# Chapter 6

# Channels, modulation, and demodulation

## 6.1   Introduction

Digital modulation (or channel encoding) is the process of converting an input sequence of bits into a waveform suitable for transmission over a communication channel. Demodulation (channel decoding) is the corresponding process at the receiver of converting the received waveform into a (perhaps noisy) replica of the input bit sequence. Chapter 1 discussed the reasons for using a bit sequence as the interface between an arbitrary source and an arbitrary channel, and Chapters 2 and 3 discussed how to encode the source output into a bit sequence.

Chapters 4 and 5 developed the signal-space view of waveforms. As explained there, the source and channel waveforms of interest can be represented as real or complex[1] $\mathcal{L}_2$ vectors. Any such vector can be viewed as a conventional function of time, $x(t)$. Given an orthonormal basis $\{\phi_1(t), \phi_2(t), \dots, \}$ of $\mathcal{L}_2$, any such $x(t)$ can be represented as

$$x(t) = \sum_j x_j \phi_j(t). \tag{6.1}$$

Each $x_j$ in (6.1) can be uniquely calculated from $x(t)$, and the above series converges in $\mathcal{L}_2$ to $x(t)$. Moreover, starting from any sequence satisfying $\sum_j |x_j|^2 < \infty$ there is an $\mathcal{L}_2$ function $x(t)$ satisfying (6.1) with $\mathcal{L}_2$ convergence. This provides a simple and generic way of going back and forth between functions of time and sequences of numbers. The basic parts of a modulator will then turn out to be a procedure for mapping a sequence of binary digits into a sequence of real or complex numbers, followed by the above approach for mapping a sequence of numbers into a waveform.

In most cases of modulation, the set of waveforms $\phi_1(t), \phi_2(t), \dots,$ in (6.1) will be chosen not as a basis for $\mathcal{L}_2$ but as a basis for some subspace[2] of $\mathcal{L}_2$ such as the set of functions that are baseband limited to some frequency $W$ or passband limited to some range of frequencies. In some cases, it will also be desirable to use a sequence of waveforms that are not orthonormal.

---

[1]As explained later, the actual transmitted waveforms are real. However, they are usually bandpass real waveforms that are conveniently represented as complex baseband waveforms.

[2]Equivalently, $\phi_1(t), \phi_2(t), \dots,$ can be chosen as a basis of $\mathcal{L}_2$ but the set of indices for which $x_j$ is allowed to be nonzero can be restricted.

We can view the mapping from bits to numerical signals and the conversion of signals to a waveform as separate layers. The demodulator then maps the received waveform to a sequence of received signals, which is then mapped to a bit sequence, hopefully equal to the input bit sequence. A major objective in designing the modulator and demodulator is to maximize the rate at which bits enter the encoder, subject to the need to retrieve the original bit stream with a suitably small error rate. Usually this must be done subject to constraints on the transmitted power and bandwidth. In practice there are also constraints on delay, complexity, compatibility with standards, etc., but these need not be a major focus here.

**Example 6.1.1.** As a particularly simple example, suppose a sequence of binary symbols enters the encoder at $T$-spaced instants of time. These symbols can be mapped into real numbers using the mapping $0 \rightarrow +1$ and $1 \rightarrow -1$. The resulting sequence $u_1, u_2, \ldots$, of real numbers is then mapped into the transmitted waveform

$$u(t) = \sum_k u_k \operatorname{sinc} \left( \frac{t}{T} - k \right). \tag{6.2}$$

At the receiver, in the absence of noise, attenuation, and other imperfections, the received waveform is $u(t)$. This can be sampled at times $T_1, T_2, \ldots$, to retrieve $u_1, u_2, \ldots$, which can be decoded into the original binary symbols.

The above example contains rudimentary forms of the two layers discussed above. The first is the mapping of binary symbols into numerical signals[3] and the second is the conversion of the sequence of signals into a waveform. In general, the set of $T$-spaced sinc functions in (6.2) can be replaced by any other set of orthogonal functions (or even non-orthogonal functions). Also, the mapping $0 \rightarrow +1$, $1 \rightarrow -1$ can be generalized by segmenting the binary stream into $b$-tuples of binary symbols, which can then be mapped into $n$-tuples of real or complex numbers. The set of $2^b$ possible $n$-tuples resulting from this mapping is called a *signal constellation*.

Modulators usually include a third layer, which maps a baseband encoded waveform, such as $u(t)$ in (6.2), into a passband waveform $x(t) = \Re\{u(t)e^{2\pi i f_c t}\}$ centered on a given carrier frequency $f_c$. At the decoder this passband waveform is mapped back to baseband before performing the other components of decoding. This frequency conversion operation at encoder and decoder is often referred to as modulation and demodulation, but it is more common today to use the word modulation for the entire process of mapping bits to waveforms. Figure 6.1 illustrates these three layers.

We have illustrated the channel above as a one way device going from source to destination. Usually, however, communication goes both ways, so that a physical location can send data to another location and also receive data from that remote location. A physical device that both encodes data going out over a channel and also decodes oppositely directed data coming in from the channel is called a *modem* (for <u>mo</u>dulator/<u>dem</u>odulator). As described in Chapter 1, feedback on the reverse channel can be used to request retransmissions on the forward channel, but in practice, this is usually done as part of an <u>a</u>utomatic <u>r</u>etransmission re<u>q</u>uest (ARQ) strategy in the data link control layer. Combining coding with more sophisticated feedback strategies than

---

[3]The word *signal* is often used in the communication literature to refer to symbols, vectors, waveforms, or almost anything else. Here we use it only to refer to real or complex numbers (or $n$-tuples of numbers) in situations where the numerical properties are important. For example, in (6.2) the *signals* (numerical values) $u_1, u_2, \ldots$ determine the real valued waveform $u(t)$, whereas the binary input *symbols* could be 'Alice' and 'Bob' as easily as 0 and 1.
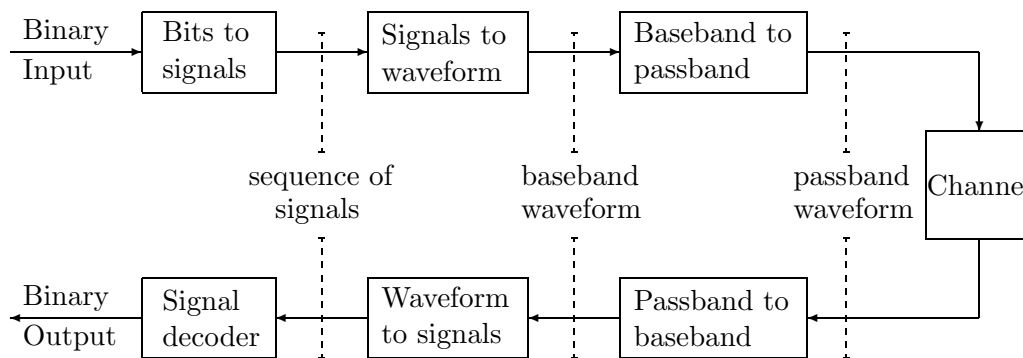
Figure 6.1: The layers of a modulator (channel encoder) and demodulator (channel decoder).

ARQ has always been an active area of communication and information theoretic research, but it will not be discussed here for the following reasons:

- It is important to understand communication in a single direction before addressing the complexities of two directions.
- Feedback does not increase channel capacity for typical channels (see [28]).
- Simple error detection and retransmission is best viewed as a topic in data networks.

There is an interesting analogy between analog source coding and digital modulation. With analog source coding, an analog waveform is first mapped into a sequence of real or complex numbers (*e.g.*, the coefficients in an orthogonal expansion). This sequence of signals is then quantized into a sequence of symbols from a discrete alphabet, and finally the symbols are encoded into a binary sequence. With modulation, a sequence of bits is encoded into a sequence of signals from a signal constellation. The elements of this constellation are real or complex points in one or several dimensions. This sequence of signal points is then mapped into a waveform by the inverse of the process for converting waveforms into sequences.

## 6.2 Pulse amplitude modulation (PAM)

*Pulse amplitude modulation*[4] (PAM) is probably the the simplest type of modulation. The incoming binary symbols are first segmented into $b$-bit blocks. There is a mapping from the set of $M = 2^b$ possible blocks into a signal constellation $\mathcal{A} = \{a_1, a_2, \ldots, a_M\}$ of real numbers. Let $R$ be the rate of incoming binary symbols in bits per second. Then the sequence of $b$-bit blocks, and the corresponding sequence, $u_1, u_2, \ldots$, of $M$-ary signals, has a rate of $R_s = R/b$ signals per second. The sequence of signals is then mapped into a waveform $u(t)$ by the use of time shifts of a basic pulse waveform $p(t)$, *i.e.*,

$$u(t) = \sum_k u_k \, p(t - kT), \tag{6.3}$$

where $T = 1/R_s$ is the interval between successive signals. The special case where $b = 1$ is called *binary* PAM and the case $b > 1$ is called *multilevel* PAM. Example 6.1.1 is an example

---

[4]The terminology comes from analog amplitude modulation, where a baseband waveform is modulated up to some passband for communication. For digital communication, the more interesting problem is turning a bit stream into a waveform at baseband.

of binary PAM where the basic pulse shape $p(t)$ is a sinc function. Comparing (6.1) with (6.3), we see that PAM is a special case of digital modulation in which the underlying set of functions $\phi_1(t), \phi_2(t), \dots$, is replaced by functions that are $T$-spaced time shifts of a basic function $p(t)$.

The following two subsections discuss the signal constellation (*i.e.*, the outer layer in Figure 6.1) and the subsequent two discuss the choice of pulse waveform $p(t)$ (*i.e.*, the middle layer in Figure 6.1). In most cases[5], the pulse waveform $p(t)$ is a baseband waveform and the resulting modulated waveform $u(t)$ is then modulated up to some passband (*i.e.*, the inner layer in Figure 6.1). Section 6.4 discusses modulation from baseband to passband and back.

### 6.2.1    Signal constellations

A *standard* $M$-PAM signal constellation $\mathcal{A}$ (see Figure 6.2) consists of $M = 2^b$ $d$-spaced real numbers located symmetrically about the origin; *i.e.*,

$$\mathcal{A} = \{\frac{-d(M-1)}{2}, \dots, \frac{-d}{2}, \frac{d}{2}, \dots, \frac{d(M-1)}{2}\}.$$

In other words, the signal points are the same as the representation points of a symmetric $M$-point uniform scalar quantizer.
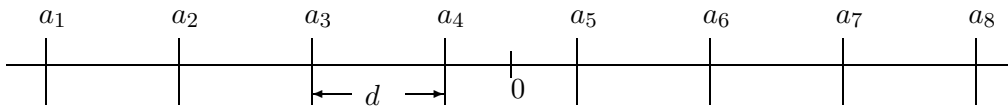


Figure 6.2: An 8-PAM signal set.

If the incoming bits are independent equiprobable random symbols (which is well approximated by effective source coding), then each signal $u_k$ is a sample value of a random variable $U_k$ that is equiprobable over the constellation (alphabet) $\mathcal{A}$. Also the sequence $U_1, U_2, \dots$, is independent and identically distributed (iid). As derived in Exercise 6.1, the mean squared signal value, or "energy per signal" $E_s = \mathsf{E}[U_k^2]$ is then given by

$$E_s = \frac{d^2(M^2 - 1)}{12} = \frac{d^2(2^{2b} - 1)}{12}. \tag{6.4}$$

For example, for $M = 2, 4$ and $8$, we have $E_s = d^2/4$, $5d^2/4$ and $21d^2/4$, respectively.

For $b$ greater than 2, $2^{2b} - 1$ is approximately $2^{2b}$, so we see that each unit increase in $b$ increases $E_s$ by a factor of 4. Thus increasing the rate $R$ by increasing $b$ requires impractically large energy for large $b$.

Before explaining why standard $M$-PAM is a good choice for PAM and what factors affect the choice of constellation size $M$ and distance $d$, a brief introduction to channel imperfections is required.

---

[5]Ultra-wide-band modulation (UAW) is an interesting modulation technique where the transmitted waveform is essentially a baseband PAM system over a 'baseband' of multiple gigahertz. This is discussed briefly in Chapter 9.

### 6.2.2 Channel imperfections: a preliminary view

Physical waveform channels are always subject to propagation delay, attenuation, and noise. Many wireline channels can be reasonably modeled using only these degradations, whereas wireless channels are subject to other degrations discussed in Chapter 9. This subsection provides a preliminary look at delay, then attenuation, and finally noise.

The time reference at a communication receiver is conventionally delayed relative to that at the transmitter. A waveform $u(t)$ at the transmitter is subject to propagation delay plus various filter delays in the modulator and demodulator. Thus $u(t)$, according to the transmitter clock, appears as $u(t-\tau)$ at the receiver, where $\tau$ is the overall delay. By delaying the receiver clock by $\tau$ from the transmitter clock, the received waveform, according to the receiver clock, is $u(t)$. With this convention, the channel can be modeled as having no delay, and all equations will be greatly simplified. This explains why communication engineers often model filters in the modulator and demodulator as being noncausal, since responses before time 0 can be added to the difference between the two clocks. *Estimating* the above fixed delay at the receiver is a significant problem called timing recovery, but is largely separable from the problem of recovering the transmitted data.

The magnitude of delay in a communication system is often important. It is one of the parameters often referred to as *quality of service* in a communication system. Delay is important for voice communication and often critically important when the communication is in the feedback loop of a real time control system. In addition to the fixed delay in time reference between modulator and demodulator, there is also delay in source encoding and decoding. Coding for error correction adds additional delay, which might or might not be counted as part of the modulator/demodulator delay. Either way, the delays in the source coding and error-correction coding are often much larger than that in the modulator/demodulator proper. Thus this latter delay can be significant, but is usually not of primary significance. Also, as channel speeds increase, the filtering delays in the modulator/demodulator become even less significant.

Amplitudes are usually measured on a different scale at transmitter and receiver. The actual power attenuation suffered in transmission is a product of amplifier gain, antenna coupling losses, antenna directional gain, propagation losses, etc. The process of finding all these gains and losses (and perhaps changing them) is called "the link budget." Such gains and losses are invariably calculated in decibels (dB). The number of decibels corresponding to a power gain $\alpha$ is defined to be $10 \log_{10} \alpha$. Thus power losses correspond to negative dB and power gains to positive dB. The use of a logarithmic measure of gain allows the various components of gain to be added rather than multiplied.

The use of decibels rather than some other logarithmic measure such as natural logs or logs to the base 2 is partly motivated by the ease of doing rough mental calculations. A factor of 2 is $10 \log_{10} 2 = 3.010 \cdots$ dB, approximated as 3 dB. Thus $4 = 2^2$ is 6 dB and 8 is 9 dB. Since 10 is 10 dB, we also see that 5 is 10/2 or 7 dB. We can just as easily see that 20 is 13 dB and so forth.

It is important to remember that the gains expressed in dB are *power* gains. Thus if there is a multiplicative gain of $g$ in a signal, this corresponds to a gain $g^2$ in power, which corresponds to $20 \log_{10} g$ dB.

The link budget in a communication system is largely separable from other issues, so the amplitude scale at the transmitter is usually normalized to that at the receiver.

By treating attenuation and delay as issues largely separable from modulation, we obtain a model of the channel in which a baseband waveform $u(t)$ is converted to passband and transmitted. At the receiver, after conversion back to baseband, a waveform $v(t) = u(t) + z(t)$ is received where $z(t)$ is noise. This noise is a fundamental limitation to communication and arises from a variety of causes, including thermal effects and unwanted radiation impinging on the receiver. Chapter 7 is largely devoted to understanding noise waveforms by modeling them as sample values of random processes. Chapter 8 then explains how best to decode signals in the presence of noise. These issues are briefly summarized here to see how they affect the choice of signal constellation.

If $p(t)$ is orthogonal to all its shifts by multiples of $T$, then, in the absence of noise, the transmitted signals $u_1, u_2, \dots$, can be retrieved from the baseband waveform $u(t)$ by the inner product operation,

$$u_k = \int u(t)\, p(t - kT)\, dt.$$

In the presence of noise, this same operation can be performed, yielding

$$v_k = \int v(t)\, p(t - kT)\, dt = u_k + z_k, \tag{6.5}$$

where $z_k = \int z(t)\, p(t - kT)\, dt$ is the projection of $z(t)$ on the shifted pulse $p(t - kT)$.

The most common (and often the most appropriate) model for noise on channels is called the additive white Gaussian noise model. As shown in Chapters 7 and 8, the above coefficients $\{z_k; k \in \mathbb{Z}\}$ in this model are the sample values of zero-mean, iid Gaussian random variables $\{Z_k; k \in \mathbb{Z}\}$. This is true no matter how the orthonormal functions $\{p(t-kT); k \in \mathbb{Z}\}$ are chosen, and these random variables are also independent of the signal random variables $\{U_k; k \in \mathbb{Z}\}$. Chapter 8 also shows that the operation in (6.5) is the appropriate operation to go from waveform to signal sequence in the layered demodulator of Figure 6.1.

Now consider the effect of the noise on the choice of $M$ and $d$ in a PAM modulator. Since the transmitted signal reappears at the receiver with a zero-mean Gaussian random variable added to it, any attempt to directly retrieve $U_k$ from $V_k$ with reasonably small probability of error[6] will require $d$ to exceed several standard deviations of the noise. Thus the noise determines how large $d$ must be, and this, combined with the power constraint, determines $M$.

The relation between error probability and signal-point spacing also helps explain why multi-level PAM systems almost invariably use a standard $M$-PAM signal set. Because the Gaussian density drops off so fast with increasing distance, the error probability due to confusion of nearest neighbors drops off equally fast. Thus error probability is dominated by the points in the constellation that are closest together. If the signal points are constrained to have some minimum distance $d$ between points, it can be seen that the minimum energy $E_s$ for a given number of points $M$ is achieved by the standard $M$-PAM set.[7]

To be more specific about the relationship between $M, d$ and the variance $\sigma^2$ of the noise $Z_k$, suppose that $d$ is selected to be $\alpha\sigma$, where $\alpha$ is chosen to make the detection sufficiently reliable. Then with $M = 2^b$, where $b$ is the number of bits encoded into each PAM signal, (6.4) becomes

$$E_s = \frac{\alpha^2 \sigma^2 (2^{2b} - 1)}{12}; \qquad b = \frac{1}{2} \log\left(1 + \frac{12 E_s}{\alpha^2 \sigma^2}\right). \tag{6.6}$$

---

[6]If error-correction coding is used with PAM, then $d$ can be smaller, but for any given error-correction code, $d$ still depends on the standard deviation of $Z_k$.

[7]On the other hand, if we choose a set of $M$ signal points to minimize $E_s$ for a given error probability, then the standard $M$-PAM signal set is not quite optimal (see Exercise 6.3).

This expression looks strikingly similar to Shannon's capacity formula for additive white Gaussian noise, which says that for the appropriate PAM bandwidth, the capacity per signal is $C = \frac{1}{2} \log(1 + \frac{E_s}{\sigma^2})$. The important difference is that in (6.6), $\alpha$ must be increased, thus decreasing $b$, in order to decrease error probability. Shannon's result, on the other hand, says that error probability can be made arbitrarily small for any number of bits per signal less than $C$. Both equations, however, show the same basic form of relationship between bits per signal and the signal to noise ratio $E_s/\sigma^2$. Both equations also say that if there is no noise ($\sigma^2 = 0$, then the the number of transmitted bits per signal can be infinitely large (*i.e.*, the distance $d$ between signal points can be made infinitesimally small). Thus both equations say that noise is a fundamental limitation on communication.

### 6.2.3  Choice of the modulation pulse

As defined in (6.3), the baseband transmitted waveform, $u(t) = \sum_k u_k \, p(t - kT)$, for a PAM modulator is determined by the signal constellation $\mathcal{A}$, the signal interval $T$ and the real $\mathcal{L}_2$ modulation pulse $p(t)$.

It may be helpful to visualize $p(t)$ as the impulse response of a linear time-invariant filter. Then $u(t)$ is the response of that filter to a sequence of $T$-spaced impulses $\{u_k \delta(t - kT)\}$. The problem of choosing $p(t)$ for a given $T$ turns out to be largely separable from that of choosing $\mathcal{A}$. The choice of $p(t)$ is also the more challenging and interesting problem.

The following objectives contribute to the choice of $p(t)$.

- $p(t)$ must be 0 for $t < -\tau$ for some finite $\tau$. To see this, assume that the $k$th input signal at the modulator is generated at time $Tk - \tau$. The contribution of $u_k$ to the transmitted waveform $u(t)$ cannot start until $kT - \tau$, which implies $p(t) = 0$ for $t < -\tau$ as stated. This rules out $\text{sinc}(t/T)$ as a choice for $p(t)$ (although $\text{sinc}(t/T)$ could be truncated at $t = -\tau$ to satisfy the condition).

- In most situations, $\hat{p}(f)$ should be essentially baseband limited to some bandwidth $B_b$ slightly larger than $\frac{1}{2T}$. We will see shortly that it cannot be baseband limited to less than $\frac{1}{2T}$. There is usually an upper limit on $B_b$ because of regulatory constraints at bandpass or to allow for other transmission channels in neighboring bands. If this limit were much larger than $\frac{1}{2T}$, then $T$ could be increased, increasing the rate of transmission.

- The retrieval of the sequence $\{u_k; k \in \mathbb{Z}\}$ from the noisy received waveform should be simple and relatively reliable. In the absence of noise, $\{u_k; k \in \mathbb{Z}\}$ should be uniquely specified by the received waveform.

The first condition above makes it somewhat tricky to satisfy the second condition. In particular, the Paley-Wiener theorem [20] states that a necessary and sufficient condition for a nonzero $\mathcal{L}_2$ function $p(t)$ to be zero for all $t < 0$ is that its Fourier transform satisfy

$$\int_{-\infty}^{\infty} \frac{|\ln |\hat{p}(f)||}{1 + f^2} \, df \; < \; \infty. \tag{6.7}$$

Combining this with the shift condition for Fourier transforms, it says that any $\mathcal{L}_2$ function that is 0 for all $t < -\tau$ for any finite delay $\tau$ must also satisfy (6.7). This is a particularly strong statement of the fact that functions cannot be both time and frequency limited. One consequence of (6.7) is that if $p(t) = 0$ for $t < -\tau$, then $\hat{p}(f)$ must be nonzero except on a set of

measure 0. Another consequence is that $\hat{p}(f)$ must go to 0 with increasing $f$ more slowly than exponentially.

The Paley-Wiener condition turns out to be useless as a tool for choosing $p(t)$. First, it distinguishes whether the above delay $\tau$ is finite or infinite, but gives no indication of its value when finite. Second, if an $\mathcal{L}_2$ function $p(t)$ is chosen with no concern for (6.7), it can then be truncated to be 0 for $t < -\tau$. The resulting $\mathcal{L}_2$ error caused by truncation can be made arbitrarily small by choosing $\tau$ sufficiently large. The tradeoff between truncation error and delay is clearly improved by choosing $p(t)$ to approach 0 rapidly as $t \to -\infty$.

In summary, we will replace the first objective above with the objective of choosing $p(t)$ to approach 0 rapidly as $t \to -\infty$. The resulting $p(t)$ will then be truncated to satisfy the original objective. Thus $p(t) \leftrightarrow \hat{p}(f)$ will be an approximation to the transmit pulse in what follows. This also means that $\hat{p}(f)$ can be strictly bandlimited to a frequency slightly larger than $\frac{1}{2T}$.

We next turn to the third objective, particularly that of easily retrieving the sequence $u_1, u_2, \ldots$, from $u(t)$ in the absence of noise. This problem was first analyzed in 1928 in a classic paper by Harry Nyquist [19]. Before looking at Nyquist's results, however, we must consider the demodulator.

### 6.2.4 PAM demodulation

For the time being, ignore the channel noise. Assume that the time reference and the amplitude scaling at the receiver have been selected so that the received baseband waveform is the same as the transmitted baseband waveform $u(t)$. This also assumes that no noise has been introduced by the channel.

The problem at the demodulator is then to retrieve the transmitted signals $u_1, u_2, \ldots$ from the received waveform $u(t) = \sum_k u_k p(t-kT)$. The middle layer of a PAM demodulator is defined by a signal interval $T$ (the same as at the modulator) and a real $\mathcal{L}_2$ waveform $q(t)$. The demodulator first filters the received waveform using a filter with impulse response $q(t)$. It then samples the output at $T$-spaced sample times. That is, the received filtered waveform is

$$r(t) = \int_{-\infty}^{\infty} u(\tau)q(t-\tau)\, d\tau, \tag{6.8}$$

and the received samples are $r(T), r(2T), \ldots, ,$.

Our objective is to choose $p(t)$ and $q(t)$ so that $r(kT) = u_k$ for each $k$. If this objective is met for all choices of $u_1, u_2, \ldots$, then the PAM system involving $p(t)$ and $q(t)$ is said to have *no intersymbol interference*. Otherwise, intersymbol interference is said to exist. The reader should verify that $p(t) = q(t) = \frac{1}{\sqrt{T}}\text{sinc}(\frac{t}{T})$ is one solution.

This problem of choosing filters to avoid intersymbol interference at first appears to be somewhat artificial. First, the form of the receiver is restricted to be a filter followed by a sampler. Exercise 6.4 shows that if the detection of each signal is restricted to a linear operation on the received waveform, then there is no real loss of generality in further restricting the operation to be a filter followed by a $T$-spaced sampler. This does not explain the restriction to linear operations, however.

The second artificiality is neglecting the noise, thus neglecting the fundamental limitation on the bit rate. The reason for posing this artificial problem is, first, that avoiding intersymbol interference is significant in choosing $p(t)$, and, second, that there is a simple and elegant solution

to this problem. This solution also provides part of the solution when noise is brought into the picture.

Recall that $u(t) = \sum_k u_k p(t - kT)$; thus from (6.8)

$$r(t) = \int_{-\infty}^{\infty} \sum_k u_k p(\tau - kT) q(t - \tau) \, d\tau. \tag{6.9}$$

Let $g(t)$ be the convolution $g(t) = p(t) * q(t) = \int p(\tau) q(t - \tau) \, d\tau$ and assume[8] that $g(t)$ is $\mathcal{L}_2$. We can then simplify (6.9) to

$$r(t) = \sum_k u_k g(t - kT). \tag{6.10}$$

This should not be surprising. The filters $p(t)$ and $q(t)$ are in cascade with each other. Thus $r(t)$ does not depend on which part of the filtering is done in one and which in the other; it is only the convolution $g(t)$ that determines $r(t)$. Later, when channel noise is added, the individual choice of $p(t)$ and $q(t)$ will become important.

There is no intersymbol interference if $r(kT) = u_k$ for each integer $k$, and from (6.10) this is satisfied if $g(0) = 1$ and $g(kT) = 0$ for each nonzero integer $k$. Waveforms with this property are said to be *ideal Nyquist* or, more precisely, *ideal Nyquist with interval $T$*.

Even though the clock at the receiver is delayed by some finite amount relative to that at the transmitter, and each signal $u_k$ can be generated at the transmitter at some finite time before $kT$, $g(t)$ must still have the property that $g(t) = 0$ for $t < -\tau$ for some finite $\tau$. As before with the transmit pulse $p(t)$, this finite delay constraint will be replaced with the objective that $g(t)$ should approach 0 rapidly as $|t| \to \infty$. Thus the function $\text{sinc}(\frac{t}{T})$ is ideal Nyquist with interval $T$, but is unsuitable because of the slow approach to 0 as $|t| \to \infty$.

As another simple example, the function $\text{rect}(t/T)$ is ideal Nyquist with interval $T$ and can be generated with finite delay, but is not remotely close to being baseband limited.

In summary, we want to find functions $g(t)$ that are ideal Nyquist but are approximately baseband limited and approximately time limited. The Nyquist criterion, discussed in the next section, provides a useful frequency characterization of functions that are ideal Nyquist. This characterization will then be used to study ideal Nyquist functions that are approximately baseband limited and approximately time limited.

## 6.3 The Nyquist criterion

The ideal Nyquist property is determined solely by the $T$-spaced samples of the waveform $g(t)$. This suggests that the results about aliasing should be relevant. Let $s(t)$ be the baseband-limited waveform generated by the samples of $g(t)$, *i.e.*,

$$s(t) = \sum_k g(kT) \, \text{sinc}(\frac{t}{T} - k). \tag{6.11}$$

---

[8]By looking at the frequency domain, it is not difficult to construct a $g(t)$ of infinite energy from $\mathcal{L}_2$ functions $p(t)$ and $q(t)$. When we study noise, however, we find that there is no point in constructing such a $g(t)$, so we ignore the possibility.

If $g(t)$ is ideal Nyquist, then all the above terms except $k = 0$ disappear and $s(t) = \text{sinc}(\frac{t}{T})$. Conversely, if $s(t) = \text{sinc}(\frac{t}{T})$, then $g(t)$ must be ideal Nyquist. Taking the Fourier transform of (6.11) shows that $g(t)$ is ideal Nyquist if and only if

$$\hat{s}(f) = T \, \text{rect}(fT). \tag{6.12}$$

From the aliasing theorem,

$$\hat{s}(f) = \text{l.i.m.} \sum_m \hat{g}(f + \frac{m}{T}) \, \text{rect}(fT). \tag{6.13}$$

The result of combining (6.12) and (6.13) is the Nyquist criterion:

**Theorem 6.3.1 (Nyquist criterion).** *Let   $\hat{g}(f)$   be   $\mathcal{L}_2$   and   satisfy   the   condition* $\lim_{|f| \to \infty} \hat{g}(f)|f|^{1+\varepsilon} = 0$ *for some $\varepsilon > 0$. Then the inverse transform, $g(t)$, of $\hat{g}(f)$ is ideal Nyquist with interval $T$ if and only if $\hat{g}(f)$ satisfies the <u>Nyquist criterion</u> for $T$, defined as*[9]

$$\text{l.i.m.} \sum_m \hat{g}(f + m/T) \, \text{rect}(fT) = T \, \text{rect}(fT). \tag{6.14}$$

**Proof:** From the aliasing theorem, the baseband approximation $s(t)$ in (6.11) converges pointwise and is $\mathcal{L}_2$. Similarly, the Fourier transform $\hat{s}(f)$ satisfies (6.13). If $g(t)$ is ideal Nyquist, then $s(t) = \text{sinc}(\frac{t}{T})$. This implies that $\hat{s}(f)$ is $\mathcal{L}_2$ equivalent to $T \, \text{rect}(fT)$, which in turn implies (6.14). Conversely, satisfaction of the Nyquist criterion (6.14) implies that $\hat{s}(f) = T \, \text{rect}(fT)$. This implies $s(t) = \text{sinc}(\frac{t}{T})$ implying that $g(t)$ is ideal Nyquist.    □

There are many choices for $\hat{g}(f)$ that satisfy (6.14), but the ones of major interest are those that are approximately both bandlimited and time limited. We look specifically at cases where $\hat{g}(f)$ is strictly bandlimited, which, as we have seen, means that $g(t)$ is not strictly time limited. Before these filters can be used, of course, they must be truncated to be strictly time limited. It is strange to look for strictly bandlimited and approximately time-limited functions when it is the opposite that is required, but the reason is that the frequency constraint is the more important. The time constraint is usually more flexible and can be imposed as an approximation.

### 6.3.1   Band-edge symmetry

The *nominal or Nyquist band* associated with a PAM pulse $g(t)$ with signal interval $T$ is defined to be $W_b = 1/(2T)$. The actual baseband bandwidth[10] $B_b$ is defined as the smallest number $B_b$ such that $\hat{g}(f) = 0$ for $|f| > B_b$. Note that if $\hat{g}(f) = 0$ for $|f| > W_b$, then the left side of (6.14) is zero except for $m = 0$, so $\hat{g}(f) = T \, \text{rect}(fT)$. This means that $B_b \geq W_b$ and equality holds if and only if $g(t) = \text{sinc}(t/T)$.

As discussed above, if $W_b$ is much smaller than $B_b$, then $W_b$ can be increased, thus increasing the rate $R_s$ at which signals can be transmitted. Thus $g(t)$ should be chosen in such a way that

---

[9]It can be seen that $\sum_m \hat{g}(f + m/T)$ is periodic and thus the $\text{rect}(fT)$ could be essentially omitted from both sides of (6.14). Doing this, however, would make the limit in the mean meaningless and would also complicate the intuitive understanding of the theorem.

[10]It might be better to call this the design bandwidth, since after the truncation necessary for finite delay, the resulting frequency function is nonzero almost everywhere. However, if the delay is large enough, the energy outside of $B_b$ is negligible. On the other hand, Exercise 6.9 shows that these approximations must be handled with great care.

$B_b$ exceeds $W_b$ by a relatively small amount. In particular, we now focus on the case where $W_b \leq B_b < 2W_b$.

The assumption $B_b < 2W_b$ means that $\hat{g}(f) = 0$ for $|f| \geq 2W_b$. Thus for $0 \leq f \leq W_b$, $\hat{g}(f + 2mW_b)$ can be nonzero only for $m = 0$ and $m = -1$. Thus the Nyquist criterion (6.14) in this positive frequency interval becomes

$$\hat{g}(f) + \hat{g}(f - 2W_b) = T \qquad \text{for } 0 \leq f \leq W_b. \tag{6.15}$$

Since $p(t)$ and $q(t)$ are real, $g(t)$ is also real, so $\hat{g}(f-2W_b) = \hat{g}^*(2W_b-f)$. Substituting this in (6.15) and letting $\Delta = f - W_b$, (6.15) becomes

$$T - \hat{g}(W_b+\Delta) = \hat{g}^*(W_b-\Delta). \tag{6.16}$$

This is sketched and interpreted in Figure 6.3. The figure assumes the typical situation in which $\hat{g}(f)$ is real. In the general case, the figure illustrates the real part of $\hat{g}(f)$ and the imaginary part satisfies $\Im\{\hat{g}(W_b+\Delta)\} = \Im\{\hat{g}(W_b-\Delta)\}$.
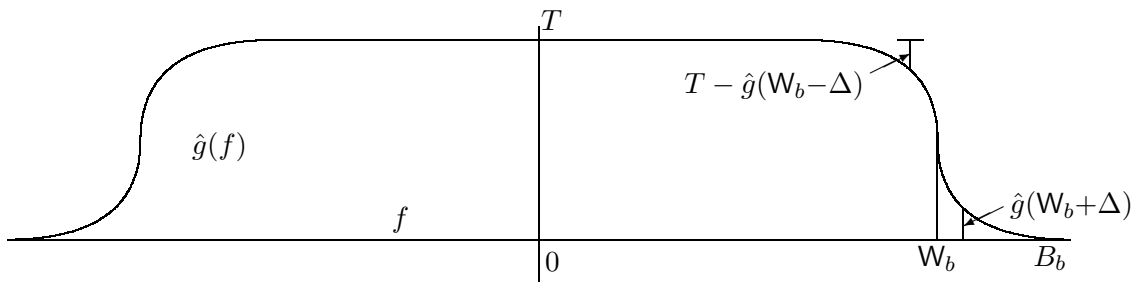


Figure 6.3: Band edge symmetry illustrated for real $\hat{g}(f)$: For each $\Delta$, $0 \leq \Delta \leq W_b$, $\hat{g}(W_b+\Delta) = T - \hat{g}(W_b-\Delta)$. The portion of the curve for $f \geq W_b$, rotated by $180^o$ around the point $(W_b, T/2)$, is equal to the portion of the curve for $f \leq W_b$.

Figure 6.3 makes it particularly clear that $B_b$ must satisfy $B_b \geq W_b$ to avoid intersymbol interference. We then see that the choice of $\hat{g}(f)$ involves a tradeoff between making $\hat{g}(f)$ smooth, so as to avoid a slow time decay in $g(t)$, and reducing the excess of $B_b$ over the Nyquist bandwidth $W_b$. This excess is expressed as a *rolloff factor*[11], defined to be $(B_b/W_b) - 1$, usually expressed as a percentage. Thus $\hat{g}(f)$ in the figure has about a 30% rolloff.

PAM filters in practice often have *raised cosine* transforms. The raised cosine frequency function, for any given rolloff $\alpha$ between 0 and 1, is defined by

$$\hat{g}_\alpha(f) = \begin{cases} T, & 0 \leq |f| \leq \frac{1-\alpha}{2T}; \\ T \cos^2\left[\frac{\pi T}{2\alpha}(|f| - \frac{1-\alpha}{2T})\right], & \frac{1-\alpha}{2T} \leq |f| \leq \frac{1+\alpha}{2T}; \\ 0, & |f| \geq \frac{1+\alpha}{2T}. \end{cases} \tag{6.17}$$

---

[11] The requirement for a small rolloff actually arises from a requirement on the transmitted pulse $p(t)$, *i.e.*, on the actual bandwidth of the transmitted channel waveform, rather than on the cascade $g(t) = p(t) * q(t)$. The tacit assumption here is that $\hat{p}(f) = 0$ when $\hat{g}(f) = 0$. One reason for this is that it is silly to transmit energy in a part of the spectrum that is going to be completely filtered out at the receiver. We see later that $\hat{p}(f)$ and $\hat{q}(f)$ are usually chosen to have the same magnitude, ensuring that $\hat{p}(f)$ and $\hat{g}(f)$ have the same rolloff.

The inverse transform of $\hat{g}_\alpha(f)$ can be shown to be (see Exercise 6.8)

$$g_\alpha(t) = \text{sinc}(\frac{t}{T}) \, \frac{\cos(\pi\alpha t/T)}{1 - 4\alpha^2 t^2/T^2} \, , \tag{6.18}$$

which decays asymptotically as $1/t^3$, compared to $1/t$ for $\text{sinc}(\frac{t}{T})$. In particular, for a rolloff $\alpha = 1$, $\hat{g}_\alpha(f)$ is nonzero from $-2\mathsf{W}_b = -1/T$ to $2\mathsf{W}_b = 1/T$ and $g_\alpha(t)$ has most of its energy between $-T$ and $T$. Rolloffs as sharp as 5–10% are used in current practice. The resulting $g_\alpha(t)$ goes to 0 with increasing $|t|$ much faster than $\text{sinc}(t/T)$, but the ratio of $g_\alpha(t)$ to $\text{sinc}(t/T)$ is a function of $\alpha t/T$ and reaches its first zero at $t = 1.5T/\alpha$. In other words, the required filtering delay is proportional to $1/\alpha$.

The motivation for the raised cosine shape is that $\hat{g}(f)$ should be smooth in order for $g(t)$ to decay quickly in time, but $\hat{g}(f)$ must decrease from $T$ at $\mathsf{W}_b(1 - \alpha)$ to 0 at $\mathsf{W}_b(1 + \alpha)$; as seen in Figure 6.3, the raised cosine function simply rounds off the step discontinuity in $\text{rect}(\frac{f}{2\mathsf{W}_b})$ in such a way as to maintain the Nyquist criterion while making $\hat{g}(f)$ continuous with a continuous derivitive, thus guaranteeing that $g(t)$ decays asympototically with $1/t^3$.

### 6.3.2   Choosing $\{p(t-kT); k \in \mathbb{Z}\}$ as an orthonormal set

The above subsection describes the choice of $\hat{g}(f)$ as a compromise between rolloff and smoothness, subject to band edge symmetry. As illustrated in figure 6.3, it is not a serious additional constraint to restrict $\hat{g}(f)$ to be real and nonnegative (why let $\hat{g}(f)$ go negative or imaginary in making a smooth transition from $T$ to 0?). After choosing $\hat{g}(f) \geq 0$, however, the question remains of choosing the transmit filter $p(t)$ and the receive filter $q(t)$ subject to $\hat{p}(f)\hat{q}(f) = \hat{g}(f)$. When studying white Gaussian noise later, we will find that $\hat{q}(f)$ should be chosen to equal $\hat{p}^*(f)$. Thus[12],

$$|\hat{p}(f)| = |\hat{q}(f)| = \sqrt{\hat{g}(f)} \, . \tag{6.19}$$

The phase of $\hat{p}(f)$ can be chosen in an arbitrary way, but this determines the phase of $\hat{q}(f) = \hat{p}^*(f)$. The requirement that $\hat{p}(f)\hat{q}(f) = \hat{g}(f) \geq 0$ means that $\hat{q}(f) = \hat{p}^*(f)$. In addition, if $p(t)$ is real then $\hat{p}(-f) = \hat{p}^*(f)$, which determines the phase for negative $f$ in terms of an arbitrary phase for $f > 0$. It is convenient here, however, to be slightly more general and allow $p(t)$ to be complex. We will prove the following important theorem:

**Theorem 6.3.2 (Orthonormal shifts).** *Let $p(t)$ be an $\mathcal{L}_2$ function such that $\hat{g}(f) = |\hat{p}(f)|^2$ satisfies the Nyquist criterion for $T$. Then $\{p(t-kT); k \in \mathbb{Z}\}$ is a set of orthonormal functions. Conversely, if $\{p(t-kT); k \in \mathbb{Z}\}$ is a set of orthonormal functions, then $|\hat{p}(f)|^2$ satisfies the Nyquist criterion.*

**Proof:** Let $q(t) = p^*(-t)$. Then $g(t) = p(t) * q(t)$ so that

$$g(kT) = \int_{-\infty}^{\infty} p(\tau)q(kT - \tau) \, d\tau = \int_{-\infty}^{\infty} p(\tau)p^*(\tau - kT) \, d\tau. \tag{6.20}$$

If $\hat{g}(f)$ satisfies the Nyquist criterion, then $g(t)$ is ideal Nyquist and (6.20) has the value 0 for each integer $k \neq 0$ and has the value 1 for $k = 0$. By shifting the variable of integration by

---

[12]A function $p(t)$ satisfying (6.19) is often called square root of Nyquist, although it is the magnitude of the transform that is the square root of the transform of an ideal Nyquist pulse.

$jT$ for any integer $j$ in (6.20), we see also that $\int p(\tau - jT)p^*(\tau - (k+j)T)\, d\tau = 0$ for $k \neq 0$ and 1 for $k = 0$. Thus $\{p(t - kT); k \in \mathbb{Z}\}$ is an orthonormal set. Conversely, assume that $\{p(t - kT); k \in \mathbb{Z}\}$ is an orthonormal set. Then (6.20) has the value 0 for integer $k \neq 0$ and 1 for $k = 0$. Thus $g(t)$ is ideal Nyquist and $\hat{g}(f)$ satisfies the Nyquist criterion. $\qquad\square$

Given this orthonormal shift property for $p(t)$, the PAM transmitted waveform $u(t) = \sum_k u_k p(t - kT)$ is simply an orthonormal expansion. Retrieving the coefficient $u_k$ then corresponds to projecting $u(t)$ onto the one dimensional subspace spanned by $\boldsymbol{p}_k$. Note that this projection is accomplished by filtering $u(t)$ by $q(t)$ and then sampling at time $kT$. The filter $q(t)$ is called the *matched filter* to $p(t)$. We discuss these filters later when noise is introduced into the picture.

Note that we have restricted the pulse $p(t)$ to have unit energy. There is no loss of generality here, since the input signals $\{u_k\}$ can be scaled arbitrarily and there is no point in having an arbitrary scale factor in both places.

For $|\hat{p}(f)|^2 = \hat{g}(f)$, the actual bandwidth of $\hat{p}(f), \hat{q}(f)$, and $\hat{g}(f)$ are the same, say $B_b$. Thus if $B_b < \infty$, we see that $p(t)$ and $q(t)$ can be realized only with infinite delay, which means that both must be truncated. Since $q(t) = p^*(-t)$, they must be truncated for both positive and negative $t$. We assume that they are truncated at such a large value of delay that the truncation error is negligible. Note that the delay generated by both the transmitter and receiver filter (*i.e.*, from the time that $u_k p(t - kT)$ starts to be formed at the transmitter to the time when $u_k$ is sampled at the receiver) is twice the duration of $p(t)$.

### 6.3.3 Relation between PAM and analog source coding

The main emphasis in PAM modulation has been that of converting a sequence of $T$-spaced signals into a waveform. Similarly, the first part of analog source coding is often to convert a waveform into a $T$-spaced sequence of samples. The major difference is that with PAM modulation, we have control over the PAM pulse $p(t)$ and thus some control over the class of waveforms. With source coding, we are stuck with whatever class of waveforms describes the source of interest.

For both systems the nominal bandwidth is $\mathsf{W}_b = 1/(2T)$ and $B_b$ can be defined as the actual baseband bandwidth of the waveforms. In the case of source coding, $B_b \leq \mathsf{W}_b$ is a necessary condition for the sampling appoximation $\sum_k u(kT)\operatorname{sinc}(\frac{t}{T} - k)$ to perfectly recreate the waveform $u(t)$. The aliasing theorem and the $T$-spaced sinc weighted sinusoid expansion were used to analyze the squared error if $B_b > \mathsf{W}_b$.

For PAM, on the other hand, the necessary condition for the PAM demodulator to recreate the initial PAM sequence is $B_b \geq \mathsf{W}_b$. With $B_b > \mathsf{W}_b$, aliasing can be used to advantage, creating an aggregate pulse $g(t)$ that is ideal Nyquist. There is considerable choice in such a pulse, and it is chosen by using contributions from both $f < \mathsf{W}_b$ and $f > \mathsf{W}_b$. Finally we saw that the transmission pulse $p(t)$ for PAM can be chosen so that its $T$-spaced shifts form an orthonormal set. The sinc functions have this property, but many other waveforms with slightly greater bandwidth have the same property but decay much faster with $t$.

## 6.4    Modulation: baseband to passband and back

The discussion of PAM in the previous 2 sections focussed on converting a $T$-spaced sequence of real signals into a real waveform of bandwidth $B_b$ slightly larger than the Nyquist bandwidth $W_b = \frac{1}{2T}$. This section focuses on converting that baseband waveform into a passband waveform appropriate for the physical medium, regulatory constraints, and avoiding other transmission bands.

### 6.4.1    Double-sideband amplitude modulation

The objective of modulating a baseband PAM waveform $u(t)$ to some high frequency passband around some carrier $f_c$ is to simply shift $\hat{u}(f)$ up in frequency to $\hat{u}(f)e^{2\pi i f_c t}$. Thus if $\hat{u}(f)$ is zero except for $-B_b \leq f \leq B_b$, then the shifted version would be zero except for $f_c - B_b \leq f \leq f_c + B_b$. This does not quite work since it results in a complex waveform, whereas only real waveforms can actually be transmitted. Thus $u(t)$ is also multiplied by the complex conjugate of $e^{2\pi i f_c t}$, i.e., $e^{-2\pi i f_c t}$, resulting in the following passband waveform:

$$x(t) \quad = \quad u(t)[e^{2\pi i f_c t} + e^{-2\pi i f_c t}] = 2u(t)\cos(2\pi f_c t), \qquad (6.21)$$
$$\hat{x}(f) \quad = \quad \hat{u}(f - f_c) + \hat{u}(f + f_c). \qquad (6.22)$$

As illustrated in Figure 6.4, $u(t)$ is both translated up in frequency by $f_c$ and also translated down by $f_c$. Since $x(t)$ must be real, $\hat{x}(f) = \hat{x}^*(-f)$, and the negative frequencies cannot be avoided. Note that the entire set of frequencies in $[-B_b, B_b]$ is both translated up to $[-B_b + f_c, B_b + f_c]$ and down to $[-B_b - f_c, \ B_b - f_c]$. Thus (assuming $f_c > B_b$) the range of nonzero frequencies occupied by $x(t)$ is twice as large as that occupied by $u(t)$.
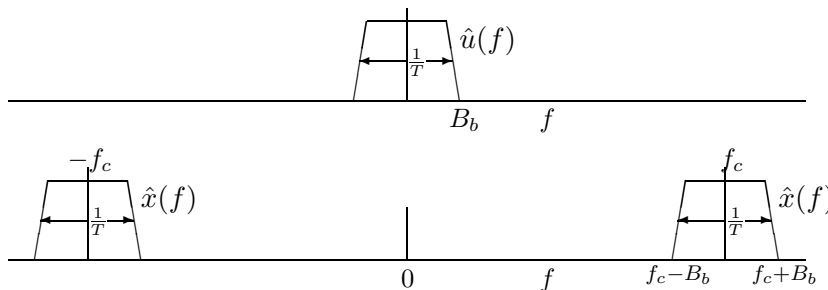


Figure 6.4:  Frequency domain representation of a baseband waveform $u(t)$ shifted up to a passband around the carrier $f_c$. Note that the baseband bandwidth $B_b$ of $u(t)$ has been doubled to the passband bandwidth $B = 2B_b$ of $x(t)$.

In the communication field, the *bandwidth* of a system is universally defined as the range of *positive* frequencies used in transmission. Since transmitted waveforms are real, the negative frequency part of those waveforms is determined by the positive part and is not counted. This is consistent with our earlier baseband usage, where $B_b$ is the bandwidth of the baseband waveform $u(t)$ in Figure 6.4, and with our new usage for passband waveforms where $B = 2B_b$ is the bandwidth of $\hat{x}(f)$.

The passband modulation scheme described by (6.21) is called *double-sideband amplitude modulation*. The terminology comes not from the negative frequency band around $-f_c$ and the

positive band around $f_c$, but rather from viewing $[f_c-B_b,\ f_c+B_b]$ as two sidebands, the upper, $[f_c,\ f_c+B_b]$, coming from the positive frequency components of $u(t)$ and the lower, $[f_c-B_b,\ f_c]$ from its negative components. Since $u(t)$ is real, these two bands are redundant and either could be reconstructed from the other.

Double-sideband modulation is quite wasteful of bandwidth since half of the band is redundant. Redundancy is often useful for added protection against noise, but such redundancy is usually better achieved through digital coding.

The simplest and most widely employed solution for using this wasted bandwidth[13] is *quadrature amplitude modulation* (QAM), which is described in the next section. PAM at passband is appropriately viewed as a special case of QAM, and thus the demodulation of PAM from passband to baseband is discussed at the same time as the demodulation of QAM.

## 6.5 Quadrature amplitude modulation (QAM)

QAM is very similar to PAM except that with QAM the baseband waveform $u(t)$ is chosen to be complex. The complex QAM waveform $u(t)$ is then shifted up to passband as $u(t)e^{2\pi i f_c t}$. This waveform is complex and is converted into a real waveform for transmission by adding its complex conjugate. The resulting real passband waveform is then

$$x(t) = u(t)e^{2\pi i f_c t} + u^*(t)e^{-2\pi i f_c t}\,. \tag{6.23}$$

Note that the passband waveform for PAM in (6.21) is a special case of this in which $u(t)$ is real. The passband waveform $x(t)$ in (6.23) can also be written in the following equivalent ways:

$$\begin{aligned}
x(t) &= 2\Re\{u(t)e^{2\pi i f_c t}\} \tag{6.24}\\
&= 2\Re\{u(t)\}\cos(2\pi f_c t) - 2\Im\{u(t)\}\sin(2\pi f_c t)\,. \tag{6.25}
\end{aligned}$$

The factor of 2 in (6.24) and (6.25) is an arbitrary scale factor. Some authors leave it out, (thus requiring a factor of $1/2$ in (6.23)) and others replace it by $\sqrt{2}$ (requiring a factor of $1/\sqrt{2}$ in (6.23)). This scale factor (however chosen) causes additional confusion when we look at the energy in the waveforms. With the scaling here, $\|\boldsymbol{x}\|^2 = 2\|\boldsymbol{u}\|^2$. Using the scale factor $\sqrt{2}$ solves this problem, but introduces many other problems, not least of which is an extraordinary number of $\sqrt{2}$'s in equations. At one level, scaling is a trivial matter, but although the literature is inconsistent, we have tried to be consistent here. One intuitive advantage of the convention here, as illustrated in Figure 6.4 is that the positive frequency part of $x(t)$ is simply $u(t)$ shifted up by $f_c$.

The remainder of this section provides a more detailed explanation of QAM, and thus also of a number of issues about PAM. A QAM modulator (see figure 6.5) has the same 3 layers as a PAM modulator, *i.e.*, first mapping a sequence of bits to a sequence of complex signals, then mapping the complex sequence to a complex baseband waveform, and finally mapping the complex baseband waveform to a real passband waveform.

The demodulator, not surprisingly, performs the inverse of these operations in reverse order, first mapping the received bandpass waveform into a baseband waveform, then recovering the

---

[13]An alternate approach is single-sideband modulation. Here either the positive or negative sideband of a double-sideband waveform is filtered out, thus reducing the transmitted bandwidth by a factor of 2. This used to be quite popular for analog communication but is harder to implement for digital communication than QAM.

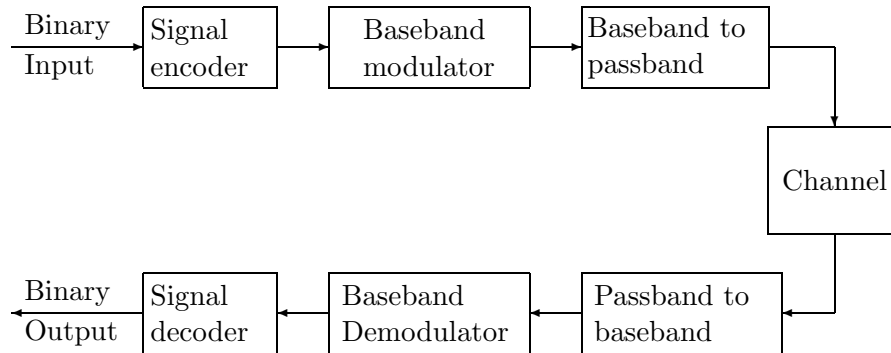sequence of signals, and finally recovering the binary digits. Each of these layers is discussed in turn.



Figure 6.5: QAM modulator and demodulator.

### 6.5.1   QAM signal set

The input bit sequence arrives at a rate of $R$ b/s and is converted, $b$ bits at a time, into a sequence of complex signals $u_k$ chosen from a *signal set* (alphabet, constellation) $\mathcal{A}$ of size $M = |\mathcal{A}| = 2^b$. The *signal rate* is thus $R_s = R/b$ signals per second, and the *signal interval* is $T = 1/R_s = b/R$ sec.

In the case of QAM, the transmitted signals $u_k$ are complex numbers $u_k \in \mathbb{C}$, rather than real numbers. Alternatively, we may think of each signal as a real 2-tuple in $\mathbb{R}^2$.
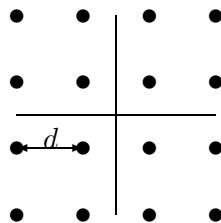
A *standard $(M' \times M')$-QAM signal set*, where $M = (M')^2$ is the Cartesian product of two $M'$-PAM sets; *i.e.*,

$$\mathcal{A} = \{(a' + ia'') \mid a' \in \mathcal{A}', a'' \in \mathcal{A}'\},$$

where

$$\mathcal{A}' = \{-d(M'-1)/2, \dots, -d/2, d/2, \dots, d(M'-1)/2\}.$$

The signal set $\mathcal{A}$ thus consists of a square array of $M = (M')^2 = 2^b$ signal points located symmetrically about the origin, as illustrated below for $M = 16$.



The minimum distance between the two-dimensional points is denoted by $d$. Also the average energy per two-dimensional signal, which is denoted $E_s$, is simply twice the average energy per

dimension:

$$E_s = \frac{d^2[(M')^2 - 1]}{6} = \frac{d^2[M-1]}{6}.$$

In the case of QAM there are clearly many ways to arrange the signal points other than on a square grid as above. For example, in an $M$-PSK (phase-shift keyed) signal set, the signal points consist of $M$ equally spaced points on a circle centered on the origin. Thus 4-PSK = 4-QAM. For large $M$ it can be seen that the signal points become very close to each other on a circle so that PSK is rarely used for large $M$. On the other hand, PSK has some practical advantages because of the uniform signal magnitudes.

As with PAM, the probability of decoding error is primarily a function of the minimum distance $d$. Not surprisingly, $E_s$ is linear in the signal power of the passband waveform. In wireless systems the signal power is limited both to conserve battery power and to meet regulatory requirements. In wired systems, the power is limited both to avoid crosstalk between adjacent wires and adjacent frequencies, and also to avoid nonlinear effects.

For all of these reasons, it is desirable to choose signal constellations that approximately minimize $E_s$ for a given $d$ and $M$. One simple result here is that a hexagonal grid of signal points achieves smaller $E_s$ than a square grid for very large $M$ and fixed minimum distance. Unfortunately, finding the optimal signal set to minimize $E_s$ for practical values of $M$ is a messy and ugly problem, and the minima have few interesting properties or symmetries. We will not spend further time on this other than a few exercises and will usually simply assume a standard $(M' \times M')$-QAM signal set, which is almost universally used in practice.

The standard $(M' \times M')$-QAM signal set is almost universally used in practice and will be assumed in what follows.

### 6.5.2 QAM baseband modulation and demodulation

A QAM baseband modulator is determined by the signal interval $T$ and a complex $\mathcal{L}_2$ waveform $p(t)$. The discrete-time sequence $\{u_k\}$ of complex signal points modulates the amplitudes of a sequence of time shifts $\{p(t-kT)\}$ of the basic pulse $p(t)$ to create a complex transmitted signal $u(t)$ as follows:

$$u(t) = \sum_{k \in \mathbb{Z}} u_k p(t-kT). \tag{6.26}$$

As in the PAM case, we could choose $p(t)$ to be $\mathrm{sinc}(\frac{t}{T})$, but for the same reasons as before, $p(t)$ should decay with increasing $|t|$ faster than the sinc function. This means that $\hat{p}(f)$ should be a continuous function that goes to zero rapidly but not instantaneously as $f$ increases beyond $1/(2T)$. As with PAM, we define $\mathsf{W}_b = \frac{1}{2T}$ to be the nominal baseband bandwidth of the QAM modulator and $B_b$ to be the actual design bandwidth.

Assume for the moment that the process of conversion to passband, channel transmission, and conversion back to baseband, is ideal, recreating the baseband modulator output $u(t)$ at the input to the baseband demodulator. The baseband demodulator is determined by the interval $T$ (the same as at the modulator) and an $\mathcal{L}_2$ waveform $q(t)$. The demodulator filters $u(t)$ by $q(t)$ and samples the output at $T$-spaced sample times. Denoting the filtered output by

$$r(t) = \int_{-\infty}^{\infty} u(\tau)q(t-\tau)\,d\tau,$$

we see that the received samples are $r(T), r(2T), \ldots$. Note that this is the same as the PAM demodulator except that real signals have been replaced by complex signals. As before, the output $r(t)$ can be represented as

$$r(t) = \sum_k u_k g(t - kT),$$

where $g(t)$ is the convolution of $p(t)$ and $q(t)$. As before, $r(kT) = u_k$ if $g(t)$ is ideal Nyquist, namely if $g(0) = 1$ and $g(kT) = 0$ for all nonzero integer $k$.

The proof of the Nyquist criterion, Theorem 6.3.1, is valid whether or not $g(t)$ is real. For the reasons explained earlier, however, $\hat{g}(f)$ is usually real and symmetric (as with the raised cosine functions) and this implies that $g(t)$ is also real and symmetric.

Finally, as discussed with PAM, $\hat{p}(f)$ is usually chosen to satisfy $|\hat{p}(f)| = \sqrt{\hat{g}(f)}$. Choosing $\hat{p}(f)$ in this way does not specify the phase of $\hat{p}(f)$, and thus $\hat{p}(f)$ might be real or complex. However $\hat{p}(f)$ is chosen, subject to $|\hat{g}(f)|^2$ satisfying the Nyquist criterion, the set of time shifts $\{p(t-kT)\}$ form an orthonormal set of functions. With this choice also, the baseband bandwidth of $u(t)$, $p(t)$, and $g(t)$ are all the same. Each has a nominal baseband bandwidth given by $\frac{1}{2T}$ and each has an actual baseband bandwidth that exceeds $\frac{1}{2T}$ by some small rolloff factor. As with PAM, $p(t)$ and $q(t)$ must be truncated in time to allow finite delay. The resulting filters are then not quite bandlimited, but is viewed as a negligible implementation error.

In summary, QAM baseband modulation is virtually the same as PAM baseband modulation. The signal set for QAM is of course complex, and the modulating pulse $p(t)$ can be complex, but the Nyquist results about avoiding intersymbol interference are unchanged.

### 6.5.3   QAM: baseband to passband and back

Next we discuss modulating the complex QAM baseband waveform $u(t)$ to the passband waveform $x(t)$. Alternative expressions for $x(t)$ are given by (6.23), (6.24). and (6.25) and the frequency representation is illustrated in Figure 6.4.

As with PAM, $u(t)$ has a nominal baseband bandwidth $\mathsf{W}_b = \frac{1}{2T}$. The actual baseband bandwidth $B_b$ exceeds $\mathsf{W}_b$ by some small rolloff factor. The corresponding passband waveform $x(t)$ has a nominal passband bandwidth $\mathsf{W} = 2\mathsf{W}_b = \frac{1}{T}$ and an actual passband bandwidth $B = 2B_b$. We will assume in everything to follow that $B/2 < f_c$. Recall that $u(t)$ and $x(t)$ are idealized approximations of the true baseband and transmitted waveforms. These true baseband and transmitted waveforms must have finite delay and thus infinite bandwidth, but it is assumed that the delay is large enough that the approximation error is negligible. The assumption[14] $B/2 < f_c$ implies that $u(t)e^{2\pi i f_c t}$ is constrained to positive frequencies and $u(t)e^{-2\pi i f_c t}$ to negative frequencies. Thus the Fourier transform $\hat{u}(f - f_c)$ does not overlap with $\hat{u}(f + f_c)$.

As with PAM, the modulation from baseband to passband is viewed as a two step process. First $u(t)$ is translated up in frequency by an amount $f_c$, resulting in a complex passband waveform $x^+(t) = u(t)e^{2\pi i f_c t}$. Next $x^+(t)$ is converted to the real passband waveform $x(t) = [x^+(t)]^* + x^+(t)$.

---

[14]Exercise 6.11 shows that when this assumption is violated, $u(t)$ can not be perfectly retrieved from $x(t)$, even in the absence of noise. The negligible frequency components of the truncated version of $u(t)$ outside of $B/2$ are assumed to cause negligible error in demodulation.

Assume for now that $x(t)$ is transmitted to the receiver with no noise and no delay. In principle, the received $x(t)$ can be modulated back down to baseband by the reverse of the two steps used in going from baseband to passband. That is, $x(t)$ must first be converted back to the complex positive passband waveform $x^+(t)$, and then $x^+(t)$ must be shifted down in frequency by $f_c$.

Mathematically, $x^+(t)$ can be retrieved from $x(t)$ simply by filtering $x(t)$ by a complex filter $h(t)$ such that $\hat{h}(f) = 0$ for $f < 0$ and $\hat{h}(f) = 1$ for $f > 0$. This filter is called a *Hilbert filter*. Note that $h(t)$ is not an $\mathcal{L}_2$ function, but it can be converted to $\mathcal{L}_2$ by making $\hat{h}(f)$ have the value 0 except in the positive passband $[\frac{-B}{2}+f_c, \frac{B}{2}+f_c]$ where it has the value 1. We can then easily retrieve $u(t)$ from $x^+(t)$ simply by a frequency shift. Figure 6.6 illustrates the sequence of operations from $u(t)$ to $x(t)$ and back again.
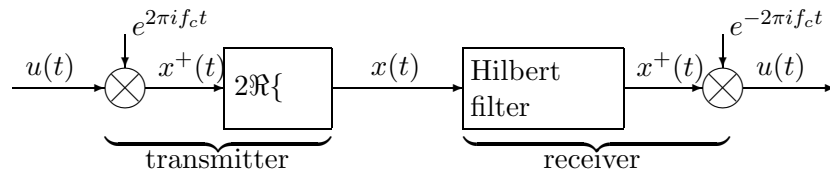


Figure 6.6: Baseband to passband and back.

### 6.5.4   Implementation of QAM

From an implementation standpoint, the baseband waveform $u(t)$ is usually implemented as two real waveforms, $\Re\{u(t)\}$ and $\Im\{u(t)\}$. These are then modulated up to passband using multiplication by in-phase and out-of-phase carriers as in (6.25), *i.e.*,

$$x(t) = 2\Re\{u(t)\} \cos(2\pi f_c t) - 2\Im\{u(t)\} \sin(2\pi f_c t).$$

There are many other possible implementations, however, such as starting with $u(t)$ given as magnitude and phase. The positive frequency expression $x^+(t) = u(t)e^{2\pi i f_c t}$ is a complex multiplication of complex waveforms which requires 4 real multiplications rather than the two above used to form $x(t)$ directly. Thus going from $u(t)$ to $x^+(t)$ to $x(t)$ provides insight but not ease of implementation.

The baseband waveforms $\Re\{u(t)\}$ and $\Im\{u(t)\}$ are easier to generate and visualize if the modulating pulse $p(t)$ is also real. From the discussion of the Nyquist criterion, this is not a fundamental limitation, and there are few reasons for desiring a complex $p(t)$. For real $p(t)$,

$$\Re\{u(t)\} = \sum_k \Re\{u_k\}\, p(\frac{t}{T}-k),$$

$$\Im\{u(t)\} = \sum_k \Im\{u_k\}\, p(\frac{t}{T}-k).$$

Letting $u'_k = \Re\{u_k\}$ and $u''_k = \Im\{u_k\}$, the transmitted passband waveform becomes

$$x(t) = 2\cos(2\pi f_c t)\left(\sum_k u'_k p(t-kT)\right) - 2\sin(2\pi f_c t)\left(\sum_k u''_k p(t-kT)\right). \qquad (6.27)$$

If the QAM signal set is a standard QAM set, then $\sum_k u'_k p(t-kT)$ and $\sum_k u''_k p(t-kT)$ are parallel baseband PAM systems. They are modulated to passband using "double-sideband"

modulation by "quadrature carriers" $\cos 2\pi f_c t$ and $-\sin 2\pi f_c t$. These are then summed (with the usual factor of 2), as shown in Figure 6.7. This realization of QAM is called *double-sideband quadrature-carrier* (DSB-QC) modulation[15].
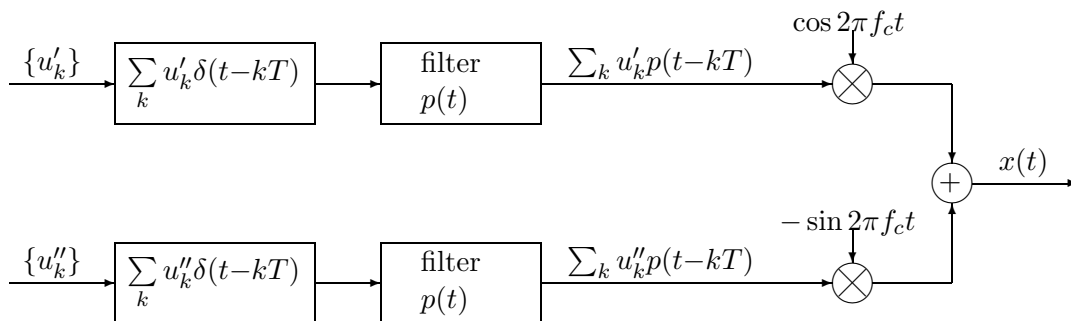


Figure 6.7: DSB-QC modulation

We have seen that $u(t)$ can be recovered from $x(t)$ by a Hilbert filter followed by shifting down in frequency. A more easily implemented but equivalent procedure starts by multiplying $x(t)$ both by $\cos(2\pi f_c t)$ and by $-\sin(2\pi f_c t)$.

Using the trigonometric identities $2\cos^2(\alpha) = 1 + \cos(2\alpha)$, $2\sin(\alpha)\cos(\alpha) = \sin(2\alpha)$, and $2\sin^2(\alpha) = 1 - \cos(2\alpha)$, these terms can be written as

$$x(t)\cos(2\pi f_c t) = \Re\{u(t)\} + \Re\{u(t)\}\cos(4\pi f_c t) + \Im\{u(t)\}\sin(4\pi f_c t), \qquad (6.28)$$

$$-x(t)\sin(2\pi f_c t) = \Im\{u(t)\} - \Re\{u(t)\}\sin(4\pi f_c t) + \Im\{u(t)\}\cos(4\pi f_c t). \qquad (6.29)$$

To interpret this, note that multiplying by $\cos(2\pi f_c t) = \frac{1}{2}e^{2\pi i f_c t} + \frac{1}{2}e^{-2\pi i f_c t}$ both shifts $x(t)$ up[16] and down in frequency by $f_c$. Thus the positive frequency part of $x(t)$ gives rise to a baseband term and a term around $2f_c$, and the negative frequency part gives rise to a baseband term and a term at $-2f_c$. Filtering out the double frequency terms then yields $\Re\{u(t)\}$. The interpretation of the sine multiplication is similar.

As another interpretation, recall that $x(t)$ is real and consists of one band of frquencies around $f_c$ and another around $-f_c$. Note also that (6.28) and (6.29) are the real and imaginary parts of $x(t)e^{-2\pi i f_c t}$, which shifts the positive frequency part of $x(t)$ down to baseband and shifts the negative frequency part down to a band around $-2f_c$. In the Hilbert filter approach, the lower band is filtered out before the frequency shift, and in the approach here, it is filtered out after the frequency shift. Clearly the two are equivalent.

It has been assumed throughout that $f_c$ is greater than the baseband bandwidth of $u(t)$. If this is not true, then, as shown in Exercise 6.11, $u(t)$ can not be retrieved from $x(t)$ by any approach.

Now assume that the baseband modulation filter $p(t)$ is real and a standard QAM signal set is used. Then $\Re\{u(t)\} = \sum u'_k p(t-kT)$ and $\Im\{u(t)\} = \sum u''_k p(t-kT)$ are parallel baseband PAM

---

[15]The terminology comes from analog modulation where two real analog waveforms are modulated respectively onto cosine and sine carriers. For analog modulation, it is customary to transmit an additional component of carrier from which timing and phase can be recovered. As we see shortly, no such additional carrier is necessary here.

[16]This shift up in frequency is a little confusing, since $x(t)e^{-2\pi i f_c t} = x(t)\cos(2\pi f_c t) - ix(t)\sin(2\pi f_c t)$ is only a shift down in frequency. What is happening is that $x(t)\cos(2\pi f_c t)$ is the real part of $x(t)e^{-2\pi i f_c t}$ and thus needs positive frequency terms to balance the negative frequency terms.

modulations. Assume also that a receiver filter $q(t)$ is chosen so that $\hat{g}(f) = \hat{p}(f)\hat{q}(f)$ satisfies the Nyquist criterion and all the filters have the common bandwidth $B_b < f_c$. Then, from (6.28), if $x(t)\cos(2\pi f_c t)$ is filtered by $q(t)$, it can be seen that $q(t)$ will filter out the component around $2f_c$. The output from the remaining component, $\Re\{u(t)\}$ can then be sampled to retrieve the real signal sequence $u'_1, u'_2, \ldots$. This plus the corresponding analysis of $-x(t)\sin(2\pi f_c t)$ is illustrated in the DSB-QC receiver in Figure 6.8. Note that the use of the filter $q(t)$ eliminates the need for either filtering out the double frequency terms or using a Hilbert filter.
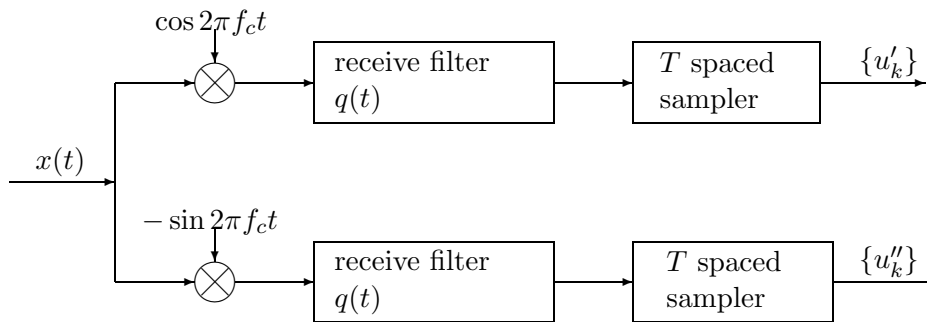


Figure 6.8: DSB-QC demodulation

The above description of demodulation ignores the noise. As explained in Section 6.3.2, however, if $p(t)$ is chosen so that $\{p(t-kT); k \in \mathbb{Z}\}$ is an orthonormal set (*i.e.*, so that $|\hat{p}(f)|^2$ satisfies the Nyquist criterion), then the receiver filter should satisfy $q(t) = p(-t)$. It will be shown later that in the presence of white Gaussian noise, this is the optimal thing to do (in a sense to be described later).

## 6.6   Signal space and degrees of freedom

Using PAM, real signals can be generated at $T$-spaced intervals and transmitted in a baseband bandwidth arbitrarily little more than $\mathsf{W}_b = \frac{1}{2T}$. Thus, over an asymptotically long interval $T_0$, and in a baseband bandwidth asymptotically close to $\mathsf{W}_b$, $2\mathsf{W}_b T_0$ real signals can be transmitted using PAM.

Using QAM, complex signals can be generated at $T$-spaced intervals and transmitted in a passband bandwidth arbitrarily little more than $\mathsf{W} = \frac{1}{T}$. Thus, over an asymptotically long interval $T_0$, and in a passband bandwidth asymptotically close to $\mathsf{W}$, $\mathsf{W}T_0$ complex signals, and thus $2\mathsf{W}T_0$ real signals can be transmitted using QAM.

The above description described PAM at baseband and QAM at passband. To get a better comparison of the two, consider an overall large baseband bandwidth $\mathsf{W}_0$ broken into $m$ passbands each of bandwidth $\mathsf{W}_0/m$. Using QAM in each band, we can asymptotically transmit $2\mathsf{W}_0 T_0$ real signals in a long interval $T_0$. With PAM used over the entire band $\mathsf{W}_0$, we again asyptotically send $2\mathsf{W}_0 T_0$ real signals in a duration $T_0$. We see that in principle, QAM and baseband PAM with the same overall bandwidth are equivalent in terms of the number of degrees of freedom that can be used to transmit real signals. As pointed out earlier, however, PAM when modulated up to passband uses only half the available degrees of freedom. Also, QAM offers considerably

more flexibility since it can be used over an arbitrary selection of frequency bands.

Recall that when we were looking at $T$-spaced truncated sinusoids and $T$-spaced sinc weighted sinusoids, we argued that the class of real waveforms occupying a time interval $(-T_0/2, T_0/2)$ and a frequency interval $(-W_0, W_0)$ has about $2T_0W_0$ degrees of freedom for large $W_0, T_0$. What we see now is that baseband PAM and passband QAM each employ about $2T_0W_0$ degrees of freedom. In other words, these simple techniques essentially use all the degrees of freedom available in the given bands.

The use of Nyquist theory here has added to our understanding of waveforms that are "essentially" time and frequency limited. That is, we can start with a family of functions that are bandlimited within a rolloff factor and then look at asymptotically small rolloffs. The discussion of noise in the next two chapters will provide a still better understanding of degrees of freedom subject to essential time and frequency limits.

### 6.6.1   Distance and orthogonality

Previous sections have shown how to modulate a complex QAM baseband waveform $u(t)$ up to a real passband waveform $x(t)$ and how to retrieve $u(t)$ from $x(t)$ at the receiver. They have also discussed signal constellations that minimize energy for given minimum distance. Finally, the use of a modulation waveform $p(t)$ with orthonormal shifts, has connected the energy difference between two baseband signal waveforms, say $u(t) = \sum u_k p(t - kT)$ and $v(t) = \sum_k v_k p(t - kt)$ and the energy difference in the signal points by

$$\|\boldsymbol{u} - \boldsymbol{v}\|^2 = \sum_k |u_k - v_k|^2.$$

Now consider this energy difference at passband. The energy $\|\boldsymbol{x}\|^2$ in the passband waveform $x(t)$ is twice that in the corresponding baseband waveform $u(t)$. Next suppose that $x(t)$ and $y(t)$ are the passband waveforms arising from the baseband waveforms $u(t)$ and $v(t)$ respectively. Then

$$x(t) - y(t) = 2\Re\{u(t)e^{2\pi i f_c t}\} - 2\Re\{v(t)e^{2\pi i f_c t}\} = 2\Re\{[u(t)-v(t)]e^{2\pi i f_c t}\}.$$

Thus $x(t) - y(t)$ is the passband waveform corresponding to $u(t) - v(t)$, so

$$\|x(t) - y(t)\|^2 = 2\|u(t) - v(t)\|^2 .$$

This says that for QAM and PAM, distances between waveforms are preserved (aside from the scale factor of 2 in energy or $\sqrt{2}$ in distance) in going from baseband to passband. Thus distances are preserved in going from signals to baseband waveforms to passband waveforms and back. We will see later that the error probability caused by noise is essentially determined by the distances between the set of passband source waveforms. This error probability is then simply related to the choice of signal constellation and the discrete coding that precedes the mapping of data into signals.

This preservation of distance through the modulation to passband and back is a crucial aspect of the signal space viewpoint of digital communication. It provides a practical focus to viewing waveforms at baseband and passband as elements of related $\mathcal{L}_2$ inner product spaces.

There is unfortunately a mathematical problem in this very nice story. The set of baseband waveforms forms a complex inner product space whereas the set of passband waveforms constitutes a real inner product space. The transformation $x(t) = \Re\{u(t)e^{2\pi i f_c t}\}$ is not linear, since,

for example, $iu(t)$ does not map into $ix(t)$ for $u(t) \neq 0$). In fact, the notion of a linear transformation does not make much sense, since the transformation goes from complex $\mathcal{L}_2$ to real $\mathcal{L}_2$ and the scalars are different in the two spaces.

**Example 6.6.1.** As an important example, suppose the QAM modulation pulse is a real waveform $p(t)$ with orthonormal $T$-spaced shifts. The set of complex baseband waveforms spanned by the orthonormal set $\{p(t-kT); k \in \mathbb{Z}\}$ has the form $\sum_k u_k p(t-kT)$ where each $u_k$ is complex. As in (6.27), this is transformed at passband to

$$\sum_k u_k p(t-kT) \rightarrow \sum_k 2\Re\{u_k\}p(t-kT)\cos(2\pi ft) - 2\sum_k \Im\{u_k\}p(t-kT)\sin(2\pi ft).$$

Each baseband function $p(t-kT)$ is modulated to the passband waveform is $2p(t-kT)\cos(2\pi f_c t)$. The set of functions $\{p(t-kT)\cos(2\pi f_c t); k \in \mathbb{Z}\}$ is not enough to span the space of modulated waveforms, however. It is necessary to add the additional set $\{p(t-kT)\sin(2\pi f_c t); k \in \mathbb{Z}\}$. As shown in Exercise 6.15, This combined set of waveforms is an orthogonal set, each of energy 2.

Another way to look at this example is to observe that modulating the baseband function $u(t)$ into the positive passband function $x^+(t) = u(t)e^{2\pi i f_c t}$ is somewhat easier to understand in that the orthonormal set $\{p(t-kT); k \in \mathbb{Z}\}$ is modulated to the orthonormal set $\{p(t-kT)e^{2\pi i f_c t}; k \in \mathbb{Z}\}$, which can be seen to span the space of complex positive frequency passband source waveforms. The additional set of orthonormal waveforms $\{p(t-kT)e^{-2\pi i f_c t}; k \in \mathbb{Z}\}$ is then needed to span the real passband source waveforms. We then see that the sine, cosine series is simply another way to express this. In the sine, cosine formulation all the coefficients in the series are real, whereas in the complex exponential formulation, there is a real and complex coefficient for each term, but they are pairwise dependent. It will be easier to understand the effects of noise in the sine, cosine formulation.

In the above example, we have seen that each orthonormal function at baseband gives rise to two real orthonormal functions at passband. It can be seen from a degrees of freedom argument that this is inevitable no matter what set of orthonormal functions are used at baseband. For a nominal passband bandwidth W, there are 2W real degrees of freedom per second in the baseband complex source waveform, which means there 2 real degrees of freedom for each orthonormal baseband waveform. At passband, we have the same 2W degrees of freedom per second, but with a real orthonormal expansion, there is only one real degree of freedom for each orthonormal waveform. Thus there must be two passband real orthonormal waveforms for each baseband complex orthonormal waveform.

The sine, cosine expansion above generalizes in a nice way to an arbitrary set of complex orthonormal baseband functions. Each complex function in this baseband set generates two real functions in an orthogonal passband set. This is expressed precisely in the following theorem which is proven in Exercise 6.16.

**Theorem 6.6.1.** *Let* $\{\theta_k(t) : k \in \mathbb{Z}\}$ *be an orthonormal set limited to the frequency band* $[-B/2, B/2]$. *Let* $f_c$ *be greater than* $B/2$ *and for each* $k \in \mathbb{Z}$ *let*

$$\psi_{k,1}(t) = \Re\left\{2\theta_k(t)\, e^{2\pi i f_c t}\right\},$$
$$\psi_{k,2}(t) = \Im\left\{-2\theta_k(t)\, e^{2\pi i f_c t}\right\}.$$

*The set* $\{\psi_{k,j}; k \in \mathbb{Z}, j \in \{1,2\}\}$ *is an orthogonal set of functions, each of energy 2. Furthermore, if* $u(t) = \sum_k u_k \theta_k(t)$, *then the corresponding passband function* $x(t) = 2\Re\{u(t)e^{2\pi i f_c t}\}$ *is given*

by

$$x(t) = \sum_k \Re\{u_k\}\, \psi_{k,1}(t) + \Im\{u_k\}\, \psi_{k,2}(t).$$

This gives us a very general way to map any orthonormal set at baseband into a related orthonormal set at passband, with two real orthonormal functions at passband corresponding to each orthonormal function at baseband. It is not limited to any particular type of modulation, and thus will allow us to make general statements about signal space at baseband and passband.

## 6.7   Carrier and phase recovery in QAM systems

Consider a QAM receiver and visualize the passband-to-baseband conversion as multiplying the positive frequency passband by the complex sinusoid $e^{-2\pi i f_c t}$. If the receiver has a phase error $\phi(t)$ in its estimate of the phase of the transmitted carrier, then it will instead multiply the incoming waveform by $e^{-2\pi i f_c t + i\phi(t)}$. We assume in this analysis that the time reference at the receiver is perfectly known, so that the sampling of the filtered output is done at the correct time. Thus the assumption is that the oscillator at the receiver is not quite in phase with the oscillator at the transmitter. Note that the carrier frequency is usually orders of magnitude higher than the baseband bandwidth, and thus a small error in timing is significant in terms of carrier phase but not in terms of sampling. The carrier phase error will rotate the correct complex baseband signal $u(t)$ by $\phi(t)$; *i.e.*, the actual received baseband signal $r(t)$ will be

$$r(t) = e^{i\phi(t)} u(t).$$

If $\phi(t)$ is slowly time-varying relative to the response $q(t)$ of the receiver filter, then the samples $\{r(kT)\}$ of the filter output will be

$$r(kT) \approx e^{i\phi(kT)} u_k,$$

as illustrated in Figure 6.9. The phase error $\phi(t)$ is said to come through *coherently*. This phase coherence makes carrier recovery easy in QAM systems.
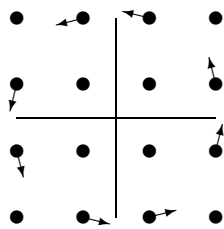


Figure 6.9: Rotation of constellation points by phase error

As can be seen from the figure, if the phase error is small enough, and the set of points in the constellation are well enough separated, then the phase error can be simply corrected by moving to the closest signal point and adjusting the phase of the demodulating carrier accordingly.

There are two complicating factors here. The first is that we have not taken noise into account yet. When the received signal $y(t)$ is $x(t) + n(t)$, then the output of the $T$ spaced sampler is not

the original signals $\{u_k\}$, but rather a noise corrupted version of them. The second problem is that if a large phase error ever occurs, it can not be corrected. For example, in Figure 6.9, if $\phi(t) = \pi/2$, then even in the absence of noise, the received samples always line up with signals from the constellation (but of course not the transmitted signals).

### 6.7.1  Tracking phase in the presence of noise

The problem of *deciding on* or *detecting* the signals $\{u_k\}$ from the received samples $\{r(kT)\}$ in the presence of noise is a major topic of Chapter 8. Here, however, we have the added complication of both detecting the transmitted signals and also tracking and eliminating the phase error.

Fortunately, the problem of decision making and that of phase tracking are largely separable. The oscillators used to generate the modulating and demodulating carriers are relatively stable and have phases which change quite slowly relative to each other. Thus the phase error with any kind of reasonable tracking will be quite small, and thus the data signals can be detected from the received samples almost as if the phase error were zero. The difference between the received sample and the detected data signal will still be nonzero, mostly due to noise but partly due to phase error. However, the noise has zero mean (as we understand later) and thus tends to average out over many sample times. Thus the general approach is to make decisions on the data signals as if the phase error is zero, and then to make slow changes to the phase based on averaging over many sample times. This approach is called *decision directed carrier recovery*. Note that if we track the phase as phase errors occur, we are also tracking the carrier, in both frequency and phase.

In a decision directed scheme, assume that the received sample $r(kT)$ is used to make a decision $d_k$ on the transmitted signal point $u_k$. Also assume that $d_k = u_k$ with very high probability. The apparent phase error for the $k$th sample is then the difference between the phase of $r(kT)$ and the phase of $d_k$. Any method for feeding back the apparent phase error to the generator of the sinusoid $e^{-2\pi i f_c t + i\phi(t)}$ in such a way as to slowly reduce the apparent phase error will tend to produce a robust carrier recovery system.

In one popular method, the feedback signal is taken as the imaginary part of $r(kT)d_k^*$. If the phase angle from $d_k$ to $r(kT)$ is $\phi_k$, then

$$r(kT)d_k^* = |r(kT)||d_k| \ e^{i\phi_k},$$

so the imaginary part is $|r(kT)||d_k| \sin \phi_k \approx |r(kT)||d_k|\phi_k$, when $\phi_k$ is small. Decision-directed carrier recovery based on such a feedback signal can be extremely robust even in the presence of substantial distortion and large initial phase errors. With a second-order phase-locked carrier recovery loop, it turns out that the carrier frequency $f_c$ can be recovered as well.

### 6.7.2  Large phase errors

A problem with decision-directed carrier recovery and with many other approaches is that the recovered phase may settle into any value for which the received eye pattern (*i.e.*, the pattern of a long string of received samples as viewed on a scope) "looks OK." With $(M \times M)$-QAM signal sets, as in Figure 6.9, the signal set has four-fold symmetry, and phase errors of $90°, 180°,$ or $270°$ are not detectable. Simple differential coding methods that transmit the "phase" (quadrantal)

part of the signal information as a change of phase from the previous signal rather than as an absolute phase can easily overcome this problem. Another approach is to resynchronize the system frequently by sending some known pattern of signals. This latter approach is frequently used in wireless systems where fading sometimes causes a loss of phase synchronization.

## 6.8    Summary of modulation and demodulation

This chapter has used the signal space developed in Chapters 4 and 5 to study the mapping of binary input sequences at a modulator into the waveforms to be transmitted over the channel. Figure 6.1 summarized this process, mapping bits to signals, then signals to baseband waveforms, and then baseband waveforms to passband waveforms. The demodulator goes through the inverse process, going from passband waveforms to baseband waveforms to signals to bits. This breaks the modulation process into three layers that can be studied more or less independently.

The development used PAM and QAM throughout, both as widely used systems, and as convenient ways to bring out the principles that can be applied more widely.

The mapping from binary digits to signals segments the incoming binary sequence into $b$-tuples of bits and then maps the set of $M = 2^b$ $n$-tuples into a constellation of $M$ signal points in $\mathbb{R}^m$ or $C^m$ for some convenient $m$. Since the $m$ components of these signal points are going to be used as coefficients in an orthogonal expansion to generate the waveforms, the objectives are to choose a signal constellation with small average energy but with a large distance between each pair of points. PAM is an example where the signal space is $\mathbb{R}^1$ and QAM is an example where the signal space is $\mathbb{C}^1$. For both of these, the standard mapping is the same as the representation points of a uniform quantizer. These are not quite optimal in terms of minimizing the average energy for a given minimum point spacing, but they are almost universally used because of the near-optimality and the simplicity.

The mapping of signals into baseband waveforms for PAM chooses a fixed waveform, $p(t)$ and modulates the sequence of signals $u_1, u_2, \ldots$ into the baseband waveform $\sum_j u_j p(t - jT)$. One of the objectives in choosing $p(t)$ is to be able to retrieve the sequence $u_1, u_2, \ldots$, from the received waveform. This involves an output filter $q(t)$ which is sampled each $T$ seconds to retrieve $u_1, u_2, \ldots$. The Nyquist criterion was derived, specifying the properties that the product $\hat{g}(f) = \hat{p}(f)\hat{q}(f)$ must satisfy to avoid intersymbol interference. The objective in choosing $\hat{g}(f)$ is a trade off between the closeness of $\hat{g}(f)$ to $T \operatorname{rect}(fT)$ and the time duration of $g(t)$, subject to satisfying the Nyquist criterion. The raised cosine functions are widely used as a good compromise between these dual objectives. For a given real $\hat{g}(f)$, the choice of $\hat{p}(f)$ usually satisfies $\hat{g}(f) = |\hat{p}(f)|^2$, and in this case $\{p(t - kT); k \in \mathbb{Z}\}$ is a set of orthonormal functions.

Most of the remainder of the chapter discussed modulation from baseband to passband. This was primarily a topic in the manipulation of Fourier transforms, and need not be summarized here.

## 6.E    Exercises

6.1. (PAM) Consider standard $M$-PAM and assume that the signals are used with equal probability. Show that the average energy per signal $E_s = \overline{U_k^2}$ is equal to the average energy $\overline{U^2} = d^2 M^2/12$ of a uniform continuous distribution over the interval $[-dM/2, dM/2]$, minus the average energy $\overline{(U - U_k)^2} = d^2/12$ of a uniform continuous distribution over the interval $[-d/2, d/2]$:

$$E_s = \frac{d^2(M^2 - 1)}{12}.$$

This establishes (6.4). Verify the formula for $M = 4$ and $M = 8$.

6.2. (PAM) A discrete memoryless source emits binary equiprobable symbols at a rate of 1000 symbols per second. The symbols from a one second interval are grouped into pairs and sent over a bandlimited channel using a standard 4-PAM signal set. The modulation uses a signal interval 0.002 and pulse $p(t) = \text{sinc}(t/T)$.

(a) Suppose that a sample sequence $u_1, \ldots, u_{500}$ of transmitted signals includes 115 appearances of $3d/2$, 130 appearances of $d/2$, 120 appearances of $-d/2$, and 135 appearances of $-3d/2$. Find the energy in the corresponding transmitted waveform $u(t) = \sum_{k=1}^{500} u_k \, \text{sinc}(\frac{t}{T} - k)$ as a function of $d$.

(b) What is the bandwidth of the waveform $u(t)$ in part (a)?

(c) Find $\mathsf{E}\left[\int U^2(t)\, dt\right]$ where $U(t)$ is the random waveform $\sum_{k=1}^{500} U_k \, \text{sinc}(\frac{t}{T} - k)$.

(d) Now suppose that the binary source is not memoryless, but is instead generated by a Markov chain where

$$\Pr(X_i{=}1 \mid X_{i-1}{=}1) = \Pr(X_i{=}0 \mid X_{i-1}{=}0) = 0.9.$$

Assume the Markov chain starts in steady state with $\Pr(X_1{=}1) = 1/2$. Using the mapping $(00 \to a_1), (01 \to a_2), (10 \to a_3), (11 \to a_4)$, find $\mathsf{E}[U_k^2]$ for $1 \le k \le 500$.

(e) Find $\mathsf{E}\left[\int U^2(t)\, dt\right]$ for this new source.

(f) For the above Markov chain, explain how we could change the above mapping to reduce the expected energy without changing the separation between signal points.

6.3. (a) Assume that the received signal in a 4-PAM system is $V_k = U_k + Z_k$ where $U_k$ is the transmitted 4-PAM signal at time $k$. Let $Z_k$ be independent of $U_k$ and Gaussian with density $f_Z(z) = \sqrt{\frac{1}{2\pi}} \exp\left\{-\frac{z^2}{2}\right\}$. Assume that the receiver chooses the signal $\tilde{U}_k$ closest to $V_k$. (It is shown in Chapter 8 that this detection rule minimizes $P_e$ for equiprobable signals.) Find the probability $P_e$ (in terms of Gaussian integrals) that $U_k \ne \tilde{U}_k$.

(b) Evaluate the partial derivitive of $P_e$ with respect to the third signal point $a_3$ (i.e., the positive inner signal point) at the point where $a_3$ is equal to its value $d/2$ in standard 4-PAM and all other signal points are kept at their 4-PAM values. Hint: This doesn't require any calculation.

(c) Evaluate the partial derivitive of the signal energy $E_s$ with respect to $a_3$.

(d) Argue from this that the minimum error probability signal constellation for 4 equiprobable signal points is not 4-PAM, but rather a constellation where the distance between the inner points is smaller than the distance from inner point to outer point on either side. (This is quite surprising intuitively to the author.)

6.4. (Nyquist) Suppose that the PAM modulated baseband waveform $u(t) = \sum_{k=-\infty}^{\infty} u_k p(t-kT)$ is received. That is, $u(t)$ is known, $T$ is known, and $p(t)$ is known. We want to determine the signals $\{u_k\}$ from $u(t)$. We assume we must use only linear operations. That is, we wish to find some waveform $d_k(t)$ for each integer $k$ such that $\int_{-\infty}^{\infty} u(t)d_k(t)\, dt = u_k$.

(a) What properites must be satisfied by $d_k(t)$ such that the above equation is satisfied no matter what values are taken by the other signals, $\dots, u_{k-2}, u_{k-1}, u_{k+1}, u_{k+2}, \dots$? These properties should take the form of constraints on the inner products $\langle p(t-kT), d_j(t)\rangle$. Do not worry about convergence, interchange of limits, etc.

(b) Suppose you find a function $d_0(t)$ that satisfies these constraints for $k=0$. Show that for each $k$, a function $d_k(t)$ satisfying these constraints can be found simply in terms of $d_0(t)$.

(c) What is the relationship between $d_0(t)$ and a function $q(t)$ that avoids intersymbol interference in the approach taken in Section 6.3 (*i.e.*, a function $q(t)$ such that $p(t) * q(t)$ is ideal Nyquist).

You have shown that the filter/sample approach in Section 6.3 is no less general than the arbitrary linear operation approach here. Note that, in the absence of noise and with a known signal constellation, it might be possible to retrieve the signals from the waveform using nonlinear operations even in the presence of intersymbol interference.

6.5. (Nyquist) Let $v(t)$ be a continuous $\mathcal{L}_2$ waveform with $v(0) = 1$ and define $g(t) = v(t)\operatorname{sinc}(\frac{t}{T})$.

(a) Show that $g(t)$ is ideal Nyquist with interval $T$.

(b) Find $\hat{g}(f)$ as a function of $\hat{v}(f)$.

(c) Give a direct demonstration that $\hat{g}(f)$ satisfies the Nyquist criterion.

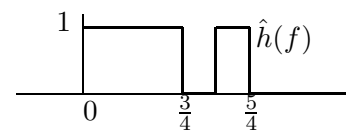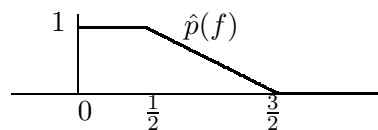(d) If $v(t)$ is baseband limited to $B_b$, what is $g(t)$ baseband limited to?

Note: The usual form of the Nyquist criterion helps in choosing waveforms that avoid intersymbol interference with prescribed rolloff properties in frequency. The approach above show how to avoid intersymbol interference with prescribed attenuation in time and in frequency.

6.6. (Nyquist) Consider a PAM baseband system in which the modulator is defined by a signal interval $T$ and a wveform $p(t)$, the channel is defined by a filter $h(t)$, and the receiver is defined by a filter $q(t)$ which is sampled at $T$-spaced intervals. The received waveform, after the receive filter $q(t)$, is then given by $r(t) = \sum_k u_k g(t-kT)$ where $g(t) = p(t) * h(t) * q(t)$.

(a) What property must $g(t)$ have so that $r(kT) = u_k$ for all $k$ and for all choices of input $\{u_k\}$? What is the Nyquist criterion for $\hat{g}(f)$?

(b) Now assume that $T = 1/2$ and that $p(t), h(t), q(t)$ and all their Fourier transforms are restricted to be real. Assume further that $\hat{p}(f)$ and $\hat{h}(f)$ are given by

$$\hat{p}(f) = \begin{cases} 1, & |f| \le 0.5; \\ 1.5 - t, & 0.5 < |f| \le 1.5 \\ 0, & |f| > 1.5 \end{cases} \qquad \hat{h}(f) = \begin{cases} 1, & |f| \le 0.75; \\ 0, & 0.75 < |f| \le 1 \\ 1, & 1 < |f| \le 1.25 \\ 0, & |f| > 1.25 \end{cases}$$

Is it possible to choose a receive filter transform $\hat{q}(f)$ so that there is no intersymbol interference? If so, give such a $\hat{q}(f)$ and indicate the regions in which your solution is nonunique.

(c) Redo part (b) with the modification that now $\hat{h}(f) = 1$ for $|f| \leq 0.75$ and $\hat{h}(f) = 0$ for $|f| > 0.75$.

(d) Explain the conditions on $\hat{p}(f)\hat{h}(f)$ under which intersymbol interference can be avoided by proper choice of $\hat{q}(f)$ (you may assume, as above, that $\hat{p}(f), \hat{h}(f), p(t),$ and $h(t)$ are all real).

6.7. (Nyquist) Recall that the $\text{rect}(t/T)$ function has the very special property that it, plus its time and frequency shifts by $kT$ and $j/T$ respectively, form an orthogonal set of functions. The function $\text{sinc}(t/T)$ has this same property. This problem is about some other functions that are generalizations of $\text{rect}(t/T)$ and which, as you will show in parts (a) to (d), have this same interesting property. For simplicity, choose $T$ to be 1.

These functions take only the values 0 and 1 and are allowed to be nonzero only over [-1, 1] rather than $[-1/2, 1/2]$ as with $\text{rect}(t)$. Explicitly, the functions considered here satisfy the following constraints:
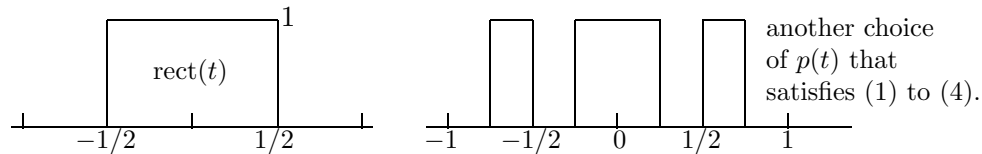
$$
\begin{aligned}
p(t) &= p^2(t) && \text{for all } t \quad \text{(0/1 property)} && (6.30) \\
p(t) &= 0 && \text{for } |t| > 1 && (6.31) \\
p(t) &= p(-t) && \text{for all } t \quad \text{(symmetry)} && (6.32) \\
p(t) &= 1 - p(t-1) && \text{for } 0 \leq t < 1/2. && (6.33)
\end{aligned}
$$

Note: Because of property (6.32), condition (6.33) also holds for $1/2 < t \leq 1$. Note also that $p(t)$ at the single points $t = \pm 1/2$ does not effect any orthogonality properties, so you are free to ignore these points in your arguments.



(a) Show that $p(t)$ is orthogonal to $p(t-1)$. Hint: evaluate $p(t)p(t-1)$ for each $t \in [0, 1]$ other than $t = 1/2$.

(b) Show that $p(t)$ is orthogonal to $p(t-k)$ for all integer $k \neq 0$.

(c) Show that $p(t)$ is orthogonal to $p(t-k)e^{2\pi imt}$ for integer $m \neq 0$ and $k \neq 0$.

(d) Show that $p(t)$ is orthogonal to $p(t)e^{2\pi imt}$ for integer $m \neq 0$. Hint: Evaluate $p(t)e^{-2\pi imt} + p(t-1)e^{-2\pi im(t-1)}$.

(e) Let $h(t) = \hat{p}(t)$ where $\hat{p}(f)$ is the Fourier transform of $p(t)$. If $p(t)$ satisfies properties (1) to (4), does it follow that $h(t)$ has the property that it is orthogonal to $h(t-k)e^{2\pi imt}$ whenever either the integer $k$ or $m$ is nonzero?
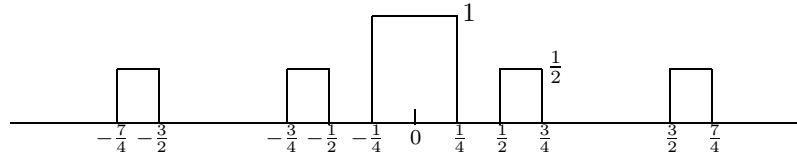
Note: Almost no calculation is required in this exercise.

6.8. (Nyquist) (a) For the special case $\alpha = 1, T = 1$, verify the formula in (6.18) for $\hat{g}_1(f)$ given $g_1(t)$ in (6.17). Hint: As an intermediate step, verify that $g_1(t) = \text{sinc}(2t) + \frac{1}{2}\text{sinc}(2t+1) + \frac{1}{2}\text{sinc}(2t-1)$. Sketch $g_1(t)$, in particular showing its value at $mT/2$ for each $m \geq 0$.

(b) For the general case $0 < \alpha < 1$, $T = 1$, show that $\hat{g}_\alpha(f)$ is the convolution of rect $f$ with a half cycle of $\beta \cos \pi \alpha f$ and specify the required value of $\beta$.

(c) Verify (6.18) for $0 < \alpha < 1$, $T = 1$ and then verify for arbitrary $T > 0$.

6.9. (Approximate Nyquist)This exercise shows that approximations to the Nyquist criterion must be treated with great care. Define $\hat{g}_k(f)$, for integer $k \geq 0$ as in the diagram below for $k = 2$. For arbitrary $k$, there are $k$ small pulses on each side of the main pulse, each of height $\frac{1}{k}$.



(a) Show that $\hat{g}_k(f)$ satisfies the Nyquist criterion for $T = 1$ and for each $k \geq 1$.

(b) Show that l.i.m.$_{k \to \infty}$ $\hat{g}_k(f)$ is simply the central pulse above. That is, this $\mathcal{L}_2$ limit satisfies the Nyquist criterion for $T = \frac{1}{2}$. To put it another way, $\hat{g}_k(f)$, for large $k$, satisfies the Nyquist criterion for $T = 1$ using 'approximately' the bandwidth $\frac{1}{4}$ rather than the necessary bandwidth $\frac{1}{2}$. The problem is that the $\mathcal{L}_2$ notion of approximation (done carefully here as a limit in the mean of a sequence of approximations) is not always appropriate, and it is often inappropriate with sampling issues.

6.10. (Nyquist) (a) Assume that $\hat{p}(f) = \hat{q}^*(f)$ and $\hat{g}(f) = \hat{p}(f)\hat{q}(f)$. Show that if $p(t)$ is real, then $\hat{g}(f) = \hat{g}(-f)$ for all $f$.

(b) Under the same assumptions, find an example where $p(t)$ is not real but $\hat{g}(f) \neq \hat{g}(-f)$ and $\hat{g}(f)$ satisifes the Nyquist criterion. Hint: Show that $\hat{g}(f) = 1$ for $0 \leq f \leq 1$ and $\hat{g}(f) = 0$ elsewhere satisfies the Nyquist criterion for $T = 1$ and find the corresponding $p(t)$.

6.11. (Passband) (a) Let $u_k(t) = \exp(2\pi i f_k t)$ for $k = 1, 2$ and let $x_k(t) = 2\Re\{u_k(t)\exp(2\pi i f_c t)\}$. Assume $f_1 > -f_c$ and find the $f_2 \neq f_1$ such that $x_1(t) = x_2(t)$.

(b) Explain that what you have done is to show that, without the assumption that the bandwidth of $u(t)$ is less than $f_c$, it is impossible to always retrieve $u(t)$ from $x(t)$, even in the absence of noise.

(c) Let $y(t)$ be a real $\mathcal{L}_2$ function. Show that the result in part (a) remains valid if $u_k(t) = y(t)\exp(2\pi i f_k t)$ (i.e., show that the result in part (a) is valid with a restriction to $\mathcal{L}_2$ functions.

(d) Show that if $u(t)$ is restricted to be real, then $u(t)$ can be retrieved almost everywhere from $x(t) = 2\Re\{u(t)\exp(2\pi i f_c t)\}$. Hint: express $x(t)$ in terms of $\cos(2\pi f_c t)$.

(e) Show that if the bandwidth of $u(t)$ exceeds $f_c$, then neither Figure 6.6 nor Figure 6.8 work correctly, even when $u(t)$ is real.

6.12. (QAM) (a) Let $\theta_1(t)$ and $\theta_2(t)$ be orthonormal complex waveforms. Let $\phi_j(t) = \theta_j(t)e^{2\pi i f_c t}$ for $j = 1, 2$. Show that $\phi_1(t)$ and $\phi_2(t)$ are orthonormal for any $f_c$.

(b) Suppose that $\theta_2(t) = \theta_1(t - T)$. Show that $\phi_2(t) = \phi_1(t - T)$ if $f_c$ is an integer multiple of $1/T$.

6.13. (QAM) (a) Assume $B/2 < f_c$. Let $u(t)$ be a real function and let $v(t)$ be an imaginary function, both baseband limited to $B/2$. Show that the corresponding passband functions, $\Re\{u(t)e^{2\pi i f_c t}\}$ and $\Re\{v(t)e^{2\pi i f_c t}\}$ are orthogonal.

(b) Give an example where the functions in part (a) are not orthogonal if $B/2 > f_c$.

6.14. (a) Derive (6.28) and (6.29) using trigonometric identities.

(b) View the left side of (6.28) and (6.29) as the real and imaginary part respectively of $x(t)e^{-2\pi i f_c t}$. Rederive (6.28) and (6.29) using complex exponentials. (Note how much easier this is than part (a).

6.15. (Passband expansions) Assume that $\{p(t-kT) : k\in\mathbb{Z}\}$ is a set of orthonormal functions. Assume that $\hat{p}(f) = 0$ for $|f| \geq f_c)$.

(a) Show that $\{\sqrt{2}p(t-kT)\cos(2\pi f_c t); k\in\mathbb{Z}\}$ is an orthonormal set.

(b) Show that $\{\sqrt{2}p(t-kT)\sin(2\pi f_c t); k\in\mathbb{Z}\}$ is an orthonormal set and that each function in it is orthonormal to the cosine set in part (a).

6.16. (Passband expansions) Prove Theorem 6.6.1. Hint: First show that the set of functions $\{\hat{\psi}_{k,1}(f)\}$ and $\{\hat{\psi}_{k,2}(f)\}$ are orthogonal with energy 2 by comparing the integral over negative frequencies with that over positive frequencies. Indicate explicitly why you need $f_c > B/2$.

6.17. (Phase and envelope modulation) This exercise shows that any real passband waveform can be viewed as a combination of phase and amplitude modulation. Let $x(t)$ be an $\mathcal{L}_2$ real passband waveform of bandwidth $B$ around a carrier frequency $f_c > B/2$. Let $x^+(t)$ be the positive frequency part of $x(t)$ and let $u(t) = x^+(t)\exp\{-2\pi i f_c t\}$.

(a) Express $x(t)$ in terms of $\Re\{u(t)\}, \Im\{u(t)\}, \cos[2\pi f_c t]$, and $\sin[2\pi f_c t]$.

(b) Define $\phi(t)$ implicitly by $e^{i\phi(t)} = \frac{u(t)}{|u(t)|}$. Show that $x(t)$ can be expressed as $x(t) = 2|u(t)|\cos[2\pi f_c t + \phi(t)]$. Draw a sketch illustrating that $2|u(t)|$ is a baseband waveform upper-bounding $x(t)$ and touching $x(t)$ roughly once per cycle. Either by sketch or words, illustrate that $\phi(t)$ is a phase modulation on the carrier.

(c) Define the *envelope* of a passband waveform $x(t)$ as twice the magnitude of its positive frequency part, *i.e.*, as $2|x^+(t)|$. Without changing the waveform $x(t)$ (or $x^+(t)$) from that before, change the carrier frequency from $f_c$ to some other frequency $f_c'$. Thus $u'(t) = x^+(t)\exp\{-2\pi i f_c't\}$. Show that $|x^+(t)| = |u(t)| = |u'(t)|$. Note that you have shown that the envelope does not depend on the assumed carrier frequency, but has the interpretation of part (b).

(d) Show the relationship of the phase $\phi'(t)$ for the carrier $f_c'$ to that for the carrier $f_c$.

(e) Let $p(t) = |x(t)|^2$ be the power in $x(t)$. Show that if $p(t)$ is lowpass filtered to bandwidth $B$, the result is $2|u(t)|^2$. Interpret this filtering as a short-term time average over $|x(t)|^2$ to interpret why the envelope squared is twice the short-term average power (and thus why the envelope is $\sqrt{2}$ times the short-term root-mean-squared amplitude).

6.18. (Carrierless amplitude-phase modulation (CAP)) We have seen how to modulate a baseband QAM waveform up to passband and then demodulate it by shifting down to baseband, followed by filtering and sampling. This exercise explores the interesting concept of eliminating the baseband operations by modulating and demodulating directly at passband. This approach is used in one of the North American standards for Asymmetrical Digital Subscriber Loop (ADSL)

(a) Let $\{u_k\}$ be a complex data sequence and let $u(t) = \sum_k u_k\, p(t-kT)$ be the corresponding modulated output. Let $\hat{p}(f)$ be equal to $\sqrt{T}$ over $f \in [3/(2T), 5/(2T)]$ and be equal to 0 elsewhere. At the receiver, $u(t)$ is filtered using $p(t)$ and the output $y(t)$ is then T-space sampled at time instants $kT$. Show that $y(kT) = u_k$ for all $k \in \mathbb{Z}$. Don't worry about the fact that the transmitted waveform $u(t)$ is complex.

(b) Now suppose that $\hat{p}(f) = \sqrt{T}\,\text{rect}(T(f - f_c)]$ for some arbitrary $f_c$ rather than $f_c = 2/T$ as in part (a). For what values of $f_c$ does the scheme still work?

(c) Suppose that $\Re\{u(t)\}$ is now sent over a communication channel. Suppose that the received waveform is filtered by a Hilbert filter before going through the demodulation procedure above. Does the scheme still work?