

# Coding and Cryptography

February 12, 2018

## Contents

<b>0</b>	<b>Miscellaneous</b>	<b>3</b>
<b>1</b>	<b>Introduction to communication channels and coding</b>	<b>4</b>
1.1	Noiseless coding . . . . .	5
1.2	Mathematical entropy . . . . .	7
<b>2</b>	<b>Error control codes</b>	<b>14</b>
2.1	Bound on codes . . . . .	17
2.2	Channel Capacity . . . . .	18
2.3	Conditional Entropy . . . . .	19

## 0 Miscellaneous

# 1 Introduction to communication channels and coding

For example, given a message  $M = \text{"Callme!"}$  which we wish to send by email. We first encode it as binary strings using ASCII. So  $f(C) = 1000011$ ,  $f(a) = 1100001$ ,  $f^*(M) = 10000111100001...0100001$ .

The message goes from the source to the receiver after encoded by the source and decoded by the receiver via a channel, where errors could occur. The basic problem is, given a source and a channel (described probabilistically, we aim to design an encoder and a decoder in order to transmit information economically, reliably, and preserving privacy (secretly).

Some examples of each aspect:

*economically*: Morse code, where common letters have shorter codewords;

*reliability*: every book has an ISBN of form  $a_1...a_{10}$  where  $a_i \in \{0, 1, ..., 9\}$  for  $1 \leq i \leq 9$  and  $a_{10} \in \{0, 1, ..., 9, X\}$ , s.t.  $10a_1 + 9a_2 + ... + a_{10} \equiv 0 \pmod{11}$ , where we treat  $X$  as 10. In this way errors can be detected, although not corrected. There is another version of ISBN which is 13 digit;

*preserve privacy* RSA.

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

A channel is determined by the probabilities  $P(y_1, ..., y_k \text{ received} | x_1, ..., x_k \text{ sent})$ .

**Definition.** A *discrete memoryless channel* (DMC) is a channel for which  $P_{ij} = P(b_j \text{ received} | a_i \text{ sent})$  is the same each time the channel is used, and is independent of all past and future. The channel matrix is the  $r \times s$  matrix with entries  $p_{ij}$ . Note the rows sum to 1.

**Example.** (Binary Symmetric Channel, BSC)

BSC has  $\Sigma_1 = \Sigma_2 = \{0, 1\}$ ,  $0 \leq p \leq 1$ . It has channel matrix  $\begin{pmatrix} 1-p & p \\ p & 1-p \end{pmatrix}$ , i.e.  $p$  is the probability symbol is mistransmitted.

**Example.** (Binary Erasure Channel)

$\Sigma_1 = \{0, 1\}$ ,  $\Sigma_2 = \{0, 1, *\}$ ,  $0 \leq p \leq 1$ . Then the channel matrix is  $\begin{pmatrix} 1-p & p & 0 \\ 0 & p & 1-p \end{pmatrix}$ , i.e.  $p$  is the probability that a symbol can't be read.

Informal definition: A channel's capacity is the highest rate at which information can be reliably transmitted over the channel. Here rate means the units of information per unit time (we want that high), and reliably means arbitrarily small error probability.

There are 3 sections:

- 1) Noiseless coding (data compression);
- 2) Error control codes;
- 3) Cryptography.

## 1.1 Noiseless coding

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

If  $x = x_1 \dots x_r$ ,  $y = y_1 \dots y_s$  are strings from  $\Sigma$ , 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 string.

**Definition.** Let  $\Sigma_1, \Sigma_2$  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.** (Greek five code)

$\Sigma_1 = \{\alpha, \beta, \dots, \omega\}$  (24 letters);  $\Sigma_2 = \{1, 2, 3, 4, 5\}$  (more used). Now let  $\alpha \rightarrow 11, \beta \rightarrow 12, \dots, \omega \rightarrow 54$ .

**Example.**  $\Sigma_1 = \{\text{all words in the dictionary}\}$ ,  $\Sigma_2 = \{A, B, \dots, \text{space}\}$ . Then  $f = \text{'spell the word and a space.'}$

We sent a message  $x_1, \dots, x_n \in \Sigma_1^*$  as  $f(x_1)f(x_2)\dots f(x_n) \in \Sigma_2^*$ , i.e. extend  $f$  to  $f^* : \Sigma_1^* \rightarrow \Sigma_2^*$ .

**Definition.** A code  $f$  is *decipherable* if  $f^*$  is injective, i.e. every string from  $\Sigma_2$  arises from at most one message.

Note that  $f$  being injective is not enough. See this example:

**Example.**  $\Sigma_1 = \{1, 2, 3, 4\}$ ,  $\Sigma_2 = \{0, 1\}$ ,  $f(1) = 0$ ,  $f(2) = 1$ ,  $f(3) = 00$ ,  $f(4) = 01$ . Then  $f$  is injective, but  $f^*(312) = 0001 = f^*(114)$  so  $f^*$  is not decipherable.

**Notation.** If  $|\Sigma_1| = m$ ,  $|\Sigma_2| = a$ , then we say  $f$  is an  $a$ -ary code of size  $m$  (in particular, if  $a = 2$  we use the word binary).

Our aim is to construct decipherable codes with short word lengths.

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

- (1) A *Block code* is a code with all codewords of the same length (eg Greek fire code);
- (2) In a *comma code* we reserve one letter from  $\Sigma_2$  that is only used to signal the end of the codeword (example 2);
- (3) 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.** (1) and (2) are special cases of (3).

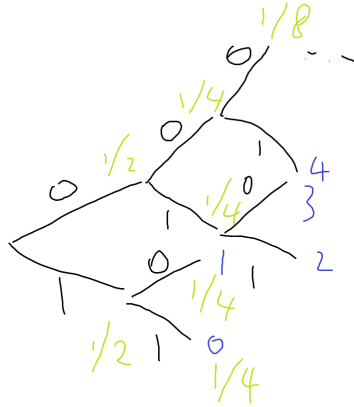
Prefix-free codes are also known as *instantaneous codes* (i.e. a word can be recognised as soon as its complete), or *self-punctuating codes*.

**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  $s_1, \dots, s_m$  exist iff

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

*Proof.* First we prove forward implication. Consider an infinite tree where each has a descendent, labelled by the elements of  $\Sigma_2$ . Each codeword corresponds to a node, the path from the root to this node spelling at the codeword. 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 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  $\leq 1$ .



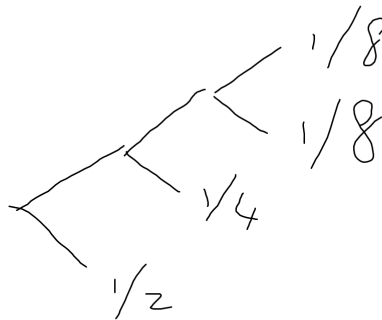
So we can restrict our attention to prefix-free codes.

## 1.2 Mathematical entropy

Entropy is a measure of 'randomness' or 'certainty'. 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.** Suppose  $p_1 = p_2 = p_3 = p_4 = \frac{1}{4}$ . We identify  $\{x_1, \dots, x_4\}$  with  $\{HH, HT, TH, TT\}$ , so  $H(x) = 2$ .

**Example.**  $(p_1, p_2, p_3, p_4) = (1/2, 1/4, 1/8, 1/8)$ . Then  $H(x) = 1/2 + 1/4 \times 2 + 1/8 \times 3 + 1/8 \times 3 = 7/4$ . So the entropy here is smaller.



So in some sense, there is more randomness in the first example than the second.

**Definition.** (Entropy)

The entropy of  $X$ ,

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

where, as most of the time in this course,  $\log = \log_2$ .

**Remark.** (1) If  $p_i = 0$ , we define  $p_i \log p_i = 0$ ;  
(2)  $H(X) \geq 0$ .

**Example.** (3)

We toss a biased coin with  $\mathbb{P}(\text{heads}) = p$ . Write  $H(p) = H(p, 1-p) = -p \log p -$

## 1 INTRODUCTION TO COMMUNICATION CHANNELS AND CODINGS

$(1-p)\log(1-p)$ . If  $p = 0$  or  $1$ , the outcome is certain and so entropy is  $0$ . Entropy is maximal where  $p = 1/2$  (check), a fair coin.

Note the entropy can also be viewed as the expected value of the information of  $X$ , where information is given by  $I(X = x) = -\log \mathbb{P}(X = x)$ . For example, if a coin always lands heads, we gain on information from tossing the coin. This 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.* It suffices to prove the inequality with  $\log$  replaced by  $\ln$ . Note  $\ln x \leq x-1$  with equality iff  $x = 1$ . Let  $I = \{1 \leq i \leq n : p_i \neq 0\}$ . Then  $\ln \frac{q_i}{p_i} \leq \frac{q_i}{p_i} - 1 \forall i \in I$ . So  $\sum_{i \in I} p_i \ln \frac{q_i}{p_i} \leq \sum_{i \in I} q_i - \underbrace{\sum_{i \in I} p_i}_{=1} \leq 0$ . Rearranging we get  $-\sum_{i \in I} p_i \ln p_i \leq -\sum_{i \in I} p_i \ln q_i$ , so  $-\sum_{i=1}^n p_i \ln p_i \leq -\sum_{i=1}^n p_i \ln q_i$ . If equality holds then  $\frac{q_i}{p_i} = 1 \forall i \in I$ , so  $\sum_{i \in I} q_i = 1 \implies p_i = q_i$  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 previous lemma.  $\square$

Two alphabets  $\Sigma_1, \Sigma_2$  with  $|\Sigma_1| = m, |\Sigma_2| = a$  ( $m \geq 2, a \geq 2$ ). We model the source as a sequence of random variables  $X_1, X_2, \dots$  taking values in  $\Sigma_1$ .

**Definition.** A *Bernoulli* or *memoryless* source is a sequence of independently, identically distributed random variables, i.e. for each  $\mu \in \Sigma_i$ , the probability of  $X_i = \mu$  is independent of  $i$  and independent of all past and future symbols emitted. Thus

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

Let  $\Sigma_1 = \{\mu_1, \dots, \mu_n\}$ ,  $p_i = \mathbb{P}(X = \mu_i)$  and assume  $p_i > 0$ . 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.** A code  $f : \Sigma_1 \rightarrow \Sigma_2^*$  is *optimal* if it has the shortest possible expected word length among decipherable code.

**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 and Kraft inequalities. Let  $q_i = \frac{a^{-s_i}}{c}$  where  $c = \sum a^{-s_i} \leq 1$  by Kraft's inequality. Note  $\sum q_i = 1$ . Now

$$\begin{aligned} H(X) &= -\sum p_i \log p_i \leq -\sum_i p_i \log q_i \\ &= \sum_i p_i (s_i \log a + \log c) \\ &= (\sum_i p_i s_i) \log a + \log c \\ &\leq E(s) \log a \end{aligned}$$

We get equality if and only if  $p_i = a^{-s_i}$  for some integers  $s_i$ .

For the upper bound, put  $s_i = \lceil -\log_a p_i \rceil$ . We have  $\log_a p_i \leq s_i < -\log_a p_i + 1$ , so  $a^{-s_i} \leq p_i \implies \sum a^{-s_i} \leq \sum p_i \leq 1$ . So by theorem (1.), there 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.

Shannon-Fano Coding (as in Goldie & Pinch):

This follows from the proof above. 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. The Kraft inequality ensures there is enough room.

**Example.**  $\mu_1, \dots, \mu_5$  emitted with probabilities 0.4, 0.2, 0.2, 0.1, 0.1. We try to construct Shannon-Fano code (with  $a = 2$ ,  $\Sigma_2 = \{0, 1\}$ ):

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

The expected word length is  $2 \times 0.4 + 3 \times 0.2 + 3 \times 0.2 + 4 \times 0.1 + 4 \times 0.1 = 2.8$ . As a comparison,  $H(X) \approx 2.12$ , which is consistent with our previous inequality.

**Definition.** (Huffman coding)

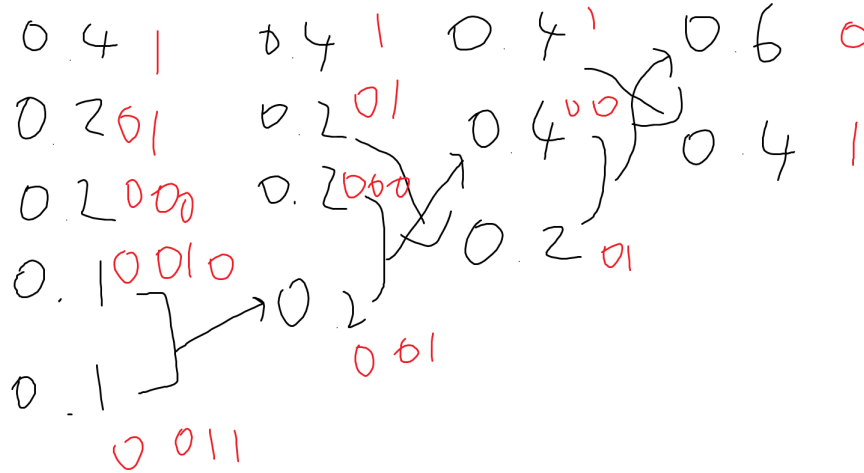
For simplicity we take  $a = 2$ . Let  $\Sigma_1 = \{\mu_1, \dots, \mu_m\}$ ,  $p_i = \mathbb{P}(X = \mu_i)$ . WLOG  $p_1 \geq p_2 \geq \dots \geq p_m$ . Huffman coding is defined inductively:

If  $m = 2$ , assign codewords 0 and 1;

If  $m > 2$ , find a Huffman coding in the case of messages  $\mu_1, \mu_2, \dots, \nu$  with probabilities  $p_1, p_2, \dots, p_{m-1} + p_m$ ; append 0 (1 respectively) to the codeword for  $\nu$  to be a codeword for  $\mu_{m-1}$  ( $\mu_m$  respectively).

**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.** We look at the previous example:



So we get  $\{1, 01, 000, 0010, 0011\}$  as the prefix-free code constructed. The expected word length is 2., which is better than the Shannon-Fano coding.

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

**Lemma.** (1.6)  
Suppose  $\mu_1, \dots, \mu_m \in \Sigma_1$  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) there exists 2 codewords of maximal length which are equal up to the last digit.

*Proof.* (i) Otherwise, swap codewords  $i$  and  $j$  to reduce the expected word length.  
(ii) If not, 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, delete the last digit of each codeword of maximal length. This maintains the prefix free condition, contradicting with  $f$  being optimal.  $\square$

*Proof.* (of 1.5 ( $a = 2$ ))  
We show, by induction on  $m$ , that any Huffman code of size  $m$  is optimal.  
 $m = 2$ : codewords 0, 1 optimal;

$m > 2$ : source  $X_m$  emits  $\mu_1, \dots, \mu_m$  with probabilities  $p_1 \geq p_2 \geq \dots \geq p_m$ . 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  $E(S_m) = E(S_{m-1}) + p_{m-1} + p_m$ . Let  $f'_m$  be an optimal code for  $X_m$ , wlog  $f'_m$  prefix free. Lemma (1.6) shows that shuffling codewords we may assume that the last two codewords of  $f'_m$  are of maximal length and differ only in the last digit, say  $y_0$  and  $y'$  (for some string  $y$ ). We define a code  $f'_{m-1}$  for  $X_{m-1}$  with  $f'_{m-1}(\mu_i) = f'_m(\mu_i) \forall 1 \leq i \leq m-2$ ,  $f'_{m-1}(\nu) = y$ . Then  $f'_{m-1}$  is a prefix free code and the expected word length satisfies

$$E(S'_m) = E(S'_{m-1}) + p_{m-1} + p_m$$

By induction hypothesis,  $f_{m-1}$  is optimal, so  $E(S_{m-1}) \leq E(S'_{m-1}) \implies E(S_m) \leq E(S'_m)$ , so  $f_m$  is optimal.  $\square$

**Remark.** Not all optimal codes are Huffman. For example, if  $p = 0.3, 0.3, 0.2, 0.2$ , we could use code 00, 10, 01, 11 which is not Huffman.

Nevertheless, the previous result says if we have a prefix-free optimal code with word lengths  $s_1, \dots, s_m$  associated with probabilities  $p_1, \dots, p_m$ , there exists a Huffman code with those word lengths.

**Definition.** The *joint entropy* of  $X$  and  $Y$  is

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

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

*Proof.* Let  $\Sigma_1 = \{x_1, \dots, x_m\}$ ,  $\Sigma_2 = \{y_1, \dots, y_n\}$   $p_{ij} = \mathbb{P}(X = x_i, Y = y_j)$ ,  $p_i = \mathbb{P}(X = x_i)$ ,  $q_j = \mathbb{P}(Y = y_j)$ . Apply Gibbs inequality with  $p_{ij}$  and  $p_i q_j$  we get

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

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

Suppose we have a source  $\Omega$  which produces a string  $X_1, X_2, \dots$  of random variables with values in  $\Sigma$ . The probability mass function (pmf) of  $X^{(n)} = (X_1, \dots, X_n)$  is given by

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

Now  $p_n : \Sigma^n \rightarrow \mathbb{R}$  by  $X^{(n)} : \Omega \rightarrow \Sigma^n$ . We can form  $p(X^{(n)} : \Omega \xrightarrow{X^{(n)}} \Sigma^n \xrightarrow{p_n} \mathbb{R}$  by  $w \mapsto p_n(X^{(n)} = X^{(n)}(w))$  a random variable.

For example, let  $\Sigma = \{A, B, C\}$ . Then  $X^{(2)} = AB(0.3), AC(0.1), BA(0.1), BA(0.2), CA(0.25), CB(0.05)$  So  $p_2(AB) = 0.3$ , etc.. And  $p_2(X^{(2)})$  takes values 0.3 with probability 0.3, 0.1 with probability 0.2, 0.2 with probability 0.2, 0.25 with probability 0.25, 0.05 with probability 0.05.

**Definition.** A sequence of random variables  $X_1, \dots$  converges in probability to  $c \in \mathbb{R}$ , written  $X_n \xrightarrow{p} c$  as  $n \rightarrow \infty$ , if  $\forall \varepsilon > 0$ ,  $P(|x_n - c| \leq \varepsilon) \rightarrow 1$  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)

Let  $X_1, X_2, \dots$  be an i.i.d. sequence of random variables with finite expected value  $\mu$ , then

$$\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p} \mu$$

as  $n \rightarrow \infty$ .

Application: Suppose  $X_1, X_2, \dots$  is a Bernoulli source. Then  $p(X_1), p(X_2), \dots$  are i.i.d. random variables, and  $p(X_1, \dots, X_n) = p(X_1) \dots p(X_n)$ . Take log of both sides we get

$$-\frac{1}{n} \log p(X_1, \dots, X_n) = -\frac{1}{n} \sum_{i=1}^n \log p(X_i) \xrightarrow{p} \mathbb{E}(-\log p(X_1)) = H(X_1)$$

as  $n \rightarrow \infty$ .

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

$$-\frac{1}{n} \log p(X_1, \dots, X_n) \xrightarrow{p} H$$

as  $n \rightarrow \infty$ .

Motivating example: suppose we have a coin with  $p(H) = p$ . If coin tossed  $N$  times, expect approximately  $pN$  heads and  $(1-p)N$  tails. The probability of a particular sequence of  $pN$  heads and  $(1-p)N$  tails equals  $p^{pN} (1-p)^{(1-p)N} = 2^{N(p \log p + (1-p) \log(1-p))} = 2^{-NH(A)}$ , where  $A$  is the result of 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  $X_1, X_2, \dots$  satisfies AEP iff it satisfies the following condition (\*):

$\forall \varepsilon < 0, \exists n_0(\varepsilon) \text{ s.t. } \forall n \geq n_0(\varepsilon) \exists \text{ a typical set } T_n \subset \Sigma^n \text{ s.t.}$

(i)  $P((X_1, \dots, X_n) \in T_n) > 1 - \varepsilon$ ;

(ii)  $2^{-n(H+\varepsilon)} \leq p(x_1, \dots, x_n) \leq 2^{-n(H-\varepsilon)}$  for all  $(x_1, \dots, x_n) \in T_n$ .

*Proof.* (sketch)

AEP  $\implies$  (\*):

Take  $T_n = \{(x_1, \dots, x_n) \in \Sigma^n : |-\frac{1}{n} \log p(x_1, \dots, x_n) - H| < \varepsilon\} = \{(x_1, \dots, x_n) \in \Sigma^n : 2^{-n(H+\varepsilon)} \leq p(x_1, \dots, x_n) \leq 2^{-n(H-\varepsilon)}\}$ .

(\*)  $\implies$  AEP:

$P(|-\frac{1}{n} \log p(X_1, \dots, X_n) - H| < \varepsilon) \geq P(T_n) \rightarrow 1$  as  $n \rightarrow \infty$ . □

**Definition.** A source  $X_1, X_2, \dots$  is *reliably encodable* at rate  $r$  if there exists  $A_n \subset \Sigma^n$  for each  $n$  s.t.:

- (i)  $\frac{\log |A_n|}{n} \rightarrow r$  as  $n \rightarrow \infty$ ;
- (ii)  $p((X_1, \dots, X_n) \in A_n) \rightarrow 1$  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, we have  $27 \approx 2^{4.755}$ . So we can encode at a rate of 4.76 bits/letter. But much lower rates suffice, there is a lot of redundancy in the english language. Hence the following definition:

**Definition.** The *information rate*,  $H$ , of a source, is the infimum of all rates at which it is reliably encodable.

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

**Theorem.** (1.9, Shannons' first encoding theorem)

If a source satisfies AEP with same constant  $H$ , then the source has information rate  $H$ .

*Proof.* Let  $\varepsilon > 0$  and let  $T_n \subset \Sigma^n$  be typical sets. Then for sufficiently large  $n \geq n_0(\varepsilon)$ ,  $p(x_1, \dots, x_n) \geq 2^{-n(H+\varepsilon)} \forall (x_1, \dots, x_n) \in T_n$ .

So,  $\mathbb{P}((X_1, \dots, X_n) \in T_n) \geq 2^{-n(H+\varepsilon)} |T_n| \implies |T_n| \leq 2^{n(H+\varepsilon)}$ . So the source is reliably encodable at rate  $H + \varepsilon$ .

Conversely, if  $H = 0$  we are done. Otherwise, pick  $0 < \varepsilon < \frac{H}{2}$ . Suppose the source is reliably encodable at rate  $H - 2\varepsilon$ , say with sets  $A_n$ . Then  $p(x_1, \dots, x_n) \leq 2^{-n(H-\varepsilon)} \forall (x_1, \dots, x_n) \in T_n$ . This implies  $P((x_1, \dots, x_n) \in A_n \cap T_n) \leq 2^{-n(H-\varepsilon)} |A_n|$ , so  $\frac{\log P(A_n \cap T_n)}{n} \leq -(H-\varepsilon) + \frac{\log |A_n|}{n} \rightarrow -(H-\varepsilon) + (H-2\varepsilon) = -\varepsilon$  as  $n \rightarrow \infty$ . (??) So  $\log p(A_n \cap T_n) \rightarrow -\infty$ , i.e.  $p(A_n \cap T_n) \rightarrow 0$ . But  $p(T_n) \rightarrow 1$  as  $n \rightarrow \infty$  and  $P(A_n) \rightarrow 1$  as  $n \rightarrow \infty$ , contradiction(??). So the source is not reliably encodable at rate  $H - 2\varepsilon$ . So the information rate is  $H$ .  $\square$

**Corollary.** A Bernoulli source  $X_1, X_2, \dots$  has information rate  $H(X_1)$ .

## 2 Error control codes

**Definition.** 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.

Note: so this is a block codes.

We use a  $[n, m]$  code to send one of  $m$  possible messages through a BSC making use of the channel  $n$  times. Suppose that we have a probability of  $p$  that a digit get mistransmitted, i.e. 0 becomes 1 or the other way.

**Definition.** The *information rate* of  $C$  is  $\rho(C) = \frac{\log m}{n}$  (as usual,  $\log$  is base 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 large information rate and small error rate. Apparently, these two are contradicting.

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  $\mathbb{P}(c \text{ sent} | x \text{ received})$ ;
- (ii) The maximum likelihood decoding rule decodes  $x \in \{0, 1\}^n$  as  $c \in C$  maximising  $\mathbb{P}(x \text{ received} | c \text{ sent})$ ;
- (iii) The minimum distance decoding rule decodes  $x \in \{0, 1\}^n$  as  $c \in C$  minimising the number of  $\{1 \leq i \leq n : x_i \neq c_i\}$ , i.e. the manhattan distance.

**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 were equally likely, then (i) and (ii) agree.

*Proof.* By Bayes rule,

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

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

**Definition.** Let  $x, y \in \{0, 1\}^n$ . Then *Hamming distance* between  $x$  and  $y$  is  $d(x, y) = |\{1 \leq i \leq n : x_i \neq y_i\}|$ .

**Lemma.** (2.2)

If  $p < \frac{1}{2}$  then (ii) and (iii) agree.

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

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

since  $p < 1/2$ ,  $\frac{p}{1-p} < 1$ . So choosing  $c$  to maximise the probability is the same as choosing  $r$  to minimise the distance  $d(x, c)$ .  $\square$

Note that  $p < \frac{1}{2}$  is really a reasonable assumption (else just revert everything we received).

**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 then we know by minimum distance to decode it as 111. By lemma (2.2), the maximum likelihood would give the same choice as well.

However, for ideal observer, we get  $\mathbb{P}(000 \text{ sent} | 110 \text{ received}) = \frac{3}{4}$  (after some calculation). So ideal observer would decode as 000 (that's why it's ideal).

**Remark.** Ideal observer rule is also known as minimum-error rule. However, it does rely on knowing the probability of the codewords sent.

From now on we use minimum distance decoding.

**Definition.** (i)  $C$  is  $d$ -error detecting if changing at mod  $d$  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.

For example, the repetition code of length  $n$ ,  $C = \{000\dots 0, 111\dots 1\}$  is  $[n, 2]$ -code. It is  $n - 1$  error-detecting, and  $\lfloor \frac{n-1}{2} \rfloor$ -error correcting. But its information rate is only  $\frac{1}{n}$ .

Another example: the simple parity check code of length  $n$  (also known as paper tape code): we view  $\{0, 1\} = \mathbb{Z}/2\mathbb{Z}$ , and  $C = \{(x_1, \dots, x_n) \in \{0, 1\}^n, \sum_{i=1}^n x_i = 0\}$ .

This is a  $[n, 2^{n-1}]$  code. It is 1-error detecting, but cannot correct any errors. Its information rate is  $\frac{n-1}{n}$ .

**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 etc.). We say such codes are equivalent.

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

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

since  $p < 1/2$ ,  $\frac{p}{1-p} < 1$ . So choosing  $c$  to maximise the probability is the same as choosing  $r$  to minimise the distance  $d(x, c)$ .

**Example.** (3, Hamming's original code, 1950)

Let  $C \subseteq F_2^7$  be defined by

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

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

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

since we can choose  $c_3, c_5, c_6, c_7$  freely, then  $c_1, c_2, c_4$  is uniquely determined. We get  $|C| = 2^4$ , so  $C$  is a  $[7, 4]$  code.

Information rate =  $\frac{\log m}{n} = \frac{4}{7}$ .

Suppose we receive  $x \in F_2^7$ . We form the syndrome  $z_x = (z_1, z_2, z_4)$  where

$$\begin{aligned} z_1 &= x_1 + x_3 + x_5 + x_7 \\ z_2 &= x_2 + x_3 + x_6 + x_7 \\ z_4 &= x_4 + x_5 + x_6 + x_7 \end{aligned}$$

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 mod 2). So this is also 1-error correcting.

Since if  $x = c + e_1$  where  $e_i = 0 \dots 010 \dots 0$  where 1 is in the  $i$ th position, then the syndrome of  $x$  is the syndrome of  $e_1$ , and for example syndrome of  $e_3$  is  $(1, 1, 0)$ , the binary expansion of 3. True for each  $1 \leq i \leq 7$ .

Recall  $d(x, y)$  is the number of differences between  $x, y$ .

**Lemma.** (2.3)

The hamming distance is a metric on  $F_2^n$ .

Don't really want to copy the proof – this is obvious.

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

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

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

Note:  $m \leq 2^n$ , with equality if  $C = F_2^n$ , this is called the trivial code. Also  $d \leq n$ , with equality in the case of the repetition code.

**Lemma.** (2.4)

Let  $C$  be a code with minimum distance  $d$ .

(i)  $C$  is  $(d - 1)$ -error detecting, but cannot detect all sets of  $d$  errors.

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

*Proof.* (i) If  $x \in F_2^n$  and  $c \in C$  with  $1 \leq d(x, c) \leq d - 1$ , then  $x \notin C$ , so errors detected. The second part is obvious.

(ii) Let  $e = \lfloor \frac{d-1}{2} \rfloor$ , so  $e \leq \frac{d-1}{2} < e + 1$ , i.e.  $2e < d \leq 2(e + 1)$ . Let  $x \in F_2^n$ . If  $\exists c_1 \in C$  with  $d(x, c_1) \leq e$ , we want to show  $d(x, c_2) > e \forall c_2 \in C, c_2 \neq c_1$ , which can be done by triangle inequality. So  $C$  is  $e$ -error correcting. The second part is trivial as well.  $\square$

**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, i.e.  $d \geq 3$ . Since 0000000 and 1110000 are both in the code we get  $d = 3$ , i.e.  $[7, 16, 3]$ -code.



New codes from old:

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

(i) The parity extension of  $C$  is

$$\bar{C} = \{(c, \dots, c_n, \sum c_i) : (c_1, \dots, c_n) \in C\}$$

It is a  $[n+1, m, d']$ -code where  $d \leq d' \leq d+1$ .

(ii) Fix  $1 \leq i \leq n$ . Deleting the  $i^{\text{th}}$  letter from each codeword gives a punctured code. Assuming  $d \geq 2$ ,  $[n-1, m, d'']$  where  $d-1 \leq d'' \leq d$ .

(iii) Fix  $1 \leq i \leq n$  and  $a \in \{0, 1\}$ . The shortened code is  $\{(c_1, \dots, c_{i-1}, c_{i+1}, \dots, c_n) : (c_1, \dots, c_{i-1}, a, c_{i+1}, \dots, c_n) \in C\}$ . It is a  $[n-1, m', d']$ -code where  $d' \geq d$  and some choice of  $a$  gives  $m' \geq \frac{m}{2}$ .

## 2.1 Bound on codes

**Definition.** Let  $X \in F_2^n$  and  $r \geq 0$ . The (closed) Hamming ball with centre  $x$  and radius  $r$  is just what we think it will be (obvious). We denote it by  $B(x, r)$ . Note that the volume  $V(n, r) = |B(x, r)| = \sum_{i=0}^r \binom{n}{i}$  is independent of  $x$ .

**Lemma.** (2.5, Hamming's bound)

If  $C \subset F_2^n$  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 done.  $\square$

**Definition.** A  $[n, m]$ -code which can correct  $e$ -errors is called *perfect* if  $m = \frac{2^n}{V(n, e)}$ .

Note that 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. We can check that it is a perfect code.

Note, however, that a perfect  $e$ -error correcting code will always incorrectly decode  $e+1$  errors.

**Definition.**  $A(n, d) = \max\{m : \exists [n, m, d]\text{-code}\}$ , i.e. size of largest code with parameters  $n$  and  $d$ . For example,  $A(n, 1) = 2^n$ ,  $A(n, n) = 2$ .

**Proposition.** (2.6)

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

The lower and upper bounds have alternative names: GSV-bounds and Hamming's bound.

*Proof.* We've already done the upper bound. For the lower bound, 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  $d(x, c) \geq d \forall c \in C$ , otherwise we can place  $C$  with  $C \cup \{x\}$ . So  $\mathbb{F}_2^n \subseteq \cup_{c \in C} B(c, d-1)$ , i.e.  $|C| \geq \frac{2^n}{V(n, d-1)}$ .  $\square$

**Example.** Consider  $n = 10, d = 3$ . Then  $V(n, 1) = 1 + 10 = 11$ .  $V(n, 2) = 1 + 10 + \binom{10}{2} = 56$ . 2.6 gives  $19 \leq A(10, 3) \leq 93$ . The exact value is 72, only known in 1999.

There exist asymptotic versions of GSV and Hamming's bound: Let  $\alpha(\delta) = \limsup \frac{1}{n} \log A(n, \delta n)$ ,  $0 \leq \delta \leq 1$ .

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

Asymptotic GSV bound:  $\alpha(\delta) \geq 1 - H(\delta)$  for  $0 < \delta < 1/2$ ;

Asymptotic Hamming:  $\alpha(\delta) \leq 1 - H(\delta/2)$ .

We'll prove the asymptotic GSV bound.

**Proposition.** (2.7)

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

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

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

*Proof.* (i)  $\implies$  (ii): GSV bound implies  $A(n, \lfloor n\delta \rfloor) \geq \frac{2^n}{V(n, \lfloor n\delta \rfloor)}$ , so  $\log A(n, \lfloor n\delta \rfloor) \geq n - \log V(n, \lfloor n\delta \rfloor) \geq n - nH(\delta)$  by (i). So  $\frac{\log A(n, \lfloor n\delta \rfloor)}{n} \geq 1 - H(\delta)$ ;

Proof of (i):

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

Taking log, we get

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

$\square$

## 2.2 Channel Capacity

Let  $|\Sigma| = q$ . A code of length  $n$  is a subset of  $\Sigma^n$  (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.

We define  $\hat{e}(C) = \max c \in C \mathbb{P}(\text{error} | c \text{ sent})$ , the maximum error probability.

**Definition.** A channel can transmit reliably at rate  $R$  if there exist 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  $\hat{e}(C_n) \rightarrow 0$ .

**Lemma.** (2.8)

Let  $\varepsilon > 0$ . A BSC with error probability  $p$  is used to send  $n$  digits. Then the probability of the BSC making at least  $n(p + \varepsilon)$  errors  $\rightarrow 0$  as  $n \rightarrow \infty$ .

*Proof.* Let  $\mu_i = 1$  if digit mistransmitted, 0 otherwise. Then  $\mu_1, \dots$  are iid rvs. Also  $P(\mu_i = 1) = p$ , so  $E(\mu_i) = p$ . So the require probability is  $P(\sum_{i=1}^n \mu_i \geq n(p + \varepsilon)) \leq \mathbb{P}(|\frac{1}{n} \sum \mu_i - p| \geq \varepsilon) \rightarrow 0$  as  $n \rightarrow \infty$  by WLLN.  $\square$

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

**Proposition.** (2.9)

The capacity of a BSC with error probability  $p < 1/4$  is not zero.

*Proof.* Choose  $\delta$  with  $2p < \delta < 1/2$ . We prove reliably encoding at rate  $R = 1 - H(\delta) > 0$ . Let  $C_n$  be the largest code of length  $n$  and minimum distance  $\lfloor n\delta \rfloor$ . So  $|C_n| = A(n, \lfloor n\delta \rfloor) \geq 2^{n(1-H(\delta))} = 2^{nR}$  by (2.7). 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,

$$\begin{aligned} \hat{e}(C_n) &\leq P(\text{BSC makes } \geq \lfloor \frac{n\delta - 1}{2} \rfloor + 1 \text{ errors}) \\ &\leq P(\text{BSC makes } \geq \frac{n\delta - 1}{2} \text{ errors}) \end{aligned}$$

Pick  $\varepsilon > 0$  s.t.  $p + \varepsilon < \frac{\delta}{2}$ . Then  $\frac{n\delta - 1}{2} = n(\frac{\delta}{2} - \frac{1}{2n}) > n(p + \varepsilon)$  for  $n$  sufficiently large. So  $\hat{e}(C_n) \leq P(\text{BSC makes } \geq n(p + \varepsilon) \text{ errors}) \rightarrow 0$  as  $n \rightarrow \infty$  by (2.8).  $\square$

## 2.3 Conditional Entropy

Let  $X$  and  $Y$  be random variables taking values in  $\Sigma_1$  and  $\Sigma_2$ . We define

$$\begin{aligned} H(X|Y = y) &= - \sum_{x \in \Sigma_1} P(X = x|Y = y) \log P(X = x|Y = y), \\ H(X|y) &= \sum_{y \in \Sigma_2} P(Y = y) H(X|Y = y) \end{aligned}$$

**Lemma.** (2.10)

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

*Proof.*

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

□

Example: A fair six sided dice is thrown.  $X$  is the value on the dice,  $y = X \bmod 2$ . So  $H(X, Y) = H(X) = \log 6$ ,  $H(Y) = \log 2 = 1$ ,  $H(X|Y) = H(X, Y) - H(Y) = \log 3$ , and  $H(Y|X) = 0$ .

**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 " $H(X, Y) \leq H(X) + H(Y)$  with equality iff  $X$  and  $Y$  independent, which is true by (1.7). □

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

**Lemma.** (2.11)

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

*Proof.* We expand  $H(X, Y, Z)$  using (2.10) in two different ways:  $H(X, Y, Z) = H(Z|X, Y) + H(X|Y) + H(Y)$  and  $H(X, Y, Z) = H(X|Y, Z) + H(Z|Y) + H(Y)$ . Since  $H(Z|X, Y) \geq 0$ , we get  $H(X|Y) \leq H(X|Y, Z) + H(Z|Y) \leq H(X|Y, Z) + H(Z)$  by corollary. □

**Lemma.** (2.12, Fano's inequality)

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

*Proof.* Let  $Z = 1$  if  $X \neq Y$ , and 0 otherwise. Then  $P(Z = 1) = p$ . So by (2.11),

$$H(X|Y) \leq H(X|Y, Z) + \underbrace{H(Z)}_{=H(p)} (*)$$

Now  $H(X|Y = y, Z = 0) = 0$  (must have  $X = y$ ),  $H(X|Y = y, Z = 1) \leq \log(m - 1)$  since  $m - 1$  choices for  $X$  remain. So,  $H(X|Y, Z) = \sum_{y,z} P(Y = y, Z = z) H(X|Y = y, Z = z) \leq \sum_y P(Y = y, Z = 1) \log(m - 1)$ , where the sum is just  $P(Z = 1) = p$  and by (\*) we get the result. □

**Definition.** 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 (1.7) and (2.10), we have  $I(X, Y) = H(X) + H(Y) - H(X, Y) \geq 0$  is symmetric in  $X$  and  $Y$ . We get equality iff  $X$  and  $Y$  are independent.

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

**Definition.** The *information channel capacity* is  $\max_X I(x; y)$ .

**Remark.** The max is over all choices of  $p_1, \dots, p_m$ ; max is always attained since  $I$  is continuous on a compact set.

The information capacity only depends on the channel matrix.

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

We'll show  $\leq$  in general, and  $\geq$  for a BSC.

We now compute the capacity of certain channels using Shannon's second coding theorem (error probability  $p$ ):

Input  $X$ ,  $P(X = 0) = 1 - \alpha$ ,  $P(X = 1) = \alpha$ , output  $Y$ ,  $P(Y = 0) = (1 - \alpha)(1 - p) + \alpha p$ ,  $P(Y = 1) = \alpha(1 - p) + (1 - \alpha)p$ .

Capacity is

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

Here we denote  $H(p) = H(p, 1 - p)$ .

Capacity of Binary Erasure Channel (erasure probability  $p$ ):

In this model we have each bit having a probability of  $p$  being erased to become a  $*$ .

Input  $X$ ,  $P(X = 0) = 1 - \alpha$ ,  $P(X = 1) = \alpha$ ; output  $Y$ ,  $P(Y = 0) = (1 - \alpha)(1 - p)$ ,  $P(Y = 1) = \alpha(1 - p)$ ,  $P(Y = *) = p$ . Now  $H(X|Y = 0) = 0$ ,  $H(X|Y = 1) = 0$ .

The only interesting case is

$$\begin{aligned} H(X|Y = *) &= - \sum_x P(X = x|Y = *) \log P(X = x|Y = *) \\ &= H(\alpha) \end{aligned}$$

where we work out

$$P(X = 0|Y = *) = 1 - \alpha, P(X = 1|Y = *) = \alpha$$

Capacity is

$$\begin{aligned}
C &= \max_{\alpha} I(X, Y) \\
&= \max_{\alpha} (H(X) - H(X|Y)) \\
&= \max_{\alpha} (H(\alpha) - pH(\alpha)) \\
&= (1 - p) \max_{\alpha} H(\alpha) \\
&= 1 - p
\end{aligned}$$

attained when  $\alpha = 1/2$ .

We model using a channel  $n$  times as the  $n^{th}$  extension, i.e. we replace input and output alphabets  $\Sigma_1$  and  $\Sigma_2$  by  $\Sigma_1^n$  and  $\Sigma_2^n$ . Channel probabilities:  $P(y_1 \dots y_n \text{ received} | x_1 \dots x_n \text{ sent}) = \prod_{i=1}^n P(y_i \text{ received} | x_i \text{ sent})$ .

**Lemma.** (2.14)

The  $n^{th}$  extension of a DMC with information capacity  $C$  has information capacity  $nC$ .

*Proof.* We take r.v. input  $X_1, \dots, X_n = X$  producing r.v. output  $Y_1, \dots, Y_n = Y$ . Now

$$H(Y|X) = \sum_x P(X = x) H(Y|X = x)$$

Since channel is memoryless,

$$H(Y|X = x) = \sum_i H(Y_i|X = x) = \sum_i H(Y_i|X_i = x_i)$$

So

$$\begin{aligned}
H(Y|X) &= \sum_x P(X = x) \sum_i H(Y_i|X_i = x_i) \\
&= \sum_i \sum_n H(Y_i|X_i = n) P(X_i = n) \\
&= \sum_i H(X_i|Y_i)
\end{aligned}$$

Now  $H(Y) \leq H(Y_1) + \dots + H(Y_n)$  by (1.7). So  $I(X, Y) = H(Y) - H(Y|X) \leq \sum_i (H(Y_i) - H(Y_i|X_i)) = \sum_{i=1}^n I(X_i, Y_i) \leq nC$  by definition of info capacity.

For equality, we 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 s.t.  $I(X_i, Y_i) = C$ .  $\square$

**Proposition.** (2.15)

For a DMC, operational capacity  $\leq$  information capacity.

*Proof.* Let  $C$  be the information capacity. Suppose we can transmit reliably at rate  $R > C$ , i.e. there is a sequence of codes  $(C_n)_{n \geq 1}$  with  $C_n$  of length  $n$  and size  $\lfloor 2^{nR} \rfloor$  such that  $\hat{e}(C_n) \rightarrow 0$  as  $n \rightarrow \infty$ . Then  $\hat{e}(C_n) = \max_{c \in C_n} P(\text{error} | c \text{ sent})$ ,  $e(C_n) = \frac{1}{|C_n|} \sum_{c \in C_n} P(\text{error} | c \text{ sent})$ . Clearly  $e(C_n) \leq \hat{e}(C_n)$ , so  $e(C_n) \rightarrow 0$  as  $n \rightarrow \infty$ . Take r.v. input  $X$ , equidistributed over  $C_n$ . Let  $Y$  be the r.v. output

when  $X$  is transmitted and decoded. So  $e(C_n) = P(x \neq y) = p$  say.

Now  $H(X) = \log(|C_n|) \geq nR - 1$  for  $n$  sufficiently large,  $H(X|Y) \leq H(p) + p \log(|C_n| - 1) \leq 1 + pnR$  (Fano's inequality),  $I(X, Y) = H(X) - H(X|Y)$ ,  $nC \geq nR - 1 - (1 + pnR)$  (2.14), so  $pnR \geq n(R - C) - 2$ . So we get  $p \geq \frac{n(R - C) - 2}{nR} \not\rightarrow 0$  as  $n \rightarrow \infty$ .

Thus our sequence of codes cannot exist.  $\square$

**Proposition.** (2.16)

Consider a BSC, error probability  $p$ . Let  $R < 1 - H(p)$ . Then there exists a sequence of codes  $(C_n)_{n \geq 1}$  of length  $n$  and size  $\lfloor 2^{nR} \rfloor$  such that  $e(C_n) \rightarrow 0$  as  $n \rightarrow \infty$ .

*Proof.* The idea is to construct codes by picking codewords at random. WLOG let  $p < 1/2$ , so  $\exists \varepsilon > 0$  s.t.  $R < 1 - H(p + \varepsilon)$ . We use minimum distance decoding (in case of tie, make arbitrary choice).

Let  $m = \lfloor 2^{nR} \rfloor$ . We pick a  $[n, m]$ -code  $C$  at random (i.e. each with probability  $\frac{1}{\binom{n}{m}}$ ), say  $C = \{c_1, \dots, c_m\}$ . Choose  $1 \leq i \leq m$  at random (i.e. each with probability  $1/m$ ). We send  $c_1$  through the channel and get output  $Y$ .

Then  $P(Y \text{ is not decoded as } c)$  is the average value of  $e(C)$  as  $C$  runs over all  $[n, m]$ -codes. We can pick  $C_n$  a  $[n, m]$ -code with  $e(C_n)$  at most this average. So it will suffice to show that

$$P(Y \text{ is not decoded as } c_i) \rightarrow 0 \text{ as } n \rightarrow \infty$$

Let  $r = \lfloor n(p + \varepsilon) \rfloor$ . Then the above probability  $\leq P(C_i \notin B(Y, r)) + P(B(Y, r) \cap C \supsetneq \{c_i\})$ , i.e. either  $c_i$  is not in the ball, or it is in the ball but is not the only one. We consider the two probabilities separately:

(i)  $P(d(c_i, y) > r) = P(\text{BSC makes } > r \text{ errors}) = P(\text{BSC makes } > n(p + \varepsilon) \text{ errors}) \rightarrow 0$  as  $n \rightarrow \infty$  by WLLN.

(ii) if  $j \neq i$ ,  $P(c_j \in B(Y, r) | c_i \in B(Y, r)) = \frac{V(n, r) - 1}{2^n - 1} \leq \frac{V(n, r)}{2^n}$ . So,

$$\begin{aligned} P(B(Y, r) \cap C \supsetneq \{c_i\}) &\leq (m - 1) \frac{V(n, r)}{2^n} \\ &\leq 2^{nR} 2^{nH(p + \varepsilon)} 2^{-n} \\ &= 2^{n(R - (1 - H(p + \varepsilon)))} \rightarrow 0 \end{aligned}$$

as  $n \rightarrow \infty$ , since  $R < 1 - H(p + \varepsilon)$ .  $\square$

**Proposition.** (2.17)

Consider a BSC with error probability  $p$ . Let  $R < 1 - H(p)$ . Then there exists a sequence of codes  $(C_n)_{n \geq 1}$  with  $C_n$  of length  $n$ , size  $\lfloor 2^{nR} \rfloor$  and  $\hat{e}(C_n) \rightarrow 0$  as  $n \rightarrow \infty$ .

*Proof.* Pick  $R'$  s.t.  $R < R' < 1 - H(p)$ . By (2.16), we construct a sequence of codes  $(C'_n)_{n \geq 1}$  with  $C'_n$  of length  $n$ , size  $\lfloor 2^{nR'} \rfloor$  and  $e(C'_n) \rightarrow 0$  as  $n \rightarrow \infty$ .

Throwing at the worst half of the codewords in  $C'_n$ , gives a code  $C_n$  with  $\hat{e}(C_n) \leq 2e(C'_n)$ . So  $\hat{e}(C_n) \rightarrow 0$  as  $n \rightarrow \infty$ .

Note  $C_n$  has length  $n$  and size  $\lfloor 2^{nR-1} \rfloor$ , but  $2^{nR-1} = 2^{n(R' - \lfloor \text{frac} 1n \rfloor)} \geq 2^{nR}$  for  $n$  sufficiently large.

We can replace  $C_n$  by a subcode of size  $\lfloor 2^{n+R} \rfloor$  and still get  $\hat{e}(C_n) \rightarrow 0$  as  $n \rightarrow \infty$ .  $\square$

Conclusion: a BSC with error probability  $p$  has operational capacity  $1 - H(p)$ .

Remark: (i) How does it work? say capacity is 0.8, and we have a message a string of 0's and 1's. Let  $R = 0.75 < 0.8$ , Then  $\exists$  a set of  $2^{0.75n}$  codewords of length  $n$  that have error probability below some prescribed threshold, hence to encode message:

- (a) break message into block's of size  $3\lceil \frac{n}{4} \rceil = m$  sufficiently large;
- (b) encode these  $m$ -blocks into  $C_n$  using codewords of length  $\frac{4}{3}m$  each  $m$  block;
- (c) transmit new message through channel.

The theorem shows good codes exists. But the proof does not construct them for us.