# Telecommunications and RF

Semester 2, 2023

*Dr Jacob Coetzee*

Tarang Janawalkar

# Contents

# 1   Telecommunications Systems

Telecommunication is the transmission of information over a distance through some technology. Telecommunications systems are designed to transmit information with as little **deterioration** as possible while satisfying design constraints, such as allowable transmittable energy and signal bandwidth.

Signal deterioration commonly measures:

- For Analog Systems: **Signal-to-Noise Ratio** (SNR) at the receiver output — the ratio of the signal power to the noise power

- For Digital Systems: **Bit Error Rate** (BER) at the receiver output — the ratio of the number of bits received in error to the total number of bits transmitted

## 1.1   Analog and Digital Signals

A system is either **analog** or **digital** based on the possible amplitudes of waveforms it can handle.

- **Analog Information Sources** produce values defined on a continuum:

    - Human voice and other sounds

- **Digital Information Sources** produce a finite set of possible symbols:

    - Computer data
    - MP3 encoder output

## 1.2   Communications System

A communications system can be summarised by the following block diagram:



where

- $m\,(t)$ is the information message signal (prior to conditioning for transmission)

- $s\,(t)$ is the conditioned signal for transmission

- $n\,(t)$ contains channel noise and interference

- $r\,(t)$ is the received signal

- $\hat{m}\,(t)$ is the reconstructed received message signal (where $\hat{m}\,(t)$ usually approximates $m\,(t)$)

3

### 1.2.1  Transmitter

The transmitter carries out **signal conditioning** by transforming the signal to a more appropriate form before transmission through the channel. Some examples of signal conditioning techniques are prodived below.

- **Low Pass Filtering** (LPF) restricts signal bandwidth to avoid wasting signal power on frequencies that are not transmitted by the channel and to avoid interference with other signals

- **Analog to Digital Conversion** (ADC) produces a digital word which represents a sample of the analog message waveform

- **Carrier Modulation** transfers the signal to a frequency band that is suitable for transmission through the channel

### 1.2.2  Channel

A communication channel refers to the physical medium that carries the signal from the transmitter to the receiver. There are two types of channels:

- **Wired**: twisted pair copper telephone lines, waveguides, coaxial cables, fibre-optic cables

- **Wireless**: air, vacuums, sea water, optical fibres

General principles of communications always apply regardless of the type of channel. However, certain conditioning methods are better suited to certain channels.
Channels often **attenuate** signals (reduce their amplitude or strength) through

- random noise

- interference from other sources

and therefore it is a key consideration in the design of a communications system.

### 1.2.3  Receiver

The receiver acts as the inverse of a transmitter. The receiver:

1. **Demodulates** the received signal by stripping the carrier from the received signal $r\left(t\right)$.

2. **Filters** out noise and interference from the demodulated signal.

3. **Reconstructs** an estimate of the original message signal $\hat{m}\left(t\right)$.

Due to the finite nature of the SNR, the estimated output of an analog signal can never be exactly equal to the original signal[1]. However, it is often possible to reconstruct a digital signal exactly using error detection and correction techniques at the receiver.

---

[1] A perfect reconstruction requires an infinite SNR which is impractical.

### 1.2.4   Information Sources

As discussed before, an information source can be classified as either **analog** or **digital**. Analog signals can be modulated or transmitted directly, or converted to digital data and transmitted using digital modulation techniques.

An analogue signal to be transmitted is called the **message signal** and is denoted $m(t)$. The spectral components of this signal lie within a finite bandwidth $W$, such that $M(f) = 0$ for $|f| > W$, where $M(f)$ is the Fourier Transform of $m(t)$. This signal's bandwidth is limited to prevent interference with other signals.

Many kinds of message signals can be considered:

- Audio

- Video

- Computer data

- Telemetry (measurements)

- Soundings (RADAR, SONAR)

- A mixure of the above (i.e., data over voice)

## 1.3   Modulation

The process of modulation produces a signal that is suitable for transmission through the channel by transforming the message signal $m(t)$ to a new signal $s(t)$.

Modulation is often performed with respect to another signal, called the **carrier** signal $c(t)$. Here the message *modulates* the carrier to produce the transmitted signal $s(t)$.

### 1.3.1   Benefits of Modulation

- Modulation shifts the spectral content of a message onto a suitable band. As the size of an antenna is related to the wavelength of a signal, higher carrier frequencies require smaller antennas.

- Modulation facilitates multiplexing, where multiple signals are transmitted over the same spectrum. This is because the frequency spectrum can be divided into non-overlapping frequency bands, each of which can carry a separate signal.

- Modulation provides some control over noise and interference by choosing a bandwidth that is smaller than the allocated channel bandwidth.

### 1.3.2   Convolution and Modulation

Consider two time domain signals $m(t)$ and $c(t)$ and their Fourier Transforms $M(f)$ and $C(f)$ respectively.

The **Convolution Property** demonstrates:

$$m(t) * c(t) = y(t) \qquad \text{Convolution in Time Domain}$$
$$M(f)\, C(f) = Y(f) \qquad \text{Multiplication in Frequency Domain}$$

The **Modulation property** demonstrates:

$$m(t)\,c(t) = y(t) \qquad \text{Multiplication in Time Domain}$$
$$M(f) * C(f) = Y(f) \qquad \text{Convolution in Frequency Domain}$$

### 1.3.3 Carrier Signal

The carrier signal $c(t)$ is a sinusoidal signal of the form:

$$c(t) = A_c \cos(2\pi f_c t + \phi)$$

where $A_c$ is the amplitude, $f_c$ is the frequency and $\phi$ is the phase of the carrier signal.

### 1.3.4 Signal Properties

The **energy** of a signal $x(t)$ is defined as:

$$E_x = \int_{-\infty}^{\infty} |x(t)|^2 \, dt$$

The rate at which energy is transmitted is called the **power** of the signal, and is defined as:

$$P_x = \lim_{T \to \infty} \frac{1}{T} \int_{-T/2}^{T/2} |x(t)|^2 \, dt$$

where $T$ is the time period over which the power is measured.

## 1.4 Amplitude Modulation Schemes

A message signal $m(t)$ with bandwidth $W$ can be amplitude modulated by mixing (multiplying) the signal with a carrier signal $c(t)$:

$$c(t) = A_c \cos(2\pi f_c t)$$

with amplitude $A_c$ and carrier frequency $f_c \gg W$.
This section will consider:

- Double Side Band — Suppressed Carrier (DSB-SC) modulation

- Double Side Band — Full Carrier (DSB-FC) modulation, or simply conventional Amplitude Modulation (AM)

- Single Side Band (SSB) modulation

### 1.4.1   Double Side Band — Suppressed Carrier (DSB-SC)

The transmitted signal is defined as:

$$s\left(t\right) = m\left(t\right)c\left(t\right) = A_c m\left(t\right)\cos\left(2\pi f_c t\right).$$

In the frequency domain, the message signal is shifted to the carrier frequency at $\pm f_c$ and the magnitude is halved.

$$S\left(f\right) = \frac{1}{2}A_c\left[M\left(f - f_c\right) + M\left(f + f_c\right)\right]$$

Assuming no channel noise, the received signal can be represented as:

$$r\left(t\right) = s\left(t\right) = A_c m\left(t\right)\cos\left(2\pi f_c t\right)$$

This signal can be demodulated coherently by multiplying a sinusoidal signal of the same frequency as the carrier. Suppose we generate the sinusoidal signal $\cos\left(2\pi f_c t + \phi\right)$ where $\phi$ is the phase of the sinusoid.

$$\begin{aligned}
y\left(t\right) &= r\left(t\right)\cos\left(2\pi f_c t + \phi\right)\\
&= A_c m\left(t\right)\cos\left(2\pi f_c t\right)\cos\left(2\pi f_c t + \phi\right)\\
&= \frac{1}{2}A_c m\left(t\right)\left[\cos\left(4\pi f_c t + \phi\right) + \cos\left(-\phi\right)\right]\\
&= \underbrace{\frac{1}{2}A_c m\left(t\right)\cos\left(\phi\right)}_{\text{filter out}} + \underbrace{\frac{1}{2}A_c m\left(t\right)\cos\left(4\pi f_c t + \phi\right)}_{\text{reject}}
\end{aligned}$$

As the frequency content of the message signal $m\left(t\right)$ is limited to $W$, a lowpass filter can be used to remove the high frequency component centred at $2f_c$. The output of this filter is then

$$\widehat{m}\left(t\right) = \frac{1}{2}A_c m\left(t\right)\cos\left(\phi\right).$$

Note that $m\left(t\right)$ is multiplied by $\cos\left(\phi\right)$, meaning the power of the message signal is reduced by a factor of $\cos^2\left(\phi\right)$. It is therefore important to choose $\phi$ such that $\cos\left(\phi\right)$ is as close to 1 as possible. This demonstrates a need for a phase-coherent or synchronous demodulator, i.e., the phase of the locally generated sinusoid should be identical to the received carrier signal.

### 1.4.2   Double Side Band — Full Carrier (DSB-FC)

This scheme is similar to DSB-SC, however the carrier is sent along with the modulated signal. We begin by defining an envelope signal $g\left(t\right)$ by amplifying and biasing the message signal, so that:

$$g\left(t\right) = 1 + \mu m_n\left(t\right)$$

where the **modulation index** $0 < \mu \leqslant 1$ is chosen to ensure $g\left(t\right) > 0$. $m_n\left(t\right)$ is the normalised message signal defined as

$$m_n\left(t\right) = \frac{m\left(t\right)}{\max\left|m\left(t\right)\right|}$$

such that $|m_n(t)| \leqslant 1$. The transmitted signal is then defined as:

$$s(t) = g(t)c(t) = A_c[1 + \mu m_n(t)]\cos(2\pi f_c t).$$

In the frequency domain, we observe similar behaviour, with impulses at $\pm f_c$.

$$S(f) = \frac{1}{2}A_c\mu[M(f - f_c) + M(f + f_c)] + \frac{1}{2}A_c[\delta(f - f_c) + \delta(f + f_c)]$$

The modulation index can be determined from the AM signal using the following equation:

$$\mu = \frac{A_{\max} - A_{\min}}{A_{\max} + A_{\min}} = \frac{\max|m(t)|}{\max|c(t)|}$$

where $A_{\max}$ and $A_{\min}$ are the maximum and minimum amplitudes of the envelope signal $g(t)$. This value is chosen to be as close to 1 as possible to maximise the power of the transmitted signal. Increasing $\mu$ beyond 1 overmodulates the signal, resulting in a full-wave rectified version of the signal.

Modulation efficiency is the percentage of the total power of the modulated signal that conveys information

$$\eta = \frac{\text{sideband power}}{\text{total power}} \times 100\%$$

As the total power comprises of both the sideband power and carrier power, the DSB-FC scheme is very inefficient. This is because the carrier signal does not contain any useful information, and therefore the power of the carrier is wasted. This can be seen when transmitting a sinusoidal message signal,

- Carrier power: $\frac{1}{2}A_c^2$

- Sideband power: $\frac{1}{4}A_c^2\mu^2$

- Total power: $\frac{1}{2}A_c^2\left(1 + \frac{1}{2}\mu^2\right)$

- Modulation efficiency: $\frac{\frac{1}{4}A_c^2\mu^2}{\frac{1}{2}A_c^2\left(1 + \frac{1}{2}\mu^2\right)} = \frac{\mu^2}{2 + \mu^2}$

The maximum efficiency is achieved when $\mu = 1$, i.e., which is 33%. This scheme is often used in AM medium-wave (MW) systems.

To demodulate a DSB-FC signal, we can use an envelope detector. In this circuit, we need to calculate the minimum and maximum frequency that the tuned circuit can respond to:

$$f_0 = \frac{1}{2\pi\sqrt{LC_1}}$$

and calculate the break frequency of the lowpass filter:

$$f_b = \frac{1}{2\pi RC_2}$$

### 1.4.3   Single Side Band (SSB)

As the double side band schemes produce a signal with twice the bandwidth of the message signal, consider either the lower or upper sidebands of the modulated signal. We can attempt to use a filter to remove the unwanted sideband, however there are two issues with this approach:

1. The filter is particularly difficult to implement when $m(t)$ has a large concentration of power close to $f = 0$

2. The filter must have a very sharp cutoff in the vicinity of the carrier frequency

Instead, we can use a phase shift method through the use of a Hilbert transform. The Hilbert transform of a signal $m(t)$ is defined as

$$\hat{m}(t) = m(t) * \frac{1}{\pi t}$$

The result of this operation is more easily understood in the frequency domain.

$$\hat{M}(f) = -j \operatorname{sgn}(f) M(f) = \begin{cases} -jM(f) & f > 0 \\ jM(f) & f < 0 \end{cases}$$

so that negative frequency components of $m$ are phase shifted by $+90°$ ($+\pi/2$) and positive frequency components are phase shifted by $-90°$ ($-\pi/2$). The magnitude remains unchanged.
Using this transform, we can define SSB signals as:

$$s_{\text{USB}}(t) = A_c \left[ m(t) \cos(2\pi f_c t) - \hat{m}(t) \sin(2\pi f_c t) \right]$$
$$s_{\text{LSB}}(t) = A_c \left[ m(t) \cos(2\pi f_c t) + \hat{m}(t) \sin(2\pi f_c t) \right]$$

or compactly,

$$s(t) = m(t) c(t) \pm \hat{m}(t) \hat{c}(t)$$

These signals can be demodulated coherently using the same process as the DSB-SC scheme. Again assuming that there is no channel noise,

$$r(t) = s(t) = A_c m(t) \cos(2\pi f_c t) \pm A_c \hat{m}(t) \sin(2\pi f_c t)$$

$$
\begin{aligned}
y(t) &= r(t) \cos(2\pi f_c t + \phi) \\
&= \left[ A_c m(t) \cos(2\pi f_c t) \pm A_c \hat{m}(t) \sin(2\pi f_c t) \right] \cos(2\pi f_c t + \phi) \\
&= A_c m(t) \cos(2\pi f_c t) \cos(2\pi f_c t + \phi) \pm A_c \hat{m}(t) \sin(2\pi f_c t) \cos(2\pi f_c t + \phi) \\
&= \frac{1}{2} A_c m(t) \left[ \cos(\phi) + \cos(4\pi f_c t + \phi) \right] \pm \frac{1}{2} A_c \hat{m}(t) \left[ -\sin(\phi) + \sin(4\pi f_c t + \phi) \right] \\
&= \underbrace{\frac{1}{2} A_c \left[ m(t) \cos(\phi) \mp \hat{m}(t) \sin(\phi) \right]}_{\text{filter out}} + \underbrace{\frac{1}{2} A_c \left[ m(t) \cos(4\pi f_c t + \phi) \pm \hat{m}(t) \sin(4\pi f_c t + \phi) \right]}_{\text{reject}}
\end{aligned}
$$

This gives the output

$$\hat{m}(t) = \frac{1}{2} A_c \left[ m(t) \cos(\phi) \mp \hat{m}(t) \sin(\phi) \right]$$

Note the $\hat{m}(t)$ term on the right hand side of the equation is referring to the Hilbert transform of $m(t)$, not the output of the demodulator.

## 2   Angle Modulation

AM suffers from poor noise performance, as amplitude variations in the received signal cannot be removed from the demodulated signal. Angle modulation schemes overcome this issue by modulating the phase or frequency of the carrier signal.

### 2.1   Carrier Signal

In these schemes, the message signal is modulated onto the angle $\theta(t)$ of the carrier signal:

$$s(t) = A_c \cos(\theta(t))$$
$$\theta(t) = 2\pi f_c t + \phi(t)$$

In **phase modulation** (PM), variations in the message signal are encoded into the phase:

$$\phi(t) = k_p m(t)$$

where $k_p$ is the **phase deviation** constant.
**Frequency modulation** (FM) considers the frequency deviation from the modulation frequency $f_c$:

$$f_i - f_c = k_f m(t)$$

where $f_i$ is the instantaneous frequency of the carrier signal and $k_f$ is the **frequency deviation** constant. To determine the phase $\phi$, consider the time derivative of the angle $\theta(t)$ with arbitrary frequency $f$

$$\frac{\mathrm{d}\theta(t)}{\mathrm{d}t} = \frac{\mathrm{d}}{\mathrm{d}t}(2\pi f t + \phi)$$
$$\frac{\mathrm{d}\theta(t)}{\mathrm{d}t} = 2\pi f_i$$
$$\theta(t) = 2\pi \int_0^t f_i \, \mathrm{d}\tau$$
$$\theta(t) = 2\pi \int_0^t f_c + k_f m(\tau) \, \mathrm{d}\tau$$
$$\theta(t) = 2\pi f_c t + 2\pi k_f \int_0^t m(\tau) \, \mathrm{d}\tau$$

#### 2.1.1   Phase Modulation

A phase modulated signal is defined as:

$$s(t) = A_c \cos\left(2\pi f_c t + k_p m(t)\right)$$

with

$$\theta(t) = 2\pi f_c t + k_p m(t)$$

where $k_p$ is the **phase deviation** constant.

### 2.1.2   Frequency Modulation

A frequency modulated signal is defined as:

$$s\left(t\right) = A_c \cos\left(2\pi f_c t + 2\pi k_f \int_0^t m\left(\tau\right) \mathrm{d}\tau\right)$$

with

$$\theta\left(t\right) = 2\pi f_c t + 2\pi k_f \int_0^t m\left(\tau\right) \mathrm{d}\tau$$

where $k_f$ is the **frequency deviation** constant.

### 2.1.3   Instantaneous Frequency

The instantaneous frequency $f_i$ may be determined from these signals by considering the time derivative of $\theta$.
In PM,

$$\begin{aligned}
\frac{\mathrm{d}\theta}{\mathrm{d}t} &= 2\pi f_c + k_p \frac{\mathrm{d}}{\mathrm{d}t} m\left(t\right) \\
&= 2\pi\left(f_c + \frac{k_p}{2\pi}\frac{\mathrm{d}}{\mathrm{d}t} m\left(t\right)\right) \implies f_i = f_c + \frac{k_p}{2\pi}\frac{\mathrm{d}}{\mathrm{d}t} m\left(t\right)
\end{aligned}$$

In FM,

$$\begin{aligned}
\frac{\mathrm{d}\theta}{\mathrm{d}t} &= 2\pi f_c + 2\pi k_f m\left(t\right) \\
&= 2\pi\left(f_c + k_f m\left(t\right)\right) \implies f_i = f_c + k_f m\left(t\right)
\end{aligned}$$

### 2.1.4   Modulation Index

The modulation index $\beta$ is defined as the ratio of the maximum frequency deviation to the modulating frequency of the message signal. In PM, if we rewrite the phase $\phi\left(t\right)$ as,

$$\begin{aligned}
\phi\left(t\right) &= k_p \max\left|m\left(t\right)\right| \frac{m\left(t\right)}{\max\left|m\left(t\right)\right|} \\
&= \beta_p m_n\left(t\right)
\end{aligned}$$

then $\beta_p$ is the modulation index for PM. By definition, FM has a modulation index of $\beta_f = \frac{k_f \max\left|m\left(t\right)\right|}{B}$ where $B$ is the bandwidth of the message signal.

### 2.1.5   Maximum Phase and Frequency Deviation

The maximum phase deviation $\Delta p$ and maximum frequency deviation $\Delta f$ are defined as:

$$\Delta p = k_p \max\left|m\left(t\right)\right| \qquad \Delta f = k_f \max\left|m\left(t\right)\right|$$

### 2.1.6 Carson's Rule

Phase and frequency modulation generally expand the bandwidth of a signal. The modulated signals bandwidth $W$ may be approximated by Carson's rule:

$$W = 2B\left(\beta + 1\right)$$

### 2.1.7 Summary of Angle Modulation Definitions

|  | Phase Modulation | Frequency Modulation |
| --- | --- | --- |
| **Modulated Signal** $s\left(t\right)$ | $A_c \cos\left(2\pi f_c t + k_p m\left(t\right)\right)$ | $A_c \cos\left(2\pi f_c t + 2\pi k_f \int_0^t m\left(\tau\right)\mathrm{d}\tau\right)$ |
| **Phase** $\phi\left(t\right)$ | $k_p m\left(t\right)$ | $2\pi k_f \int_0^t m\left(\tau\right)\mathrm{d}\tau$ |
| **Instantaneous Frequency** $f_i$ | $f_c + \dfrac{k_p}{2\pi}\dfrac{\mathrm{d}}{\mathrm{d}t}m\left(t\right)$ | $f_c + k_f m\left(t\right)$ |
| **Modulation Index** | $\beta_p = k_p \max\left\lvert m\left(t\right)\right\rvert$ | $\beta_f = \dfrac{k_f \max\left\lvert m\left(t\right)\right\rvert}{B}$ |
| **Maximum Deviation** | $\Delta p = k_p \max\left\lvert m\left(t\right)\right\rvert$ | $\Delta f = k_f \max\left\lvert m\left(t\right)\right\rvert$ |

## 2.2 Narrowband Angle Modulation

When $\beta \ll 1$, angle modulation is referred to as narrowband angle modulation. The modulated signal can be approximated as:

$$\begin{aligned}
s\left(t\right) &= A_c \cos\left(2\pi f_c t + \phi\left(t\right)\right) \\
&= A_c \cos\left(2\pi f_c t\right)\cos\left(\phi\left(t\right)\right) - A_c \sin\left(2\pi f_c t\right)\sin\left(\phi\left(t\right)\right) \\
&= A_c \cos\left(2\pi f_c t\right) - A_c \phi\left(t\right)\sin\left(2\pi f_c t\right)
\end{aligned}$$

This is equivalent to a DSB-FC signal with a phase modulated carrier, as the message signal is now modulated onto a sine carrier instead of cosine. The modulated signal bandwidth $W$ is approximately $2B$.

Narrowband angle modulation (low-index angle modulation) does not provide better noise immunity than AM, and is seldom used on its own in practical communication systems.

## 2.3 Wideband Angle Modulation

Assuming the message signal is a sinusoidal signal, $\phi\left(t\right) = \beta \sin\left(2\pi f_m t\right)$,

$$\begin{aligned}
s\left(t\right) &= A_c \cos\left(2\pi f_c t + \beta \sin\left(2\pi f_m t\right)\right) \\
&= A_c \cos\left(2\pi f_c t\right)\cos\left(\beta \sin\left(2\pi f_m t\right)\right) - A_c \sin\left(2\pi f_c t\right)\sin\left(\beta \sin\left(2\pi f_m t\right)\right)
\end{aligned}$$

To simplify this expression, we can use Bessel functions. Notably

$$\cos\left(\beta\sin\left(2\pi f_m t\right)\right) = J_0\left(\beta\right) + 2\sum_{n=1}^{\infty} J_{2n}\left(\beta\right)\cos\left(2\pi\left(2n\right)f_m t\right)$$

$$\sin\left(\beta\sin\left(2\pi f_m t\right)\right) = 2\sum_{n=1}^{\infty} J_{2n-1}\left(\beta\right)\sin\left(2\pi\left(2n-1\right)f_m t\right)$$

where $J_n\left(\beta\right)$ are Bessel functions of the first kind and order $n$, evaluated at $\beta$.
These results can be used to show that a wideband FM modulated signal can be expressed as:

$$s\left(t\right) = \sum_{n=-\infty}^{\infty} A_c J_n\left(\beta\right)\cos\left(2\pi\left(f_c + n f_m\right)t\right)$$

so that the carrier frequency is located at $2\pi f_c$ with magnitude $A_c J_0\left(\beta\right)$, with an infinite number of sidebands at $2\pi\left(f_c \pm nfm\right)$ with magnitudes $A_c J_{\pm n}\left(\beta\right)$.
When deciding the number of sidebands $n$ to include, consider

$$\left|J_{\pm n}\left(\beta\right)\right| \geqslant 0.1$$

For large values of $\beta$, $n \approx \beta$ is sufficient.
According to Carson's rule, $W = 2f_m\left(\beta + 1\right)$ contains at least 98% of the signal power.

### 2.3.1  Bessel Function Properties

The Bessel functions of the first kind $J_n\left(x\right)$ are defined as the solutions to the Bessel differential equation

$$x^2 \frac{\mathrm{d}^2 y}{\mathrm{d}x^2} + x\frac{\mathrm{d}y}{\mathrm{d}x} + \left(x^2 - n^2\right)y = 0$$

- $J_n\left(\beta\right)$ is a real function

- $J_n\left(\beta\right) = J_{-n}\left(\beta\right)$

- $\sum_{n=-\infty}^{\infty} J_n^2\left(\beta\right) = 1$

For small values of $\beta$,

- $J_0\left(\beta\right) = 1$

- $J_1\left(\beta\right) = \beta/2$

- $J_n\left(\beta\right) = 0$ for $n > 2$

## 2.4  Effect of Bandwidth

Rewriting Carson's rule using $\beta$,

$$W = 2B\left(\beta + 1\right) = \begin{cases} 2B\left(k_p \max\left|m\left(t\right)\right| + 1\right) & \text{PM} \\ 2\left(k_f \max\left|m\left(t\right)\right| + B\right) & \text{FM} \end{cases}$$

From these equations

- Increasing the amplitude of the modulating signal has the same effect in both PM and FM.

- Increase the message signal bandwidth $B$ has a greater effect on the bandwidth of a PM signal than for FM.

## 2.5 Angle Modulator Implementation

FM can be generated using a Voltage Controlled Oscillator (VCO). A varactor diode is a capacitor whose capacitance changes with applied voltage. This capacitor can be used in the tuned circuit of an oscillator. If the message signal is applied to the varactor, the output frequency of the oscillator will change in accordance with the message signal.

The time varying capacitance of the varactor diode is given by

$$C_v(t) = C_a + k_0 m(t)$$

- When $m(t) = 0$, the frequency of the tuned circuit is given by

$$f_i = f_c = \frac{1}{2\pi\sqrt{LC_a}}$$

- When $m(t) \neq 0$, the frequency of the tuned circuit is given by

$$f_i = \frac{1}{2\pi\sqrt{L\left(C_a + k_0 m(t)\right)}} = \frac{1}{2\pi\sqrt{LC_a\left(1 + k_0 m(t)/C_a\right)}} = f_c \frac{1}{\sqrt{1 + k_0 m(t)/C_a}}$$

Narrowband FM can be converted to wideband using a narrowband-to-wideband convertor by multiplying the frequencies of the narrowhand signal.

## 2.6 Demodulating Angle Modulated Signals

FM demodulators are implemented by generating an AM signal whose amplitude is proportional to the instantaneous frequency of the FM signal. We can then use an AM demodulator to recover the message signal.

In such a circuit, the LTI system is a differentiator whose frequency response is approximately a straight line in the frequency band of the FM signal. $|H| = 2\pi f$.

## 2.7 AM vs. FM

- FM capture effect: When two FM signals are received, the stronger signal is demodulated and the weaker signal is ignored.

  - The complete suppression of the weaker signal occurs at the receiver limited, where it is treated as noise and rejected.
  - When both signals are of equal strength, the receiver may switch between the two signals.

- FM requires a higher bandwidth $W_{\mathrm{AM}} < W_{\mathrm{FM}}$.

- FM rejects amplitude noise cause by lightning and other man-made noise.

- AM demodulators: envelope detector, product detector.

- FM demodulators: PLL, ratio detector, frequency discriminator, slope detector.

# 3    Digital Modulation

A source encoder converts an analog signal into a sequence of binary digits (bits). Digital modulation is used to transmit digital information over an additive white Gaussian noise (AWGN) channel.

**Definition 3.1** (Baseband Channel)**.** A baseband channel is a channel that includes the zero frequency $f = 0$, and there is no need to transmit a carrier signal.

**Definition 3.2** (Bandpass Channel)**.** A communication channel that is far from the zero frequency is known as a bandpass channel, and such a signal is impressed on a sinusoidal carrier, that shifts the signal to the desired frequency band that is *passed* by the channel.

Pulse code modulation (PCM) is a method used to digitally represent a sampled analog signal. To transmit these digits, we can use electrical pulses to represent binary digits.

**Definition 3.3** (Bit)**.** A bit is a binary digit, i.e., a digit that can take on one of two values; 0 or 1.

**Definition 3.4** (Digital Symbol)**.** A digital symbol refers to a sequence of several bits. For a sequence of $k$ bits, there are $M = 2^k$ unique symbols.

## 3.1    Waveform Representation of Digital Signals

**Symbol rate** $R_s$ is the number of symbols transmitted per second, measured in symbols/s.
**Bit rate** $R_b$ is the number of bits transmitted per second, measured in bits/s.
Symbol rate and bit rate are related by

$$R_b = kR_s$$

where $k$ is the number of bits per symbol. When a symbol contains only one bit $(k = 1)$, $R_b = R_s$. The duration of a single bit, or symbol can be determined using the reciprocal of each quantity:

$$T_s = \frac{1}{R_s} \qquad T_b = \frac{1}{R_b}$$

These quantities are related by

$$T_b = \frac{1}{k}T_s$$

To effectively transmit digital signals through an analog channel, each pulse must carry sufficient energy either by increasing the amplitude or duration of each pulse.

### 3.1.1    Antipodal and Orthogonal Signalling

Antipodal signalling is a method of transmitting digital signals where the information bit 1 is represented by a pulse $p(t)$ of duration $T_b$, and the information bit 0 is represented by $-p(t)$.
In orthogonal signalling, information is represented by two orthogonal pulses $p_1(t)$ and $p_2(t)$ such that

$$\int_0^{T_b} p_1(t) p_2(t) \, \mathrm{d}t = 0$$

Each of the above methods require precise synchronisation between the receiver and transmitter.

15

### 3.1.2   Not-Return To Zero

A non-return to zero (NRZ) **line code** is a binary code in which bits are represented by two different levels of voltage.
For bipolar NRZ, the bit 1 is represented by a positive voltage, and the bit 0 is represented by a negative voltage.
For unipolar NRZ, the bit 1 is represented by a positive voltage, and the bit 0 is represented by a zero voltage.

### 3.1.3   Pulse Amplitude Modulation

In pulse amplitude modulation (PAM), the amplitude of a pulse is varied between various voltage levels to represent digital symbols of $k$ bits. The number of voltage levels is given by $M = 2^k$.

## 3.2   Signal Bandwidth

The bandwidth of a signal provides a measure of the extent of significant spectral content of the signal for positive frequencies. There are various definitions of bandwidth,

- **Null-to-null bandwidth**: range of frequencies between zeros in the magnitude spectrum

- **3-dB bandwidth**: range of frequencies where the magnitude spectrum falls no lower than $1/\sqrt{2}$ of the maximum value of the magnitude spectrum

- **Equivalent noise bandwidth**: width of a fictitious rectangular spectrum such that the power in the rectangular band is equal to the power associated with the actual spectrum over positive frequencies

### 3.2.1   Pulse Dilemma

The pulse dilemma is the trade-off between the bandwidth and the duration of a pulse. A narrow pulse has a large bandwidth, and a wide pulse has a small bandwidth.

- A perfectly band-limited pulse implies an infinite duration, which is not realisable.

- A precisely duration-limited pulse implies an infinite bandwidth, which is also not realisable.

When pulses are filtered by a communications system to reduce bandwidth, they are **spread** in time. This causes pulses of adjacent symbols to overlap in time, causing **intersymbol interference** (ISI).
In the following sections, we will consider various pulse shapes and analyse their bandwidths.

**Definition 3.5** (Spectral Efficiency)**.** The spectral efficiency of a digital signal is the number of bits per second that can be supported by a hertz of bandwidth:

$$\eta = \frac{R_b}{B} \, \text{bits} \, \text{s}^{-1} \, \text{Hz}^{-1}$$

### 3.2.2   Rectangular Pulse

If a symbol is represented by a rectangular pulse, it's spectrum is the sinc function:

$$\Pi\left(\frac{t}{T_s}\right) \overset{\mathscr{F}}{\Longleftrightarrow} T_s \operatorname{sinc}\left(fT_s\right)$$

By using the null-to-null bandwidth definition, the bandwidth of this pulse is

$$B = \frac{1}{T_s} = R_s.$$

For binary modulation using bipolar NRZ, the spectral efficiency is given by

$$\eta = \frac{R_b}{R_s} = \frac{R_b}{B} = \frac{R_s}{R_s} = 1$$

### 3.2.3   Nyquist Pulse

A Nyquist pulse is a pulse that satisfies the **Nyquist criterion**, which states that the pulse must be zero at the sampling times of adjacent pulses. Assuming a sampling frequency of $R_s = 1/T_s$, the Nyquist pulse is defined

$$p\left(t\right) = \frac{\sin\left(\pi R_s t\right)}{\pi R_s t}$$

where each pulse is sampled at $t = nT_s$. From this definition, it can be seen that

$$p\left(nT_s\right) = \begin{cases} 1 & n = 0 \\ 0 & n \neq 0 \end{cases}$$

The spectrum of this function is a rectangular function,

$$P\left(f\right) = \frac{1}{R_s}\Pi\left(\frac{f}{R_s}\right)$$

with bandwidth:

$$B \geqslant \frac{R_s}{2}$$

For binary modulation, the spectral efficiency of this pulse is given by

$$\eta = \frac{R_b}{B} = \frac{R_s}{R_s/2} = 2$$

In practise, this pulse shape is not physically realisable due to the following:

- infinite pulse in time domain

- sharp roll-off in frequency domain

- requirement for perfect synchronisation to ensure that the pulse is zero at the sampling times of adjacent pulses

### 3.2.4   Raised Cosine Pulse

If we relax the filtering and clock timing requirements of the Nyquist pulse, we can use a raised cosine (RC) pulse. The RC pulse has the form

$$p\left(t\right) = \frac{\sin\left(\pi R_s t\right)}{\pi R_s t} \frac{\cos\left(\pi r R_s t\right)}{1 - \left(2r R_s t\right)^2}$$

where $r$ is the roll-off factor. The spectrum of this function is a finite bandwidth piecewise function of the form,

$$P\left(f\right) = \begin{cases} 1 & |f| < \frac{1-r}{2T_s} \\ \frac{1}{2}\left[1 + \cos\left(\frac{\pi T_s}{r}\left[|f| - \frac{1-r}{2T_s}\right]\right)\right] & \frac{1-r}{2T_s} < f < \frac{1+r}{2T_s} \\ 0 & |f| \geqslant \frac{1+r}{2T_s} \end{cases}$$

The bandwidth of this pulse is

$$B = \frac{R_s}{2}\left(1 + R\right)$$

The spectral efficiency of this pulse is given by

$$\eta = \frac{R_b}{B} = \frac{R_s}{R_s/2\left(1 + R\right)} = \frac{2}{1 + R}$$

## 3.3   Energy of a Digital Symbol

The symbol energy $E$ is the instantaneous power integrated over the duration of one pulse $T_s$. For a symbol waveform $s\left(t\right)$:

$$E_s = \int_0^{T_s} s^2\left(t\right) \mathrm{d}t$$

such that $E$ is the area under the squared symbol waveform.

If there are $M$ symbols, each with energy $E_i$, the average symbol energy (assuming all $M$ symbols are equiprobable) is the average of the symbol energies:

$$E_s = \frac{1}{M}\sum_{i=1}^{M} E_i$$

## 3.4   Additive White Gaussian Noise

Additive white Gaussian noise (AWGN) is a model for the effect of external noise on a signal. The probability density function of the AWGN process $n$ is given by

$$f_N\left(n\right) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{\left(n - \mu\right)^2}{2\sigma^2}\right]$$

where $\mu$ is the mean of the process, and $\sigma^2$ is the variance of the process.

The term *white* denotes a process in which all frequency components appear with equal power, i.e., the power spectral density is constant.

This is a result of the Wiener-Khinchin theorem, which states that the power spectral density of a wide sense stationary random process is the Fourier transform of the autocorrelation function of the process. For the above AWGN process, the autocorrelation function is given by

$$R_N\left(\tau\right) = \frac{N_0}{2}\delta\left(\tau\right)$$

where $\delta\left(\tau\right)$ is the Dirac delta function. The power spectral density is then

$$S_N\left(f\right) = \frac{N_0}{2}$$

$N_0$ refers to the thermal noise power, which is defined

$$N_0 = kT$$

where $k = 1.380\,649 \times 10^{-23}\,\mathrm{J\,K^{-1}}$ is Boltzmann's constant and $T$ is the temperature in Kelvin. At room temperature, $20°C$, $N_0 = 4 \times 10^{-21}\,\mathrm{W} = -174\,\mathrm{dBm}$.

## 3.5   Binary Transceiver

In a simple binary transciever, bits are represented by impulses, where each pulse is one bit duration $T_b$ apart. These pulses are then converted into an appropriate waveform for transmission, i.e., as a rectangular, since, raised cosine, Gaussian, etc., pulse shape.
The receiver can then detect these pulses by sampling the received signal at the center of each bit period, where the voltage of the signal is either positive or negative, corresponding to a 1 or 0. However, this fails in the presence of noise, and we must consider other approaches.

## 3.6   Optimum Receiver

The **correlator** is an optimum receiver for binary signals in an AWGN channel. The correlator performs waveform recovery in preparation for the next step of detection.

$$\int_0^{T_s} s_i\left(t\right) s_i\left(t\right) \mathrm{d}t$$

The correlator integrates the product of the received signal and a replica of the transmitted symbol (called the reference signal) for $T_s$ seconds (symbol period), and dumps the result to the decision circuit by closing the switch every $T_s$ seconds.

$$
\begin{aligned}
T_s &= T_b & &\text{for binary modulation} \\
T_s &= T_b \log_2\left(M\right) & &\text{for M-ary modulation}
\end{aligned}
$$

A decision is then made regarding the symbol that was transmitted, based on the test statistic $r$. As an example, the received signal may have the following form:

$$r\left(t\right) = s_i\left(t\right) + n\left(t\right)$$

where we transmit symbol $i$. By performing integration with $s_j$, we obtain

$$r_j(t) = \begin{cases} \displaystyle\int_0^{T_s} s_j^2(t)\,\mathrm{d}t + \int_0^{T_s} n(t)\,s_j(t)\,\mathrm{d}t & j = i \\ \displaystyle\int_0^{T_s} s_i(t)\,s_j(t)\,\mathrm{d}t + \int_0^{T_s} n(t)\,s_j(t)\,\mathrm{d}t & j \neq i \end{cases}$$

$$= \begin{cases} E_j + \hat{n}(t) & j = i \\ \displaystyle\int_0^{T_s} s_i(t)\,s_j(t)\,\mathrm{d}t + \hat{n}(t) & j \neq i \end{cases}$$

where $\hat{n}(t)$ is also a Gaussian random process.

- When the correlator is matched to the transmitted symbol, the output is distributed about the autocorrelation of the transmitted symbol.

- When the correlator is matched to the wrong symbol, the output is distributed about the cross-correlation of the transmitted and reference symbol.

Using the test statistic $r$, we can then make a decision regarding the transmitted symbol, by comparing $r$ with a particular threshold.

## 3.7   Matched Filter

The filter corresponding to the corellator is called the matched filter. The matched filter is a linear filter that maximises the SNR at a particular sampling time (typically at the end of the bit period).

$$\left(\frac{S}{N}\right) = \frac{|x(t_0)|^2}{\mathrm{E}\left[N_0^2(t)\right]}$$

where $x(t_0)$ is the signal at the sampling time $t_0$ and $N_0(t)$ is the noise process. The transfer function that achieves this is given by

$$H(\omega) = \frac{1}{2\pi C}\frac{X^*(\omega)}{S_N(\omega)}e^{-j\omega t_o}$$

where $X^*(\omega)$ is the complex conjugate of the Fourier transform of the signal, $S_N(\omega)$ is the power spectral density of the noise, and $C$ is an arbitrary real constant.
As the noise is assumed to be AWGN, $S_N(\omega) = N_0/2$. The impulse response then becomes

$$h_{\mathrm{opt}}(t) = Kx^*(t_0 - t)$$

where $K$ is a constant, and $x^*(t)$ is the complex conjugate of the signal.
For a symbol shape $s_i(t)$, this tells us that the impulse response of the matched filter is $s_i(T_s - t)$.

### 3.7.1   Error Performance

The bit error rate is given by the ratio of the number of bit errors to the total number of bits transmitted. The theoretical bit error rate is a function of the bit energy $E_b$ and the noise power

spectral density $N_0$. For binary signals

$$P_b = \frac{1}{2} \operatorname{erfc} \left( \sqrt{\frac{d_{01}^2}{N_0}} \right)$$

where $d_{01}$ is the Euclidean distance between the two symbols,

$$d_{01}^2 = \int_0^{T_b} (s_0 - s_1)^2 \, \mathrm{d}t$$

# 4   Information Theory

Information is a measure of the uncertainty of an event.

- An event that is 100% probably is unsurprising and yields no information.

- A less probable event is more surprising and therefore yields more information.

The study of information theory therefore applies to random events.

**Definition 4.1** (Alphabet). An alphabet $\mathcal{A}$ is a finite set of $M$ symbols $\{a_1, a_2, ..., a_M\}$ used to represent information from a source. For a source $\{X_n\}_{n=-\infty}^{\infty}$, the alphabet $\mathcal{X} = \{x_1, x_2, ..., x_M\}$ represents the set of all possible values that $X_n$ can take.

**Definition 4.2** (Probability of a symbol). The probability of a symbol $x \in \mathcal{X}$ is the probability that the random variable $X$ takes the value $x$:

$$p(x) = \Pr(X = x).$$

When all $M$ symbols in this alphabet are equiprobable, $p(x) = 1/M$ for all $x \in \mathcal{X}$.

## 4.1   Measure of Information

The informational content associated with a symbol $x$ is inversely proportional to the probability of the occurrence of that symbol:

$$I(x) = \log_2 \left( \frac{1}{p(x)} \right) = -\log_2 (p(x)).$$

This suggests that an event with low probability carries more information.

**Definition 4.3** (Entropy). The entropy of a discrete random variable $X$ is the average information carried by each symbol:

$$H(X) = \mathrm{E}[I(x)] = \mathrm{E}[-\log(p(X))] = -\sum_{i=1}^{M} p(x_i) \log_2 (p(x_i)).$$

The entropy of a source $\{X_n\}_{n=-\infty}^{\infty}$ reflects the average amount of uncertainty that is resolved by using the alphabet.

For a binary source, entropy is measured in bits per symbol (bit/symbol). The bit rate $R_b$ is given by the product of the entropy $H(X)$ and the symbol rate $R_s$:

$$R_b = H(X) R_s$$
$$\mathrm{bit/s} = \mathrm{bit/symbol} \times \mathrm{symbol/s}$$

## 4.2   Multiple Sources of Information

Given two discrete random variables $X$ and $Y$ with alphabets $\mathcal{X}$ and $\mathcal{Y}$, the joint probability of $x \in \mathcal{X}$ and $y \in \mathcal{Y}$ is given by

$$p\left(x,\,y\right) = \Pr\left(X = x \cap Y = y\right).$$

When $X$ and $Y$ are independent, the joint probability is given by the product of the individual (marginal) probabilities:

$$p\left(x,\,y\right) = p\left(x\right)p\left(y\right).$$

The conditional probability of $x$ given $y$ is given by

$$p\left(x \mid y\right) = \frac{p\left(x,\,y\right)}{p\left(y\right)} \implies p\left(x,\,y\right) = p\left(x \mid y\right)p\left(y\right)$$

and the conditional probability of $y$ given $x$ is given by

$$p\left(y \mid x\right) = \frac{p\left(x,\,y\right)}{p\left(x\right)} \implies p\left(x,\,y\right) = p\left(y \mid x\right)p\left(x\right).$$

### 4.2.1   Joint Entropy

The joint entropy of $X$ and $Y$ is given by

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

Joint entropy measures the amount of uncertainty associated with the joint distribution of $X$ and $Y$.

**Theorem 4.2.1** (Chain Rule for Entropy)**.** *The above result can be expressed as*

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

*This is known as the chain rule for entropy.*

*Proof.* We can prove this result by expanding the joint entropy:

$$
\begin{aligned}
H(X,\,Y) &= -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p\,(x,\,y) \log_2 \left( p\,(x,\,y) \right) \\
&= -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p\,(x,\,y) \log_2 \left( p\,(x)\,p\,(y \mid x) \right) \\
&= -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p\,(x,\,y) \left[ \log_2 \left( p\,(x) \right) + \log_2 \left( p\,(y \mid x) \right) \right] \\
&= -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p\,(x,\,y) \log_2 \left( p\,(x) \right) - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p\,(x,\,y) \log_2 \left( p\,(y \mid x) \right) \\
&= -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p\,(x)\,p\,(y \mid x) \log_2 \left( p\,(x) \right) - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p\,(x)\,p\,(y \mid x) \log_2 \left( p\,(y \mid x) \right) \\
&= -\sum_{x \in \mathcal{X}} p\,(x) \log_2 \left( p\,(x) \right) \sum_{y \in \mathcal{Y}} p\,(y \mid x) - \sum_{x \in \mathcal{X}} p\,(x) \sum_{y \in \mathcal{Y}} p\,(y \mid x) \log_2 \left( p\,(y \mid x) \right) \\
&= H(X) + \sum_{x \in \mathcal{X}} p\,(x)\,H\,(Y \mid X) \\
&= H(X) + \mathrm{E}_X \left[ H\,(Y \mid X) \right] \\
&= H(X) + H\,(Y \mid X)
\end{aligned}
$$

A similar proof can be used to show that $H(X,\,Y) = H(Y) + H(X \mid Y)$. $\qquad \square$

## 4.3   Conditional Entropy

The conditional entropy of $Y$ given $X$ is given by

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

Conditional entropy measures the amount of information needed to describe the uncertainty of $Y$ given $X = x$.

*Proof.* As $H(Y \mid X)$ is a function of $Y$, we can write

$$
H(Y \mid X) = \sum_{x \in \mathcal{X}} p\,(y \mid x) \log_2 \left( p\,(y \mid x) \right).
$$

Additionally, the expected value of $H(Y \mid X)$ with respect to $X$ is given by

$$
\mathrm{E}_X \left[ H\,(Y \mid X) \right] = \sum_{x \in \mathcal{X}} p\,(x)\,H\,(Y \mid X).
$$

which is equivalent to the definition of $H(Y \mid X)$. Therefore, by combining the two equations, we

obtain

$$
\begin{aligned}
H\left(Y \mid X\right) &= \mathrm{E}_X\left[H\left(Y \mid X\right)\right] \\
&= \sum_{x \in \mathcal{X}} p\left(x\right) H\left(Y \mid X\right) \\
&= -\sum_{x \in \mathcal{X}} p\left(x\right) \sum_{y \in \mathcal{Y}} p\left(y \mid x\right) \log_2\left(p\left(y \mid x\right)\right) \\
&= -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p\left(x\right) p\left(y \mid x\right) \log_2\left(p\left(y \mid x\right)\right) \\
&= -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p\left(x,\, y\right) \log_2\left(p\left(y \mid x\right)\right)
\end{aligned}
$$

$\square$

## 4.4   Mutual Information

The mutual information between $X$ and $Y$ is given by

$$
\begin{aligned}
I\left(X,\, Y\right) &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p\left(x,\, y\right) \log_2\left(\frac{p\left(x,\, y\right)}{p\left(x\right) p\left(y\right)}\right) \\
&= H\left(X\right) - H\left(X \mid Y\right) \\
&= H\left(Y\right) - H\left(Y \mid X\right)
\end{aligned}
$$

Mutual information measures the amount of information that $X$ and $Y$ share. It is also the reduction in uncertainty of $X$ when $Y$ is known, or the reduction in uncertainty of $Y$ when $X$ is known.

# 5   Source Coding

Given a source $X$ with $n$ symbols, the minimum number of bits per symbol required to transmit the source is given by the entropy $H\left(X\right)$.

## 5.1   Source Coding Theorem

Given a source $X$, let $C : \mathcal{X} \to \mathcal{A}$ be any coding function that maps the source alphabet $\mathcal{X} = \{x_1,\, x_2,\, ...,\, x_M\}$ to the set of sequences $\mathcal{A} = \{a_1,\, a_2,\, ...,\, a_M\}$, i.e.,

$$
C\left(x_i\right) = a_i.
$$

Then, the smallest possible expected code word length $L = |C\left(X\right)|$ is the entropy of $X$.

$$
\mathrm{E}\left[L\right] = \sum_{x=1}^{M} L_i p_i \geqslant H\left(X\right).
$$

For a sequence of $n$ symbols, the expected code word length is at least $nH\left(X\right)$.

Therefore, for **lossless encoding**, a source with entropy $H(X)$ can be encoded with an arbitrarily small error probability at a rate $R$ bits per symbol, as long as

$$R > H(X).$$

When all symbols in the source alphabet are equiprobable, the entropy is maximised, and the minimum number of bits per symbol is

$$R > \log_2(M).$$

This is used in Pulse Code Modulation (PCM).

## 5.2   Designing a Code

When designing the optimal code $C$, the following criteria must be satisfied:

1. **Uniqueness**: each code word must be uniquely decodable, i.e., there must be no ambiguity in decoding a code word from a binary sequence.

2. **Instantaneous**: each code word must be recognisable without the need to look ahead at the next code word in the sequence.

3. **Prefix**: no code word can be a prefix of another code word.

Additionally, the code must be **efficient**, i.e., the average code word length $\mathrm{E}[L]$, must be minimised (close to the entropy of the source). The efficiency of a code is given by

$$\eta = \frac{H(X)}{\mathrm{E}[L]} = \frac{H(X)}{\sum_{i=1}^{M} L_i p_i}.$$

## 5.3   Shannon-Fano Coding

The Shannon-Fano coding algorithm is an algorithm that constructs a prefix code by recursively splitting the source alphabet into two sets, such that the entropy of each set is almost equal.
The steps for the algorithm are as follows:

1. Order the source alphabet $\mathcal{X} = \{x_1, x_2, ..., x_M\}$ in descending order of probability.

2. Split the alphabet into two sets $A$ and $B$ such that

$$\sum_{x \in A} p(x) \approx \sum_{x \in B} p(x).$$

3. Assign a 0 to all symbols in set $A$ and a 1 to all symbols

4. Repeat steps 2 and 3 for each set until each set contains only one symbol.

5. The code word for each symbol $x_i$ is the sequence of 0s and 1s assigned to that symbol.

This algorithm does not always generate an optimal code however. Instead, we can use the Huffman coding algorithm, which generates codes in reverse, i.e., from the leaf to the root.

25

## 5.4   Huffman Coding

The Huffman coding algorithm is an algorithm that constructs a prefix-free code by recursively combining the two least probable symbols into a single symbol, until only one symbol remains. The steps for the algorithm are as follows:

1. Order the source alphabet $\mathcal{X} = \{x_1, x_2, ..., x_M\}$ in descending order of probability.

2. Combine the two least probable symbols $x_i$ and $x_j$ into a single symbol $x_k$, where $k > j > i$.

3. Assign a 0 to $x_i$ and a 1 to $x_j$.

4. Repeat steps 2 and 3 until only one symbol remains.

5. The code word for each symbol $x_i$ is the sequence of 0s and 1s assigned to that symbol.

In this algorithm, step 3 may be performed in multiple ways, for example,

- $x_i$ is assigned 1 and $x_j$ is assigned 0

- $x_i$ is assigned 0 and $x_j$ is assigned 1

- the assignment changes at each iteration

and therefore the Huffman algorithm is not unique. Additionally, when the least probable symbol $x_i$ has a probability $p_i$, and the next least probable symbol $x_j$ shares a probability $p_j$ with another symbol, the algorithm is ambiguous.

## 5.5   Variance of Code Word Length

When two codes share the same average code word length, the code with the smaller variance of code length is preferred. The variance of a code length is given by

$$\sigma^2 = \sum_{i=1}^{M} \left(L_i - \mathrm{E}\left[L\right]\right)^2 p_i.$$

A high variance typically indicates that the number of bits per symbol is high. Therefore, it is beneficial to minimise the variance of a code word length.