}_2 . These are elements of \mathbb{F}_2^n for some length n , where we take addition to be element-wise. That is, an element of \mathbb{F}_2^n is an n -long vector with all entries being 0 or 1. We use this to define a linear code.

Definition 3.2 (Linear Code)

A code $C \subseteq \mathbb{F}_2^n$ is *linear* if $\mathbf{0} = (0, \dots, 0) \in C$ and for all x and y in C , $x + y \in C$.

Equivalently, $C \subseteq \mathbb{F}_2^n$ is linear if and only if it is a vector space over \mathbb{F}_2 . The *rank* of a code C is its dimension as such a vector space.

Note: A code of length n and rank k is called an (n, k) code.

Corollary: If C is an (n, k) code, it has a basis v_1, \dots, v_k . Then $C = \{\sum \lambda_i v_i : \lambda_i \in \mathbb{F}_2\}$. So in fact $|C| = 2^k$: an (n, k) code is an $[n, 2^k]$ code, and has information rate k/n .

Definition 3.3 (Dot Product)

For x and y in \mathbb{F}_2^n , we define the dot product $x \cdot y$ to be:

$$\sum_{i=1}^n x_i y_i \in \mathbb{F}_2.$$

By symmetry, $x \cdot y = y \cdot x$. This is also bilinear: $x \cdot (y + z) = x \cdot y + x \cdot z$.

Note: $x \cdot x = 0$ does not mean that $x = \mathbf{0}$, just that x has an even number of 1s.

Proposition 3.4 (Linear Code Construction)

Let $P \subseteq \mathbb{F}_2^n$ be any subset. Then $C = \{x \in \mathbb{F}_2^n : (p \cdot x = 0) \forall p \in P\}$ is a linear code.

Proof: $\mathbf{0} \in C$, since $p \cdot \mathbf{0} = 0$ for all $p \in P$. Also, if x and y are in C , then $p \cdot (x + y) = p \cdot x + p \cdot y$ by linearity, and this is 0, so $x + y \in C$. \square

Note: P is then called a *set of parity checks*, and C is a *parity check code* over P .

Definition 3.5 (Dual Code)

Let $C \subseteq \mathbb{F}_2^n$ be a linear code. The *dual code* C^\perp is defined to be

$$C^\perp = \{x \in \mathbb{F}_2^n : x \cdot y = 0 \forall y \in C\}.$$

Note: Dual codes are also linear codes by Proposition 3.4, but it is possible that they intersect non-trivially with their original code C .

Take $V = \mathbb{F}_2^n$, and V^* is the set of linear maps from $V \rightarrow \mathbb{F}_2$. Then consider $\phi : V \rightarrow V^*$, which sends $x \mapsto \theta_x$, with $\theta_x : y \mapsto x \cdot y$ is a linear map in V^* .

Then ϕ is a linear map! Suppose $x \in \ker \phi$. Then $x \cdot y = 0$ for all $y \in V$. Taking $y = e_i$, which is the vector with all entries 0 except entry i , we get $x_i = 0$. But this is true for all i , so in fact $x = \mathbf{0}$, and thus the kernel is trivial.

But since $\dim V = \dim V^*$, ϕ must be an isomorphism. So $\phi(C^\perp) = \{\theta \in V^* : \theta(x) = 0 \ \forall x \in C\}$, which is the “annihilator of C ” C^0 .

This means that $\dim C + \dim \phi(C^\perp) = \dim C + \dim C^\perp = n$.

Corollary: Any linear code is a parity check code.