

Part II – Coding and Cryptography (Rough)

Based on lectures by Dr R. Camina

Notes taken by Bhavik Mehta

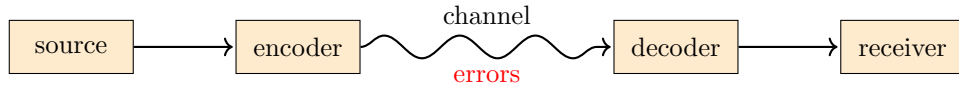
Lent 2017

Contents

1	Noiseless Coding	4
1.1	Mathematical Entropy	6
1.2	Shannon-Fano coding	9
1.3	Huffman coding	10
2	Error Control Codes	15
2.1	New codes from old	19
2.2	Bound on codes	19
2.3	Channel Capacity	21
2.4	Conditional Entropy	23
2.5	New codes from old	29

Introduction to communication channels and coding

For example, given a message m = ‘Call me!’ which we wish to send by email, first encode as binary strings using ASCII. So, $f(C) = 1000011$, $f(a) = 1100001$, and $f^*(m) = 1000011\ 1100001 \dots 0100001$.



Basic problem: Given a source and a channel (described probabilistically) we aim to design an encoder and a decoder in order to transmit information both *economically* and *reliably* (coding) and maybe also to *preserve privacy* (cryptography).

Example.

- ‘Economically’: In Morse code, common letters have shorter codewords:

$$A = .- \quad E = . \quad Q = --.-$$

- ‘Reliably’: Every book has an ISBN of form $a_1a_2 \dots a_{10}$ where $a_i \in \{0, 1, \dots, 9\}$ for $1 \leq i \leq 9$, $a_{10} \in \{0, 1, \dots, 9, X\}$ such that

$$10a_1 + 9a_2 + \dots + a_{10} \equiv 0 \pmod{11}$$

so errors can be detected (but not corrected). Similarly a 13-digit ISBN has

$$x_1 + 3x_2 + x_3 + 3x_4 + \dots + 3x_{12} + x_{13} \equiv 0 \pmod{10}$$

for $0 \leq x_i \leq 10$, doesn’t necessarily spot transpositions.

- ‘Preserve privacy’ e.g. RSA.

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

A channel is determined by the probabilities

$$\mathbb{P}(y_1 \dots y_k \text{ received} \mid x_1 \dots x_k \text{ sent})$$

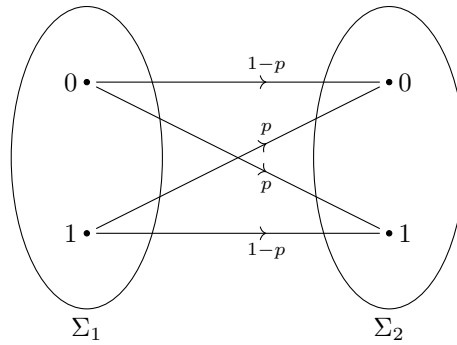
Definition (Discrete memoryless channel). A **discrete memoryless channel** (DMC) is a channel for which

$$P_{ij} = \mathbb{P}(b_j \text{ received} \mid a_i \text{ sent})$$

is the same each time the channel is used and is independent of all past and future uses.

The channel matrix is the $r \times s$ matrix with entries p_{ij} (note the rows sum to 1).

Example. Binary Symmetric Channel (BSC) has $\Sigma_1 = \Sigma_2 = \{0, 1\}$, $0 \leq p \leq 1$:

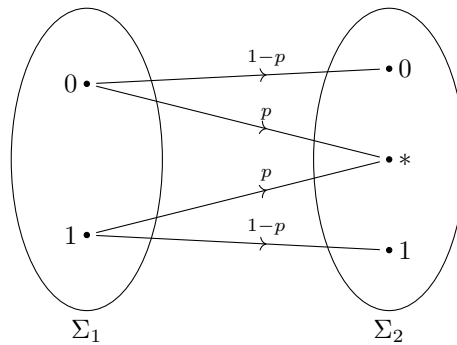


with channel matrix

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

i.e. p is the probability a symbol is mistransmitted.

Another example is given by the Binary Erasure channel, $\Sigma_1\{0,1\}$, $\Sigma_2 = \{0,1,*\}$ and $0 \leq p \leq 1$.



with channel matrix

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

i.e. p is the probability a symbol can't be read.

Informally, a channel's capacity is the highest rate at which information can be reliably transmitted over the channel. Rate refers to units of information per unit time, which we want to be high. Similarly, reliably means we want an arbitrarily small error probability.

1 Noiseless Coding

Notation. For Σ an alphabet, let $\Sigma^* = \bigcup_{n \geq 0} \Sigma^n$ be the set of all finite strings of elements of Σ .

If $x = x_1 \dots x_r$, $y = y_1 \dots y_s$ are strings from Σ , write xy for the concatenation $x_1 \dots x_r y_1 \dots y_s$. Further, $|x_1 \dots x_r y_1 \dots y_s| = r + s$, the length of the string.

Definition (Code). Let Σ_1, Σ_2 be two alphabets. A **code** is a function $f : \Sigma_1 \rightarrow \Sigma_2^*$. The strings $f(x)$ for $x \in \Sigma_1$ are called **codewords**.

Example.

- 1) Greek fire **code**:

$$\begin{aligned}\Sigma_1 &= \{\alpha, \beta, \gamma, \dots, \omega\} && 24 \text{ letters} \\ \Sigma_2 &= \{1, 2, 3, 4, 5\}\end{aligned}$$

so, $\alpha \mapsto 11, \beta \mapsto 12, \dots, \omega \mapsto 54$.

- 2) $\Sigma_1 = \{\text{all words in the dictionary}\}$, and $\Sigma_2 = \{A, B, \dots, Z, [\text{space}]\}$ and $f = \text{'spell the word and a space'}$.

Send a message $x_1 \dots x_n \in \Sigma_1^*$ as $f(x_1) \dots f(x_n) \in \Sigma_2^*$ i.e. extend f to $f^* : \Sigma_1^* \rightarrow \Sigma_2^*$.

Definition (Decipherable). A **code** f is **decipherable** if f^* is injective, i.e. every string from Σ_2^* arises from at most one message. Clearly we need f injective, but this is not enough.

Example. Take $\Sigma_1 = \{1, 2, 3, 4\}$, $\Sigma_2 = \{0, 1\}$ with

$$f(1) = 0, f(2) = 1, f(3) = 00, f(4) = 01.$$

f injective but $f^*(312) = 0001 = f^*(114)$ so f^* not **decipherable**.

Notation. If $|\Sigma_1| = m$, $|\Sigma_2| = a$, then we say f is an a -ary code of size m . (If $a = 2$ we say binary).

Aim. Construct **decipherable codes** with short word lengths.

Provided $f : \Sigma_1 \rightarrow \Sigma_2^*$ is injective, the following codes are always decipherable.

- (i) A **block code** is a code with all codewords of the same length (e.g. **Greek fire code**).
- (ii) In a **comma code**, we reserve one letter from Σ_2 that is only used to signal the end of the codeword (e.g. **Example 2 above**).
- (iii) A **prefix-free code** is a code where no codeword is a prefix of another (if $x, y \in \Sigma_2^*$, x is a prefix of y if $y = xz$ for some $z \in \Sigma_2^*$).

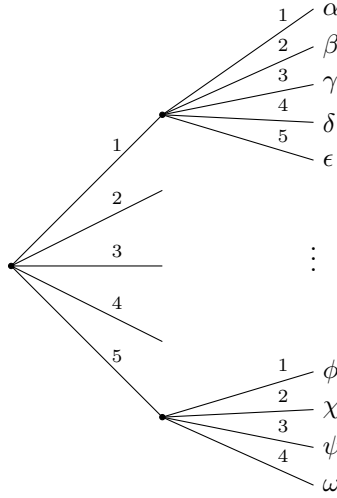
Remark. (i) and (ii) are special cases of (iii).

Prefix-free codes are also known as **instantaneous codes** (i.e. a word can be recognised as soon as it is complete) or **self-punctuating codes**.

Theorem 1.1 (Kraft's inequality). Let $|\Sigma_1| = m, |\Sigma_2| = a$. A **prefix-free code** $f : \Sigma_1 \rightarrow \Sigma_2^*$ with word lengths s_1, \dots, s_m exists iff

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

Proof. (\Rightarrow) Consider an infinite tree where each node has a descendant, labelled by the elements of Σ_2 . Each **codeword** corresponds to a node, the path from the root to this node spelling out the codeword. For example,



Assuming f is **prefix-free**, no codeword is the ancestor of any other. Now view the tree as a network with water being pumped in at a constant rate and dividing the flow equally at each node.

The total amount of water we can extract at the codewords is $\sum_{i=1}^m a^{-s_i}$, which is therefore ≤ 1 .

(\Leftarrow) Conversely, suppose we can construct a prefix-free code with word lengths s_1, \dots, s_m , wlog $s_1 \leq s_2 \leq \dots \leq s_m$. We pick codewords of lengths s_1, s_2, \dots sequentially ensuring previous codewords are not prefixes. Suppose there is no valid choice for the r th codeword. Then reconstructing the tree as above gives $\sum_{i=1}^{r-1} a^{-s_i} = 1$, contradicting our assumption. So we can construct a prefix-free code. (There is a more algebraic proof in Welsh.) \square

Theorem 1.2 (McMillan). Every **decipherable code** satisfies **Kraft's inequality**.

Proof. (Karush) Let $f : \Sigma_1 \rightarrow \Sigma_2^*$ be a **decipherable code** with word lengths s_1, \dots, s_m , let $s = \max_{1 \leq i \leq m} s_i$. Let $r \in \mathbb{N}$,

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

where b_l is the # of ways of choosing r codewords of total length l . f decipherable $\implies b_l \leq |\Sigma_2|^l = a^l$.

Thus

$$\begin{aligned} \left(\sum_{i=1}^m a^{-s_i} \right)^r &\leq \sum_{l=1}^{rs} a^l a^{-l} = rs \\ \Rightarrow \sum_{i=1}^m a^{-s_i} &\leq (rs)^{\frac{1}{r}} \rightarrow 1 \text{ as } r \rightarrow \infty. \end{aligned}$$

(As $\frac{\log r + \log s}{r} \rightarrow 0$ as $r \rightarrow \infty$).

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

□

Corollary. A decipherable code with prescribed word lengths exists iff there exists a prefix-free code with the same word lengths.

So we can restrict our attention to prefix-free codes.

1.1 Mathematical Entropy

Definition (Entropy). The entropy of X :

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

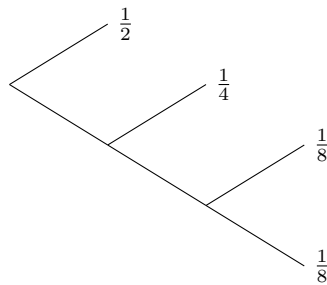
where, in this course, $\log = \log_2$.

Remark.

- (i) If $p_i = 0$, we take $p_i \log p_i = 0$.
- (ii) $H(x) \geq 0$.

Example.

1. Suppose $p_1 = p_2 = p_3 = p_4 = \frac{1}{4}$. We identify $\{x_1, x_2, x_3, x_4\}$ with $\{\text{HT}, \text{HT}, \text{TH}, \text{TT}\}$. Then $H(X) = 2$.
2. Take $(p_1, p_2, p_3, p_4) = (\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8})$.

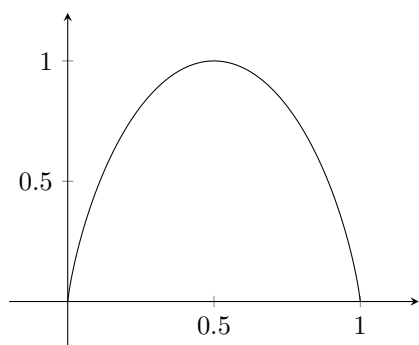


$$H(X) = \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{8} \times 3 = \frac{7}{4}.$$

So example 1 is more random than example 2.

Entropy is a measure of ‘randomness’ or ‘uncertainty’. Consider a random variable X taking values x_1, \dots, x_n with probability p_1, \dots, p_n ($\sum p_i = 1, 0 \leq p_i \leq 1$). The entropy $H(X)$ is roughly speaking the expected number of tosses of a fair coin needed to simulate X (or the expected number of yes/no questions we need to ask in order to establish the value of X).

Example. We toss a biased coin, $\mathbb{P}(\text{heads}) = p, \mathbb{P}(\text{tails}) = 1-p$. Write $H(p) = H(p, 1-p) = -p \log p - (1-p) \log(1-p)$. If $p = 0$ or 1 , the outcome is certain and so $H(p) = 0$. **Entropy** is maximal where $p = \frac{1}{2}$, i.e. a fair coin.



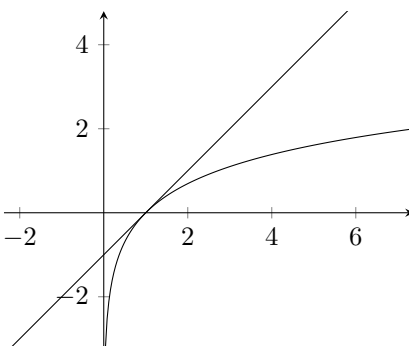
Note the **entropy** can also be viewed as the expected value of the information of X , where information is given by $I(X = x) = -\log_2 \mathbb{P}(X = x)$. For example, if a coin always lands heads we gain no information from tossing the coin. The entropy is the average amount of information conveyed by a random variable X .

Lemma 1.3 (Gibbs’ Inequality). Let p_1, \dots, p_n and q_1, \dots, q_n be probability distributions. Then

$$-\sum p_i \log p_i \leq -\sum p_i \log q_i$$

with equality iff $p_i = q_i$.

Proof. Since $\log x = \frac{\ln x}{\ln 2}$ it suffices to prove the inequality with \log replaced with \ln . Note $\ln x \leq x - 1$, equality iff $x = 1$.



Let $I = \{1 \leq i \leq n \mid p_i \neq 0\}$

$$\begin{aligned}
& \ln \frac{q_i}{p_i} \leq \frac{q_i}{p_i} - 1 \quad \forall i \in I \\
& \sum_{i \in I} p_i \ln \frac{q_i}{p_i} \leq \sum_{i \in I} q_i - \underbrace{\sum_{i \in I} p_i}_{=1} \leq 0 \\
& \implies -\sum_{i \in I} p_i \ln p_i \leq -\sum_{i \in I} p_i \ln q_i \\
& \implies -\sum_{i=1}^n p_i \ln p_i \leq -\sum_{i=1}^n p_i \ln q_i
\end{aligned}$$

If equality holds then $\frac{q_i}{p_i} = 1 \quad \forall i \in I$. So, $\sum_{i \in I} q_i = 1$ and hence $p_i = q_i$ for $1 \leq i \leq n$. \square

Corollary. $H(p_1, \dots, p_n) \leq \log n$ with equality iff $p_1 = p_2 = \dots = p_n = \frac{1}{n}$.

Proof. Take $q_1 = q_2 = \dots = q_n = \frac{1}{n}$ in [Gibbs' Inequality](#). \square

Suppose we have two alphabets Σ_1, Σ_2 with $|\Sigma_1| = m$ and $|\Sigma_2| = a$, for $m \geq 2$ and $a \geq 2$. We model the source as a sequence of random variables X_1, X_2, \dots taking values in Σ_1 .

Definition (Memoryless source). A **Bernoulli** or **memoryless** source is a sequence of independently, identically distributed random variables.

That is, for each $\mu \in \Sigma_1$, $\mathbb{P}(X_i = \mu)$ is independent of i and independent of all past and future symbols emitted. Thus

$$\mathbb{P}(X_1 = x_1, X_2 = x_2, \dots, X_k = x_k) = \prod_{i=1}^k \mathbb{P}(X_i = x_i).$$

Let $\Sigma_1 = \{\mu_1, \dots, \mu_n\}$, $p_i = \mathbb{P}(X = \mu_i)$ (assume $p_i > 0$).

Definition (Expected word length). The **expected word length** of a code $f : \Sigma_1 \rightarrow \Sigma_2^*$ with word lengths s_1, \dots, s_m is $E(S) = \sum_{i=1}^m p_i s_i$.

Definition (Optimal code). A [code](#) $f : \Sigma_1 \rightarrow \Sigma_2^*$ is **optimal** if it has the shortest possible [expected word length](#) among [decipherable](#) codes.

Theorem 1.4 (Shannon's Noiseless Coding Theorem). The minimum [expected word length](#) of a [decipherable code](#) $f : \Sigma_1 \rightarrow \Sigma_2^*$ satisfies

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

Proof. The lower bound is given by combining [Gibbs' Inequality](#) and [Kraft's inequality](#). Let $q_i = \frac{a^{-s_i}}{c}$ where $c = \sum a^{-s_i} \leq 1$ by Kraft's inequality. Note $\sum q_i = 1$.

$$\begin{aligned}
H(X) &= -\sum p_i \log p_i \leq -\sum_i p_i \log q_i \\
&= \sum p_i (s_i \log a + \log c) \\
&= \left(\sum p_i s_i \right) \log a + \underbrace{\log c}_{\leq 0} \leq E(S) \log a \\
\implies \frac{H(X)}{\log a} &\leq E(S)
\end{aligned}$$

We get equality $\iff p_i = a^{-s_i}$ for some integers s_i . For the upper bound put

$$s_i = \lceil -\log_a p_i \rceil$$

where $\lceil x \rceil$ means least integer $\geq x$.

We have

$$\begin{aligned}
-\log_a p_i &\leq s_i < -\log_a p_i + 1 \\
\implies a^{-s_i} &\leq p_i \implies \sum a^{-s_i} \leq \sum p_i \leq 1.
\end{aligned}$$

So by [Theorem 1.1](#), \exists a prefix-free code with word lengths s_1, \dots, s_m . Also,

$$\begin{aligned}
E(S) &= \sum p_i s_i \\
&< \sum p_i (-\log_a p_i + 1) \\
&= \frac{H(X)}{\log a} + 1
\end{aligned}$$

□

Remark. The lower bound holds for all [decipherable codes](#).

1.2 Shannon-Fano coding

Follows from above proof. Set $s_i = \lceil -\log_a p_i \rceil$ and construct a [prefix-free code](#) with word lengths s_1, \dots, s_m by taking the s_i in increasing order ensuring that previous [codewords](#) are not prefixes. [Kraft's inequality](#) ensures there is enough room.

Example. Suppose μ_1, \dots, μ_5 are emitted with probabilities 0.4, 0.2, 0.2, 0.1, 0.1.

A possible [Shannon-Fano](#) code (with $a = 2$, $\Sigma_2 = \{0, 1\}$) has

p_i	$\lceil -\log_2 p_i \rceil$	
0.4	2	00
0.2	3	010
0.2	3	100
0.1	4	1100
0.1	4	1110

This has [expected word length](#)

$$\begin{aligned}
&= 2 \times 0.4 + 3 \times 0.2 + 3 \times 0.2 + 4 \times 0.1 + 4 \times 0.1 \\
&= 2.8.
\end{aligned}$$

compare $H(X) \approx 2.12$.

1.3 Huffman coding

For simplicity, take $a = 2$. Take $\Sigma_1 = \{\mu_1, \dots, \mu_m\}$ with $p_i = \mathbb{P}(X = \mu_i)$. Without loss of generality, $p_1 \geq p_2 \geq \dots \geq p_m$. Huffman coding is defined inductively.

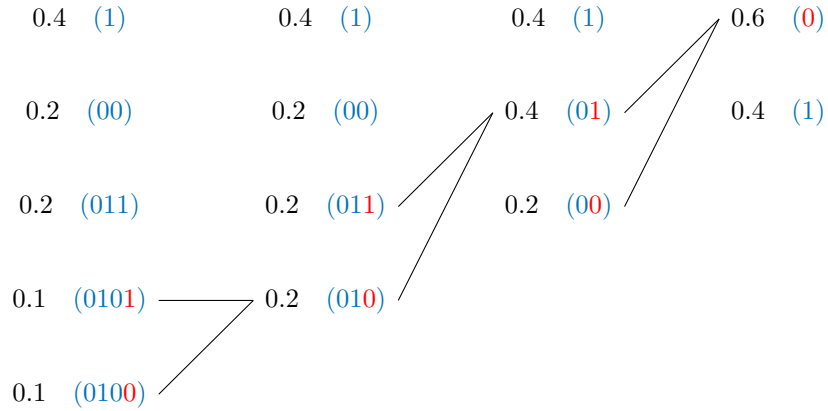
If $m = 2$, assign codewords 0 and 1. If $m > 2$, find a Huffman coding in the case of messages $\mu_1, \mu_2, \dots, \mu_m$, with probabilities $p_1, p_2, \dots, p_{m-1} + p_m$.

Append 0 (resp, 1) to the codeword for μ_{m-1} (resp, μ_m).

Remark.

- i) This construction gives a **prefix-free** code.
- ii) We exercise some choice when some of the p_i are equal. So **Huffman codes** are not unique.

Example. Use the same **example probabilities** as earlier.



So $\{1, 00, 011, 0101, 0100\}$ is the **prefix-free code** constructed. The **expected word length** is:

$$\begin{aligned}
 &= 1 \times 0.4 + 2 \times 0.2 + 2 \times 0.2 + 4 \times 0.1 + 4 \times 0.1 \\
 &= 0.4 + 0.4 + 0.6 + 0.4 + 0.4 \\
 &= 2.2.
 \end{aligned}$$

This is better than **Shannon-Fano**, which gave 2.8.

Theorem 1.5. **Huffman coding** is **optimal**.

Lemma 1.6. Suppose we have $\mu_1, \dots, \mu_m \in \Sigma_1$ 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

- i) If $p_i > p_j$, then $s_i \leq s_j$.
- ii) \exists two **codewords** of maximal length which are equal up to the last digit.

Proof.

- i) If not, then swap the i th and j th **codewords**. This decreases the **expected word length**, contradicting f **optimal**.

- ii) If not, then either only one codeword of maximal length, or any two codewords of maximal length differ before the last digit. In either case, delete the last digit of each codeword of maximal length. This maintains the [prefix-free](#) condition, contradicting f optimal. \square

Proof of Theorem 1.5. We only take the case $a = 2$. We show by induction on m that any [Huffman code](#) of size m is [optimal](#).

For $m = 2$, [codewords](#) 0, 1 are optimal.

For $m > 2$, say source X_m emits μ_1, \dots, μ_m with probabilities $p_1 \geq p_2 \geq \dots \geq p_m$ while source X_{m-1} emits $\mu_1, \dots, \mu_{m-2}, \nu$ with probabilities $p_1, \dots, p_{m-2}, p_{m-1} + p_m$.

We construct a Huffman coding f_{m-1} for X_{m-1} and extend to a Huffman coding for X_m . Then the [expected codeword length](#) satisfies:

$$E(s_m) = E(s_{m-1}) + p_{m-1} + p_m \quad (\dagger)$$

Let f'_m be an optimal code for X_m , WLOG f'_m prefix-free. [Lemma 1.6](#) tells us that by shuffling codewords, we may assume that the last two codewords are of maximal length and differ only in the last digit, say y_0 and y_1 for some string y .

We define a code f'_{m-1} for X_{m-1} with

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

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

$$E(s'_m) = E(s'_{m-1}) + p_{m-1} + p_m \quad (\ddagger)$$

Induction hypothesis tells us f_{m-1} is optimal, so $E(s_{m-1}) \leq E(s'_{m-1})$. Hence $E(s_m) \leq E(s'_m)$ by (\dagger) , (\ddagger) . So f_m optimal. \square

Remark. Not all [optimal](#) codes are [Huffman](#). For instance, take $m = 4$, and probabilities 0.3, 0.3, 0.2, 0.2. An optimal code is given by 00, 01, 10, 11, but this is not Huffman.

But, the previous result says that if we have a prefix-free optimal code with word lengths s_1, \dots, s_m and associated probabilities p_1, \dots, p_m , \exists a Huffman code with these word lengths.

Definition (Joint entropy). The **joint entropy** of X and Y is

$$H(X, Y) = - \sum_{x \in \Sigma_1} \sum_{y \in \Sigma_2} \mathbb{P}(X = x, Y = y) \log \mathbb{P}(X = x, Y = y).$$

Lemma 1.7. $H(X, Y) \leq H(X) + H(Y)$, with equality $\iff X, Y$ independent.

Proof. Take $\Sigma_1 = \{x_1, \dots, x_m\}$, $\Sigma_2 = \{y_1, \dots, y_n\}$. Let $p_{ij} = \mathbb{P}(X = x_i, Y = y_j)$, as well as $p_i = \mathbb{P}(X = x_i)$, $q_j = \mathbb{P}(Y = y_j)$. Apply [Gibbs' Inequality](#) to the distributions p_{ij} and $p_i q_j$:

$$\begin{aligned} - \sum p_{ij} \log(p_{ij}) &\leq - \sum p_{ij} \log(p_i q_j) \\ &= - \sum_i \left(\sum_j p_{ij} \log p_i \right) - \sum_j \left(\sum_i p_{ij} \log q_j \right) \\ &= - \sum p_i \log p_i - \sum q_j \log q_j \end{aligned}$$

That is, $H(X, Y) \leq H(X) + H(Y)$. Equality $\iff p_{ij} = p_i q_j \forall i, j \iff X, Y$ independent. \square

Suppose a source Ω produces a stream X_1, X_2, \dots of random variables with values in Σ . The probability mass function (p.m.f.) of $X^{(n)} = (X_1, \dots, X_n)$ is given by

$$p_n(x_1, \dots, x_n) = \mathbb{P}(X_1, \dots, X_n = x_1, \dots, x_n) \quad \forall x_1, \dots, x_n \in \Sigma^n$$

Now,

$$\begin{aligned} p_n &: \Sigma^n \rightarrow \mathbb{R} \\ X^{(n)} &: \Omega \rightarrow \Sigma^n \end{aligned}$$

can form

$$p_n(X^{(n)}) : \Sigma \xrightarrow{X^{(n)}} \Sigma^n \xrightarrow{p_n} \mathbb{R}$$

a random variable sending $\omega \mapsto p_n(X^{(n)} = X^{(n)}(\omega))$.

Example. Take $\Sigma = \{A, B, C\}$, with

$$X^{(2)} = \begin{cases} AB & p = 0.3 \\ AC & p = 0.1 \\ BC & p = 0.1 \\ BA & p = 0.2 \\ CA & p = 0.25 \\ CB & p = 0.05 \end{cases}$$

So, $p_2(AB) = 0.3$, etc, and $p_2(X^{(2)})$ takes values

$$p_2(X^{(2)}) = \begin{cases} 0.3 & p = 0.3 \\ 0.1 & p = 0.2 \\ 0.2 & p = 0.2 \\ 0.25 & p = 0.25 \\ 0.05 & p = 0.05 \end{cases}$$

Definition (Convergence in probability). A sequence of random variables X_1, X_2, \dots **converges in probability** to $c \in \mathbb{R}$, written $X_n \xrightarrow{p} c$ as $n \rightarrow \infty$, if

$$\forall \epsilon > 0 \quad \mathbb{P}(|X_n - c| \leq \epsilon) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

So, X_n and c can take very different values for large n , but only on a set with small probability.

Theorem (Weak law of large numbers). X_1, X_2, \dots an independent, identically distributed sequence of random variables with finite expected value μ , then

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

Example (Application). Take X_1, X_2, \dots a **Bernoulli source**. Then $p(X_1), p(X_2), \dots$ are i.i.d. random variables

$$\begin{aligned} p(X_1, \dots, X_n) &= p(X_1) \dots p(X_n) \\ -\frac{1}{n} \log p(X_1, \dots, X_n) &= -\frac{1}{n} \sum_{i=1}^n \log p(X_i) \xrightarrow{p} E(-\log p(X_1)) = H(X_1) \quad \text{as } n \rightarrow \infty. \end{aligned}$$

Definition (Asymptotic Equipartition Property).

A source X_1, X_2, \dots satisfies the **Asymptotic Equipartition Property** (AEP) if for some $H \geq 0$ we have

$$-\frac{1}{n} \log p(X_1, \dots, X_n) \xrightarrow{p} H \quad \text{as } n \rightarrow \infty.$$

Example. Consider a coin, $p(H) = p$. If coin tossed N times, expect approximately pN heads and $(1-p)N$ tails.

$$\begin{aligned} & \mathbb{P}(\text{particular sequence of } pN \text{ heads and } (1-p)N \text{ tails}) \\ &= p^{pN} (1-p)^{(1-p)N} \\ &= 2^{N(p \log p + (1-p) \log(1-p))} = 2^{-NH(A)} \end{aligned}$$

where A is the result of an independent coin toss. So, with high probability we will get a typical sequence, and its probability will be close to $2^{-NH(A)}$.

Lemma 1.8. A source X_1, \dots, X_n satisfies [AEP](#) iff it satisfies the following.

$\forall \epsilon > 0, \exists n_0(\epsilon)$ such that $\forall n \geq n_0(\epsilon) \exists$ a ‘typical set’ $T_n \subset \Sigma^n$ with

$$\begin{aligned} & \mathbb{P}((X_1, \dots, X_n) \in T_n) > 1 - \epsilon \\ & 2^{-n(H+\epsilon)} \leq p(x_1, \dots, x_n) \leq 2^{-n(H-\epsilon)} \quad \forall (x_1, \dots, x_n) \in T_n \end{aligned} \quad (*)$$

Sketch proof. [AEP](#) \Rightarrow $(*)$. Take

$$\begin{aligned} T_n &= \left\{ (x_1, \dots, x_n) \in \Sigma^n \mid \left| -\frac{1}{n} \log p(x_1, \dots, x_n) - H \right| < \epsilon \right\} \\ &= \left\{ (x_1, \dots, x_n) \in \Sigma^n \mid 2^{-n(H+\epsilon)} \leq p(x_1, \dots, x_n) \leq 2^{-n(H-\epsilon)} \right\} \end{aligned}$$

$(*) \Rightarrow$ [AEP](#)

$$\mathbb{P} \left(\left| -\frac{1}{n} \log p(X_1, \dots, X_n) - H \right| < \epsilon \right) \geq \mathbb{P}(T_n) \rightarrow 1 \quad \text{as } n \rightarrow \infty. \quad \square$$

Definition (Reliably encodable). A source X_1, X_2, \dots is **reliably encodable** at rate r if $\exists A_n \subset \Sigma^n$ for each n such that

- (i) $\frac{\log |A_n|}{n} \rightarrow r \quad \text{as } n \rightarrow \infty$
- (ii) $\mathbb{P}((X_1, \dots, X_n) \in A_n) \rightarrow 1 \quad \text{as } n \rightarrow \infty$.

So, in principle you can encode at rate almost r with negligible error for long enough strings.

So, if $|\Sigma| = a$, you can [reliably encode](#) at rate $\log a$. However you can often do better. For example, consider telegraph English with 26 letters and a space. $27 \approx 2^{4.756}$, so can encode at rate of 4.76 bits/letter. But much lower rates suffice, as there is a lot of redundancy in the English language. Hence the following definition.

Definition (Information rate). The **information rate** H of a source is the infimum of all values at which it is [reliably encodable](#).

Roughly, nH is the number of bits required to encode (X_1, \dots, X_n) .

Theorem 1.9 (Shannon's first coding theorem). If a source satisfies [AEP](#) with some constant H , then the source has information rate H .

Proof. Let $\epsilon > 0$ and let $T_n \subset \Sigma^n$ be [typical sets](#). Then for sufficiently large $n \geq n_0(\epsilon)$,

$$p(x_1, \dots, x_n) \geq 2^{-n(H+\epsilon)} \quad \forall (x_1, \dots, x_n) \in T_n.$$

So,

$$\begin{aligned} \mathbb{P}((X_1, \dots, X_n) \in T_n) &\geq 2^{-n(H+\epsilon)} |T_n| \\ \implies |T_n| &\leq 2^{n(H+\epsilon)} \end{aligned}$$

therefore the source is [reliably encodable](#) at rate $H + \epsilon$.

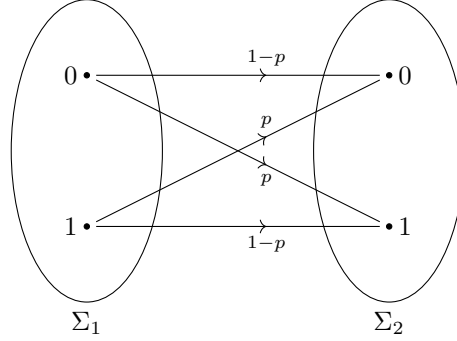
Conversely, if $H = 0$, done. Otherwise, pick $0 < \epsilon < \frac{H}{2}$. Suppose the source is reliably encodable at rate $H - 2\epsilon$, say with sets A_n .

$$\begin{aligned} p(x_1, \dots, x_n) &\leq 2^{-n(H-\epsilon)} \quad \forall (x_1, \dots, x_n) \in T_n \\ \implies \mathbb{P}((x_1, \dots, x_n) \in A_n \cap T_n) &\leq 2^{-n(H-\epsilon)} |A_n| \\ \implies \frac{\log \mathbb{P}(A_n \cap T_n)}{n} &\leq -(H - \epsilon) + \frac{\log |A_n|}{n} \\ &\rightarrow -(H - \epsilon) + (H - 2\epsilon) = -\epsilon \quad \text{as } n \rightarrow \infty. \quad \square \end{aligned}$$

2 Error Control Codes

Definition (Binary code). A $[n, m]$ **binary code** is a subset $C \subset \{0, 1\}^n$ of size $|C| = m$. We say C has length n . The elements of C are called **codewords**.

We use a $[n, m]$ -code to send one of m possible messages through a **BSC** making use of the channel n times.



Definition (Information rate). The **information rate** of C is

$$\rho(C) = \frac{\log m}{n}$$

where we continue to use $\log = \log_2$.

Note since $C \subset \{0, 1\}^n$, $\rho(C) \leq 1$, with equality iff $C = \{0, 1\}^n$. A **code** with size $m = 1$ has **information rate** 0.

We aim to design codes with both a large information rate and a small error rate, which are contradicting aims. The error rate depends on the decoding rule. We consider 3 possible rules.

- (i) The **ideal observer** decoding rule decodes $x \in \{0, 1\}^n$ as the **codeword** c maximising $\mathbb{P}(c \text{ sent} \mid x \text{ received})$.
- (ii) The **maximum likelihood** decoding rule decodes $x \in \{0, 1\}^n$ as the codeword c maximising $\mathbb{P}(x \text{ received} \mid c \text{ sent})$.
- (iii) The **minimum distance** decoding rule decodes $x \in \{0, 1\}^n$ as the codeword c minimising $\#\{1 \leq i \leq n \mid x_i \neq c_i\}$.

Remark. Some convention should be agreed in the case of a ‘tie’, e.g. choose at random, or ask for message to be sent again.

Lemma 2.1. If all messages are equally likely, then the **ideal observer** and **maximum likelihood rule** agree.

Proof. By Bayes’ rule,

$$\begin{aligned} \mathbb{P}(c \text{ sent} \mid x \text{ received}) &= \frac{\mathbb{P}(c \text{ sent}, x \text{ received})}{\mathbb{P}(x \text{ received})} \\ &= \frac{\mathbb{P}(c \text{ sent})}{\mathbb{P}(x \text{ received})} \mathbb{P}(x \text{ received} \mid c \text{ sent}). \end{aligned}$$

We suppose $\mathbb{P}(c \text{ sent})$ is independent of c . So for fixed x , maximising $\mathbb{P}(c \text{ sent} \mid x \text{ received})$ is the same as maximising $\mathbb{P}(x \text{ received} \mid c \text{ sent})$. \square

Definition (Hamming distance). Let $x, y \in \{0, 1\}^n$. Then **Hamming distance** between x and y is

$$d(x, y) = \# \{ 1 \leq i \leq n \mid x_i \neq y_i \}.$$

Lemma 2.2. If $p < \frac{1}{2}$, then **maximum likelihood** and **minimum distance** agree.

Proof. Suppose $d(x, c) = r$,

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

Since $p < \frac{1}{2}$, $\frac{p}{1-p} < 1$. So choosing c to maximise $\mathbb{P}(x \text{ received} \mid c \text{ sent})$ is the same as choosing c to minimise $d(x, c)$. \square

Note we assumed $p < \frac{1}{2}$, which is not unreasonable.

Example. Suppose codewords 000 and 111 are sent with probabilities $\alpha = \frac{9}{10}$ and $1-\alpha = \frac{1}{10}$ respectively. We use a **BSC** with $p = \frac{1}{4}$. If we receive 110, how should it be decoded? Clearly **minimum distance** and therefore **maximum likelihood** (by **Lemma 2.2**) say decode as 111. For **ideal observer**:

$$\begin{aligned} \mathbb{P}(000 \text{ sent} \mid 110 \text{ received}) &= \frac{\mathbb{P}(000 \text{ sent}, 110 \text{ received})}{\mathbb{P}(110 \text{ received})} \\ &= \frac{\alpha p^2 (1-p)}{\alpha p^2 (1-p) + (1-\alpha) p (1-p)^2} \\ &= \frac{\alpha p}{\alpha p + (1-\alpha)(1-p)} \\ &= \frac{9/40}{9/40 + 3/40} = \frac{3}{4} \\ \mathbb{P}(000 \text{ sent} \mid 110 \text{ received}) &= \frac{(1-\alpha)p^2(1-p)}{(1-\alpha)p^2(1-p) + \alpha p(1-p)^2} \\ &= \frac{(1-\alpha)(1-p)}{(1-\alpha)(1-p) + \alpha p} \\ &= \frac{3/40}{9/40 + 3/40} = \frac{1}{4}. \end{aligned}$$

So the ideal observer rule says decode as 000.

Remark. The **ideal observer** rule is also known as the minimum-error rule. But it does rely on knowing the probabilities of the codewords sent.

From now on, we use **minimum distance** decoding.

Definition. For a **binary code** C ,

- (i) C is **d -error detecting** if changing at most d letters of a codeword cannot give another codeword.

- (ii) C is **e -error correcting** if knowing that the string received has at most e errors is sufficient to determine which codeword was sent.

Example.

1. The repetition code of length n has:

$$C = \{\underbrace{0 \cdots 0}_n, \underbrace{1 \cdots 1}_n\}$$

This is an $[n, 2]$ -code. It is $n - 1$ error detecting, and $\lfloor \frac{n-1}{2} \rfloor$ error correcting. But it has information rate $= \frac{1}{n}$.

2. The simple parity check code of length n (also known as paper tape code). We view $\{0, 1\} = \mathbb{F}_2 = \mathbb{Z}/2\mathbb{Z}$ (i.e. do arithmetic modulo 2).

$$C = \left\{ (x_1, \dots, x_n) \in \{0, 1\}^n \mid \sum_{i=1}^n x_i = 0 \right\}.$$

This is a $[n, 2^{n-1}]$ -code. It is 1-error detecting and 0-error correcting. It has information rate $\frac{n-1}{n}$.

3. Hamming's original code (1950). Let $C \subseteq \mathbb{F}_2^7$ be defined by

$$c_1 + c_3 + c_5 + c_7 = 0$$

$$c_2 + c_3 + c_6 + c_7 = 0$$

$$c_4 + c_5 + c_6 + c_7 = 0$$

where all arithmetic is modulo 2. Since we may choose c_3, c_5, c_6, c_7 freely, the c_1, c_2, c_4 are uniquely determined. So we get $|C| = 2^4$ i.e. C is a $[7, 16]$ -code. It has information rate $\frac{\log m}{n} = \frac{4}{7}$. Note $c_{3,5,6,7}$ are free, $c_{1,2,4}$ are check digits.

Suppose we receive $x \in \mathbb{F}_2^7$. We form the syndrome $z_x = (z_1, z_2, z_4)$ where

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

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

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

If $x \in C$, then $z = (0, 0, 0)$. If $d(x, c) = 1$ for some $c \in C$, then the place where x and c differ is given by $z_1 + 2z_2 + 4z_4$ (not modulo 2). Since if $x = c + e_i$ where $e_i = 0 \dots 010 \dots 0$ (1 in the i th place), then the syndrome of x is the syndrome of e_i . For example the syndrome of e_3 is $(1, 1, 0)$, the binary expansion of 3. True in fact for each $1 \leq i \leq 7$.

Remark. Suppose we change our code $C \subset \{0, 1\}^n$ by using the same permutation to reorder each codeword. This gives a code with the same mathematical properties (e.g. information rate, error detection rate). We say such codes are equivalent.

Lemma 2.3. The Hamming distance is a metric on \mathbb{F}_2^n .

Proof. Clearly $d(x, y) \geq 0$ with equality iff $x = y$. Also $d(x, y) = d(y, x)$. Check triangle inequality, let $x, y, z \in \mathbb{F}_2^n$:

$$\begin{aligned} \{1 \leq i \leq n \mid x_i \neq z_i\} &\subseteq \{1 \leq i \leq n \mid x_i \neq y_i\} \cup \{1 \leq i \leq n \mid y_i \neq z_i\} \\ \implies d(x, z) &\leq d(x, y) + d(y, z) \end{aligned} \quad \square$$

Remark. We could also write $d(x, y) = \sum_{i=1}^n d_1(x_i, y_i)$ where d_1 is the discrete metric on $\{0, 1\}$.

Definition (Minimum distance). The **minimum distance** of a code C is the smallest **Hamming distance** between distinct codewords.

Notation. An $[n, m]$ -code with **minimum distance** d is sometimes called an $[n, m, d]$ -code.

Remark.

- $m \leq 2^n$, with equality if $C = \mathbb{F}_2^n$, this is called the trivial code.
- $d \leq n$, with equality in the case of the **repetition code**.

Lemma 2.4. Let C be a code with **minimum distance** d .

- C is $(d - 1)$ -**error detecting**, but cannot detect all sets of errors.
- C is $\lfloor \frac{d-1}{2} \rfloor$ -**error correcting**, but cannot correct all sets of $\lfloor \frac{d-1}{2} \rfloor + 1$ errors.

Proof.

- If $x \in \mathbb{F}_2^n$ and $c \in C$ with $1 \leq d(x, c) \leq d - 1$ then $x \notin C$. So errors are detected. Suppose $c_1, c_2 \in C$ with $d(c_1, c_2) = d$. Then c_1 can be corrupted to c_2 with just d errors, this set of errors will not be detected.
- Let $e = \lfloor \frac{d-1}{2} \rfloor$, so $e \leq \lfloor \frac{d-1}{2} \rfloor < e + 1$, i.e. $2e < d \leq 2(e + 1)$. Let $x \in \mathbb{F}_2^n$. If $\exists c_1 \in C$ with $d(x, c_1) \leq e$, we want to show $d(x, c_2) > e \quad \forall c_2 \in C, c_2 \neq c_1$. By the triangle inequality,

$$\begin{aligned} d(x, c_2) &\geq d(c_1, c_2) - d(x, c_1) \\ &\geq d - e \\ &> e \end{aligned}$$

so C is e -error correcting.

Let $c_1, c_2 \in C$ with $d(c_1, c_2) = d$. Let $x \in \mathbb{F}_2^n$ differ from c_1 in precisely $e + 1$ places, where c_1 and c_2 differ. Then $d(x, c_1) = e + 1$ and $d(x, c_2) = d - (e + 1) \leq e + 1$. So C cannot correct all $e + 1$ errors. \square

Example.

- 1) The **repetition code** is a $[n, 2, n]$ -code, it is $n - 1$ **error detecting** and $\lfloor \frac{n-1}{2} \rfloor$ **error correcting**.
- 2) The simple **parity check** code is a $[n, 2^{n-1}, 2]$ -code, it is 1-error detecting and 0-error correcting.
- 3) **Hamming's original** $[7, 16]$ -code is 1-error correcting $\implies d \geq 3$. Since 0000000 and 1110000 are both valid **codewords**, $d = 3$. That is, it is a $[7, 16, 3]$ -code.

2.1 New codes from old

Let C be an $[n, m, d]$ -code.

i) The parity extension of C is

$$\overline{C} = \left\{ (c_1, c_2, \dots, c_n, \sum c_i) \mid (c_1, \dots, c_n) \in C \right\}.$$

It is a $[n+1, m, d']$ code, for $d' = d$ or $d+1$, depending on d being odd or even.

ii) Fix $1 \leq i \leq n$. Deleting the i th letter from each codeword gives a punctured code. If $d \geq 2$, the new code is $[n-1, m, d'']$ for $d'' = d-1$ or d .

iii) Fix $1 \leq i \leq n$ and $a \in \{0, 1\}$. The shortened code is

$$\{(c_1, \dots, c_{i-1}, c_{i+1}, \dots, c_n) \mid (c_1, \dots, c_{i-1}, a, c_{i+1}, \dots, c_n) \in C\}.$$

It is a $[n-1, m', d']$ -code, where $d' \geq d$ and some choice of a gives $m' \geq \frac{m}{2}$.

2.2 Bound on codes

Definition (Hamming ball). Let $x \in \mathbb{F}_2^n$ and $r \geq 0$. The (closed) **Hamming ball** with centre x and radius r is

$$B(x, r) = \{y \in \mathbb{F}_2^n \mid d(x, y) \leq r\}.$$

Note the ‘volume’

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

which is independent of x .

Lemma 2.5 (Hamming’s bound). If $C \subset \mathbb{F}_2^n$ is e -error correcting, then

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

Proof. Since C is e -error correcting, the Hamming balls $B(c, e)$ are disjoint for $c \in C$. So

$$\begin{aligned} \sum_{c \in C} |B(c, e)| &\leq |\mathbb{F}_2^n| \\ \implies |C| V(n, e) &\leq 2^n \\ \implies |C| &\leq \frac{2^n}{V(n, e)}. \quad \square \end{aligned}$$

Definition (Perfect code). A $[n, m]$ -code which can correct e errors is called **perfect** if

$$m = \frac{2^n}{V(n, e)}.$$

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

Example. Hamming's original $[7, 16, 3]$ -code, can correct 1 error.

$$\frac{2^n}{V(n, e)} = \frac{2^7}{V(7, 1)} = \frac{2^7}{1 + 7} = 2^4 = m$$

so it is a perfect code.

Remark. A perfect e -error correcting code will always incorrectly decode $e + 1$ errors.

Definition (Maximum size of code).

$$A(n, d) = \max \{ m \mid \exists \text{ a } [n, m, d]\text{-code} \}$$

i.e. size of largest code with parameters n and d .

Example. $A(n, 1) = 2^n$, the trivial code $= \mathbb{F}_2^n$. $A(n, n) = 1$, for the repetition code.

Proposition 2.6 (Gilbert-Shannon-Varsharov bound).

$$\frac{2^n}{V(n, d-1)} \underset{\text{GSV}}{\leq} A(n, d) \underset{\text{Hamming}}{\leq} \frac{2^n}{V(n, \lfloor \frac{n-1}{2} \rfloor)}$$

Proof. Let C be a code of length n and minimum distance d of largest possible size. Then $\nexists x \in \mathbb{F}_2^n$ such that $d(x, c) \geq d \quad \forall c \in C$, otherwise we would replace C with $C \cup \{x\}$

$$\begin{aligned} \implies \mathbb{F}_2^n &\subseteq \bigcup_{c \in C} B(c, d-1) \\ \implies 2^n &\leq |C| V(n, d-1) \\ \implies |C| &\geq \frac{2^n}{V(n, d-1)}. \quad \square \end{aligned}$$

Example. Take $n = 10, d = 3$.

$$\begin{aligned} V(n, 1) &= 1 + 10 = 11 \\ V(n, 2) &= 1 + 10 + \binom{10}{2} = 56 \end{aligned}$$

Proposition 2.6 gives $\frac{2^{10}}{56} \leq A(10, 3) \leq \frac{2^{10}}{11}$, so $19 \leq A(10, 3) \leq 93$. The exact value $A(10, 3) = 72$ was only found in 1999.

There exist asymptotic versions of the Gilbert-Shannon-Varsharov bound and Hamming's bound. Let

$$\alpha(\delta) = \limsup \frac{1}{n} \log A(n, \delta n), \quad 0 \leq \delta \leq 1.$$

Notation. $H(\delta) = -\delta \log \delta - (1 - \delta) \log(1 - \delta)$.

The asymptotic GSV says

$$\alpha(\delta) \geq 1 - H(\delta) \quad \text{for } 0 < \delta < \frac{1}{2}$$

while the asymptotic Hamming bound says

$$\alpha(\delta) \leq 1 - H\left(\frac{\delta}{2}\right).$$

We prove the asymptotic GSV bound

Proposition 2.7. Let $0 < \delta < \frac{1}{2}$. Then

$$(i) \quad \log V(n, \lfloor n\delta \rfloor) \leq nH(\delta)$$

$$(ii) \quad \frac{\log A(n, \lfloor n\delta \rfloor)}{n} \geq 1 - H(\delta)$$

Proof (i) \Rightarrow (ii). The [Gilbert-Shannon-Varsharov bound](#) gives

$$\begin{aligned} A(n, \lfloor n\delta \rfloor) &\geq \frac{2^n}{V(n, \lfloor n\delta \rfloor) - 1} \\ &\geq \frac{2^n}{V(n, \lfloor n\delta \rfloor)} \\ \Rightarrow \log A(n, \lfloor n\delta \rfloor) &\geq n - \log V(n, \lfloor n\delta \rfloor) \\ &\geq n - nH(\delta) \quad \text{by (i)} \\ \Rightarrow \frac{\log A(n, \lfloor n\delta \rfloor)}{n} &\geq 1 - H(\delta) \quad \square \end{aligned}$$

Proof of (i).

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

2.3 Channel Capacity

Let $|\Sigma| = q$. A code of length n is a subset of Σ^n (usually we take $q = 2$).

A code is used to send messages through a [discrete memoryless channel](#) with q input letters. For each code a decoding rule is chosen.

Definition (Maximum error probability). The **maximum error probability** is

$$\hat{e}(C) = \max_{c \in C} \mathbb{P}(\text{error} \mid c \text{ sent}).$$

Definition (Reliable transmission). A channel can transmit **reliably at rate** R if there exists a sequence of codes C_1, C_2, \dots where C_n is a code of length n and size $\lfloor 2^{nR} \rfloor$ such that

$$\hat{e}(C_n) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Definition. The (operational) **channel capacity** is the supremum over all **reliable transmission** rates.

Lemma 2.8. Let $\epsilon > 0$. A **BSC** with error probability p is used to send n digits. Then

$$\mathbb{P}(\text{BSC makes } \geq n(p + \epsilon) \text{ errors}) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Proof. Let

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

$\mu_1, \mu_2, \mu_3, \dots$ are i.i.d. random variables.

$$\left. \begin{array}{l} \mathbb{P}(\mu_i = 1) = p \\ \mathbb{P}(\mu_i = 0) = 1 - p \end{array} \right\} \implies \mathbb{E}(\mu_i) = p$$

$$\begin{aligned} \mathbb{P}(\text{BSC makes } \geq n(p + \epsilon) \text{ errors}) &= \mathbb{P}\left(\sum_{i=1}^n \mu_i \geq n(p + \epsilon)\right) \\ &\leq \mathbb{P}\left(\left|\frac{1}{n} \sum \mu_i - p\right| \geq \epsilon\right) \end{aligned}$$

This $\rightarrow 0$ as $n \rightarrow \infty$ by **WLLN**. □

Remark. $\sum_{i=1}^n \mu_i$ is a binomial random variable with parameters n and p .

Proposition 2.9. The **capacity** of a **BSC** with error probability $p < \frac{1}{4}$ is $\neq 0$.

Proof. Choose δ with $2p < \delta < \frac{1}{2}$. We prove **reliably encoding** at rate $R = 1 - H(\delta) > 0$. Let C_n be the largest code of length n and **minimum distance** $\lfloor n\delta \rfloor$.

$$|C_n| = A(n, \lfloor n\delta \rfloor) \geq 2^{n(1-H(\delta))} \quad \text{by Proposition 2.7} \quad = 2^{nR}.$$

Replacing C_n by a subcode gives $|C_n| = \lfloor 2^{nR} \rfloor$ and still minimum distance $\geq \lfloor n\delta \rfloor$. Using **minimum distance decoding**,

$$\hat{e}(C_n) \leq \mathbb{P}(\text{BSC makes } \geq \lfloor \frac{\lfloor n\delta \rfloor - 1}{2} \rfloor \text{ errors}) \tag{1}$$

$$\leq \mathbb{P}(\text{BSC makes } \geq \lfloor \frac{n\delta - 1}{2} \rfloor \text{ errors}). \tag{2}$$

$$\tag{3}$$

Pick $\epsilon > 0$ such that $p + \epsilon < \frac{\delta}{2}$. Then $\frac{n\delta - 1}{2} = n(\frac{\delta}{2} - \frac{1}{2n}) > n(p + \epsilon)$ for n sufficiently large. Therefore, $\hat{e}(C_n) \leq \mathbb{P}(\text{BSC makes } \geq n(p + \epsilon) \text{ errors}) \rightarrow 0$ as $n \rightarrow \infty$ by **Lemma 2.8**. □

2.4 Conditional Entropy

Definition (Conditional Entropy). Let X and Y be random variables taking values in Σ_1 and Σ_2 . We define

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

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

Lemma 2.10. $H(X, Y) = H(X | Y) + H(Y)$

Proof.

$$\begin{aligned} H(X | Y) &= - \sum_{y \in \Sigma_2} \sum_{x \in \Sigma_1} \mathbb{P}(X = x | Y = y) \mathbb{P}(Y = y) \log \mathbb{P}(X = x | Y = y) \\ &= - \sum_{y \in \Sigma_2} \sum_{x \in \Sigma_1} \mathbb{P}(X = x, Y = y) \log \left(\frac{\mathbb{P}(X = x, Y = y)}{\mathbb{P}(Y = y)} \right) \\ &= - \sum_{y \in \Sigma_2} \sum_{x \in \Sigma_1} \mathbb{P}(X = x, Y = y) \log \mathbb{P}(X = x, Y = y) \\ &\quad + \sum_{y \in \Sigma_2} \left(\sum_{x \in \Sigma_1} \mathbb{P}(X = x, Y = y) \right) \log \mathbb{P}(Y = y) \\ &= H(X, Y) - H(Y). \end{aligned}$$

□

Example. A fair six-sided dice is thrown. X is the value on the dice,

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

$$\begin{aligned} H(X, Y) &= H(X) = \log 6 & H(Y) &= \log 2 = 1 \\ H(X | Y) &= H(X, Y) - H(Y) = \log 3 \\ H(Y, X) &= 0. \end{aligned}$$

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

Proof. Since $H(X | Y) = H(X, Y) - H(Y)$, this is equivalent to showing $H(X, Y) \leq H(X) + H(Y)$ with equality iff X and Y are independent. We showed this in [Lemma 1.7](#). □

Notation. In the definition of conditional entropy we can replace random variables X and Y with random vectors $\mathbf{X} = (X_1, \dots, X_r)$ and $\mathbf{Y} = (Y_1, \dots, Y_s)$. This defines

$$H(X_1, \dots, X_r | Y_1, \dots, Y_s) := H(\mathbf{X}, \mathbf{Y}).$$

Lemma 2.11. $H(X | Y) \leq H(X | Y, Z) + H(Z)$.

Proof. We expand $H(X, Y, Z)$ using [Lemma 2.10](#) in two different ways.

$$\begin{aligned} H(X, Y, Z) &= H(Z | X, Y) + H(X | Y) + H(Y) \\ H(X, Y, Z) &= H(X | Y, Z) + H(Z | Y) + H(Y) \end{aligned}$$

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

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

□

Lemma 2.12 (Fano's inequality). Let X, Y be random variables taking values in Σ_1 with $|\Sigma_1| = m$. Let $p = P(X \neq Y)$. Then $H(X | Y) \leq H(p) + p \log(m - 1)$.

Proof. Let

$$Z = \begin{cases} 1 & \text{if } X \neq Y \\ 0 & \text{if } X = Y \end{cases}$$

so $P(Z = 1) = p$ and $P(Z = 0) = 1 - p$. By Lemma 2.11,

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

Now,

$$\begin{aligned} H(X | Y = y, Z = 0) &= 0, & \text{as } Z = 0 \text{ implies } X = Y \\ H(X | Y = y, Z = 1) &\leq \log(m - 1) & \text{since } m - 1 \text{ choices for } X \text{ remain.} \end{aligned}$$

So,

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

Definition (Mutual information). Let X, Y be random variables. The **mutual information** is $I(X; Y) = H(X) - H(X | Y)$ i.e. the amount of information about X conveyed by Y .

By Lemma 1.7 and Lemma 2.10, $I(X; Y) = H(X) + H(Y) - H(X, Y) \geq 0$ with equality iff X and Y are independent. Note the symmetry $I(X; Y) = I(Y; X)$.

Consider a DMC. Let X take values in Σ_1 , where $|\Sigma_1| = m$ with probabilities p_1, \dots, p_m . Let Y be the random variables output when the channel is given input X .

Definition (Information channel capacity). The **information channel capacity** is

$$\max_X I(X; Y).$$

Remark.

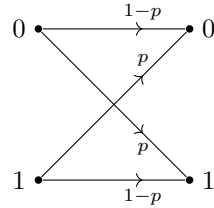
- (i) The maximum is over all choices of p_1, \dots, p_m .
- (ii) The maximum is attained, since we have a continuous function on the compact set $\{(p_1, \dots, p_m) \mid p_i \geq 0, \sum p_i = 1\} \subset \mathbb{R}^n$.
- (iii) The information capacity only depends on the channel matrix.

Theorem 2.13 (Shannon's Second Coding Theorem). **Operational capacity = information capacity.**

We will show \leq in general, \geq for a BSC. We now compute the capacity of certain channels using Theorem 2.13.

Capacity of a Binary Symmetric Channel

Take error probability p .



$$\text{Input } X: \quad \mathbb{P}(X = 0) = 1 - \alpha$$

$$\mathbb{P}(X = 1) = \alpha$$

$$\text{Output } Y: \quad \mathbb{P}(Y = 0) = (1 - \alpha)(1 - p) + \alpha p$$

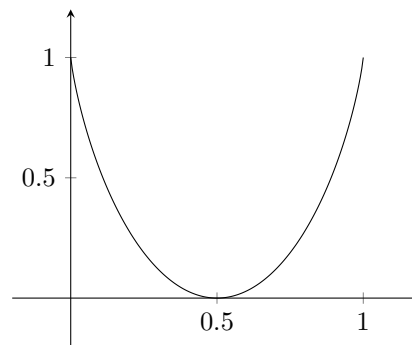
$$\mathbb{P}(Y = 1) = \alpha(1 - p) + (1 - \alpha)p$$

$$\begin{aligned} \text{Capacity is } C &= \max_{\alpha} I(X; Y) \\ &= \max_{\alpha} (H(Y) - H(Y | X)) \\ &= \max_{\alpha} (H(\alpha(1 - p) + (1 - \alpha)p) - H(p)) \end{aligned}$$

$$\text{since } H(Y | X) = \mathbb{P}(X = 0)H(p) + \mathbb{P}(X = 1)H(p) = H(p)$$

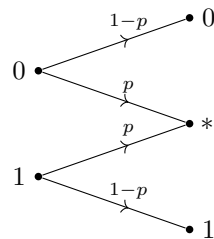
$$= 1 - H(p)$$

$$= 1 + p \log p + (1 - p) \log(1 - p)$$



Capacity of Binary Erasure Channel

Again, take error probability p .



$$\begin{aligned}
\text{Input } X: \quad & \mathbb{P}(X = 0) = 1 - \alpha \\
& \mathbb{P}(X = 1) = \alpha \\
\text{Output } Y: \quad & \mathbb{P}(Y = 0) = (1 - \alpha)(1 - p) \\
& \mathbb{P}(Y = 1) = \alpha(1 - p) \\
& \mathbb{P}(Y = *) = p
\end{aligned}$$

We have $H(X | Y = 0) = 0$, $H(X | Y = 1) = 0$.

$$\begin{aligned}
H(X | Y = *) &= - \sum_x \mathbb{P}(X = x | Y = *) \log \mathbb{P}(X = x | Y = *) \\
\mathbb{P}(X = 0 | Y = *) &= \frac{\mathbb{P}(X = 0, Y = *)}{\mathbb{P}(Y = *)} = \frac{(1 - \alpha)p}{p} = 1 - \alpha.
\end{aligned}$$

Similarly, $\mathbb{P}(X = 1 | Y = *) = \alpha$. So $H(X | Y = *) = H(\alpha)$, so $H(X | Y) = pH(\alpha)$.

$$\begin{aligned}
\text{Capacity is } C &= \max_{\alpha} I(X; Y) \\
&= \max_{\alpha} (H(X) - H(X | Y)) \\
&= \max_{\alpha} (H(\alpha) - pH(\alpha)) \\
&= (1 - p) \max_{\alpha} H(\alpha) \\
&= 1 - p, \quad \text{attained when } \alpha = \frac{1}{2}.
\end{aligned}$$

We model using a channel n times as the n th extension, i.e. replace input and output alphabets Σ_1 and Σ_2 by Σ_1^n and Σ_2^n with channel probabilities

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

Lemma 2.14. The n th extension of a DMC with information capacity C has information capacity nC .

Proof. We take random variable input $(X_1, \dots, X_n) = \mathbf{X}$ producing random variable output $(Y_1, \dots, Y_n) = \mathbf{Y}$. Now, $H(\mathbf{Y} | \mathbf{X}) = \sum_{\mathbf{x}} \mathbb{P}(\mathbf{X} = \mathbf{x}) H(\mathbf{Y} | \mathbf{X} = \mathbf{x})$. Since the channel is memoryless, $H(\mathbf{Y} | \mathbf{X} = \mathbf{x}) = \sum_i H(Y_i | \mathbf{X} = \mathbf{x}) = \sum_i H(Y_i | X_i = x_i)$. So,

$$\begin{aligned}
H(\mathbf{Y} | \mathbf{X}) &= \sum_{\mathbf{x}} \mathbb{P}(\mathbf{X} = \mathbf{x}) \sum_i H(Y_i | X_i = x_i) \\
&= \sum_i \sum_n H(Y_i | X_i = u) \mathbb{P}(X_i = u) \\
&= \sum_i H(Y_i | X_i).
\end{aligned}$$

Now $H(\mathbf{Y}) \leq H(Y_1) + \dots + H(Y_n)$ by Lemma 1.7, thus

$$\begin{aligned} I(X; Y) &= H(\mathbf{Y}) - H(\mathbf{Y} | \mathbf{X}) \\ &\leq \sum_i (H(Y_i) - H(Y_i | X_i)) \\ &= \sum_i I(X_i; Y_i) \\ &\leq nC \quad \text{by definition of information capacity.} \end{aligned}$$

For equality, need Y_1, \dots, Y_n to be independent. This can be achieved by taking X_1, \dots, X_n independent and choosing the distribution such that $I(X_i; Y_i) = c$. \square

Lemma 2.15. Let $P \subset \mathbb{F}_2^n$ be a subset. Then $C = \{x \in \mathbb{F}_2^n \mid p \cdot x = 0 \quad \forall p \in P\}$ is a linear code.

Proof.

- (i) $0 \in C$ since $p \cdot 0 = 0 \quad \forall p \in P$
- (ii) If $x, y \in C$, then $p \cdot (x + y) = p \cdot x + p \cdot y = 0 \implies x + y \in C$.

\square

P is called a set of **parity checks** and C is a **parity check code**.

Definition (Dual code). Let $C \subset \mathbb{F}_2^n$ be a linear code. The **dual code** is

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

This is a code by ??.

Lemma 2.16. $\dim C + \dim C^\perp = n$.

Warning: we can have $C \cap C^\perp \neq \{0\}$.

Proof. $V = \mathbb{F}_2^n$, $V^* = \{\text{linear maps: } V \rightarrow \mathbb{F}_2\}$. Consider

$$\begin{aligned} \phi : V &\longrightarrow V^* \\ x &\longmapsto \theta_x \quad \text{where } \theta_x : y \longmapsto x \cdot y. \end{aligned}$$

ϕ is a linear map.

Suppose $x \in \ker \phi$, then $x \cdot y = 0 \quad \forall y \in V$. Taking $y = e_i = (0, \dots, 0, 1, 0, \dots, 0)$ (with 1 in the i th place) gives $x_i = 0$. So $\ker \phi = \{0\}$. Since $\dim V = \dim V^*$, it follows that ϕ is an isomorphism. Thus

$$\begin{aligned} \theta(C^\perp) &= \{\theta \in V^* \mid \theta(x) = 0 \quad \forall x \in C\} \\ &= \text{'annihilator of } C' = C^\circ \end{aligned}$$

so $\dim C + \dim \phi(C^\perp) = \dim V$, and $\dim C + \dim C^\perp = n$. \square

Corollary. $(C^\perp)^\perp = C$ for any linear code C . In particular, any linear code is a parity check code.

Proof. Let $x \in C$. Then $x \cdot y = 0 \forall y \in C^\perp$, so $x \in (C^\perp)^\perp$, i.e. $C \subseteq (C^\perp)^\perp$. By ?? (twice), $\dim(C^\perp)^\perp = \dim C$, so $C = (C^\perp)^\perp$. \square

Definition. Let C be a (n, k) linear code.

- (i) A **generator matrix** for C is a $k \times n$ matrix whose rows are a basis for C .
- (ii) A **parity check matrix** for C is a generator matrix for C^\perp . It is a $(n - k) \times n$ matrix.

Lemma 2.17. Every (n, k) linear code is equivalent to a linear code with generator matrix $\begin{pmatrix} I_k & B \end{pmatrix}$.

Proof. We can perform row operations:

- swap 2 rows
- add one row to another

(multiplying by a scalar is not useful in \mathbb{F}_2). By Gaussian elimination, we get G , the generator matrix in row echelon form, e.g.

Permuting the columns of G gives an equivalent code, i.e. $l(i) = i \ 1 \leq i \leq k$, i.e.

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

\square

Definition. The **weight** $w(C)$ of a linear code C is the minimum weight of a non-zero codeword.

By ?? this is the same as minimum distance.

Syndrome decoding

Let C be a (n, \hat{k}) -linear code with parity check matrix H . Then $C = \{x \in \mathbb{F}_2^n : Hx = 0\}$ where x is a column vector.

Suppose we receive $y = c + e$ where $c \in C$ is a codeword and $e \in \mathbb{F}_2^n$ is an error. We compute the **syndrome** Hy . Suppose we know C is K -error correcting. Then we tabulate the syndromes He for all $e \in \mathbb{F}_2^n$ with $w(e) \leq K$. If we receive y we search for Hy in our list. If successful we get $Hy = He$ for some $e \in \mathbb{F}_2^n$ with $w(e) \leq K$. We decode y as $c = y - e$. Then $c \in C$ as $He = Hy - He = 0$, and $d(y, c) = w(e) \leq k$.

Recall [Hamming's original code](#):

$$\begin{aligned} c_1 + c_3 + c_5 + c_7 &= 0 \\ c_2 + c_3 + c_6 + c_7 &= 0 \\ c_4 + c_5 + c_6 + c_7 &= 0 \end{aligned}$$

So, $C^\perp = \langle (1 \ 0 \ 1 \ 0 \ 1 \ 0 \ 1) (0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1) (0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 1) \rangle$, giving

$$H = \begin{pmatrix} 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{pmatrix} \quad \text{and} \quad Hy = z = (z_1 \ z_2 \ z_4)$$

In general, we have

Definition (Hamming codes). Let $d \geq 1$, $n = 2^d - 1$. Let H be the $d \times n$ matrix whose columns are the non-zero elements of \mathbb{F}_2^d . The Hamming $(n, n - d)$ linear code is the linear code with parity check matrix H (original is $d = 3$).

Lemma 2.18. Let C be a linear code with parity check matrix H . Then $w(C) = d$ iff

- (i) any $(d - 1)$ -columns of H are linearly independent.
- (ii) some d columns of H are linearly dependent.

Proof. Suppose C has length n . Then $C = \{x \in \mathbb{F}_2^n \mid Hx = 0\}$. If H has columns v_1, \dots, v_n then

$$(x_1, \dots, x_n) \in C \iff \sum_{i=1}^n x_i v_i = 0$$

i.e. codewords are dependence relations between columns. □

Lemma 2.19. The Hamming $(n, n - d)$ linear code has minimum distance $d(C) = 3$ and is a perfect 1-error correcting code.

Proof. Any two columns of H are linearly independent (where H is the parity check matrix of C), but $\exists 3$ that are linearly dependent. Hence $d(C) = 3$ by Lemma 2.18. By Lemma 2.4, C is $\lfloor \frac{3-1}{2} \rfloor = 1$ error correcting.

To be perfect,

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

Here, $n = 2^d - 1$, $e = 1$ so

$$\frac{2^n}{V(n, e)} = \frac{2^n}{1 + 2^d - 1} = 2^{n-d} = |C|. \quad \square$$

2.5 New codes from old

The following construction is specific to linear codes.

Definition (Bar product). Let C_1, C_2 linear codes of length n with $C_2 \subseteq C_1$, i.e. C_2 is a subcode of C_1 . The **bar product** is

$$C_1 \mid C_2 = \{(x \mid x + y) \mid x \in C_1, y \in C_2\}$$

a linear code of length $2n$.

Lemma 2.20. Take C_1, C_2 as above.

- (i) $\text{rank}(C_1 \mid C_2) = \text{rank}(C_1) + \text{rank}(C_2)$
- (ii) $w(C_1 \mid C_2) = \min\{2w(C_1), w(C_2)\}$

Proof. Let x_1, \dots, x_k a basis for C_1 . Let y_1, \dots, y_l be a basis for C_2 . Then $\{(x_i \mid x_i) : 1 \leq i \leq k\} \cup \{(0 \mid y_j)\}$

□