# ECS452 2016/2    Part I.1    Dr.Prapun

## 1 Elements of a Digital Communication System

**1.1.** Figure 1 illustrates the functional diagram and the basic elements of a digital communication system.
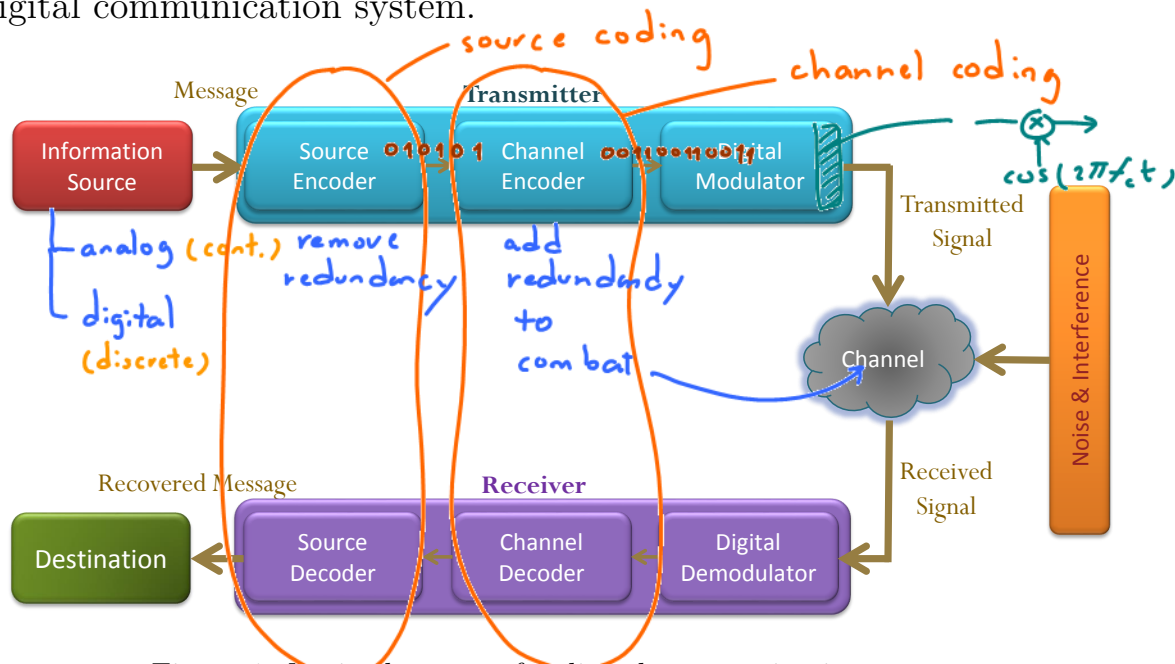


Figure 1: Basic elements of a digital communication system

**1.2.** The source output may be either

- an **analog** signal, such as an audio or video signal,

or

- a **digital** signal, such as the output of a computer, that is discrete in time and has a finite number of output characters.

4

In ECS 332, we talked about the conversion from analog signal to bits
① PCM→ via the process of sampling and quantization ⟨ uniform
② DPCM                                                   non-uniform
③ DM

**1.3.** Source Coding: The messages produced by the source are converted into a sequence of bits.

- This process is called source coding or source encoding.

- For this course, we want to also represent the source output (message) by as few bits as possible.

  ○ In other words, we seek an efficient representation of the source output that results in little or no redundancy.

  ○ Therefore, source coding may be referred to as **data compression**. (Huffman coding)

**1.4. Channel Coding**:

- Introduce, in a controlled manner, some *redundancy* in the binary information sequence that can be used at the receiver to overcome the effects of noise and interference encountered in the transmission of the signal through the channel.

  ○ The added redundancy serves to increase the reliability of the received data and improves the fidelity of the received signal.

- See Examples 1.5 and 1.6.

**Example 1.5.** Trivial channel coding: Repeat each binary digit $n$ times, where $n$ is some positive integer.

**Example 1.6.** More sophisticated channel coding: Taking $k$ information bits at a time and mapping each $k$-bit sequence into a unique $n$-bit sequence, called a **codeword**.

- The amount of redundancy introduced by encoding the data in this manner is measured by the ratio $n/k$. The reciprocal of this ratio, namely $k/n$, is called the rate of the code or, simply, the **code rate**.
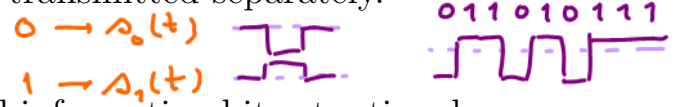
**1.7.** The binary sequence at the output of the channel encoder is passed to the **digital modulator**, which serves as the interface to the physical (analog) communication channel.

- Since nearly all the communication channels encountered in practice are capable of transmitting electrical signals (waveforms), the primary purpose of the digital modulator is to map the binary information sequence into signal waveforms.

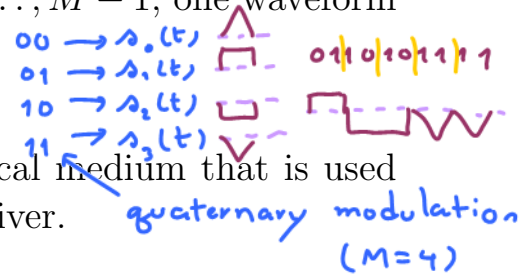In ECS332, we studied line coding.

5

- The digital modulator may simply map the binary digit 0 into a waveform $s_0(t)$ and the binary digit 1 into a waveform $s_1(t)$. In this manner, each bit from the channel encoder is transmitted separately. We call this **binary modulation**.

- The modulator may transmit $b$ coded information bits at a time by using $M = 2^b$ distinct waveforms $s_i(t), i = 0, 1, \ldots, M - 1$, one waveform for each of the $2^b$ possible $b$-bit sequences. We call this $M$-**ary modulation** ($M > 2$).

**1.8.** The **communication channel** is the physical medium that is used to send the signal from the transmitter to the receiver.

- The physical channel may be

  ○ a pair of wires that carry the electrical signal, or

  ○ an optical fiber that carries the information on a modulated light beam, or

  ○ an underwater ocean channel in which the information is transmitted acoustically, or

  ○ free space over which the information-bearing signal is radiated by use of an antenna.

  Other media that can be characterized as communication channels are data storage media, such as magnetic tape, magnetic disks, and optical disks.

- Whatever the physical medium used for transmission of the information, the essential feature is that the transmitted signal is corrupted in a random manner by a variety of possible mechanisms, such as additive **thermal noise** generated by electronic devices; man-made noise, e.g., automobile ignition noise; and **atmospheric noise**, e.g., electrical lightning discharges during thunderstorms.

- Other channel impairments including noise, attenuation, distortion, fading, and interference (such as interference from other users of the channel).

**1.9.** At the receiving end of a digital communication system, the digital demodulator processes the channel-corrupted transmitted waveform and reduces the waveforms to a sequence of numbers that represent estimates of the transmitted data symbols (binary or $M$-ary).

This sequence of numbers is passed to the channel decoder, which attempts to reconstruct the original information sequence from knowledge of the code used by the channel encoder and the redundancy contained in the received data.

- A measure of how well the demodulator and decoder perform is the frequency with which errors occur in the decoded sequence. More precisely, the average probability of a bit-error at the output of the decoder is a measure of the performance of the demodulator-decoder combination.

- In general, the probability of error is a function of the code characteristics, the types of waveforms used to transmit the information over the channel, the transmitter power, the characteristics of the channel (i.e., the amount of noise, the nature of the interference), and the method of demodulation and decoding.

**1.10.** As a final step, when an analog output is desired, the source decoder accepts the output sequence from the channel decoder and, from knowledge of the source encoding method used, attempts to reconstruct the original signal from the source.

- Because of channel decoding errors and possible distortion introduced by the source encoder, and perhaps, the source decoder, the signal at the output of the source decoder is an approximation to the original source output.

- The difference or some function of the difference between the original signal and the reconstructed signal is a measure of the distortion introduced by the digital communication system.

distortion

*Lossless*

# 2  Source Coding

In this section, we look at the "source encoder" part of the system. This part removes redundancy from the message stream or sequence. We will focus only on binary source coding.

*so that each transmitted bit (symbol) carries max information*

**2.1.** The material in this section is based on [C & T Ch 2, 4, and 5].

*Here, assume that the message is already available in digital format.*

## 2.1  General Concepts

**Example 2.2.** Suppose your message is a paragraph of (written natural) text in English.

- Approximately 100 possibilities for characters/symbols. *source symbols*

  - For example, a character-encoding scheme called (ASCII) (American Standard Code for Information Interchange) originally[1] had 128 specified characters – the numbers 0–9, the letters a–z and A–Z, some basic punctuation symbols[2], and a blank space.

- Do we need 7 bits per characters?

**2.3.** A sentence of English–or of any other language–always has more information than you need to decipher it. The meaning of a message can remain unchanged even though parts of it are removed.

**Example 2.4.**

- "J-st tr- t- r–d th-s s-nt-nc-." [3]

- "Thanks to the redundancy of language, yxx cxn xndxrstxnd whxt x xm wrxtxng xvxn xf x rxplxcx xll thx vxwxls wxth xn 'x' (t gts lttl hrdr f y dn't vn kn whr th vwls r)." [4]

---

[1]Being American, it didn't originally support accented letters, nor any currency symbols other than the dollar. More advanced Unicode system was established in 1991.

[2]There are also some control codes that originated with Teletype machines. In fact, among the 128 characters, 33 are non-printing control characters (many now obsolete) that affect how text and space are processed and 95 printable characters, including the space.

[3]Charles Seife, Decoding the Universe. Penguin, 2007

[4]Steven Pinker, The Language Instinct: How the Mind Creates Language. William Morrow, 1994

**2.5.** It is estimated that we may only need about 1 bits per character in English text.

*Simplified model of the message source (to make the analysis easier)*

**Definition 2.6. Discrete Memoryless Sources (DMS)**: Let us be more specific about the information source. $\boxed{\text{source}} \longmapsto X_1, X_2, X_3, \ldots$

- The message that the information source produces can be represented by a **vector** of characters $X_1, X_2, \ldots, X_n$.

  ○ A perpetual message source would produce a never-ending **sequence** of characters $X_1, X_2, \ldots$.

- These $X_k$'s are random variables (at least from the perspective of the decoder; otherwise, these is no need for communication).

- For simplicity, we will assume our source to be discrete and memoryless.

  ○ Assuming a discrete source means that the random variables are all discrete; that is, they have supports which are countable. Recall that "countable" means "finite" or "countably infinite".

    * We will further assume that they all share the same support and that the support is finite. This support is called the source alphabet.

  ○ Assuming a memoryless source means that there is no dependency among the characters in the sequence.

    * More specifically,

    $$p_{X_1, X_2, \ldots, X_n}(x_1, x_2, \ldots, x_n) = p_{X_1}(x_1) \times p_{X_2}(x_2) \times \cdots \times p_{X_n}(x_n).$$
    (1)

    * This means the current model of the source is far from the source of normal English text. English text has dependency among the characters. However, this simple model provide a good starting point.

    * We will further assume that all of the random variables share the same probability mass function (pmf). We denote this shared pmf by $p_X(x)$. In which case, (1) becomes
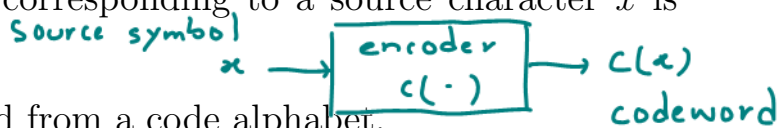
    $$p_{X_1, X_2, \ldots, X_n}(x_1, x_2, \ldots, x_n) = p_X(x_1) \times p_X(x_2) \times \cdots \times p_X(x_n). \quad (2)$$

*All RVs are i.i.d.*

· We will also assume that the pmf $p_X(x)$ is known. In practice, there is an extra step of estimating this $p_X(x)$.

· To save space, we may see the pmf $p_X(x)$ written simply as $p(x)$, i.e. without the subscript part.

* The shared support of $X$ which is usually denoted by $S_X$ becomes the source alphabet. Note that we also often see the use of $\mathcal{X}$ to denote the support of $X$.

- In conclusion, our simplified source code can be characterized by a random variable $X$. So, we only need to specify its pmf $p_X(x)$.

**Example 2.7.** See slides.

**Definition 2.8.** An **encoder** $c(\cdot)$ is a function that maps each of the character in the (source) alphabet into a corresponding (binary) codeword.

- In particular, the codeword corresponding to a source character $x$ is denoted by $c(x)$.

  Source symbol
  $x \longrightarrow$ | encoder $c(\cdot)$ | $\longrightarrow c(x)$
  codeword

- Each codeword is constructed from a code alphabet.

  ○ A binary codeword is constructed from a two-symbol alphabet, wherein the two symbols are usually taken as 0 and 1.

  ○ It is possible to consider non-binary codeword. Morse code discussed in Example 2.13 is one such example.

- Mathematically, we write

  source alphabet

  $$\text{Encoder } c : S_X \rightarrow \{0,1\}^*$$

  code alphabet

  where

  $$\{0,1\}^* = \{\varepsilon, 0, 1, 00, 01, 10, 11, 000, 001, 010, 011, \ldots\}$$

  is the set of all finite-length binary strings.

  message                                         encoded string
  | source | $\longrightarrow x_1\ x_2\ x_3 \cdots \longrightarrow$ | encoder $c(\cdot)$ | $\longrightarrow c(x_1)\ c(x_2)\ c(x_3) \cdots$
  DMS
  i.i.d. $\sim p_X(x)$

- The length of the codeword associated with source character $x$ is denoted by $\ell(x)$.

  codeword length

10

      ○ In fact, writing this as $\ell(c(x))$ may be clearer because we can see that the length depends on the choice of the encoder. However, we shall follow the notation above[5].

**Example 2.9.** $c(\text{red}) = 00$, $c(\text{blue}) = 11$ is a source code for $S_X = \{\text{red, blue}\}$.

**Example 2.10.** Suppose the message is a sequence of basic English words which happen according to the probabilities provided in the table below.

*Naive coding*

| $x$ | $p(x)$ | Codeword $c(x)$ | $\ell(x)$ |
|-----|--------|-----------------|-----------|
| Yes | 4% | 00 | 2 |
| No | 3% | 01 | 2 |
| OK | 90% | 10 | 2 |
| Thank You | 3% | 11 | 2 |

$\mathbb{E}[\ell(X)] = 2 \times 0.04 + 2 \times 0.03$
$+ 2 \times 0.9 + 2 \times 0.03$
$= 2(0.04 + 0.03 + 0.9 + 0.03)$
$= 2 \times 1 = 2$ bits

**Definition 2.11.** The <mark>expected length</mark> of a code $c(\cdot)$ for (a DMS source which is characterized by) a random variable $X$ with probability mass function $p_X(x)$ is given by

*Expected code length*

$$\mathbb{E}[\ell(X)] = \sum_{x \in S_X} p_X(x)\ell(x).$$

**Example 2.12.** Back to Example 2.10. Consider a new encoder:

| $x$ | $p(x)$ | Codeword $c(x)$ | $\ell(x)$ |
|-----|--------|-----------------|-----------|
| Yes | 4% | 01 | 2 |
| No | 3% | 001 | 3 |
| OK | 90% | 1 | 1 |
| Thank You | 3% | 0001 | 4 |

$\mathbb{E}[\ell(X)] = 2 \times 0.04 + 3 \times 0.03$
$+ 1 \times 0.9 + 4 \times 0.03$
$= \frac{1}{100}(8 + 9 + 90 + 12)$
$= 1.19$ bits

Observe the following:

- Data compression can be achieved by <mark>assigning short descriptions to the most frequent outcomes of the data source,</mark> and necessarily longer descriptions to the less frequent outcomes.

- When we calculate the expected length, we don't really use the fact that the source alphabet is the set $\{\text{Yes, No, OK, Thank You}\}$. We

---

[5]which is used by the standard textbooks in information theory.

would get the same answer if it is replaced by the set $\{1, 2, 3, 4\}$, or the set $\{a, b, c, d\}$. All that matters is that the alphabet size is 4, and the corresponding probabilities are $\{0.04, 0.03, 0.9, 0.03\}$.

Therefore, for brevity, we often find DMS source defined only by its alphabet size and the list of probabilities.

**Example 2.13.** The Morse code is a reasonably efficient code for the English alphabet using an alphabet of four symbols: a dot, a dash, a letter space, and a word space. [See Slides]

- Short sequences represent frequent letters (e.g., a single dot represents E) and long sequences represent infrequent letters (e.g., Q is represented by "dash,dash,dot,dash").

**Example 2.14.** Thought experiment: Let's consider the following code

| $x$ | $p(x)$ | Codeword $c(x)$ | $\ell(x)$ |
|---|---|---|---|
| 1 | 4% | 0 | 1 |
| 2 | 3% | 1 | 1 |
| 3 | 90% | 0 | 1 |
| 4 | 3% | 1 | 1 |

$$\mathbb{E}[\ell(x)] = 1 \quad \text{bit}$$

Not non-singular
= singular

This code is bad because we have ambiguity at the decoder. When a codeword "0" is received, we don't know whether to decode it as the source symbol "1" or the source symbol "3". If we want to have lossless source coding, this ambiguity is not allowed.

**Definition 2.15.** A code is <mark>nonsingular</mark> if <mark>every source symbol in the source alphabet has different codeword.</mark>

Equivalently, for any $x_1, x_2 \in S_x$, if $x_1 \neq x_2$, then $c(x_1) \neq c(x_2)$.

As seen from Example 2.14, nonsingularity is an important concept. However, it turns out that this property is not enough.

**Example 2.16.** Another thought experiment: Let's consider the following code

Observe that this code is non-singular.

| $x$ | $p(x)$ | Codeword $c(x)$ | $\ell(x)$ |
|---|---|---|---|
| 1 | 4% | 01 | 2 |
| 2 | 3% | 010 | 3 |
| 3 | 90% | 0 | 1 |
| 4 | 3% | 10 | 2 |

$$\mathbb{E}[\ell(x)] = 2 \times 0.04 + 3 \times 0.03 + 1 \times 0.9 + 2 \times 0.03$$
$$= 1.13 \quad \text{bits}$$

Problem:
When "010" is received at the decoder, how can we know whether it is

Not prefix-free ⟹ Not prefix code

Not UD

"010" or
"0" followed by "10" or
"01" followed by "0"

12

**2.17.** We usually wish to convey a sequence (string) of source symbols. So, we will need to consider **concatenation** of codewords; that is, if our source string is

$$X_1, X_2, X_3, \ldots$$

then the corresponding encoded string is

$$c(X_1)c(X_2)c(X_3)\cdots.$$

In such cases, to ensure decodability, we may

  (a) use fixed-length code (as in Example 2.10), or

  (b) use variable-length code and

    (i) add a special symbol (a "comma" or a "space") between any two codewords

$$0 \;\; \cancel{*} \;\; 10 \;\; \cancel{*} \;\; 010 \;\; \cancel{*} \;\; \ldots$$

    or

    (ii) use uniquely decodable codes.

**Definition 2.18.** A code is called **uniquely decodable** (UD) if any encoded string has only one possible source string producing it.

**Example 2.19.** The code used in Example 2.16 is not uniquely decodable because source string "2", source string "34", and source string "13" share the same code string "010".

**2.20.** It may not be easy to check unique decodability of a code. (See Example 2.28.) Also, even when a code is uniquely decodable, one may have to look at the entire string to determine even the first symbol in the corresponding source string. Therefore, we focus on a subset of uniquely decodable codes called prefix code.

*prefix-free code*

**Definition 2.21.** A code is called a **prefix code** if no codeword is a prefix[6] of any other codeword.

- Equivalently, a code is called a **prefix code** if you can put all the codewords into a binary tree where all of them are leaves.
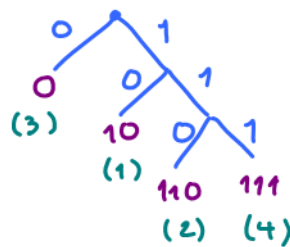
---

[6]String $s_1$ is a prefix of string $s_2$ if there exist a string $s_3$, possibly empty, such that $s_2 = s_1 s_3$.

*upside-down tree*

13

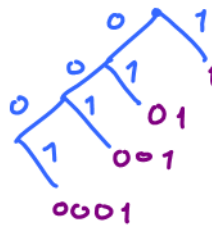- A more appropriate name would be "**prefix-free**" code.

**Example 2.22.**

| $x$ | Codeword $c(x)$ |
|---|---|
| 1 | 10 |
| 2 | 110 |
| 3 | 0 |
| 4 | 111 |



**Example 2.23.** The code used in Example 2.12 <u>is</u> a prefix code.

| $x$ | Codeword $c(x)$ |
|---|---|
| 1 | 01 |
| 2 | 001 |
| 3 | 1 |
| 4 | 0001 |



**2.24.** <mark>Any prefix code is uniquely decodable.</mark>

- <mark>The end of a codeword is immediately recognizable.</mark>

- Each source symbol <mark>can be decoded as soon as we come to the end of the codeword</mark> corresponding to it. In particular, we need not wait to see the codewords that come later.

- Therefore, another name for "prefix code" is <mark>**instantaneous code**</mark>.

**Example 2.25.** The code used in Example 2.12 (, Example 2.22, and Example 2.23) is a prefix code and hence it is uniquely decodable.

**2.26.** The nesting relationship among all the types of source codes is shown in Figure 2.

**Example 2.27.**

| $x$ | Codeword $c(x)$ |
|---|---|
| 1 | 1 |
| 2 | 10 |
| 3 | 100 |
| 4 | 1000 |

Non-singular: Yes

Prefix-free: No

$c(1)$ is a prefix of $c(2)$

UD: Yes.

Observe that a "1" signifies the beginning of every codeword.

Problem: Need to look at the first bit of the next codeword to determine the end of the current one } ⇒ not instantaneous.
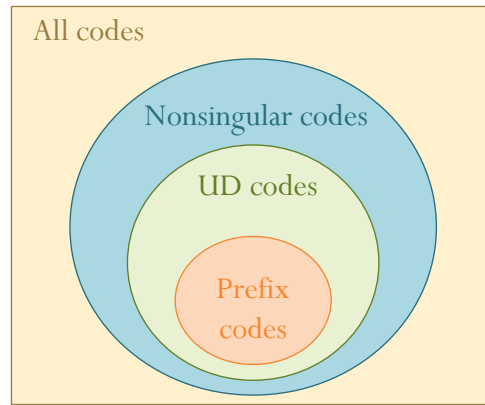
100/1000/1/10/100/1/1/1

14

Figure 2: Classes of codes

**Example 2.28.** [5, p 106–107]

| $x$ | Codeword $c(x)$ |
|---|---|
| 1 | 10 |
| 2 | 00 |
| 3 | 11 |
| 4 | 110 |

This code is not a prefix code because codeword "11" is a prefix of codeword "110".

This code is uniquely decodable. To see that it is uniquely decodable, take any code string and start from the beginning.

- If the first two bits are 00 or 10, they can be decoded immediately.

- If the first two bits are 11, we must look at the following bit(s).

  ○ If the next bit is a 1, the first source symbol is a 3.

  ○ If the next bit is a 0, we need to count how many 0s are there before 1 shows up again.

  ○ If the length of the string of 0's immediately following the 11 is even, the first source symbol is a 3.

  ○ If the length of the string of 0's immediately following the 11 is odd, the first codeword must be 110 and the first source symbol must be 4.

By repeating this argument, we can see that this code is uniquely decodable.

**Sirindhorn International Institute of Technology**

**Thammasat University**

School of Information, Computer and Communication Technology

# ECS452 2016/2  Part I.2  Dr.Prapun

**2.29.** For our present purposes, a better code is one that is uniquely decodable and has a shorter expected length than other uniquely decodable codes. We do not consider other issues of encoding/decoding complexity or of the relative advantages of block codes or variable length codes. [6, p 57]

## 2.2 Optimal Source Coding: Huffman Coding

In this section we describe a very popular source coding algorithm called the Huffman coding.

**Definition 2.30.** Given a source with known probabilities of occurrence for symbols in its alphabet, to construct a binary Huffman code, create a binary tree by repeatedly combining[7] the probabilities of the two least likely symbols.

*Steps*
*① Combine the two least-likely (combined) symbols.*
*② Repeat ① until you can not repeat more.*

*Once we have a binary tree, we know this is a prefix code.*

- Developed by David Huffman as part of a class assignment[8].

---

[7]The Huffman algorithm performs *repeated source reduction* [6, p 63]:

- At each step, two source symbols are combined into a new symbol, having a probability that is the sum of the probabilities of the two symbols being replaced, and the new reduced source now has one fewer symbol.

- At each step, the two symbols to combine into a new symbol have the two lowest probabilities.

  ○ If there are more than two such symbols, select any two.

[8]The class was the first ever in the area of information theory and was taught by Robert Fano at MIT in 1951.

  ○ Huffman wrote a term paper in lieu of taking a final examination.

  ○ It should be noted that in the late 1940s, Fano himself (and independently, also Claude Shannon) had developed a similar, but suboptimal, algorithm known today as the ShannonFano method. The difference between the two algorithms is that the ShannonFano code tree is built from the top down, while the Huffman code tree is constructed from the bottom up.

- By construction, Huffman code is a prefix code.

**Example 2.31.**

| $x$ | $p_X(x)$ | Codeword $c(x)$ | $\ell(x)$ |
|---|---|---|---|
| 1 | 0.5 | O | 1 |
| 2 | 0.25 | 1 O | 2 |
| 3 | 0.125 | 1 1 O | 3 |
| 4 | 0.125 | 1 1 1 | 3 |

(annotations left of table: 1/2, 1/4, 1/8, 1/8)

(tree annotations: 0.25, 0.5, 0, 1 labels)

$$\mathbb{E}\left[\ell(X)\right] = \underbrace{1\times 0.5 + 2\times 0.25}_{1} + \underbrace{3\times 0.125 + 3\times 0.125}_{3/4} = 1.75 \ \text{bits}$$

Note that for this particular example, the values of $2^{\ell(x)}$ from the Huffman encoding is inversely proportional to $p_X(x)$:

$$p_X(x) = \frac{1}{2^{\ell(x)}}.$$

In other words,

$$\ell(x) = \log_2 \frac{1}{p_X(x)} = -\log_2(p_X(x)).$$

Therefore,

$$\mathbb{E}\left[\ell(X)\right] = \sum_x p_X(x)\ell(x) = -\sum_x p_X(a)\log_2 p_X(a) \equiv H(X)$$

$H(x) = -\sum_x p_X(a)\log_2 p_X(a)$

$= -\frac{1}{2}\log_2\frac{1}{2} - \frac{1}{4}\log_2\frac{1}{4}$

$-\frac{1}{8}\log_2\frac{1}{8} - \frac{1}{8}\log_2\frac{1}{8}$

$= \frac{1}{2} + \frac{1}{2} + \frac{3}{4}$

$= 1 + \frac{3}{4} = 1.75$

$[\text{bits}]$

**Example 2.32.**

| $x$ | $p_X(x)$ | Codeword $c(x)$ | $\ell(x)$ |
|---|---|---|---|
| 'a' | 0.4 | O | 1 |
| 'b' | 0.3 | 1 O | 2 |
| 'c' | 0.1 | 1 1 O | 3 |
| 'd' | 0.1 | 1 1 1 O | 4 |
| 'e' | 0.06 | 1 1 1 1 O | 5 |
| 'f' | 0.04 | 1 1 1 1 1 | 5 |

(tree annotations: 0.1, 0.2, 0.3, 0.6, with 0/1 labels)

$$\mathbb{E}\left[\ell(X)\right] = 0.4\times 1 + 0.3\times 2 + 0.1\times 3 + 0.1\times 4 + 0.06\times 5 + 0.04\times 5$$

$$= 0.4 + 0.6 + 0.3 + 0.4 + 0.5 = 2.2 \ \text{bits}$$

$H(x) = 2.1435$

17

**Example 2.33.**

| $x$ | $p_X(x)$ | Codeword $c(x)$ | $\ell(x)$ |
|-----|----------|-----------------|-----------|
| 1 | 0.25 | 10 | 2 |
| 2 | 0.25 | 00 | 2 |
| 3 | 0.2 | 01 | 2 |
| 4 | 0.15 | 110 | 3 |
| 5 | 0.15 | 111 | 3 |

$\mathbb{E}\left[\ell(X)\right] = 2.3$ bits

**Example 2.34.**

| $x$ | $p_X(x)$ | Codeword $c(x)$ | $\ell(x)$ |
|-----|----------|-----------------|-----------|
| a | 1/3 | 00 | 2 |
| b | 1/3 | 01 | 2 |
| c | 1/4 | 10 | 2 |
| d | 1/12 | 11 | 2 |

$\mathbb{E}\left[\ell(X)\right] = \mathbb{E}\left[2\right] = 2$ bits

| $x$ | $p_X(x)$ | Codeword $c(x)$ | $\ell(x)$ |
|-----|----------|-----------------|-----------|
| a | 1/3 | 0 | 1 |
| b | 1/3 | 10 | 2 |
| c | 1/4 | 110 | 3 |
| d | 1/12 | 111 | 3 |

$\mathbb{E}\left[\ell(X)\right] = 1 \times \frac{1}{3} + 2 \times \frac{1}{3} + 3 \times \frac{1}{4} + 3 \times \frac{1}{12} = 2$ bits.
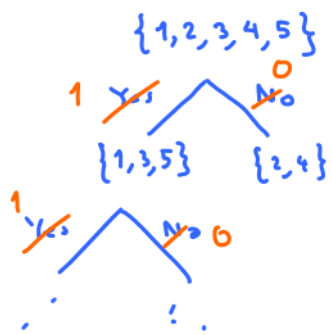
**2.35.** The set of codeword lengths for Huffman encoding is not unique. There may be more than one set of lengths but all of them will give the same value of expected length.

**Definition 2.36.** A code is **optimal** for a given source (with known pmf) if it is uniquely decodable and its corresponding expected length is the shortest among all possible uniquely decodable codes for that source.

Theorem: **2.37.** The Huffman code is optimal.

It is important to remember that there are many optimal codes: starting from an optimal code, inverting all the bits or exchanging two codewords of the same length will give another optimal code.

18

$\{1,2,3,4,5\}$

1 Yes     0 No     odd?

$\{1,3,5\}$     $\{2,4\}$

            0 1 1

1 Yes     No 0     prime?

## 2.3 Source Extension (Extension Coding)

**2.38.** One can usually (not always) do better in terms of expected length (per source symbol) by encoding blocks of several source symbols.

**Definition 2.39.** In, an **$n$-th extension** coding, $n$ successive source symbols are grouped into blocks and the encoder operates on the blocks rather than on individual symbols. [4, p. 777]

**Example 2.40.**

| $x$ | $p_X(x)$ | Codeword $c(x)$ | $\ell(x)$ |
|---|---|---|---|
| Y(es) | 0.9 | O | 1 |
| N(o) | 0.1 | 1 | 1 |

(between the rows: hand-drawn tree with 0, 1, 1 branches)

$H(X) = -0.9\log_2 0.9 - 0.1\log_2 0.1$

$= 0.4685 < 1$

(a) First-order extension:

$\mathbb{E}\left[\ell(X)\right] = \mathbb{E}[1] = 1$

0 1 1 0 0 0 1 0 0 1 1 1
YNNYYYYNYYNNN...
10 110 0 110 10 111

→ independent

For DMS, $X_1 \perp\!\!\!\perp X_2$

(b) Second-order Extension:

| $x_1 x_2$ | $p_{X_1,X_2}(x_1,x_2) = p_{X_1}(x_1)\, p_{X_2}(x_2)$ | $c(x_1,x_2)$ | $\ell(x_1,x_2)$ |
|---|---|---|---|
| YY | $0.9 \times 0.9 = 0.81$ | O | 1 |
| YN | $0.9 \times 0.1 = 0.09$ | 10 | 2 |
| NY | $0.1 \times 0.9 = 0.09$ | 110 | 3 |
| NN | $0.1 \times 0.1 = 0.01$ | 111 | 3 |

(hand-drawn tree with 0.19, 0.1, 0, 1 labels)

$\mathbb{E}\left[\ell(X_1,X_2)\right] = 0.81 \times 1 + 0.09 \times (2+3) + 0.01 \times 3 = 1.29$ bits per 2 source symbols

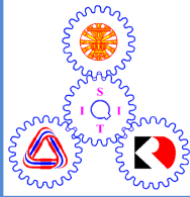$\underbrace{\phantom{0.09 \times}}_{0.45} \quad \underbrace{\phantom{0.01}}_{0.03}$

$L_2 = 0.645$ bits

(c) Third-order Extension: per source symbol

| $x_1 x_2 x_3$ | $p_{X_1,X_2,X_3}(x_1,x_2,x_3)$ | $c(x_1,x_2,x_3)$ | $\ell(x_1,x_2,x_3)$ |
|---|---|---|---|
| YYY | $0.9 \times 0.9 \times 0.9 = 0.729$ | | |
| YYN | $0.9 \times 0.9 \times 0.1 = 0.081$ | | |
| YNY | | | |
| $\vdots$ | | | |

[HW]

$\mathbb{E}\left[\ell(X_1,X_2,X_3)\right] = 1.5980$ bits per 3 source symbols

$L_3 = 0.5327$ bits per source symbol.

# ECS452 2016/2　　Part I.3　　Dr.Prapun

## 2.4　(Shannon) Entropy for Discrete Random Variables

Entropy is a measure of uncertainty of a random variable [5, p 13].

*Entropy quantifies the amount of uncertainty a RV has.*
*measures　　　　　　　　　　　　(randomness)*

It arises as the answer to a number of natural questions. One such question that will be important for us is "What is the average length of the shortest *description* of the random variable?"

**Definition 2.41.** The **entropy** $H(X)$ of a discrete random variable $X$ is defined by

*negative sign*

*Recall:* $\log_2 a = \dfrac{\ln(a)}{\ln(2)}$

*Recall*
$\mathbb{E}[g(X)]$
$= \sum_x g(x) p_x(x)$

$$H(X) = -\sum_{x \in S_X} p_X(x) \log_2 p_X(x) = -\mathbb{E}[\log_2 p_X(X)].$$

*can be omitted*

$= \mathbb{E}[i(X)]$

$i(x) = -\log_2 p_X(x)$ = the amount of info. associated with each realized value of $X$.

- The log is to the base 2 and entropy is expressed in bits (per symbol).

  ○ The base of the logarithm used in defining $H$ can be chosen to be any convenient real number $b > 1$ but if $b \neq 2$ the unit will not be in bits.

  ○ If the base of the logarithm is $e$, the entropy is measured in nats.

  ○ Unless otherwise specified, base 2 is our default base.

  $0 \log_2 0 = 0$

- Based on continuity arguments, we shall assume that $0 \ln 0 = 0.$

  *L'Hôspital's*
  *rule*

20

$$\lim_{x \to 0^+} x \ln x = \lim_{x \to 0^+} \frac{\ln x}{\frac{1}{x}} = \lim_{x \to 0^+} \frac{\frac{1}{x}}{-\frac{1}{x^2}} = \lim_{x \to 0^+} (-x) = 0$$

**Example 2.42.** The entropy of the random variable $X$ in Example 2.31 is 1.75 bits (per symbol).

**Example 2.43.** The entropy of a fair coin toss is 1 bit (per toss).

The probabilities involved are $\frac{1}{2}, \frac{1}{2}$

$$H\left(\left[\frac{1}{2} \; \frac{1}{2}\right]\right) = H(X) = -\frac{1}{2}\log_2\frac{1}{2} - \frac{1}{2}\log_2\frac{1}{2} = \frac{1}{2} + \frac{1}{2} = 1 \quad \text{bit}$$

Three ways for denoting entropy

**2.44.** Note that entropy is a functional of the (unordered) probabilities from the pmf of $X$. It does not depend on the actual values taken by the random variable $X$, Therefore, sometimes, we write ②$H(p_X)$ instead of ①$H(X)$ to emphasize this fact. Moreover, because we use only the probability values, we can use the row vector representation $\underline{p}$ of the pmf $p_X$ and simply express the entropy as $H(\underline{p})$③

In MATLAB, to calculate $H(X)$, we may define a row vector pX from the pmf $p_X$. Then, the value of the entropy is given by

$$\texttt{HX = -pX*(log2(pX))'.}$$

**Example 2.45.** The entropy of a uniform (discrete) random variable $X$ on $\{1, 2, 3, \ldots, n\}$:

$$P_X(x) = \begin{cases} 1/n, & x = 1, 2, \ldots, n, \\ 0, & \text{otherwise.} \end{cases}$$

$$\underline{p} = \left[\frac{1}{n} \; \frac{1}{n} \; \cdots \; \frac{1}{n}\right]$$
$$\underbrace{\phantom{\frac{1}{n} \; \frac{1}{n} \; \cdots \; \frac{1}{n}}}_{n \text{ values}}$$

$$H(X) = -\sum_x P_X(x)\log_2 P_X(x)$$
$$= -\sum_{x=1}^{n} \frac{1}{n}\log_2\frac{1}{n} = -n\underbrace{A}_{A} = -\cancel{n}\frac{1}{\cancel{n}}\log_2\frac{1}{n}$$
$$= \log_2 n = \log_2\left|S_X\right|$$

**Example 2.46.** The entropy of a Bernoulli random variable $X$:

$$P_X(x) = \begin{cases} p, & x = 1, \\ 1-p, & x = 0, \\ 0, & \text{otherwise} \end{cases}$$

$$\underline{p} = [1-p \;\; p]$$

Binary RV

$$H(X) = -\sum_x P_X(x)\log_2 P_X(x)$$
$$= -p\log_2 p - (1-p)\log_2(1-p)$$

Ex. $p = 0.3 \Rightarrow H(X) = 0.8813$

Ex. $p = 0.5 \Rightarrow H(X) = 1$

$$P_X(x) = \begin{cases} p, & x = a, \\ 1-p, & x = b, \\ 0, & \text{otherwise.} \end{cases}$$

21

**Definition 2.47.** Binary Entropy Function : We define $h_b(p)$, $h(p)$ or $H(p)$ to be $-p \log p - (1-p) \log (1-p)$, whose plot is shown in Figure 3.
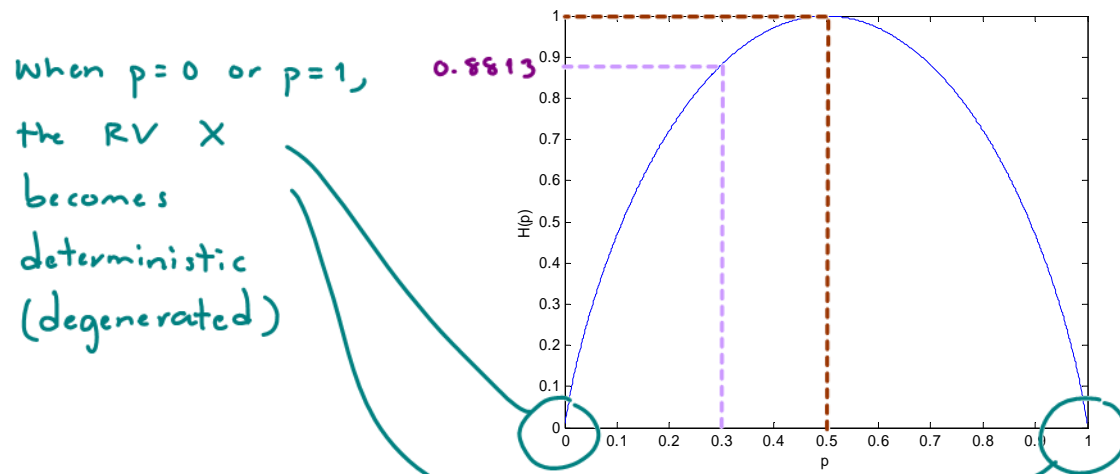
*when p = 0 or p = 1, the RV X becomes deterministic (degenerated)*

*0.8813*



Figure 3: Binary Entropy Function

**2.48.** Two important facts about entropy:

(a) $H(X) \leq \log_2 |S_X|$ with equality if and only if $X$ is a uniform random variable.

(b) $H(X) \geq 0$ with equality if and only if $X$ is not random.

In summary,

$$\underbrace{0}_{\text{deterministic}} \leq H(X) \leq \underbrace{\log_2 |S_X|}_{\text{uniform}}.$$

**Theorem 2.49.** The expected length $\mathbb{E}[\ell(X)]$ of any uniquely decodable binary code for a random variable $X$ is greater than or equal to the entropy $H(X)$; that is,

$$\mathbb{E}[\ell(X)] \geq H(X)$$

with equality if and only if $2^{-\ell(x)} = p_X(x)$. [5, Thm. 5.3.1]

**Definition 2.50.** Let $L(c, X)$ be the expected codeword length when random variable $X$ is encoded by code $c$.

Let $L^*(X)$ be the minimum possible expected codeword length when random variable $X$ is encoded by a uniquely decodable code $c$:

$$L^*(X) = \min_{\text{UD } c} L(c, X).$$

*Expected length of the optimal UD code*

22

**2.51.** Given a random variable $X$, let $c_{\text{Huffman}}$ be the Huffman code for this $X$. Then, from the optimality of Huffman code mentioned in 2.37,

*Expected length of the optimal UD code is the same as the expected length of Huffman code.*

$$L^*(X) = L(c_{\text{Huffman}}, X).$$

**Theorem 2.52.** The optimal code for a random variable $X$ has an expected length less than $H(X) + 1$:

$$L^*(X) < H(X) + 1.$$

**2.53.** Combining Theorem 2.49 and Theorem 2.52, we have
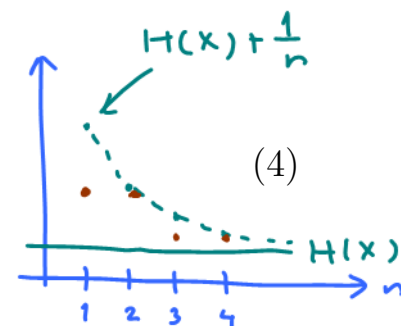
$$H(X) \le L^*(X) < H(X) + 1. \tag{3}$$

**Definition 2.54.** Let $L_n^*(X)$ be the minimum expected codeword length per symbol when the random variable $X$ is encoded with $n$-th extension uniquely decodable coding. Of course, this can be achieve by using $n$-th extension Huffman coding.

**2.55.** An extension of (3):

$$H(X) \le L_n^*(X) < H(X) + \frac{1}{n}. \tag{4}$$

*$H(X) + \frac{1}{n}$*

In particular,

$$\lim_{n \to \infty} L_n^*(X) = H(X).$$

In otherwords, by using large block length, we can achieve an expected length per source symbol that is arbitrarily close to the value of the entropy.

**2.56.** Operational meaning of entropy: Entropy of a random variable is the average length of its shortest description.

**2.57.** References

- Section 16.1 in Carlson and Crilly [4]

- Chapters 2 and 5 in Cover and Thomas [5]

- Chapter 4 in Fine [6]

- Chapter 14 in Johnson, Sethares, and Klein [8]

- Section 11.2 in Ziemer and Tranter [17]

# ECS452 2016/2     Part I.4     Dr.Prapun

## 3  An Introduction to Digital Communication Systems Over Discrete Memoryless Channel (DMC)

In this section, we keep our analysis of the communication system simple by considering purely digital systems. To do this, we assume all non-source-coding parts of the system, including the physical channel, can be combined into an "equivalent channel" which we shall simply refer to in this section as the "channel".

### 3.1  Discrete Memoryless Channel (DMC) Models

**Example 3.1.** The **binary symmetric channel (BSC)**, which is the simplest model of a channel with errors, is shown in Figure 4.
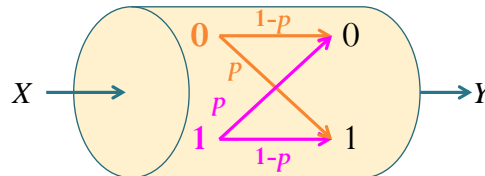


Figure 4: Binary symmetric channel and its channel diagram

- "Binary" means that the there are two possible values for the input and also two possible values for the output. We normally use the symbols 0 and 1 to represent these two values.

- Passing through this channel, the input symbols are complemented with crossover probability $p$.

switchover probability

bit-flip probability

24

- It is simple, yet it captures most of the complexity of the general problem.

**Example 3.2.** Consider a BSC whose samples of input and output are provided below

Time index: 1 2 3 4 5                                                   20
Channel input x: 1 0 1 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1
Channel output y: 1 1 1 1 1 1 1 1 1 0 1 1 0 1 0 1 1 1 1 1

Estimate the following (unconditional and conditional) probabilities by their relative frequencies.

$$\mathbf{p} = \begin{bmatrix} p(0) & p(1) \end{bmatrix} = \begin{bmatrix} \frac{3}{20} & \frac{17}{20} \end{bmatrix} \qquad x = \{0, 1\}$$

$$p(0) = P[X = 0] \approx \frac{3}{20} \qquad\qquad p(1) = P[X = 1] \approx 1 - P[x=0] = \frac{17}{20}$$

$$g(0) = P[Y = 0] \approx \frac{3}{20} \qquad\qquad g(1) = P[Y = 1] \approx 1 - \frac{3}{20} = \frac{17}{20} \qquad g = \begin{bmatrix} g(0) & g(1) \end{bmatrix}$$

$$= \begin{bmatrix} \frac{3}{20} & \frac{17}{20} \end{bmatrix}$$

$$Q(0|0) = P[Y = 0|X = 0] \approx \frac{2}{3} \qquad Q(1|0) = P[Y = 1|X = 0] \approx 1 - \frac{2}{3} = \frac{1}{3}$$

$$\begin{array}{c|cc} x\backslash Y & 0 & 1 \\\hline 0 & 2/3 & 1/3 \\ 1 & 1/17 & 16/17 \end{array}$$

$$Q(0|1) = P[Y = 0|X = 1] \approx \frac{1}{17} \qquad Q(1|1) = P[Y = 1|X = 1] \approx 1 - \frac{1}{17} = \frac{16}{17} \qquad Q =$$

**Definition 3.3.** Our model for **discrete memoryless channel (DMC)** is shown in Figure 5.
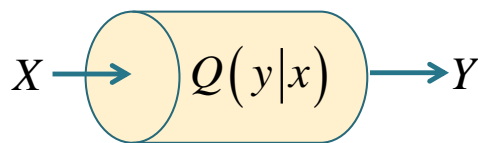
$$X \longrightarrow \boxed{Q(y|x)} \longrightarrow Y$$

Figure 5: Discrete memoryless channel

- The channel input is denoted by a random variable $X$.

  ○ The pmf $p_X(x)$ is usually denoted by simply $p(x)$ and usually expressed in the form of a row vector $\mathbf{p}$ or $\underline{\pi}$.

  ○ The support $S_X$ is often denoted by $\mathcal{X}$.
  └ channel input alphabet

25

- Similarly, the channel output is denoted by a random variable $Y$.

  ○ The pmf $p_Y(y)$ is usually denoted by simply $q(y)$ and usually expressed in the form of a row vector $\mathbf{q}$.

  ○ The support $S_Y$ is often denoted by $\mathcal{Y}$.

- The channel corrupts its input $X$ in such a way that when the input is $X = x$, its output $Y$ is randomly selected from the conditional pmf $p_{Y|X}(y|x)$. $= P\left[Y=y \mid X=x\right]$

$$Q(y|x) \equiv P\left[\underbrace{Y=y}_{A} \mid \underbrace{X=x}_{B}\right] = P(A|B) = \frac{P(A\cap B)}{P(B)} = \frac{P[Y=y, X=x]}{P[X=x]}$$

channel transition probability

  ○ This conditional pmf $p_{Y|X}(y|x)$ is denoted by $Q(y|x)$ and usually expressed in the form of a probability transition matrix $\mathbf{Q}$:
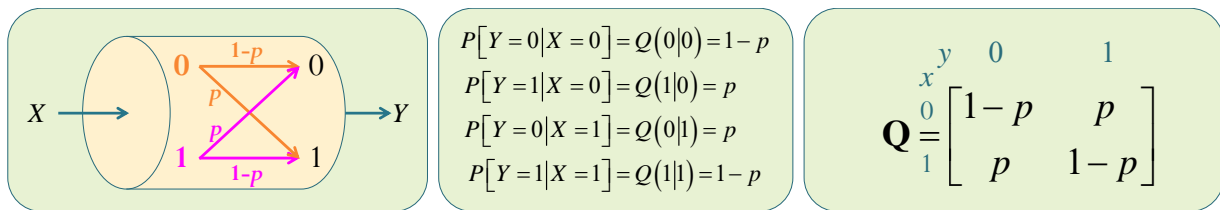
$$\begin{array}{c} y \\ x \begin{bmatrix} \ddots & & \vdots & & \iddots \\ \cdots & & \boxed{P\left[Y=y|X=x\right]} & & \cdots \\ \iddots & & \vdots & & \ddots \end{bmatrix} \end{array}$$

  ○ The channel is called memoryless[9] because its channel output at a given time is a function of the channel input at that time and is not a function of previous channel inputs.

  ○ Here, the transition probabilities are assumed constant. However, in many commonly encountered situations, the transition probabilities are time varying. An example is the wireless mobile channel in which the transmitter-receiver distance is changing with time.

---

[9]Mathematically, the condition that the channel is memoryless may be expressed as [12, Eq. 6.5-1 p. 355]

$$p_{X_1^n|Y_1^n}(x_1^n|y_1^n) = \prod_{k=1}^{n} Q(y_k|x_k).$$

**Example 3.4.** For a binary symmetric channel (BSC) defined in 3.1, we now have three equivalent ways to specify the relevant probabilities:



$$P[Y=0|X=0]=Q(0|0)=1-p$$
$$P[Y=1|X=0]=Q(1|0)=p$$
$$P[Y=0|X=1]=Q(0|1)=p$$
$$P[Y=1|X=1]=Q(1|1)=1-p$$

$$\mathbf{Q} = \begin{array}{c} \\ {\scriptstyle x} \\ 0 \\ 1 \end{array} \begin{array}{cc} {\scriptstyle y} & {\scriptstyle 0} \quad {\scriptstyle 1} \\ \left[\begin{array}{cc} 1-p & p \\ p & 1-p \end{array}\right] \end{array}$$

**Example 3.5.** Suppose, for a DMC, we have $\mathcal{X} = \{x_1, x_2\}$ and $\mathcal{Y} = \{y_1, y_2, y_3\}$. Then, its probability transition matrix $\mathbf{Q}$ is of the form

$$\mathbf{Q} = \left[ \begin{array}{ccc} Q(y_1|x_1) & Q(y_2|x_1) & Q(y_3|x_1) \\ Q(y_1|x_2) & Q(y_2|x_2) & Q(y_3|x_2) \end{array} \right].$$

You may wonder how this $\mathbf{Q}$ happens in real life. Let's suppose that the input to the channel is binary; hence, $\mathcal{X} = \{0, 1\}$ as in the BSC. However, in this case, after passing through the channel, some bits can be lost[10] (rather than corrupted). In such case, we have three possible outputs of the channel: 0, 1, e where the "e" represents the case in which the bit is erased by the channel.

**Example 3.6.** Consider a DMC whose samples of input and output are provided below

$$\overset{A \cap B}{P[X=1, Y=2]} = \frac{7}{20}$$

x: 1 1 1 1 1 1 1 1 1 1 0 0 1 1 1 0 1 1 0 1

y: 1 3 2 2 1 2 1 2 2 3 1 1 1 3 1 3 2 3 1 2

$$\overset{B \quad A}{= P[Y=2 \,|\, X=1]} \overset{A}{P[X=1]}$$

Estimate its input probability vector $\underline{\mathbf{p}}$, output probability vector $\underline{\mathbf{q}}$, and $\mathbf{Q}$ matrix.

$$= \frac{7}{16} \times \frac{16}{20} = \frac{7}{20}$$

$$20-4$$

$$\underline{\mathbf{p}} = \left[ p(0) \quad p(1) \right] \approx \left[ \frac{4}{20} \quad \frac{16}{20} \right] = \left[ 0.2 \quad 0.8 \right]$$

$$\underline{\mathbf{q}} = \left[ q(1) \quad q(2) \quad q(3) \right] \approx \left[ \frac{8}{20} \quad \frac{7}{20} \quad \frac{5}{20} \right] = \left[ 0.4 \quad 0.35 \quad 0.25 \right]$$

$$\mathbf{Q} = \begin{array}{c} \\ {\scriptstyle x} \\ 0 \\ 1 \end{array} \begin{array}{c} {\scriptstyle y} \quad 1 \qquad 2 \qquad 3 \\ \left[ \begin{array}{ccc} \frac{3}{4} & \frac{0}{4} & \frac{1}{4} \\ \frac{5}{16} & \frac{7}{16} & \frac{4}{16} \end{array} \right] \end{array}$$

$$Q(y|x) = P[Y=y \,|\, X=x]$$

---
[10]The receiver knows which bits have been erased.

27

**3.7.** Knowing the input probabilities $\underline{p}$ and the channel probability transition matrix $\mathbf{Q}$, we can calculate the output probabilities $\underline{q}$ from

$$\underline{q} = \underline{p}\mathbf{Q}.$$

To see this, recall the **total probability theorem**: If a (finite or infinitely) countable collection of events $\{B_1, B_2, \ldots\}$ is a partition of $\Omega$, then

$$P(A) = \sum_i P(A \cap B_i) = \sum_i P(A|B_i)P(B_i). \qquad (5)$$



$P(A) = P(A \cap B_1) + P(A \cap B_2)$
$\quad + P(A \cap B_3) + P(A \cap B_4) + P(A \cap B_5)$

For us, event $A$ is the event $[Y = y]$. Applying this theorem to our variables, we get

$$A = [Y=y]$$
$$B_x = [X = x]$$

$$q(y) = P[Y = y] = \sum_x P[X = x, Y = y]$$

$$= \sum_x P[Y = y|X = x]P[X = x] = \sum_x Q(y|x)p(x).$$

This is exactly the same as the matrix multiplication calculation performed to find each element of $\underline{q}$.

**Example 3.8.** For a binary symmetric channel (BSC) defined in 3.1,

first column of $Q$ matrix

$$q(0) = P[Y = 0] = P[Y = 0, X = 0] + P[Y = 0, X = 1]$$
$$= P[Y = 0|X = 0]P[X = 0] + P[Y = 0|X = 1]P[X = 1]$$
$$= Q(0|0)p(0) + Q(0|1)p(1) = (1-p) \times p_0 + p \times p_1 \quad = \begin{bmatrix} p(0) & p(1) \end{bmatrix} \begin{bmatrix} Q(0|0) \\ Q(0|1) \end{bmatrix}$$
$$q(1) = P[Y = 1] = P[Y = 1, X = 0] + P[Y = 1, X = 1]$$
$$= P[Y = 1|X = 0]P[X = 0] + P[Y = 1|X = 1]P[X = 1]$$
$$= Q(1|0)p(0) + Q(1|1)p(1) = p \times p_0 + (1-p) \times p_1 \quad = \begin{bmatrix} p(0) & p(1) \end{bmatrix} \begin{bmatrix} Q(1|0) \\ Q(1|1) \end{bmatrix}$$

second column of $Q$ matrix

$Q$

$$\underline{q} = \begin{bmatrix} q(0) & q(1) \end{bmatrix} = \begin{bmatrix} \underline{p} & \boxed{\phantom{x}} \end{bmatrix} \quad \underline{p}\begin{bmatrix} \boxed{\phantom{x}} \end{bmatrix} = \underline{p}\begin{bmatrix} \boxed{\phantom{x}} & \boxed{\phantom{x}} \end{bmatrix}$$

28

block matrix multiplication (factoring)

$$ \mathscr{g} = \begin{bmatrix} & g(y) & \end{bmatrix} = P Q = \underbrace{\begin{bmatrix} \phantom{xxxxxx} \end{bmatrix}}_{P} \begin{bmatrix} \overset{y}{\boxed{\phantom{x}}} \end{bmatrix} $$

P

Q

**3.9.** Recall, from ECS315, that there is another matrix called the **joint probability matrix P**. This is the matrix whose elements give the joint probabilities $P_{X,Y}(x,y) = P[X = x, Y = y]$:

$$\mathbf{P} = \begin{matrix} & y \\ x & \begin{bmatrix} \ddots & & \vdots & & \cdot^{\cdot^{\cdot}} \\ \cdots & & P[X = x, Y = y] & & \cdots \\ \cdot^{\cdot^{\cdot}} & & \vdots & & \ddots \end{bmatrix} \end{matrix}.$$

Recall also that we can get $p(x)$ by adding the elements of **P** in the row corresponding to $x$. Similarly, we can get $q(y)$ by adding the elements of **P** in the column corresponding to $y$.

By definition, the relationship between the conditional probability $Q(y|x)$ and the joint probability $P_{X,Y}(x,y)$ is

$$Q(y|x) = \frac{P_{X,Y}(x,y)}{p(x)}.$$

*[handwritten: $P(A|B) = \frac{P(A \cap B)}{P(B)}$    $A = [Y = y]$    $B = [X = x]$]*

Equivalently,

$$P_{X,Y}(x,y) = p(x)Q(y|x).$$

*[handwritten: $P[Y = y | X = x] = \frac{P[X = x, Y = y]}{P[X = x]}$]*
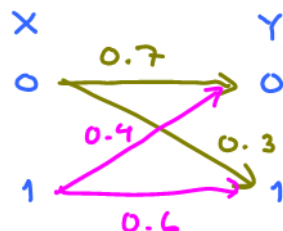
Therefore, to get the matrix **P** from matrix **Q**, we need to multiply each row of **Q** by the corresponding $p(x)$. This could be done easily in MATLAB by first constructing a diagonal matrix from the elements in $\underline{p}$ and then multiply this to the matrix **Q**:

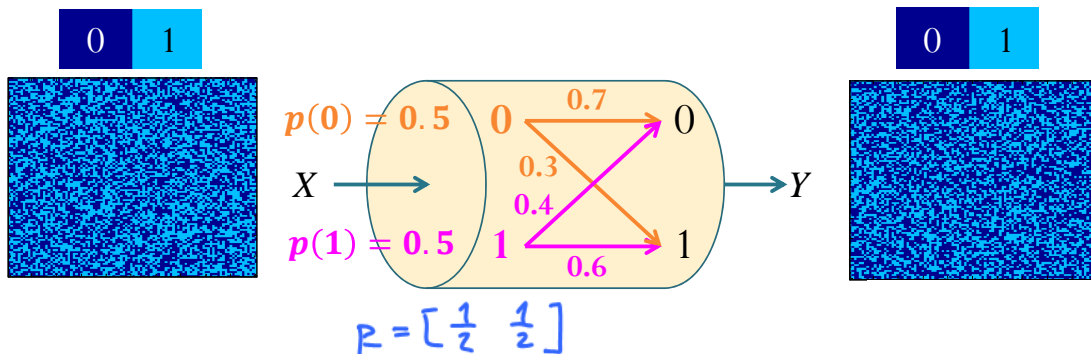$$\mathbf{P} = \left(\text{diag}\left(\underline{\mathbf{p}}\right)\right)\mathbf{Q}.$$

**Example 3.10. Binary Asymmetric Channel (BAC)**: Consider a binary input-output channel whose matrix of transition probabilities is

*[handwritten: $x \backslash y$   0   1]*

$$\mathbf{Q} = \begin{matrix} 0 \\ 1 \end{matrix}\begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix}.$$

(a) Draw the channel diagram.



29

$$p(0) = 0.5 \quad 0 \xrightarrow{0.7} 0$$

with branches 0.3, 0.4 and $p(1) = 0.5 \quad 1 \xrightarrow{0.6} 1$

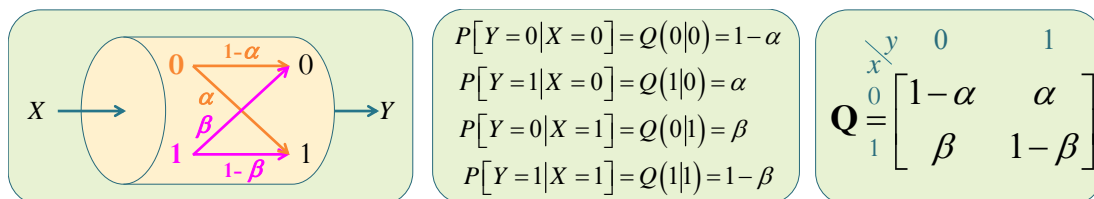$$\underline{p} = \left[ \tfrac{1}{2} \quad \tfrac{1}{2} \right]$$

(b) If the two inputs are equally likely, find the corresponding output probability vector $\underline{q}$ and the joint probability matrix $\mathbf{P}$ for this channel.

$$\underline{q} = \underline{p}\, Q = \left[ \tfrac{1}{2} \quad \tfrac{1}{2} \right] \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix} = \begin{bmatrix} 0.55 & 0.45 \end{bmatrix}$$
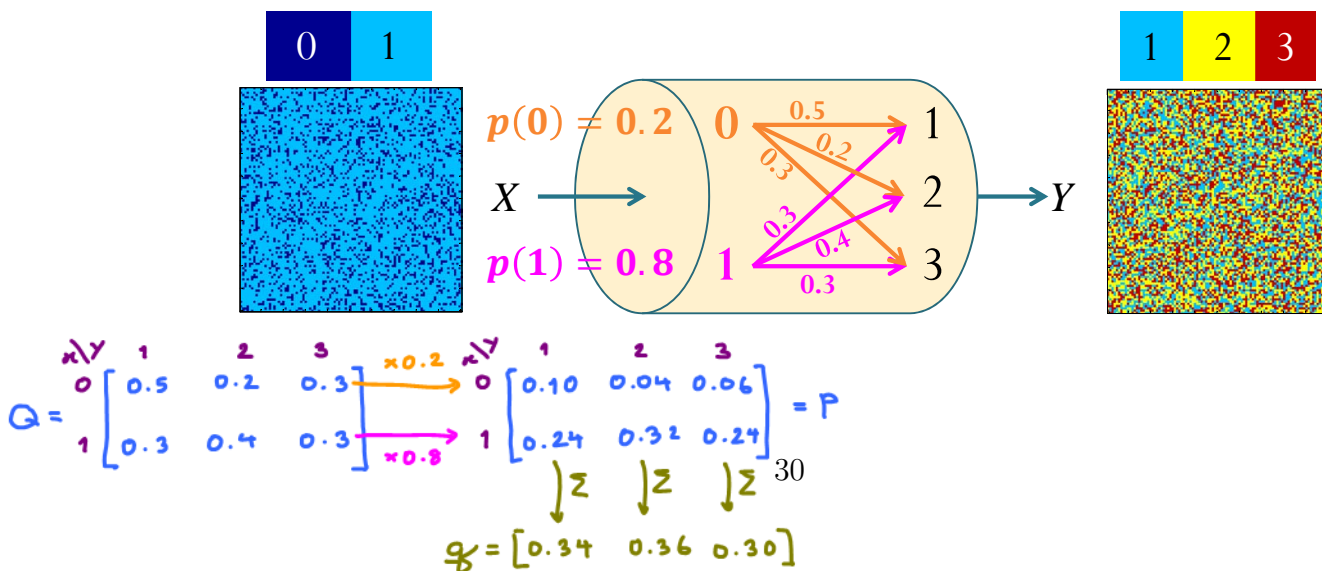
$$Q = \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix} \xrightarrow[\times p(1)]{\times p(0)} \begin{array}{c} x \backslash y \\ 0 \\ 1 \end{array} \begin{matrix} 0 & 1 \\ \begin{bmatrix} 0.35 & 0.15 \\ 0.2 & 0.3 \end{bmatrix} \end{matrix} = P$$

$$\downarrow \Sigma \qquad \downarrow \Sigma$$
$$\begin{bmatrix} 0.55 & 0.45 \end{bmatrix} = \underline{q}$$

[17, Ex. 11.3]

**Example 3.11.** Similar to Example 3.4 where we have three equivalent ways to specify BSC. We also have three different ways to describe BAC:
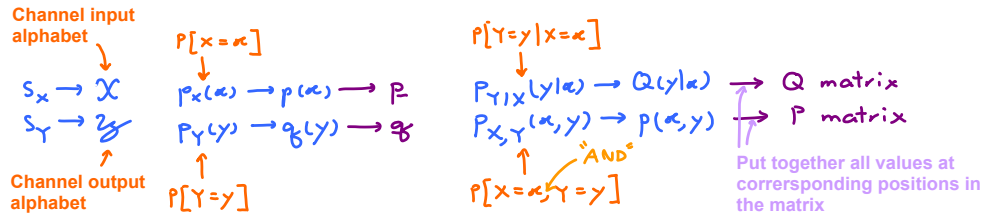


$$P[Y = 0 | X = 0] = Q(0|0) = 1 - \alpha$$
$$P[Y = 1 | X = 0] = Q(1|0) = \alpha$$
$$P[Y = 0 | X = 1] = Q(0|1) = \beta$$
$$P[Y = 1 | X = 1] = Q(1|1) = 1 - \beta$$

$$\mathbf{Q} = \begin{array}{c} x \backslash y \\ 0 \\ 1 \end{array} \begin{matrix} 0 & 1 \\ \begin{bmatrix} 1 - \alpha & \alpha \\ \beta & 1 - \beta \end{bmatrix} \end{matrix}$$

**Example 3.12.** Find the output probability vector $\underline{q}$ and the joint probability matrix $\mathbf{P}$ for the following DMC:



$$Q = \begin{array}{c} x \backslash y \\ 0 \\ 1 \end{array} \begin{matrix} 1 & 2 & 3 \\ \begin{bmatrix} 0.5 & 0.2 & 0.3 \\ 0.3 & 0.4 & 0.3 \end{bmatrix} \end{matrix} \xrightarrow[\times 0.8]{\times 0.2} \begin{array}{c} x \backslash y \\ 0 \\ 1 \end{array} \begin{matrix} 1 & 2 & 3 \\ \begin{bmatrix} 0.10 & 0.04 & 0.06 \\ 0.24 & 0.32 & 0.24 \end{bmatrix} \end{matrix} = P$$

$$\downarrow \Sigma \quad \downarrow \Sigma \quad \downarrow \Sigma$$

$$\underline{q} = \begin{bmatrix} 0.34 & 0.36 & 0.30 \end{bmatrix}$$

30

# Summary:

**DMC:** **X: Channel Input**
**Y: Channel Output**

**Notation used in digital commu. (and info. theory) class is different from probability class**

**Channel input alphabet**

$P[X=x]$

$P[Y=y|X=x]$

$S_X \to \mathcal{X}$    $P_X(x) \to p(x) \to \underline{p}$    $P_{Y|X}(y|x) \to Q(y|x) \to$ Q matrix

$S_Y \to \mathcal{Y}$    $P_Y(y) \to q(y) \to \underline{q}$    $P_{X,Y}(x,y) \to p(x,y) \to$ P matrix

**Channel output alphabet**

$P[Y=y]$

$P[X=x, Y=y]$    "AND"

**Put together all values at corresponding positions in the matrix**

**When the alphabets are lists of integers,**

we usually write $p(x)$ and $q(y)$

as $p_x$ and $q_y$ respectively.

Ex. For BSC, $\mathcal{X} = \{0,1\}$
$p_0 = p(0) = P[X=0]$
$p_1 = p(1) = P[X=1]$

**Alternatively, when the members of the alphabet(s) are explicitly indexed,**
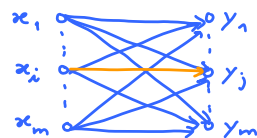
we often define $p_i \equiv p(x_i)$ and $q_j \equiv q(y_j)$

**DMC is defined by its Q matrix:**

When $\mathcal{X} = \{x_1, x_2, x_3, \ldots, x_m\}$ and $\mathcal{Y} = \{y_1, y_2, y_3, \ldots, y_n\}$,

$$Q = \begin{bmatrix} & y_1 \cdots y_j \cdots y_n \\ x_1 & \\ \vdots & \\ x_i & Q(y_i|x_i) \\ \vdots & \\ x_m & \end{bmatrix}$$

$\times p(x_1)$
$\times p(x_i)$
$\times p(x_m)$

$$\begin{bmatrix} & y_1 \cdots y_j \cdots y_n \\ x_1 & \\ \vdots & \\ x_i & P_{X,Y}(x_i, y_j) \\ \vdots & \\ x_m & \end{bmatrix} = P$$

**The P matrix can be derived from the Q matrix by scaling each row of the Q matrix by the corresponding p(x)**

$\Sigma \downarrow \quad \Sigma \downarrow \quad \Sigma \downarrow$

$[ \qquad \qquad ] = \underline{q}$

Alternatively, $\underline{q} = \underline{p}\, Q$

**Equivalently, we may define a DMC via its channel diagram (of transition probabilities):**

$x_1 \quad y_1$
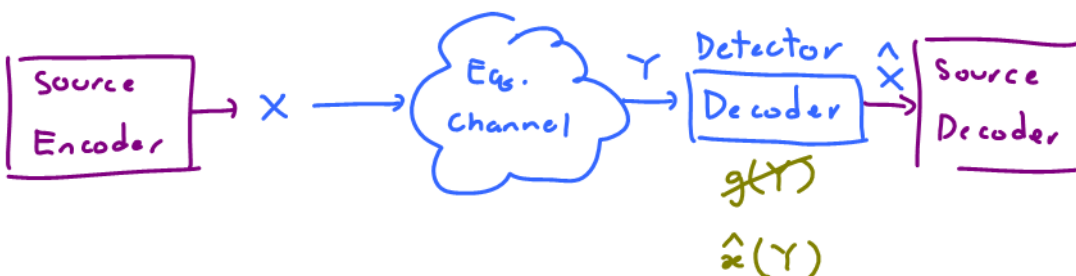$x_i \quad y_j$
$x_m \quad y_m$

**Each arrow should be labeled with**

$Q(y_j | x_i)$.

## 3.2  Decoder and Symbol Error Probability

**3.13.** Knowing the characteristics of the channel and its input, on the receiver side, we can use this information to build a "good" receiver.

We now consider a part of the receiver called the **(channel) decoder**. Its job is to guess the value of the channel input[11] $X$ from the value of the received channel output $Y$. We denote this guessed value by $\hat{X}$.



**3.14.** A "good" receiver is the one that (often) guesses correctly. So, our goal here is to

- maximize the probability of correct guessing
- minimize  "            "         "    wrong guessing

Quantitatively, to measure the performance of a decoder, we define a quantity called the (symbol) error probability.

**Definition 3.15.** The **(symbol) error probability**, denoted by $P(\mathcal{E})$, can be calculated from

$$P(\mathcal{E}) = P\left[\hat{X} \neq X\right].$$

**3.16.** A "good" detector should guess based on all the information it has obtained. Here, the only information it can observe is the value of $Y$. So, a detector is a function of $Y$, say, $g(Y)$. Therefore, $\hat{X} = g(Y)$.
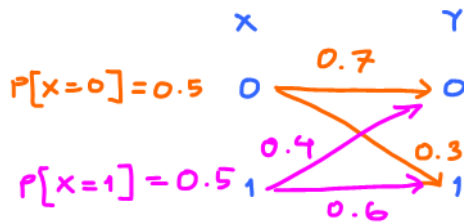
We will write $\hat{X}$ as $\hat{x}(Y)$ to emphasize that the decoded value $\hat{X}$ depends on the observed value $Y$ and that the detector is a deterministic function of the channel output $Y$; the randomness in the decoded value $\hat{X}$ comes from the randomness in $Y$.

**Definition 3.17.** A "naive" decoder is a decoder that simply sets $\hat{X} = Y$.

---

[11]To simplify the analysis, we still haven't considered the channel encoder. (It may be there but is included in the equivalent channel or it may not be in the system at all.)

**Example 3.18.** Consider the BAC channel and input probabilities specified in Example 3.10. Find $P(\mathcal{E})$ when $\hat{X} = Y$.

$$P[X=0]=0.5 \quad 0 \xrightarrow{0.7} 0$$
$$P[X=1]=0.5 \quad 1 \xrightarrow{0.6} 1$$

(arrows: $0.4$, $0.3$)

$$P(\mathcal{E}) = P[\hat{X} \neq X] = P[Y \neq X]$$
$$= 0.2 + 0.15 = 0.35$$

$$Q = \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix} \xrightarrow[\times p(1)]{\times p(0)} \begin{array}{c} Y \\ 0 \\ 1 \end{array} \begin{bmatrix} 0.35 & 0.15 \\ 0.2 & 0.3 \end{bmatrix} = P$$

**3.19.** For general DMC, the error probability of the naive decoder is

$$P(\mathcal{E}) = P\left[\hat{X} \neq X\right] = P[Y \neq X] = 1 - P[Y = X]$$
$$= 1 - \sum_x P[Y = x, X = x] = 1 - \sum_x P[Y = x \,|X = x]P[X = x]$$
$$= 1 - \sum_x Q(x\,|x)\,p(x)$$

**Example 3.20.** With the derived formula, let's revisit Example 3.18.

$$P(\mathcal{E}) = 1 - (Q(0\,|0)\,p(0) + Q(1\,|1)\,p(1)) = 1 - \left(0.7 \times \frac{1}{2} + 0.6 \times \frac{1}{2}\right) = 0.35.$$

**Example 3.21.** Find the error probability $P(\mathcal{E})$ when a naive decoder is used with a DMC channel in which $\mathcal{X} = \{0,1\}$, $\mathcal{Y} = \{1,2,3\}$, $\mathbf{Q} = \begin{bmatrix} 0.5 & 0.2 & 0.3 \\ 0.3 & 0.4 & 0.3 \end{bmatrix}$ and $\mathbf{p} = [0.2, 0.8]$.

**From Example 3.12,**

$$Q = \begin{array}{c} {}^{x\backslash y} \\ 0 \\ 1 \end{array} \begin{array}{ccc} 1 & 2 & 3 \\ \end{array} \begin{bmatrix} 0.5 & 0.2 & 0.3 \\ 0.3 & 0.4 & 0.3 \end{bmatrix} \xrightarrow[\times 0.8]{\times 0.2} \begin{array}{c} {}^{x\backslash y} \\ 0 \\ 1 \end{array} \begin{array}{ccc} 1 & 2 & 3 \\ \end{array} \begin{bmatrix} 0.10 & 0.04 & 0.06 \\ 0.24 & 0.32 & 0.24 \end{bmatrix} = P$$

$$P(\mathcal{E}) = P[\hat{X} \neq X]$$
$$= P[Y \neq X]$$
$$= 1 - 0.24 = 0.76$$

$$P(\mathcal{C}) = P[Y = x]$$
$$= 0.24$$

33

**Example 3.22.** DIY Decoder: Consider a different decoder specified in the decoding table below. Find the error probability $P(\mathcal{E})$ when such decoder is used in Example 3.21.

Decoding table →

| $y$ | $\hat{x}(y)$ |
|-----|------|
| 1   | 0    |
| 2   | 1    |
| 3   | 0    |

$\hat{x}$   0   1   0

$$\begin{array}{c|ccc} & 1 & 2 & 3 \\ \hline 0 & \boxed{0.10} & 0.04 & \boxed{0.06} \\ 1 & 0.24 & \boxed{0.32} & 0.24 \end{array} = P$$

$$P(\mathcal{C}) = P[\hat{X} = X] = 0.1 + 0.32 + 0.06 = 0.48$$

$$P(\mathcal{E}) = P[\hat{X} \ne X] = 1 - 0.48 = 0.52$$

**Example 3.23.** Repeat Example 3.22 but use the following decoder

| $y$ | $\hat{x}(y)$ |
|-----|------|
| 1   | 1    |
| 2   | 1    |
| 3   | 0    |

$\hat{x}$   1   1   0

$$\begin{array}{c|ccc} & 1 & 2 & 3 \\ \hline 0 & 0.10 & 0.04 & \boxed{0.06} \\ 1 & \boxed{0.24} & \boxed{0.32} & 0.24 \end{array} = P$$

$$P(\mathcal{C}) = 0.24 + 0.32 + 0.06 = 0.62$$

$$P(\mathcal{E}) = 1 - P(\mathcal{C}) = 1 - 0.62 = 0.38$$

Observation: For each column of the **P** matrix, we circle the probability corresponding to the row of $x$ that has the same value as $\hat{x}(y)$.

**3.24.** A recipe for finding $P(\mathcal{E})$ of any (DIY) decoder:

(a) Find the **P** matrix by scaling each row of the **Q** matrix by its corresponding $p(x)$.

(b) Write $\hat{x}(y)$ values on top of the $y$ values for the **P** matrix.

(c) For each $y$ column in the **P** matrix, circle the element whose corresponding $x$ value is the same as $\hat{x}(y)$.

(d) $P(\mathcal{C}) =$ the sum of the circled probabilities. $P(\mathcal{E}) = 1 - P(\mathcal{C})$.

## 3.3 Optimal Decoding for DMC

From the previous section, we now know how to compute the error probability for any given decoder. Here, we will attempt to find the "best" decoder. Of course, by "best", we mean "having minimum value of error probability". It is interesting to first consider the question of how many reasonable decoders we can use.

**Definition 3.25.** "Reasonable" Decoder: Recall from 3.16 that a decoder $\hat{x}(\cdot)$ is a function that map each observed valued of the channel output $y$ to the guessed value of the channel input. Therefore we can think of a decoder as a table:



We have already seen this table representation in Example 3.22 and Example 3.23. Such table has $|\mathcal{Y}|$ rows. For each value of $y$, we need to specify what is the value of $\hat{x}(y)$. To have a chance of correct guessing, any "reasonable" decoder would select the value of $\hat{x}(y)$ from $\mathcal{X}$. Therefore, there are $|\mathcal{X}|^{|\mathcal{Y}|}$ reasonable decoders.

**Example 3.26.** The naive decoder in Example 3.21 is not a reasonable decoder. The channel input $X$ is either 0 or 1. So, it does not make sense have a guess value of $\hat{x}(2) = 2$ or $\hat{x}(3) = 3$.

**Example 3.27.** "Reasonable" Decoder for BSC: For BSC in Example 3.1, any decoder has to answer two important questions:

(a) What should be the guess value of $X$ when $Y = 0$ is observed?
$$\hat{x}(0) = 0 \text{ or } \hat{x}(0) = 1$$

(b) What should be the guess value of $X$ when $Y = 1$ is observed?
$$\hat{x}(1) = 0 \text{ or } \hat{x}(1) = 1$$

Essentially, any reasonable decoder for the BSC need to complete this table:



35

So, only four reasonable decoders for BSC.

| ① | $y$ | $\hat{x}(y)$ |
|---|---|---|
| | 0 | 0 |
| | 1 | 1 |

Naive

$\hat{x}(y) = y$

| ② | $y$ | $\hat{x}(y)$ |
|---|---|---|
| | 0 | 1 |
| | 1 | 0 |

$\hat{x}(y) = 1 - y = \bar{y}$

| ③ | $y$ | $\hat{x}(y)$ |
|---|---|---|
| | 0 | 0 |
| | 1 | 0 |

$\hat{x}(y) \equiv 0$

| ④ | $y$ | $\hat{x}(y)$ |
|---|---|---|
| | 0 | 1 |
| | 1 | 1 |

$\hat{x}(y) \equiv 1$

**Example 3.28.** For the DMC defined in Example 3.21, how many reasonable decoders are there?

$$\underbrace{|\mathcal{X}| \times |\mathcal{X}| \times \cdots \times |\mathcal{X}|}_{|\mathcal{Y}| \text{ times}} = |\mathcal{X}|^{|\mathcal{Y}|}$$

| $y$ | $\hat{x}(y)$ |
|---|---|
| 1 | 0 or 1 |
| 2 | 0 or 1 |
| 3 | 0 or 1 |

$2^3 = 8$

We calculate the error probability of three decoders in Example 3.21, Example 3.22, and Example 3.23. There are still many possibilities to evaluate. Using MATLAB, we can find the error probability for all possible reasonable decoders.

**3.29.** For general DMC, it would be tedious to list all possible decoders. It is even more time-consuming to try to calculate the error probability for all of them. Therefore, in this section, we will derive a visual construction and a formula of the "optimal" decoder.

**3.30.** From the recipe 3.24 for finding $P(\mathcal{C})$ and $P(\mathcal{E})$, we see that $P(\mathcal{C})$ is the sum of our circled numbers. So, to maximize $P(\mathcal{C})$, we want to circle the largest number. For row $y$ in the decoding table, whatever the value we select for $\hat{x}(y)$ will determine which number will be circled in the column corresponding to $y$ in matrix $\mathbf{P}$. To maximize $P(\mathcal{C})$, we want to circle the largest number in the column. This means $\hat{x}(y)$ should be the same as the $x$ value that maximizes the probability value in the corresponding column.

**Example 3.31.** For the DMC and the input probablities defined in Example 3.21, the joint pmf matrix $\mathbf{P}$ was found to be

$$\mathbf{P} = \begin{array}{c|ccc} x \backslash y & 1 & 2 & 3 \\ \hline 0 & 0.1 & 0.04 & 0.06 \\ 1 & 0.24 & 0.32 & 0.24 \end{array}$$

Therefore, the optimal decoder is

① Decoder table

| $y$ | $\hat{x}(y)$ |
|---|---|
| 1 | 1 |
| 2 | 1 |
| 3 | 1 |

$P(\mathcal{C}) = 0.24 + 0.32 + 0.24 = 0.8$

$P(\mathcal{E}) = 1 - 0.8 = 0.2$

**3.32.** Mathematically, we first note that to minimize $P(\mathcal{E})$, we need to maximize $P(\mathcal{C})$. Here, we apply the total probability theorem by using the events $[Y = y]$ to partition the sample space:

$$P(\mathcal{C}) = \sum_y P(\mathcal{C}\,|[Y = y])\,P[Y = y].$$

Event $\mathcal{C}$ is the event $[\hat{X} = X]$. Therefore,

$$P(\mathcal{C}) = \sum_y P\left[\hat{X} = X\,|Y = y\right]P[Y = y].$$

Now, recall that our decoder is a function of $Y$; that is $\hat{X} = \hat{x}(Y)$. So,

$$P(\mathcal{C}) = \sum_y P[\hat{x}(Y) = X\,|Y = y]\,P[Y = y]$$

$$= \sum_y P[X = \hat{x}(y)\,|Y = y]\,P[Y = y]$$

In this form, we see[12] that for each $Y = y$, we should maximize $P[X = \hat{x}(y)\,|Y = y]$. Therefore, for each $y$, the decoder $\hat{x}(y)$ should output the value of $x$ which maximizes $P[X = x|Y = y]$:

$$\hat{x}_{\text{optimal}}(y) = \arg\max_x P[X = x\,|Y = y].$$

In other words, the *optimal* decoder is the decoder that maximizes the "a posteriori probability" $P[X = x\,|Y = y]$.

**Definition 3.33.** The optimal decoder derived in 3.32 is called the **maximum a posteriori probability (MAP) decoder**:

$$\hat{x}_{\text{MAP}}(y) = \hat{x}_{\text{optimal}}(y) = \arg\max_x P[X = x\,|Y = y].$$

Ex. $\quad 5 - x^2 = f(x)$

$\max_x f(x) = 5$

$\arg\max_x f(x) = 0$

**3.34.** After the fact, it is quite intuitive that this should be the best decoder. Recall that the decoder don't have a direct access to the $X$ value.

- Without knowing the value of $Y$, to minimize the error probability, it should guess the most likely value of $X$ which is the value of $x$ that maximize $P[X = x]$.

- Knowing $Y = y$, the decoder can update its probability about $x$ from $P[X = x]$ to $P[X = x|Y = y]$. Therefore, the decoder should guess the value of the most likely $x$ value conditioned on the fact that $Y = y$.

---

[12]We also see that any decoder that produces random results (on the support of $X$) can not be better than our optimal decoder. Outputting the value of $x$ which does not maximize the a posteriori probability reduces the contribution in the sum that gives $P(\mathcal{C})$.

**3.35.** We can "simplify" the formula for the MAP decoder even further.

.

Fist, recall "Form 1" of the Bayes' theorem:

$$P(B|A) = P(A|B)\frac{P(B)}{P(A)}.$$

Here, we set $B = [X = x]$ and $A = [Y = y]$.

$$P\left[X=x \,|\, Y=y\right] = \left(P\left[Y=y \,|\, X=x\right]\right)\frac{P\left[X=x\right]}{P\left[Y=y\right]}$$

*Does not depend on $x$ and $> 0$. So can ignore.*

Therefore,

$$\boxed{\hat{x}_{\text{MAP}}(y) = \arg\max_x Q\left(y \,|\, x\right) p\left(x\right)}. \qquad (6)$$

*Optimal*

**3.36.** A recipe for finding the MAP decoder and its corresponding error probability:

(a) Find the **P** matrix by scaling elements in each row of the **Q** matrix by their corresponding prior probability $p(x)$.

(b) Select (by circling) the maximum value in each column (for each value of $y$) in the **P** matrix.

  • If there are multiple max values in a column, select one. This won't affect the optimality of your answer.

  (i) The corresponding $x$ value is the value of $\hat{x}$ for that $y$.

  (ii) The sum of the selected values from the **P** matrix is $P(\mathcal{C})$.

(c) $P(\mathcal{E}) = 1 - P(\mathcal{C})$.

**Example 3.37.** We have applied recipe 3.36 back when we try to find the optimal decoder in Example 3.31.

**Example 3.38.** Find the MAP decoder *optimal* and its corresponding error probability for the DMC channel whose $\mathbf{Q}$ matrix is given by

$$Q = \begin{array}{c} x \backslash y \\ 0 \\ 1 \end{array} \begin{array}{ccc} 1 & 2 & 3 \\ \begin{bmatrix} 0.5 & 0.2 & 0.3 \\ 0.3 & 0.4 & 0.3 \end{bmatrix} \end{array}$$

*(handwritten annotations:)*

$\hat{x}$ : 0   1   0

$y$ : 1   2   3

$\times 0.6 \longrightarrow$
$\times 0.4 \longrightarrow$

$$\begin{array}{c} 0 \\ 1 \end{array} \begin{bmatrix} \boxed{0.3} & 0.12 & \boxed{0.18} \\ 0.12 & \boxed{0.16} & 0.12 \end{bmatrix} = P$$

and $\mathbf{p} = [0.6, 0.4]$. Note that the DMC is the same as in Example 3.21 but the input probabilities are different.

*(handwritten:)* $P(u)$ $\qquad$ $P(1)$

$P(C) = 0.3 + 0.16 + 0.18 = 0.64$

$P(\mathcal{E}) = 1 - 0.64 = 0.36$

| $y$ | $\hat{x}(y)$ |
|---|---|
| 1 | 0 |
| 2 | 1 |
| 3 | 0 |

**Definition 3.39.** In many scenarios, the MAP decoder is too complicated or the prior probabilities are unknown. In such cases, we may consider using a *suboptimal* decoder that ignores the prior probability term in (6). This decoder is called the **maximum likelihood (ML)** decoder:

$$\boxed{\hat{x}_{\mathrm{ML}}(y) = \arg\max_x Q(y\,|x)}. \tag{7}$$

*optimal*

**3.40.** ML decoder is the same as the MAP decoder when $X$ is a uniform random variable. In other words, when the prior probabilities $p(x)$ are uniform, the ML decoder is optimal.

**3.41.** A recipe for finding the ML decoder and its corresponding error probability:

(a) Select (by circling) the maximum value in each column (for each value of $y$) in the $\mathbf{Q}$ matrix.

- If there are multiple max values in a column, select one. Different choices will lead to different $P(\mathcal{E})$. However, if the information about $\mathbf{p}$ is not available at the decoder, it can not determine which choice is better anyway.

- The corresponding $x$ value is the value of $\hat{x}$ for that $y$.

(b) Find the $\mathbf{P}$ matrix by scaling elements in each row of the $\mathbf{Q}$ matrix by their corresponding prior probability $p(x)$.

(c) In the **P** matrix, select the elements corresponding to the selected positions in the **Q** matrix.

(d) The sum of the selected values from the **P** matrix is $P(\mathcal{C})$.

(e) $P(\mathcal{E}) = 1 - P(\mathcal{C})$.

**Example 3.42.** Find the ML decoder and its corresponding error probability a the DMC channel in Example 3.21 whose **Q** matrix is

$$\hat{x}_{ML} \quad 0 \quad 1 \quad 0$$

$$Q = \begin{array}{c} x \backslash y \\ 0 \\ 1 \end{array} \begin{array}{ccc} 1 & 2 & 3 \\ \boxed{0.5 \; 0.2 \; \boxed{0.3}} \\ 0.3 \; \boxed{0.4} \; 0.3 \end{array} \xrightarrow[\times 0.8]{\times 0.2} \begin{bmatrix} \boxed{0.1} & 0.04 & \boxed{0.06} \\ 0.24 & \boxed{0.32} & 0.24 \end{bmatrix} = P$$

$$\begin{array}{c|c} y & \hat{x}_{ML}(y) \\ \hline 1 & 0 \\ 2 & 1 \\ 3 & 0 \end{array}$$

and $\underline{p} = [0.2, 0.8]$.

$$P(\mathcal{C}) = 0.1 + 0.32 + 0.06 = 0.48$$

$$P(\mathcal{E}) \approx 1 - 0.48 = 0.52$$

**Example 3.43.** Find the ML decoder and the corresponding error probability for a communication over BSC with $p = 0.1$ and $p_0 = 0.8$.

$$\hat{x}_{ML}(y) \quad 0 \quad 1 \qquad\qquad \text{crossover probability}$$
$$\hat{x}_{ML}(y) \quad 0 \quad 1$$

$$Q = \begin{array}{c} x\backslash y & 0 & 1 \\ 0 & 1-p & p \\ 1 & p & 1-p \end{array} = \begin{array}{c} x\backslash y & 0 & 1 \\ 0 & \boxed{0.9} & 0.1 \\ 1 & 0.1 & \boxed{0.9} \end{array} \xrightarrow[0.2]{0.8} \begin{array}{c} x\backslash y & 0 & 1 \\ 0 & \boxed{0.72} & 0.08 \\ 1 & 0.02 & \boxed{0.18} \end{array} = P$$

$$\begin{array}{c|c} y & \hat{x}_{ML}(y) \\ \hline 0 & 0 \\ 1 & 1 \end{array} \qquad \hat{x}_{ML}(y) \equiv y$$

$$P(\mathcal{C}) = 0.72 + 0.18 = 0.9$$

$$P(\mathcal{E}) = 1 - 0.9 = 0.1$$

Note that

- the prior probabilities $p_0$ (and $p_1$) is not used when finding $\hat{x}_{ML}$.

- the ML decoder and the MAP decoder are the same in this example

  ○ ML decoder can be optimal even when the prior probabilities are not uniform.

**3.44.** In general, for BSC, it's straightforward to show that

(a) when $p < 0.5$, we have $\hat{x}_{\mathrm{ML}}(y) = y$ with corresponding $P(\mathcal{E}) = p$.

(b) when $p > 0.5$, we have $\hat{x}_{\mathrm{ML}}(y) = 1-y$ with corresponding $P(\mathcal{E}) = 1-p$.

(c) when $p = 0.5$, all four reasonable decoders have the same $P(\mathcal{E}) = 1/2$.

- In fact, when $p = 0.5$, the channel completely destroys any connection between $X$ and $Y$. In particular, in this scenario, $X \perp\!\!\!\perp Y$. So, the value of the observed $y$ is useless.

## 3.4 Optimal Block Decoding for Communications Over BSC

**3.45.** The decoding techniques (MAP and ML) discussed in the previous section can be extended to the case in which we simultaneously consider $n$ consecutive channel output symbols resulted from having $n$ input symbols.

Notation-wise, this simply means we consider an input-output vector pair $(\underline{\mathbf{X}}, \underline{\mathbf{Y}})$ instead of an input-output symbol pair $(X, Y)$

$$(X_1, X_2, \ldots, X_n) = \underline{X} \longrightarrow \boxed{DMC} \longrightarrow \underline{Y} = (Y_1, Y_2, \ldots, Y_n)$$

**3.46.** By the memoryless property of the channel,

$$Q\left(\underline{y}|\underline{x}\right) = Q(y_1|x_1) \times Q(y_2|x_2) \times \cdots \times Q(y_n|x_n).$$

**Example 3.47.** For a DMC in which $\mathcal{X} = \{0, 1\}$, $\mathcal{Y} = \{1, 2, 3\}$, $\mathbf{Q} =$

$\begin{array}{c} 0 \\ 1 \end{array} \left[ \begin{array}{ccc} 0.5 & 0.2 & 0.3 \\ 0.3 & 0.4 & 0.3 \end{array} \right]$, find

*Extended Q matrix*

(a) $Q(122|100)$ $= 0.2 \times 0.2 \times 0.3 = 0.012$

(b) $Q(333|111)$ $= (0.3)^3 = 0.027$

| $\underline{x}$ \ $\underline{y}$ | 111 | 112 | 121 | $\cdots$ | 122 | $\cdots$ | 333 |
|---|---|---|---|---|---|---|---|
| 000 | | | | | | | |
| 001 | | | | | | | |
| 010 | | | | | | | |
| $\vdots$ | | | | | | | |
| 100 | | | | | | 0.012 | |
| $\vdots$ | | | | | | | |
| 111 | | | | | | | 0.027 |

*☆ positions that $\underline{x}$ and $\underline{y}$ are the same $\nearrow^{n-d}$ different*

$d(\underline{x}, \underline{y})$

**Example 3.48.** For BSC, find

(a) $Q(101|100)$ $= p \times (1-p) \times (1-p) = (1-p)^2 p^1$

(b) $Q(111|111)$ $= (1-p)^3 p^0$

**3.49.** For BSC,

$$Q(y_i|x_i) = \begin{cases} p, & y_i \neq x_i, \\ 1 - p, & y_i = x_i. \end{cases}$$

Therefore,

$$Q\left(\underline{y}|\underline{x}\right) = p^{d\left(\underline{x},\underline{y}\right)}(1-p)^{n-d\left(\underline{x},\underline{y}\right)} = \left(\frac{p}{1-p}\right)^{d\left(\underline{x},\underline{y}\right)}(1-p)^n, \qquad (8)$$

where $d\left(\underline{\mathbf{x}}, \underline{\mathbf{y}}\right)$ is the number of coordinates in which the two blocks $\underline{\mathbf{x}}$ and $\underline{\mathbf{y}}$ differ.

*distance (difference)*

$d(101, 100) = 1$
$d(111, 111) = 0$
$d(00101, 01111) = 2$

*Hamming distance*

42

**3.50.** To recover the value of $\underline{\mathbf{x}}$ from the observed value of $\underline{\mathbf{y}}$, we can apply the vector version of what we studied about optimal decoder in the previous section.

- The optimal decoder is again given by the MAP decoder:

$$\hat{\underline{\mathbf{x}}}_{\text{MAP}}\left(\underline{\mathbf{y}}\right) = \arg\max_{\underline{\mathbf{x}}} Q\left(\underline{\mathbf{y}}\,|\underline{\mathbf{x}}\right) p\left(\underline{\mathbf{x}}\right). \tag{9}$$

- When the prior probabilities $p\left(\underline{\mathbf{x}}\right)$ is unknown or when we want simpler decoder, we may consider using the ML decoder:

$$\hat{\underline{\mathbf{x}}}_{\text{ML}}\left(\underline{\mathbf{y}}\right) = \arg\max_{\underline{\mathbf{x}}} Q\left(\underline{\mathbf{y}}\,|\underline{\mathbf{x}}\right). \tag{10}$$

*no $p(x)$ term in MLD*

Plugging-in

$$Q\left(\underline{\mathbf{y}}|\underline{\mathbf{x}}\right) = p^{d\left(\underline{\mathbf{x}},\underline{\mathbf{y}}\right)}(1-p)^{n-d\left(\underline{\mathbf{x}},\underline{\mathbf{y}}\right)} = \left(\frac{p}{1-p}\right)^{d\left(\underline{\mathbf{x}},\underline{\mathbf{y}}\right)}(1-p)^{n}, \tag{11}$$

from (8), gives

$$\hat{\underline{\mathbf{x}}}_{\text{MAP}}\left(\underline{\mathbf{y}}\right) = \arg\max_{\underline{\mathbf{x}}} \left(\frac{p}{1-p}\right)^{d\left(\underline{\mathbf{x}},\underline{\mathbf{y}}\right)}(1-p)^{n} p\left(\underline{\mathbf{x}}\right) \tag{12}$$

$$= \arg\max_{\underline{\mathbf{x}}} \left(\frac{p}{1-p}\right)^{d\left(\underline{\mathbf{x}},\underline{\mathbf{y}}\right)} p\left(\underline{\mathbf{x}}\right). \tag{13}$$

and

$$\hat{\underline{\mathbf{x}}}_{\text{ML}}\left(\underline{\mathbf{y}}\right) = \arg\max_{\underline{\mathbf{x}}} \left(\frac{p}{1-p}\right)^{d\left(\underline{\mathbf{x}},\underline{\mathbf{y}}\right)}. \tag{14}$$

**3.51.** <mark>**Minimum-distance decoder**</mark> as a ML decoding of block codes over BSC:

From (14) (or directly from (8)), note that when <mark>$p < 0.5$,</mark> which is usually the case for practical systems, we have $p < 1 - p$ and hence $0 < \frac{p}{1-p} < 1$. In which case, to maximize $Q\left(\underline{\mathbf{y}}|\underline{\mathbf{x}}\right)$, we need to minimize $d\left(\underline{\mathbf{x}},\underline{\mathbf{y}}\right)$. In other words, $\hat{\underline{\mathbf{x}}}_{\text{ML}}\left(\underline{\mathbf{y}}\right)$ should be the codeword $\underline{\mathbf{x}}$ which has the minimum distance from the observed $\underline{\mathbf{y}}$:

$$\hat{\underline{\mathbf{x}}}_{\text{ML}}\left(\underline{\mathbf{y}}\right) = \arg\min_{\underline{\mathbf{x}}} d\left(\underline{\mathbf{x}},\underline{\mathbf{y}}\right). \tag{15}$$

In conclusion, for block coding over BSC with <mark>$p < 0.5$,</mark> <mark>the ML decoder is the same as the minimum distance decoder.</mark>

## 3.5 Repetition Code for Channel Coding in Communications Over BSC

**3.52.** Recall that **channel coding** introduces, in a controlled manner, some *redundancy* in the (binary) information sequence that can be used at the receiver to overcome the effects of noise and interference encountered in the transmission of the signal through the channel.



**3.53.** **Repetition Code**: Repeat each bit $n$ times, where $n$ is some positive integer.

- Use the channel $n$ times to transmit 1 info-bit

- The (transmission) rate is $\frac{1}{n}$ [bpcu].

  ○ bpcu = bits per channel use

**3.54.** Two classes of channel codes

(a) Block codes

- To be discussed here.
- Realized by combinational/combinatorial circuit.

  *AND, OR, NOT, NAND, etc gates*
  *and*
  *wires*

(b) Convolutional codes

- Encoder has memory.
- Realized by sequential circuit. (Recall state diagram, flip-flop, etc.)

**Definition 3.55. Block Encoding**: Take $k$ (information) bits at a time and map each $k$-bit sequence into a (unique) $n$-bit sequence, called a **codeword**.



44

- The code is called $(n, k)$ code.

- Working with $k$-info-bit blocks means there are potentially $M = 2^k$ different information blocks.

  ○ The table that lists all the $2^k$ mapping from the $k$-bit info-block $\underline{\mathbf{s}}$ to the $n$-bit codeword $\underline{\mathbf{x}}$ is called the **codebook**.

  ○ The $M$ info-blocks are denoted by $\underline{\mathbf{s}}^{(1)}, \underline{\mathbf{s}}^{(2)}, \ldots, \underline{\mathbf{s}}^{(M)}$. The corresponding $M$ codewords are denoted by $\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}, \ldots, \underline{\mathbf{x}}^{(M)}$, respectively.

  | index $i$ | info-block $\underline{\mathbf{s}}$ | codeword $\underline{\mathbf{x}}$ |
  |---|---|---|
  | 1 | $\underline{\mathbf{s}}^{(1)} = 000 \ldots 0$ | $\underline{\mathbf{x}}^{(1)} =$ |
  | 2 | $\underline{\mathbf{s}}^{(2)} = 000 \ldots 1$ | $\underline{\mathbf{x}}^{(2)} =$ |
  | $\vdots$ | $\vdots$ | $\vdots$ |
  | $M$ | $\underline{\mathbf{s}}^{(M)} = 111 \ldots 1$ | $\underline{\mathbf{x}}^{(M)} =$ |

  ○ By the bijective mapping from $\underline{\mathbf{s}}$ to $\underline{\mathbf{x}}$,

$$p_i \equiv p\left(\underline{\mathbf{x}}^{(i)}\right) \equiv P\left[\underline{\mathbf{X}} = \underline{\mathbf{x}}^{(i)}\right] = P\left[\underline{\mathbf{S}} = \underline{\mathbf{s}}^{(i)}\right].$$

- To have unique codeword for each information block, we need $n \geq k$. Of course, with some redundancy added to combat the error introduced by the channel, we need $n > k$.

  ○ The amount of redundancy is measured by the ratio $\frac{n}{k}$.

  ○ The number of redundant bits is $r = n - k$.

- Here, we use the channel $n$ times to convey $k$ (information) bits.

  ○ The ratio $\frac{k}{n}$ is called the rate of the code or, simply, the **code rate**.

  ○ The (transmission) rate is $R = \frac{k}{n} = \frac{\log_2 M}{n}$ [bpcu].

**Example 3.56.** Find the codebook and code rate for the encoder which uses repetition code with $n = 5$.

codebook

| $\underline{s}$ | $\underline{x}$ |
|---|---|
| 0 | 00000 ≡ $\underline{0}$ |
| 1 | 11111 ≡ $\underline{1}$ |

45

**Example 3.57.** To get some idea about the difficulty of finding an optimal encoder, we need to consider the size of our search space. For $k = 5$ and $n = 10$, how many encoders are possible?

$$\binom{2^n}{2^k} = \binom{2^{10}}{2^5} = \binom{1024}{32} \approx 5 \times 10^{60}$$

| index | $\underline{s}$ | $\underline{x}$ |
|---|---|---|
| 1 | 00000 | - - - - - |
| 2 | 00001 | |
| ⋮ | 00010 | $n$ positions |
| ⋮ | ⋮ | |
| $2^k = 2^5 = 32$ | 11111 | |

**3.58. Decoding**: When the mapping from the information block $\underline{s}$ to the codeword $\underline{x}$ is invertible, the task of any decoder can be separated into two steps:

- First, find $\hat{\underline{x}}$ which is its guess of the $\underline{x}$ value based on the observed value of $\underline{y}$.

- Second, map $\hat{\underline{x}}$ back to the corresponding $\hat{\underline{s}}$ based on the codebook.

You may notice that it is more important to recover the index of the codeword than the codeword itself. Only its index is enough to indicate which info-block produced it.

**Example 3.59.** Repetition Code and Majority Voting: Back to Example 3.53.

Recall: ① MAP decoder is optimal ( min $P(\mathcal{E})$ )
② ML decoder is suboptimal. However, it can be optimal (the same as the MAPD) when the codewords are equally-likely.
③ ML decoder is the same as the min-distance decoder when the crossover probability $p$ of BSC is < 0.5.

Let $\underline{0}$ and $\underline{1}$ denote the $n$-dimensional row vectors $00\ldots0$ and $11\ldots1$, respectively. Observe that

$$d\left(\underline{x}, \underline{y}\right) = \begin{cases} \#1 \text{ in } \underline{y}, & \text{when } \underline{x} = \underline{0}, \\ \#0 \text{ in } \underline{y}, & \text{when } \underline{x} = \underline{1}. \end{cases}$$

Ex $\underline{y} = 01010$

Therefore, the `minimum distance decoder is`

$d(\underline{0}, \underline{y}) = 2 = \#1 \text{ in } \underline{y}$

$$\hat{\underline{x}}_{\text{ML}}\left(\underline{y}\right) = \begin{cases} \underline{0}, & \text{when } \#1 \text{ in } \underline{y} < \#0 \text{ in } \underline{y}, \\ \underline{1}, & \text{when } \#1 \text{ in } \underline{y} > \#0 \text{ in } \underline{y}. \end{cases}$$

$d(\underline{1}, \underline{y}) = 3 = \#0 \text{ in } \underline{y}$

Equivalently,

$$\hat{s}_{\text{ML}}\left(\underline{y}\right) = \begin{cases} 0, & \text{when } \#1 \text{ in } \underline{y} < \#0 \text{ in } \underline{y}, \\ 1, & \text{when } \#1 \text{ in } \underline{y} > \#0 \text{ in } \underline{y}. \end{cases}$$

This is the same as taking a majority vote among the received bit in the $\underline{\mathbf{y}}$ vector.

The corresponding error probability is

$$P\left(\mathcal{E}\right) = \sum_{c=\left\lceil \frac{n}{2} \right\rceil}^{n} \binom{n}{c} p^c (1-p)^{n-c}.$$

For example, when $p = 0.01$, we have $P(\mathcal{E}) \approx 10^{-5}$. Figure 6 compares the error probability when different values of $n$ are used.
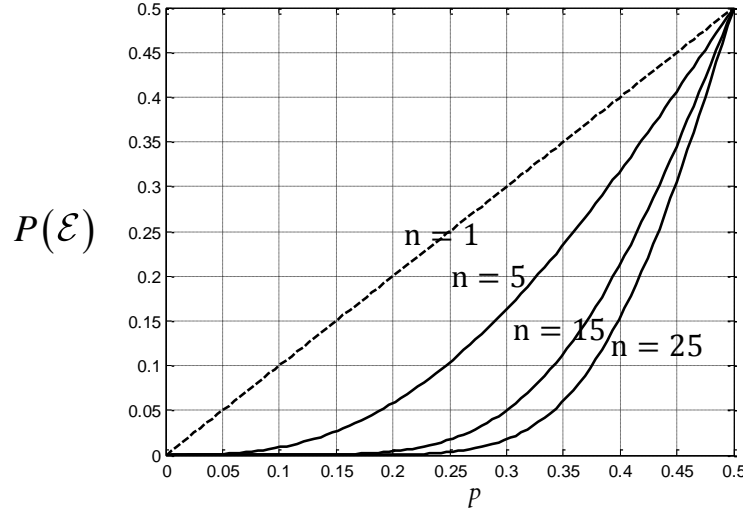


Figure 6: Error probability for a system that uses repetition code at the transmitter (repeat each info-bit $n$ times) and majority voting at the receiver. The channel is assumed to be binary symmetric with crossover probability $p$.

- Notice that the error probability decreases to 0 when $n$ is increased. It is then possible to transmit with arbitrarily low probability of error using this scheme.

- However, the (transmission) rate $R = \frac{k}{n} = \frac{1}{n}$ is also reduced as $n$ is increased.

So, in the limit, although we can have very small error probability, we suffer tiny (transmission) rate.

**3.60.** We may then ask "what is the maximum (transmission) rate of information that can be *reliably* transmitted over a communications channel?" Here, reliable communication means that the error probability can be made arbitrarily small. Shannon provided the solution to this question in his seminal work. We will revisit this question in the next section.

# 4 Mutual Information and Channel Capacity

In Section 2, we have seen the use of a quantity called entropy to measure the amount of randomness in a random variable. In this section, we introduce several more information-theoretic quantities. These quantities are important in the study of Shannon's results.

## 4.1 Information-Theoretic Quantities

**Definition 4.1.** Recall that, the **entropy** of a discrete random variable $X$ is defined in Definition 2.41 to be

*Amount of randomness in RV X*

$$\rightarrow H(X) = -\sum_{x \in S_X} p_X(x) \log_2 p_X(x) = -\mathbb{E}\left[\log_2 p_X(X)\right]. \tag{16}$$

Similarly, the entropy of a discrete random variable $Y$ is given by

$$H(Y) = -\sum_{y \in S_Y} p_Y(y) \log_2 p_Y(y) = -\mathbb{E}\left[\log_2 p_Y(Y)\right]. \tag{17}$$

In our context, the $X$ and $Y$ are input and output of a discrete memoryless channel, respectively. In such situation, we have introduced some new notations in Section 3.1:

$$
\begin{array}{c|c|c}
S_X \equiv \mathcal{X} & p_X(x) \equiv p(x) \rightsquigarrow \mathbf{p} & p_{Y|X}(y|x) \equiv Q(y|x) \rightsquigarrow \mathbf{Q} \text{ matrix} \\
S_Y \equiv \mathcal{Y} & p_Y(y) \equiv q(y) \rightsquigarrow \mathbf{q} & p_{X,Y}(x,y) \equiv p(x,y) \rightsquigarrow \mathbf{P} \text{ matrix}
\end{array}
$$

Under such notations, (16) and (17) become

$$H(X) = -\sum_{x \in \mathcal{X}} p(x) \log_2 p(x) = -\mathbb{E}\left[\log_2 p(X)\right] \tag{18}$$

and

$$H(Y) = -\sum_{y \in \mathcal{Y}} q(y) \log_2 q(y) = -\mathbb{E}\left[\log_2 q(Y)\right]. \tag{19}$$

**Definition 4.2.** The **joint entropy** for two random variables $X$ and $Y$ is given by

$$H(X,Y) = -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \log_2 p(x,y) = -\mathbb{E}\left[\log_2 p(X,Y)\right].$$

*all possible pairs (x,y)*

48

**Example 4.3.** Random variables $X$ and $Y$ have the following joint pmf matrix $\mathbf{P}$:

$$
P = \begin{array}{c c}
 & \begin{array}{cccc} Y_1 & Y_2 & Y_3 & Y_4 \end{array} \\
\begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array} &
\left[ \begin{array}{cccc}
\frac{1}{8} & \frac{1}{16} & \frac{1}{16} & \frac{1}{4} \\
\frac{1}{16} & \frac{1}{8} & \frac{1}{16} & 0 \\
\frac{1}{32} & \frac{1}{32} & \frac{1}{16} & 0 \\
\frac{1}{32} & \frac{1}{32} & \frac{1}{16} & 0
\end{array} \right]
\end{array}
\quad
\begin{array}{l} \xrightarrow{\;\Sigma\;} \; 1/2 \\ \xrightarrow{\;\Sigma\;} \; 1/4 \\ \xrightarrow{\;\Sigma\;} \; 1/8 \\ \xrightarrow{\;\Sigma\;} \; 1/8 \end{array} \; p(x)
$$

$$\Sigma\downarrow \; \Sigma\downarrow \; \Sigma\downarrow \; \Sigma\downarrow$$
$$1/4 \quad 1/4 \quad 1/4 \quad 1/4$$

Find $H(X)$, $H(Y)$ and $H(X,Y)$.

$$H(X) = -\tfrac{1}{2}\log_2 \tfrac{1}{2} - \tfrac{1}{4}\log_2 \tfrac{1}{4} - 2\times\tfrac{1}{8}\log_2\tfrac{1}{8} = \tfrac{1}{2} + \tfrac{2}{4} + \tfrac{6}{8} = \tfrac{7}{4} \;\text{ bits (per symbol)}$$

$$H(Y) = -4\times\tfrac{1}{4}\log_2\tfrac{1}{4} = 2 \;\text{ bits (per symbol)}$$

$$H(X,Y) = -\tfrac{1}{8}\log_2\tfrac{1}{8} - \tfrac{1}{16}\log_2\tfrac{1}{16} - \tfrac{1}{16}\log_2\tfrac{1}{16} \cdots \qquad (16 \text{ terms})$$

$$= \left(-\tfrac{1}{4}\log_2\tfrac{1}{4}\right) + 2\left(-\tfrac{1}{8}\log_2\tfrac{1}{8}\right) + 6\left(-\tfrac{1}{16}\log_2\tfrac{1}{16}\right) + 4\left(-\tfrac{1}{32}\log_2\tfrac{1}{32}\right)$$

$$= \tfrac{27}{8} \;\text{ bits (per symbol)}$$

**Definition 4.4.** The (conditional) entropy of $Y$ when we know $X = x$ is denoted by $H\left(Y\,|X=x\right)$ or simply $H(Y|x)$. It can be calculated from

The amount of randomness still remained in $Y$ when we know that $X = x$.

$$\boxed{H\left(Y\,|x\right)} = -\sum_{y\in\mathcal{Y}} Q\left(y\,|x\right)\log_2 Q\left(y\,|x\right)$$

- Note that the above formula is what we should expect it to be. When we want to find the entropy of $Y$, we use (19):

$$H\left(Y\right) = -\sum_{y\in\mathcal{Y}} q\left(y\right)\log_2 q\left(y\right).$$

When we have an extra piece of information that $X = x$, we should update the probability about $Y$ from the unconditional probability $q(y)$ to the conditional probability $Q(y|x)$.

- Note that when we consider $Q(y|x)$ with the value of $x$ fixed and the value of $y$ varied, we simply get the whole $x$-row from $\mathbf{Q}$ matrix. So, to

49

find $H(Y|x)$, we simply find the "usual" entropy from the probability values in the row corresponding to $x$ in the **Q** matrix.

**Example 4.5.** Consider the following DMC (actually BAC)

$H(x) = \log_2 2 = 1$

$-\frac{1}{2}\log_2\frac{1}{2} + -\frac{1}{2}\log_2\frac{1}{2}$

$p(0) = 1/2$  $0 \xrightarrow{1/2} 0$

$X \longrightarrow$   $\xrightarrow{1/2}$   $\longrightarrow Y$

$p(1) = 1/2$  $1 \xrightarrow{1} 1$

$H(X,Y) = \left(-\frac{1}{4}\log_2\frac{1}{4}\right) \times 2 - \frac{1}{2}\log_2\frac{1}{2}$

$= \frac{3}{2}$

$H(Y) = -\frac{1}{4}\log_2\frac{1}{4} + $

$-\frac{3}{4}\log_2\frac{3}{4}$

$Q = \begin{array}{c} x\backslash y \\ 0 \\ 1 \end{array} \begin{bmatrix} 0 & 1 \\ 1/2 & 1/2 \\ 0 & 1 \end{bmatrix}$

$\xrightarrow{\times 1/2}$

$\xrightarrow[\times 1/2]{}$

$\begin{array}{c} x\backslash y \\ 0 \\ 1 \end{array} \begin{bmatrix} 0 & 1 \\ 1/4 & 2/4 \\ 0 & 1/2 \end{bmatrix} = P$

$\Sigma\downarrow \qquad \downarrow\Sigma$

$1/4 \qquad 3/4$

Originally $P[Y = y] = q(y) = \begin{cases} 1/4, & y = 0, \\ 3/4, & y = 1, \\ 0, & \text{otherwise.} \end{cases}$

(a) Suppose we know that $X = 0$.

The "x = 0" row in the **Q** matrix gives $Q(y|0) = \begin{cases} 1/2, & y = 0, 1, \\ 0, & \text{otherwise;} \end{cases}$

that is, given $x = 0$, the RV $Y$ will be uniform.

$H(Y|0) = \log_2 2 = 1$

(b) Suppose we know that $X = 1$. The "x = 1" row in the **Q** matrix

gives $Q(y|1) = \begin{cases} 1, & y = 1, \\ 0, & \text{otherwise;} \end{cases}$  that is, given $x = 1$, the RV $Y$ is

degenerated (deterministic).

$H(Y|1) = 0$

**Definition 4.6. Conditional entropy**: The (average) conditional entropy of $Y$ when we know $X$ is denoted by $H(Y|X)$. It can be calculated from
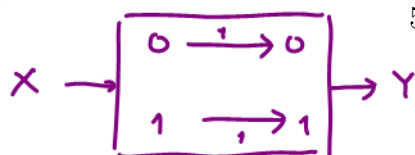
The average amount of randomness still remained in Y when we know X

$$\boxed{H(Y|X)} = \sum_{x \in \mathcal{X}} p(x) H(Y|x).$$

**Example 4.7.** In Example 4.5,

$H(Y|X) = p(0)H(Y|0) + p(1)H(Y|1) = \frac{1}{2}\times 1 + \frac{1}{2}\times 0 = \frac{1}{2}$

Easy Example

$X \longrightarrow \boxed{\begin{array}{c} 0 \xrightarrow{1} 0 \\[4pt] 1 \xrightarrow{1} 1 \end{array}} \longrightarrow Y$

50

$\left. \begin{array}{l} H(Y|0) = 0 \\ H(Y|1) = 0 \end{array} \right\} \rightarrow H(Y|X) = \sum p(x) \times 0$

$= 0.$

**4.8.** An alternative way to calculate $H(Y|X)$ can be derived by first rewriting it as

$$H(Y|X) = \sum_{x \in \mathcal{X}} p(x) H(Y|x) = -\sum_{x \in \mathcal{X}} p(x) \sum_{y \in \mathcal{Y}} Q(y|x) \log_2 Q(y|x)$$

$$= -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \log_2 Q(y|x) = -\mathbb{E}[\log_2 Q(Y|X)]$$

Note that $Q(y|x) = \frac{p(x,y)}{p(x)}$. Therefore,

Compare this with

P(A\B) = P(A∪B) - P(B)

$$\boxed{H(Y|X)} = -\mathbb{E}[\log_2 Q(Y|X)] = -\mathbb{E}\left[\log_2 \frac{p(X,Y)}{p(X)}\right]$$

$$= (-\mathbb{E}[\log_2 p(X,Y)]) - (-\mathbb{E}[\log_2 p(X)])$$

$$= \boxed{H(X,Y) - H(X)}$$

**Example 4.9.** In Example 4.5,

H(Y|x) = H(x,Y) - H(x) = $\frac{3}{2}$ - 1 = $\frac{1}{2}$ .

**Example 4.10.** Continue from Example 4.3. Random variables $X$ and $Y$ have the following joint pmf matrix $\mathbf{P}$:

$$P = \begin{bmatrix} \frac{1}{8} & \frac{1}{16} & \frac{1}{16} & \frac{1}{4} \\ \frac{1}{16} & \frac{1}{8} & \frac{1}{16} & 0 \\ \frac{1}{32} & \frac{1}{32} & \frac{1}{16} & 0 \\ \frac{1}{32} & \frac{1}{32} & \frac{1}{16} & 0 \end{bmatrix} \longrightarrow \begin{bmatrix} 1/4 & 1/2 & 1/4 & 0 \end{bmatrix} \quad Q$$

Find $H(Y|X)$ and $H(X|Y)$.

From 4.3
$$\begin{cases} H(x) = \frac{7}{4} \\ H(Y) = 2 \\ H(x,Y) = \frac{27}{8} \end{cases}$$

$$H(Y|x_2) = \left(-\frac{1}{4}\log_2\frac{1}{4}\right) \times 2 - \frac{1}{2}\log_2\frac{1}{2}$$
$$= 3/2$$

$$H(Y|X) = H(x,Y) - H(x)$$
$$= \frac{27}{8} - \frac{14}{8} = \frac{13}{8}$$

$$H(x|Y) = H(Y,x) - H(Y)$$
$$= \frac{27}{8} - 2 = \frac{11}{8}$$

$$I(X;Y) = \frac{3}{8}$$

$H(x) = \frac{14}{8}$   $H(Y) = \frac{16}{8}$

$H(Y|x) = \frac{13}{8}$

$H(x|Y) = \frac{11}{8}$

= H(x) + H(Y) - H(x,Y)
= H(x) - H(x|Y)
= H(Y) - H(Y|x)

51

# ECS 452: In-Class Exercise # _

## Instructions

1. Separate into groups of no more than three persons.
2. The group cannot be the same as your former group.
3. Only one submission is needed for each group.
4. **Write down all the steps** that you have done to obtain your answers. You may not get full credit even when your answer is correct without showing how you get your answer.
5. **Do not panic.**

1. Consider two random variables $X$ and $Y$ whose joint pmf matrix is given by

$$P = \begin{array}{c|cc} {}_X \backslash {}^y & 3 & 5 \\ \hline 1 & 1/18 & 5/18 \\ 2 & 4/9 & 2/9 \end{array} \quad \begin{array}{l} \xrightarrow{\Sigma} \frac{6}{18} = \frac{1}{3} \\ \xrightarrow{\Sigma} \frac{6}{9} = \frac{2}{3} \end{array} \quad p(x)$$

$$g(y) \quad \begin{array}{cc} \Sigma\downarrow & \Sigma\downarrow \\ \frac{9}{18} & \frac{9}{18} \end{array}$$

Calculate the following quantities. Except for part (e), all of your answers should be of the form X.XXXX.

a. $H(X,Y) = -\frac{1}{18}\log_2\frac{1}{18} - \frac{5}{18}\log_2\frac{5}{18} - \frac{4}{9}\log_2\frac{4}{9} - \frac{2}{9}\log_2\frac{2}{9}$

  $\approx 1.7472$ bits (per symbol)

b. $H(X) = H\left(\begin{bmatrix} 1/3 & 2/3 \end{bmatrix}\right) = -\frac{1}{3}\log_2\frac{1}{3} - \frac{2}{3}\log_2\frac{2}{3} \approx 0.9183$ bits (per symbol)

c. $H(Y) = H\left(\begin{bmatrix} 1/2 & 1/2 \end{bmatrix}\right) = \log_2 2 = 1.000$ bit (per symbol)

  └ For uniform RV, the entropy is simply $\log_2$ of the size of its support.

d. $H(Y|X) = H(X,Y) - H(X) \approx 1.7472 - 0.9183$

  $\approx 0.8289$ bits (per symbol)

  $H(Y|X) = p(1)H(Y|1) + p(2)H(Y|2)$
   ↑ 1/3      ↑ 2/3

  $\approx 0.8289$ bits (per symbol)

  $H(Y|1) = H\left(\begin{bmatrix} 1/6 & 5/6 \end{bmatrix}\right) \approx 0.6500$

e. **Q** matrix

$$P = \begin{bmatrix} 1/18 & 5/18 \\ 4/9 & 2/9 \end{bmatrix} \quad \begin{array}{c} \times \frac{1}{p(x)} \\ \times 3 \\ \xrightarrow{\phantom{xx}} \\ \times 3/2 \\ \xrightarrow{\phantom{xx}} \end{array} \quad \begin{array}{c|cc} {}_x\backslash{}^Y & 3 & 5 \\ \hline 1 & 1/6 & 5/6 \\ 2 & 2/3 & 1/3 \end{array} = Q$$

  $H(Y|2) = H\left(\begin{bmatrix} 2/3 & 1/3 \end{bmatrix}\right) \approx 0.9183$

**Definition 4.11.** The **mutual information**[13] $I(X;Y)$ between two random variables $X$ and $Y$ is defined as

$$I(X;Y) = H(X) - H(X|Y) \tag{20}$$
$$= H(Y) - H(Y|X) \tag{21}$$
$$= H(X) + H(Y) - H(X,Y) \tag{22}$$
$$= \mathbb{E}\left[\log_2 \frac{p(X,Y)}{p(X)q(Y)}\right] = \sum_{x \in \mathcal{X}}\sum_{y \in \mathcal{Y}} p(x,y) \log \frac{p(x,y)}{p(x)q(y)} \tag{23}$$
$$= \mathbb{E}\left[\log_2 \frac{P_{X|Y}(X|Y)}{p(X)}\right] = \mathbb{E}\left[\log_2 \frac{Q(Y|X)}{q(Y)}\right]. \tag{24}$$

*Interpretation ※1*

- Mutual information quantifies the reduction in the uncertainty of one random variable due to the knowledge of another.

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

*Amount of reduction in the randomness of X due to the knowledge of Y.*

*Amount of randomness in X*

*→ Amount of randomness still remained in X when we know Y.*

*Interpretation ※2*

- Mutual information is a measure of the amount of information one random variable contains about another [5, p 13].

*Interpretation ※3*

- It is also natural to think of $I(X;Y)$ as a measure of how far $X$ and $Y$ are from being independent.

  ∘ Technically, it is the (Kullback-Leibler) divergence between the joint and product-of-marginal distributions.

**4.12.** Some important properties

(a) $H(X,Y) = H(Y,X)$ and $I(X;Y) = I(Y;X)$.
    However, in general, $H(X|Y) \neq H(Y|X)$.

(b) $I$ and $H$ are always $\geq 0$.

(c) There is a one-to-one correspondence between Shannon's information measures and set theory. We may use an **information diagram**, which

---

[13]The name mutual information and the notation $I(X;Y)$ was introduced by [Fano, 1961, Ch 2].

is a variation of a Venn diagram, to represent relationship between Shannon's information measures. This is similar to the use of the Venn diagram to represent relationship between probability measures. These diagrams are shown in Figure 7.
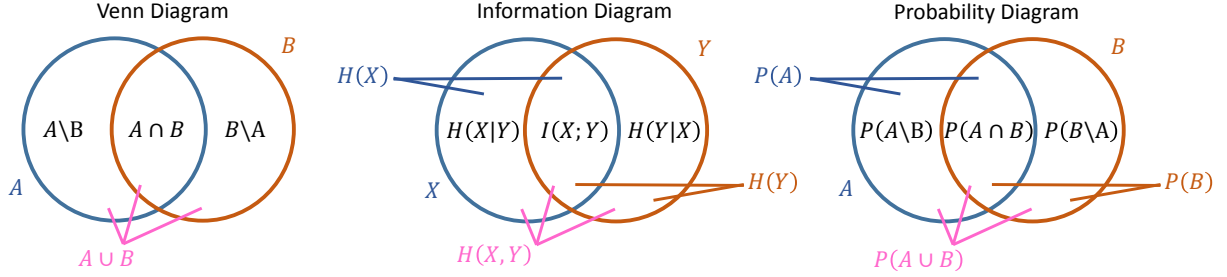


Figure 7: Venn diagram and its use to represent relationship between information measures and relationship between probabilities

- Chain rule for information measures:

$$H\left(X,Y\right) = H\left(X\right) + H\left(Y\,|X\,\right) = H\left(Y\right) + H\left(X\,|Y\,\right).$$

(d) $I(X;Y) \geq 0$ with equality if and only if $X$ and $Y$ are independent.

- When this property is applied to the information diagram (or definitions (20), (21), and (22) for $I(X,Y)$), we have

  (i) $H(X|Y) \leq H(X),$ } Conditioning only reduces entropy
  (ii) $H(Y|X) \leq H(Y),$
  (iii) $H(X,Y) \leq H(X) + H(Y)$

  Moreover, each of the inequalities above becomes equality if and only $X \perp\!\!\!\perp Y$.

(e) We have seen in Section 2.4 that

$$\underset{\text{deterministic (degenerated)}}{0} \leq H\left(X\right) \leq \underset{\text{uniform}}{\log_2 |\mathcal{X}|}. \tag{25}$$

Similarly,

$$\underset{\text{deterministic (degenerated)}}{0} \leq H\left(Y\right) \leq \underset{\text{uniform}}{\log_2 |\mathcal{Y}|}. \tag{26}$$

53

For conditional entropy, we have

$$\underset{\substack{\exists g\ Y=g(X)}}{0} \le H\left(Y\,|X\right) \le \underset{\substack{X \perp\!\!\!\perp Y}}{H\left(Y\right)} \tag{27}$$

and

$$\underset{\substack{\exists g\ X=g(Y)}}{0} \le H\left(X\,|Y\right) \le \underset{\substack{X \perp\!\!\!\perp Y}}{H\left(X\right)}. \tag{28}$$

For mutual information, we have

$$\underset{\substack{X \perp\!\!\!\perp Y}}{0} \le I\left(X;Y\right) \le \underset{\substack{\exists g\ X=g(Y)}}{H\left(X\right)} \tag{29}$$

and

$$\underset{\substack{X \perp\!\!\!\perp Y}}{0} \le I\left(X;Y\right) \le \underset{\substack{\exists g\ Y=g(X)}}{H\left(Y\right)}. \tag{30}$$

Combining 25, 26, 29, and 30, we have

$$0 \le I\left(X;Y\right) \le \min\left\{H\left(X\right), H\left(Y\right)\right\} \le \min\left\{\log_2 |\mathcal{X}|, \log_2 |\mathcal{Y}|\right\} \tag{31}$$

(f) $H\left(X\,|X\right) = 0$ and $I(X;X) = H(X)$.

**Example 4.13.** Find the mutual information $I(X;Y)$ between the two random variables $X$ and $Y$ whose joint pmf matrix is given by $\mathbf{P} = \begin{bmatrix} \frac{1}{2} & \frac{1}{4} \\ \frac{1}{4} & 0 \end{bmatrix}$

$\xrightarrow{\quad} 3/4$
$\xrightarrow{\quad} 1/4$
$\underset{2}{\downarrow}\ \underset{2}{\downarrow}$
$\frac{3}{4}\ \frac{1}{4}$

$I(X;Y) = H(X) + H(Y) - H(X,Y) \approx 0.1226$

$\llcorner -\frac{1}{2}\log_2\frac{1}{2} - \left(\frac{1}{4}\log_2\frac{1}{4}\right)\times 2 = 1.5$

$H(X) = H\left(\left[\,3/4 \quad 1/4\,\right]\right) = -\frac{3}{4}\log_2\frac{3}{4} \ -\frac{1}{4}\log_2\frac{1}{4} \approx 0.8113$

**Example 4.14.** Find the mutual information $I(X;Y)$ between the two random variables $X$ and $Y$ whose $\mathbf{p} = \left[\frac{1}{4}, \frac{3}{4}\right]$ and $\mathbf{Q} = \begin{bmatrix} \frac{1}{4} & \frac{3}{4} \\ \frac{3}{4} & \frac{1}{4} \end{bmatrix}$.

Method ① $\quad Q = \left[\quad\quad\right] \xrightarrow[\times 3/4]{\times 1/4} \left[\quad\quad\right] = P$

$\gamma = p\,Q$
$= \left[\frac{5}{8} \quad \frac{3}{8}\right]$

Method ② $\quad H(Y|x_1) = H\left(\left[1/4 \quad 3/4\right]\right) \approx 0.8113$

$H(Y|x_2) = H\left(\left[3/4 \quad 1/4\right]\right) \approx 0.8113$

$H(Y|X) = p(x_1)H(Y|x_1) + p(x_2)H(Y|x_2) \approx 0.8113$

$\quad\quad\quad\quad\quad\quad 1/4 \quad\quad 54 \quad\quad\quad 3/4$

$I(X;Y) = H(Y) - H(Y|X) \approx 0.1432$

$\quad\quad\quad 0.9544$

## 4.2 Operational Channel Capacity

**4.15.** In Section 3, we have studied how to compute the error probability $P(\mathcal{E})$ for digital communication systems over DMC. At the end of that section, we studied block encoding where the channel is used $n$ times to transmit a $k$-bit info-block.

In this section, our consideration is "reverse".

**4.16.** In this and the next subsections, we introduce a quantity called channel capacity which is crucial in benchmarking communication system. Recall that, in Section 2 where source coding was discussed, we were interested in the minimum rate (in bits per source symbol) to represent a source. Here, we are interested in the maximum rate (in bits per channel use) that can be sent through a given channel *reliably*.

**4.17.** Here, **reliable communication** means *arbitrarily* small error probability *can be achieved.*   *[positive]*

- This seems to be an impossible goal.

  - If the channel introduces errors, how can one correct them all?

    * Any correction process is also subject to error, ad infinitum.

**Definition 4.18.** Given a DMC, its **"operational" channel capacity** is the maximum rate at which *reliable* communication over the channel *is possible.*   *[by using appropriate block length $n$, channel encoder, channel decoder.]*

- The channel capacity is the maximum rate in bits per channel use at which information *can be* sent with **arbitrarily low** error probability.

**4.19.** Claude Shannon showed, in his 1948 landmark paper, that this operational channel capacity is the same as the information channel capacity which we will discuss in the next subsection. From this, we can omit the

words "operational" and "information" and simply refer to both quantities as the *channel capacity*.

**Example 4.20.** In Example 4.34, we will find that the capacity of a BSC with crossover probability $p = 0.1$ is approximately 0.531 bits per channel use. This means that for any rate $R < 0.531$ and any error probability $P(\mathcal{E})$ that we desire, as long as it is greater than 0, we can find a suitable $n$, a rate $R$ encoder, and a corresponding decoder which will yield an error probability that is at least as low as our set value.

- Usually, for very low desired value of $P(\mathcal{E})$, we may need large value of $n$.

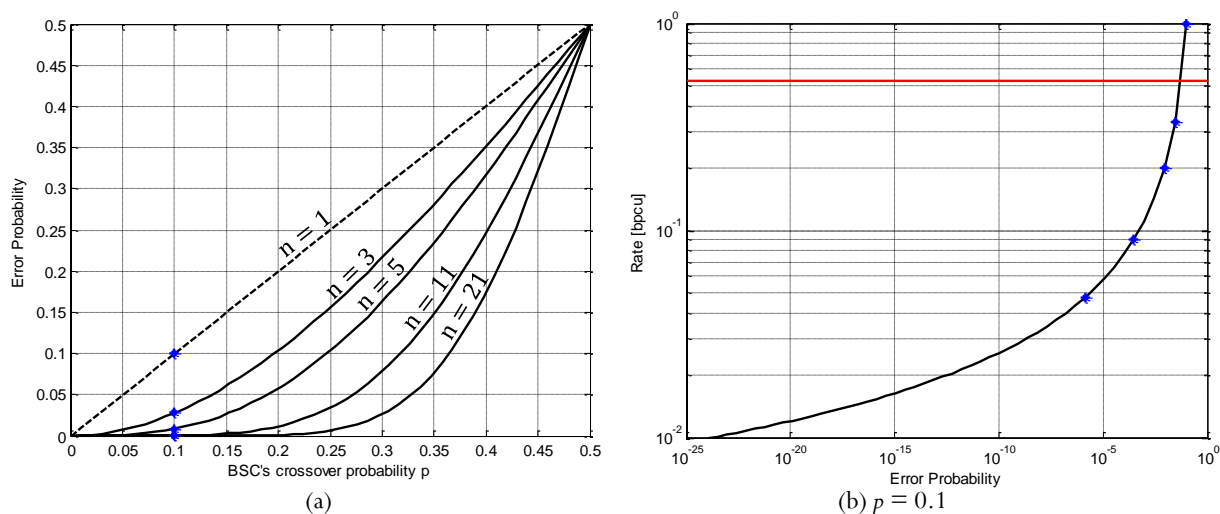**Example 4.21.** Repetition code is not good enough.



Figure 8: Performance of repetition coding with majority voting at the decoder

- From Figure 8b, it is clear that repetition coding can not reliably achieve such transmission rate. In fact, when we require error probability to be less than 0.1, the required repetition code needs $n \geq 3$. (For simplicity, let's assume only odd value of $n$ can be used here.) However, this means the rate is $\leq \frac{1}{3} \approx 0.33$ which is already significantly less than 0.531.

- In fact, for any rate $> 0$, we can see from Figure 8b that communication system based on repetition coding is not "reliable" according

56

to Definition 4.17. For example, for rate $= 0.02$ bits per channel use, repetition code can't satisfy the requirement that the error probability must be less than $10^{-15}$. In fact, Figure 8b shows that as we reduce the error probability to 0, the rate also goes to 0 as well. Therefore, there is no positive rate that works for all error probability.

- However, because the channel capacity is 0.531 [bpcu], there must exist other encoding techniques which give better error probability than repetition code.

  ○ Although Shannon's result gives us the channel capacity, it does not give us any explicit instruction on how to construct codes which can achieve that value.

## 4.3    Information Channel Capacity

**4.22.** In Section 4.1, we have studied how to compute the value of mutual information $I(X;Y)$ between two random variables $X$ and $Y$. Recall that, here, $X$ and $Y$ are the channel input and output, respectively. We have also seen, in Example 4.13, how to compute $I(X;Y)$ when the joint pmf matrix $\mathbf{P}$ is given. Furthermore, we have also worked on Example 4.14 in which the value of mutual information is computed from the prior probability vector $\underline{\mathbf{p}}$ and the channel transition probability matrix $\mathbf{Q}$. This second type of calculation is crucial in the computation of channel capacity. This kind of calculation is so important that we may write the mutual information $I(X;Y)$ as $I(\underline{\mathbf{p}}, \mathbf{Q})$.

**Definition 4.23.** Given a DMC channel, we define its "information" channel capacity as

$$C = \max_{\underline{\mathbf{p}}} I\left(X;Y\right) = \max_{\underline{\mathbf{p}}} I\left(\underline{\mathbf{p}}, \mathbf{Q}\right), \tag{32}$$

where the maximum is taken over all possible input pmfs $\underline{\mathbf{p}}$.

- Again, as mentioned in 4.19, Shannon showed that the "information" channel capacity defined here is equal to the "operational" channel capacity defined in Definition 4.18.

  ○ Thus, we may drop the word "information" in most discussions of channel capacity.

**Example 4.24.** The capacity of a BAC whose $Q(1|0) = 0.9$ and $Q(0|1) = 0.4$ can be found by first realizing that $I(X;Y)$ here is a function of a single variable: $p_0$. The plot of $I(X;Y)$ as a function of $p_0$ gives some rough estimates of the answers. One can directly solve for the optimal $p_0$ by simply taking derivative of $I(X;Y)$ and set it equal to 0. This gives the capacity value of 0.0918 bpcu which is achieved by $\underline{\mathbf{p}} = [0.5376, \ 0.4624]$.
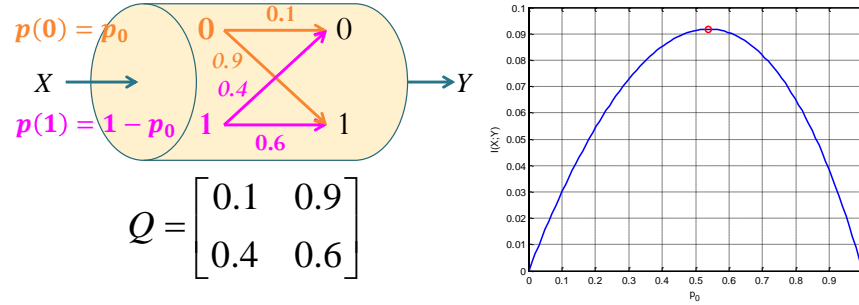


Figure 9: Maximization of mutual information to find capacity of a BAC channel. Capacity of 0.0918 bits is achieved by $\underline{\mathbf{p}} = [0.5376, \ 0.4624]$
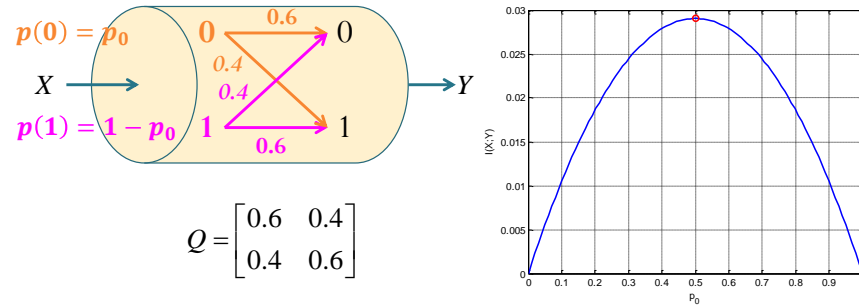


Figure 10: Maximization of mutual information to find capacity of a BSC channel. Capacity of 0.029 bits is achieved by $\underline{\mathbf{p}} = [0.5, \ 0.5]$

**4.25. Blahut-Arimoto Algorithm** [5, Section 10.8]: Alternatively, in 1972, Arimoto [1] and Blahut [2] independently developed an iterative algorithm to help us approximate the pmf $\underline{\mathbf{p}}^*$ which achieves capacity $C$. To do this, start with any (guess) input pmf $p_0(x)$, define a sequence of pmfs $p_r(x)$, $r = 0, 1, \ldots$ according to the following iterative prescription:

(a) $q_r(y) = \sum\limits_{x} p_r(x) Q(y|x)$ for all $y \in \mathcal{Y}$.

(b) $c_r(x) = 2^{\left( \sum\limits_{y} Q(y|x) \log_2 \frac{Q(y|x)}{q_r(y)} \right)}$ for all $x \in \mathcal{X}$.

58

(c) It can be shown that

$$\log_2 \left( \sum_x p_r(x) c_r(x) \right) \leq C \leq \log_2 \left( \max_x c_r(x) \right).$$

- If the lower-bound and upper-bound above are close enough. We take $p_r(x)$ as our answer and the corresponding capacity is simply the average of the two bounds.
- Otherwise, we compute the pmf

$$p_{r+1}(x) = \frac{p_r(x) c_r(x)}{\sum_x p_r(x) c_r(x)} \quad \text{for all } x \in \mathcal{X}$$

and repeat the steps above with index $r$ replaced by $r+1$.

## 4.4 Special Cases for Calculation of Channel Capacity

In this section, we study special cases of DMC whose capacity values can be found (relatively) easy.

**Example 4.26.** Find the channel capacity of a noiseless binary channel (a BSC whose crossover probability is $p = 0$).



To maximize $I(X;Y)$,
we maximize $H(X)$
by using
$P = [\tfrac{1}{2} \ \tfrac{1}{2}]$.

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

$H(X|Y) = 0 \leftarrow$ Intuitively, knowing $Y$, the value of $X$ is completely determined.

$C = \log_2 |\mathcal{X}| = \log_2 2$
$= 1$ (bpcu)
bit per channel use

$I(X;Y) = H(X)$

**Example 4.27.** Noisy Channel with Nonoverlapping Outputs: Find the channel capacity of a DMC whose

Ex. 4.27 b

$$Q = \begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \begin{bmatrix} 1/8 & 7/8 & 0 & 0 & 0 \\ 0 & 0 & 1/3 & 2/3 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$
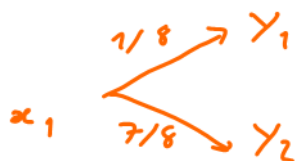
$x \backslash y \quad Y_1 \quad Y_2 \quad Y_3 \quad Y_4 \quad Y_5$

$I(X;Y) = H(X) - \underbrace{H(X|Y)}_{0}$

$\leq \log_2 |\mathcal{X}| = \log_2 2 = 1$
with equality iff $X$ is uniform
$P = [\tfrac{1}{2} \ \tfrac{1}{2}]$

$C = 1$ bpcu which happens when $P = [\tfrac{1}{2} \ \tfrac{1}{2}]$.

$C = \log_2 |\mathcal{X}|$
$= \log_2 3$ bpcu
is achieved by
$P = [\tfrac{1}{3} \ \tfrac{1}{3} \ \tfrac{1}{3}]$.



59

In this example, the channel appears to be noisy, but really is not. Even though the output of the channel is a random consequence of the input, the input can be determined from the output, and hence every transmitted bit can be recovered without error.

**4.28.** Reminder:

(a) Some definitions involving entropy

    (i) Binary entropy function: $h(p) = -p \log_2 p - (1-p) \log_2 (1-p)$

    (ii) $H(X) = -\sum_x p(x) \log_2 p(x)$

    (iii) $H(\underline{\mathbf{p}}) = -\sum_i p_i \log_2 (p_i)$

(b) A key entropy property that will be used frequently in this section is that for any random variable $X$,

<span>max</span>

$H(X) \leq \log_2 |\mathcal{X}|$ with equality iff $X$ is uniform. $[NO^2]$

**4.29.** A DMC is a **noisy channel with nonoverlapping outputs** when there is only one non-zero element in each column of its $\mathbf{Q}$ matrix. For such channel,

$$C = \log_2 |\mathcal{X}| \text{ is achieved by uniform } p(x).$$

**Definition 4.30.** A DMC is called **symmetric** if (1) all the rows of its probability transition matrix $\mathbf{Q}$ are permutations of each other and (2) so are the columns.

**Example 4.31.** For each of the following $\mathbf{Q}$, is the corresponding DMC symmetric?

$$\frac{1}{6} + \frac{3}{6} = \frac{4}{6}$$

$$\begin{bmatrix} 0.3 & 0.2 & 0.5 \\ 0.5 & 0.3 & 0.2 \\ 0.2 & 0.5 & 0.3 \end{bmatrix}, \quad \begin{bmatrix} 0.3 & 0.2 & 0.5 \\ 0.5 & 0.2 & 0.3 \\ 0.2 & 0.5 & 0.3 \end{bmatrix}, \quad \begin{bmatrix} 1/3 & 1/6 & 1/2 \\ 1/3 & 1/2 & 1/6 \end{bmatrix}, \quad \begin{bmatrix} 0.1 & 0.9 \\ 0.4 & 0.6 \end{bmatrix}$$

①✓ symmetric    ①✓ not symmetric    ① not symmetric    ①✗ not symmetric
②✓            ②✗          ②✗       ②✗

⇓            ⇑        ② ✓ 2/3   2/3   2/3

weakly symmetric   not weakly symmetric   weakly symmetric    not weakly symmetric

**4.32.** Q: Does symmetric DMC always have square $\mathbf{Q}$?

$$\begin{bmatrix} 1/5 & 1/5 & 1/5 & 1/5 & 1/5 \end{bmatrix}$$

$$\begin{bmatrix} 1/2 & 0 & 1/2 & 0 \\ 0 & 1/2 & 0 & 1/2 \end{bmatrix}$$

$$\begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix}$$

$$\begin{bmatrix} a & a & b & b \\ b & b & a & a \end{bmatrix}$$

where $a + b = \frac{1}{2}$

<span>60</span>

**Example 4.33.** Find the channel capacity of a DMC whose

$$Q = \begin{array}{c} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{array}\begin{bmatrix} 0.3 & 0.2 & 0.5 \\ 0.5 & 0.3 & 0.2 \\ 0.2 & 0.5 & 0.3 \end{bmatrix}.$$

*(handwritten: column labels $y_1$ $y_2$ $y_3$, $r$ over first row circled; diagram with $x_1, x_2, x_3 \to y_1, y_2, y_3$ with values 0.3, 0.5, 0.2, 0.5)*

Solution: First, recall that the capacity $C$ of a given DMC can be found by (32):

*(handwritten: Earlier, we used $I(X;Y) = H(X) - H(X|Y)$)*

$$C = \max_{\underline{p}} I(X;Y) = \max_{\underline{p}} I(\underline{p}, Q).$$

*(handwritten: was 0 / now, $\neq 0$.)*

*(handwritten:* $I(X;Y) = H(Y) - H(Y|X)$
$= \sum_x p(x) H(Y|x) = \sum_x p(x) H(\underline{r}) = H(\underline{r})$*)*

We see that, to maximize $I(X;Y)$, we need to maximize $H(Y)$. Of course, we know that the maximum value of $H(Y)$ is $\log_2 |\mathcal{Y}|$ which happens when $Y$ is uniform. Therefore, if we can find $\underline{p}$ which makes $Y$ uniform, then this same $\underline{p}$ will give the channel capacity.

*(handwritten left margin: first column sum / second column sum)*

*(handwritten:*
$\underline{g_5} = \underline{p} Q$

$\left[\tfrac{1}{3}\times 1, \tfrac{1}{3}\times 1, \tfrac{1}{3}\times 1\right] = \left[\tfrac{1}{3} \ \tfrac{1}{3} \ \tfrac{1}{3}\right]\begin{bmatrix} 0.3 & 0.2 & 0.5 \\ 0.5 & 0.3 & 0.2 \\ 0.2 & 0.5 & 0.3 \end{bmatrix}$

Try uniform X

$C = \log_2|\mathcal{Y}| - H(\underline{r}) = \log_2 3 - H([0.3 \ 0.2 \ 0.5]) \approx 1.5850 - 1.4855$
$= 0.0995 \ [bpcu].$*)*

Remark: If we can't find $\underline{p}$ that makes $Y$ uniform, then $C < \log_2 |\mathcal{Y}| - H(\underline{r})$ and we have to find a different technique to calculate $C$.

**Example 4.34.** Find the channel capacity of a BSC whose crossover probability is $p = 0.1$.

*(handwritten:*
$Q = \begin{array}{c} 0 \\ 1 \end{array}\begin{bmatrix} 1-p & p \\ p & 1-p \end{bmatrix} = \begin{bmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{bmatrix}$

$C = \log_2|\mathcal{Y}| - H(\underline{r})$
$= \log_2 2 - H([0.9 \ 0.1])$
$= 1 - H(0.1) \approx 1 - 0.469 = 0.531.$*)*

**Definition 4.35.** A DMC is called **weakly symmetric** if (1) all the rows of its probability transition matrix $Q$ are permutations of each other and (2) all the column sums are equal.

- It should be clear from the definition that a symmetric channel is automatically weakly symmetric.

*(handwritten right: condition (1) makes all $H(Y|x)$ the same. Therefore, $H(Y|X)$ is easy to find and does not depend on $\underline{p}$.)*

*(handwritten bottom left: Condition (2) guarantees that uniform X gives uniform Y.)*

**4.36.** For a weakly symmetric channel,

$$C = \log_2 |\mathcal{Y}| - H(\underline{\mathbf{r}}),$$

where $\underline{\mathbf{r}}$ is any row from the $\mathbf{Q}$ matrix. The capacity is achieved by a uniform pmf on the channel input.

- Important special case: For BSC, $C = 1 - H(p)$.

**4.37.** Properties of channel capacity

(a) $C \geq 0$

(b) $C \leq \min \{\log_2 |\mathcal{X}|, \log_2 |\mathcal{Y}|\}$

**Example 4.38.** Find the channel capacity of a DMC whose

$$\mathbf{Q} = \begin{bmatrix} 0.9 & 0.05 & 0.05 \\ 0.05 & 0.9 & 0.05 \\ 0.025 & 0.025 & 0.95 \end{bmatrix}$$

Suppose four choices are provided:
(a) 1.0944 (b) 1.5944 (c) 2.0944 (d) 2.5944

**Example 4.39.** Another case where capacity can be easily calculated: Find the channel capacity of a DMC of which all the rows of its $\mathbf{Q}$ matrix are the same.

$$\mathbf{Q} = \begin{bmatrix} \overset{\underline{\gamma}}{0.5} & 0.2 & 0.3 \\ 0.5 & 0.2 & 0.3 \end{bmatrix}$$

$$H(Y|x) = H(\underline{\gamma})$$

$$H(Y) = H(\underline{\gamma})$$

$$I(X;Y) = H(Y) - H(Y|X) = H(\underline{\gamma}) - H(\underline{\gamma})$$
$$= 0.$$

$$\underline{q} = \underline{p}\,\mathbf{Q} = [p_1 \quad p_2]\begin{bmatrix} 0.5 & 0.2 & 0.3 \\ 0.5 & 0.2 & 0.3 \end{bmatrix}$$

$$= \begin{bmatrix} 0.5 & 0.2 & 0.3 \end{bmatrix} = \underline{\gamma}$$

$$p_1 \times 0.5 + p_2 \times 0.5 = 0.5(p_1 + p_2)$$

$$\Rightarrow C = \max_{\underline{p}} \widetilde{I(X;Y)} = 0$$

(can use any $\underline{p}$) [bpcu]

**4.40.** In this section, we worked with "toy" examples in which finding capacity is relatively easy. In general, there is no closed-form solution for computing capacity. When we have to deal with cases that do not fit in any special family of $\mathbf{Q}$ described in the examples above, the maximum may be found by standard nonlinear optimization techniques or the Blahut-Arimoto Algorithm discussed in 4.25.

## 4.5   Shannon's Coding theorem

**4.41. Shannon's (Noisy Channel) Coding theorem** [Shannon, 1948]

(a) Reliable communication over a (discrete memoryless) channel is possible if the communication rate $R$ satisfies $R < C$, where $C$ is the channel capacity.

In particular, for any $R < C$, there exist codes (encoders and decoders) with sufficiently large $n$ such that

$$P\left(\mathcal{E}\right) \leq 2^{-n \times E(R)},$$

where $E(R)$ is

- a positive function of $R$ for $R < C$ and
- completely determined by the channel characteristics

(b) At rates higher than capacity, reliable communication is impossible.

**4.42.** Significance of Shannon's (noisy channel) coding theorem:

(a) Express the limit to reliable communication

(b) Provides a yardstick to measure the performance of communication systems.

- A system performing near capacity is a near optimal system and does not have much room for improvement.
- On the other hand a system operating far from this fundamental bound can be improved (mainly through coding techniques).

**4.43.** Shannon's nonconstructive proof for his coding theorem

- Shannon introduces a method of proof called **random coding**.
- Instead of looking for the best possible coding scheme and analyzing its performance, which is a difficult task,
  - all possible coding schemes are considered
    - ∗ by generating the code randomly with appropriate distribution
  - and the performance of the system is averaged over them.

- Then it is proved that if $R < C$, the average error probability tends to zero.

- Again, Shannon proved that

  - as long as $R < C$,
  - at any arbitrarily small (but still positive) probability of error,
  - one can find (there exist) at least one code (with sufficiently long block length $n$) that performs better than the specified probability of error.

- If we used the scheme suggested and generate a code at random, the code constructed is likely to be good for long block lengths.

- No structure in the code. Very difficult to decode

**4.44.** Practical codes:

- In addition to achieving low probabilities of error, useful codes should be "simple", so that they can be encoded and decoded efficiently.

- Shannon's theorem does not provide a practical coding scheme.

- Since Shannon's paper, a variety of techniques have been used to construct good error correcting codes.

  - The entire field of coding theory has been developed during this search.

- Turbo codes have come close to achieving capacity for Gaussian channels.