

10 Discrete Input Gaussian Channel

In the previous chapter it was seen that to reach capacity for a Gaussian channel the input variable should be Gaussian distributed. Normally, it is not possible to use a continuous distribution for the transmitted signals and instead discrete values are used. That is, the transmitted variable is discrete while the noise on the channel is Gaussian. Furthermore, in most applications the transmitted signal alternatives are considered to be equally likely. In this section a constraint capacity, in form of the mutual information, for an M-PAM modulated signal will be derived in the case when the (finite number) signal alternatives are equally likely. The loss made by using uniformly distributed inputs will be derived and addressed as the shaping gain. Finally, a parameter called the SNR gap is derived to show how far from the capacity an uncoded system is working. This latter value is derived for a certain obtained bit error rate.

10.1 M-PAM signalling

When transmitting discrete data over a channel the bits must be represented by signals such that it can be transmitted over a continuous media. It should also be possible for the receiver to decode back to the discrete form even though the signal was distorted by the channel. This process is called modulation and demodulation, and one of the basic modulation scheme is M-ary Pulse Amplitude Modulation (M-PAM). The number M is the number of signal alternatives, i.e. the number of different signals used in the scheme. Since the transmitted data is often binary, this number will here be assumed to be a power of 2, $M = 2^k$. In an M-PAM scheme a signal is built from an amplitude and a pulse form, where the amplitude is the information carrier and the pulse form common for

all signal alternatives. To minimise the average signal energy the amplitudes are centered around zero, e.g. the binary case has the amplitudes -1 and 1 . If $M = 4$ the amplitudes $-3, -1, 1$ and 3 are used. In this way the minimum difference between two amplitude values are always 2. For an arbitrary M the amplitude values can be described by

$$A_i = M - 1 - 2i, \quad i = 0, 1, 2, \dots, M - 1 \quad (10.1)$$

which holds for all positive integer M and not just powers of 2 [54]. Then, to form the signal the amplitude is applied to a pulse form $g(t)$, meaning that the general form of a signal alternative in M-PAM can be written as

$$s_i(t) = A_i g(t) \quad (10.2)$$

In Figure 10.1 a graphical view of the 2-PAM and 4-PAM signal alternatives are shown.

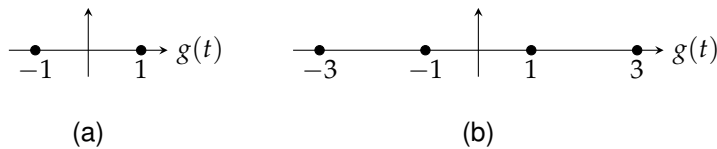


Figure 10.1: Graphical representation of (a) 2-PAM and (b) 4-PAM.

Assuming an infinite binary information sequence to be transmitted, where tuples of $k = \log M$ bits are mapped to an amplitude A_i , the transmitted signal is

$$s(t) = \sum_{\ell} A_{i_{\ell}} g(t - \ell T_s) \quad (10.3)$$

where T_s is the signalling interval.

The pulse form $g(t)$ has the energy $\int_{\mathbb{R}} g^2(t) dt = E_g$. By letting A be a random variable for the amplitude level, the symbol energy becomes,

$$E_s = E \left[\int_{\mathbb{R}} (A g(t))^2 dt \right] = E[A^2] \int_{\mathbb{R}} (g(t))^2 dt = E[A^2] E_g \quad (10.4)$$

For equally likely signal alternatives and levels according to a PAM constellations this yields

$$E_s = E[A^2] E_g = \sum_{i=0}^{M-1} \frac{1}{M} A_i^2 E_g = \frac{M^2 - 1}{3} E_g \quad (10.5)$$

EXAMPLE 10.1 [Origin of BSC] Considering a 2-PAM signal constellation used to communicate over a channel with AWGN. The signals are chosen from the

10.1. M-PAM signalling

237

two signal alternatives in Figure 10.1(a). To transmit a mapping between the information bit a and the amplitude is used according to $s_a(t) = s_a \cdot g(t)$, where $s_a = (-1)^a$ and $g(t) = \sqrt{E_g}\phi(t)$, i.e.

$$s_a(t) = \begin{cases} \sqrt{E_g}\phi(t), & a = 0 \\ -\sqrt{E_g}\phi(t), & a = 1 \end{cases} \quad (10.6)$$

The basis function $\phi(t)$ is a scaled version of $g(t)$ such that it has unit energy, $\int_{\mathbb{R}} \phi^2(t)dt = 1$. The energy per transmitted information bit for this constellation is

$$E_b = \sum_{a=0}^1 \frac{1}{2} \int_{\mathbb{R}} s_a^2(t)dt = \sum_{a=0}^1 \frac{1}{2} E_g \int_{\mathbb{R}} \phi^2(t)dt = E_g \quad (10.7)$$

On the channel white noise with density $R_{\eta}(f) = N_0/2$ is added to the signal. In the receiver, after filtering and ML detection, this means the received signal can be viewed as the point $r = s + z$ in the signal space, where $s = \pm\sqrt{E_b}$ is the transmitted signal amplitude and $z \sim N(0, \sqrt{N_0/2})$. In Figure 10.2 the probability distributions for the received value conditioned on the transmitted s is shown.

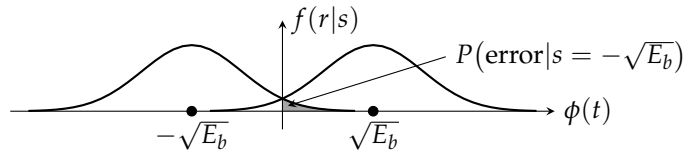


Figure 10.2: The conditional distributions at the receiver side in a 2-PAM transmission over an AWGN channel.

If the two signal alternatives are equally likely, an ML receiver follows a simple decoding rule. If the received value is positive the estimated transmitted amplitude $\sqrt{E_b}$, and if the value is negative the estimated transmitted amplitude is $-\sqrt{E_b}$. Hence the probability of erroneous estimation, conditioned on the transmitted amplitude $-\sqrt{E_b}$ is

$$\begin{aligned} P(\text{error}|s = -\sqrt{E_b}) &= P(r > 0|s = -\sqrt{E_b}) \\ &= P(z > \sqrt{E_b}) \\ &= P(z_{\text{norm}} > \sqrt{\frac{E_b}{N_0/2}}) = Q\left(\sqrt{2\frac{E_b}{N_0}}\right) \end{aligned} \quad (10.8)$$

where $z_{\text{norm}} \sim N(0, 1)$ is a normalised Gaussian variable and

$$Q(x) = \int_x^{\infty} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt \quad (10.9)$$

an error function. Similarly, the error probability conditioned on the transmitted amplitude $\sqrt{E_g}$ gets the same value. This error probability is the probability that a 1 is transmitted and a zero is received, and vice versa. That is, the channel can now be modelled as a binary symmetric channel, BSC, with the cross over probability equal to

$$\varepsilon = Q\left(\sqrt{2\frac{E_b}{N_0}}\right) \tag{10.10}$$

In Figure 10.3 the error probability ε is plotted as a function of the signal to noise ratio E_b/N_0 . With this mapping the capacity for the BSC, $C_{BSC} = 1 - h(\varepsilon)$ is plotted in Figure 10.4.

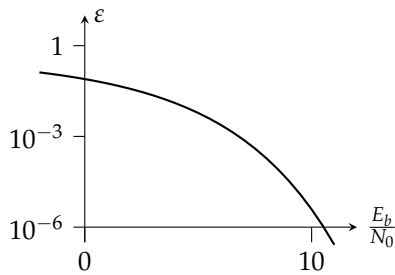


Figure 10.3: The error probability of a BSC as a function of E_b/N_0 for an AWGN channel.

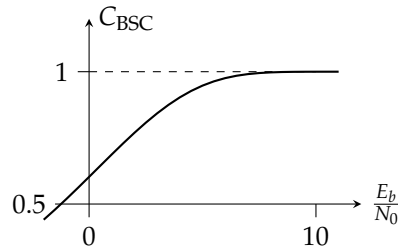


Figure 10.4: The capacity of a BSC as a function of E_b/N_0 for an AWGN channel.

In Chapter 8 it was seen that the interpretation of the mutual information is the same in the discrete and the continuous case, i.e. the amount of information achieved about one variable by observing another. The derivations to get there was based on a Riemann sum limit value. Thus, the interpretation still holds if one of the variables is discrete and the other continuous. This is an important fact since the capacity for a channel is the mutual information maximised over the input distribution. In this section the considered channel model has discrete input signals, but white noise is added on the channel. To derive the capacity for this model the mutual information should be maximised over all distributions for the input signals. However, in many cases the signals are transmitted with equal probabilities, and the counterpart of the capacity, a constraint capacity, is the mutual information for the case with discrete, equally likely inputs and white noise added in the transmission. The mutual information between discrete and continuous variables can be derived in the same manner as before by $I(X; Y) = H(Y) - H(Y|X)$. Letting $p(x)$ be the distribu-

10.1. M-PAM signalling

239

tion for the input symbols X , the conditional entropy can be written as

$$\begin{aligned}
 H(Y|X) &= - \sum_x \int_{\mathbb{R}} f(x, y) \log f(y|x) dy \\
 &= - \sum_x \int_{\mathbb{R}} f(y|x) p(x) \log f(y|x) dy \\
 &= \sum_x p(x) \left(- \int_{\mathbb{R}} f(y|x) \log f(y|x) dy \right) \\
 &= \sum_x p(x) H(Y|X = x)
 \end{aligned} \tag{10.11}$$

Except for the probabilities of X also the density function of Y conditioned on X is needed, which is the channel transition density function. Furthermore, by expressing $f(y) = \sum_x f(y|x) p(x)$, the entropy of Y can be written as

$$\begin{aligned}
 H(Y) &= - \int_{\mathbb{R}} f(y) \log f(y) dy \\
 &= - \int_{\mathbb{R}} \left(\sum_x f(y|x) p(x) \right) \log \left(\sum_x f(y|x) p(x) \right) dy
 \end{aligned} \tag{10.12}$$

For an M-PAM signal constellation with equally likely signal alternatives the density of Y becomes

$$f(y) = \frac{1}{M} \sum_x f(y|x) \tag{10.13}$$

Assuming an additive Gaussian noise with zero mean and variance $N_0/2$, the conditional density function is

$$f(y|x) = \frac{1}{\sqrt{\pi N_0}} e^{-(y-x)^2/N_0} \tag{10.14}$$

In Figure 10.5 the resulting density function for an 8-PAM constellation is shown. In the figure the contribution from the eight conditional density functions $f(y|x)$ are shown as dashed curves, while the density function for Y is shown as a solid curve. In this example the noise variance is quite high to show the behaviour. In the case of a more moderate noise the eight peaks corresponding to the signal alternatives will be more separated by deeper valleys.

To get the entropy of Y the function $-f(y) \log f(y)$ is integrated by numerical methods. The entropy conditional entropy $H(Y|X)$ can be derived from

$$H(Y|X) = \frac{1}{2} \log \pi e N_0 \tag{10.15}$$

since $[Y|X = x] \sim N(x, \sqrt{N_0/2})$.

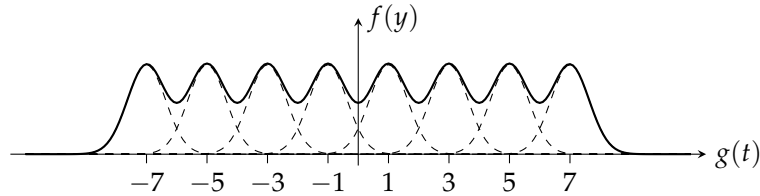


Figure 10.5: The density functions $f(y|x)$ for an 8-PAM signal transmitted over a Gaussian channel, together with the resulting $f(y)$.

In Figure 10.6 plots of the mutual information $I(X;Y) = H(Y) - H(Y|X)$ for M-PAM signalling is shown for the case of equiprobable signal alternatives and additive Gaussian noise. Here, the mutual information $I(X;Y)$ is a measure of how much information can be transmitted over the channel for each channel use, i.e. for each signal alternative sent. The plots typically flattens at the maximum transmitted bits for the number of signal alternatives as the channel becomes good. For example, six bits can be written as 64 different binary vectors and 64-PAM therefore flattens at 6 bits/channel use.

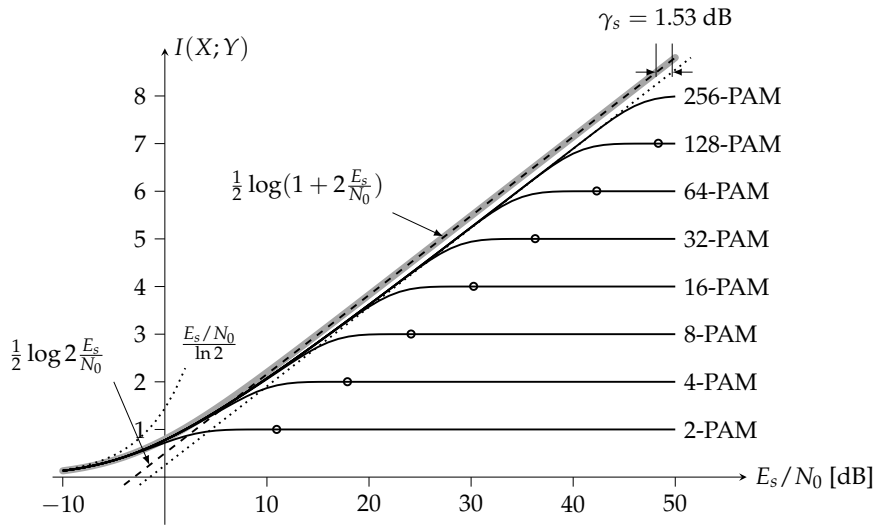


Figure 10.6: Constraint capacity for discrete uniformly distributed signal constellations, like M-PAM, transmitted over an AWGN channel. The grey shaded line in the figure is the capacity $C = \frac{1}{2} \log(1 + 2 \frac{E_s}{N_0})$. The circles on the curves mark the SNR where the uncoded M-PAM signalling gives an error probability of 10^{-6} .

10.2. A note on the dimensionality

241

By assuming Nyquist sampling, $W = 2T_s$, the capacity becomes

$$\begin{aligned} C &= \frac{1}{2} \log\left(1 + \frac{P}{N_0 W}\right) \\ &= \frac{1}{2} \log\left(1 + \frac{E_s/T_s}{N_0/2T_s}\right) \\ &= \frac{1}{2} \log\left(1 + 2\frac{E_s}{N_0}\right) \quad [\text{bit/channel use}] \end{aligned} \quad (10.16)$$

where E_s is the average signalling energy in one signal interval. In the figure the capacity is shown as the thick grey line.

For good channels, as $\frac{E_s}{N_0}$ becomes large, the one in the capacity formula can be neglected. Furthermore, for bad channels, as $\frac{E_s}{N_0}$ becomes small, the following series expansion can be used

$$\frac{1}{2} \log(1+x) = \frac{1}{2 \ln 2} \left(x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots\right) \approx \frac{x}{2 \ln 2} \quad (10.17)$$

Hence, the asymptotic behaviour of the capacity function is given by

$$C \approx \begin{cases} \frac{1}{2} \log 2 \frac{E_s}{N_0}, & \frac{E_s}{N_0} \text{ large} \\ \frac{1}{\ln 2} \frac{E_s}{N_0}, & \frac{E_s}{N_0} \text{ small} \end{cases} \quad (10.18)$$

In Figure 10.6 these two functions are shown as a dashed line and a dotted line, respectively. There is also one more dotted line, located 1.53 dB to the right of the asymptotic capacity. This shows the asymptotic slope of the achievable bit rate for equiprobable M-PAM signalling. The 1.53 dB gap shows the possible gain by not restricting to equiprobable signal alternatives. This quantity is called the *shaping gain* and will be further elaborated in Section 10.3.

Assuming the Nyquist sampling rate $F_s = 2W$, the signal period is $1/2W$. With P as the average power the energy per signal becomes $E_s = P/2W$. Hence, the signal to noise ratio in the capacity formula can be derived as

$$\text{SNR} = \frac{P}{WN_0} = 2\frac{E_s}{N_0} \quad (10.19)$$

which explains the capacity formula in (10.16).

10.2 A note on the dimensionality

Figure 10.6 contains a lot of information about how a real system can be expected to behave compared to what the capacity limit promise. The mutual

information plotted for different M-PAM constellations in the figure shows how practical systems behave at the best, using equally likely signal alternatives. It also shows the asymptotic loss made by using uniform distribution instead of Gaussian. In the plot the the unit of the x -axis is often expressed as *bits/transmission per dimension*, where *transmission* means the transmission of a signal alternative. The extra added term *per dimension* can be interpreted in a variety of ways, and it is worth to give a special note on this. In an information theoretical view the *per dimension* can be any dimension and it is not strictly coupled to the time series of signal alternatives or the dimensionality of the signal constellation.

To get a better understanding of how the capacity relates to the dimensionality of the signal consider an N -dimensional signal and introduce N orthonormal basis function, $\phi_i(t), i = 1, 2, \dots, N$. The orthonormality requirement means

$$\int_{\mathbb{R}} \phi_i(t)\phi_j(t)dt = \delta_{i-j} = \begin{cases} 1, & i = j, \\ 0, & i \neq j \end{cases} \quad (10.20)$$

The basis functions used represent the span of the signal in different dimensions. However, it does not say how they are differentiated. A PAM signal can be seen as N consecutive signal alternatives separated in time, and then the base pulses $g(t - nT_s)$ and $g(t - kT_s)$ are orthogonal if $n \neq k$ and the pulse duration is T . In a QAM constellation it is often used that the basis functions $\phi_1(t) = \sqrt{2}\cos(2\pi f_c t)$ and $\phi_2(t) = \sqrt{2}\sin(2\pi f_c t)$ have an orthogonal behaviour for $f_c \gg 1/T_s$. So, the dimensions here can be seen as the dimensionality of the signal constellation, but it has also a more general interpretation and can e.g. be seen as separation in time. The signal is constructed from N (real) signal amplitudes, $s_i, i = 1, 2, \dots, N$, as

$$s(t) = \sum_{i=1}^N s_i\phi_i(t) \quad (10.21)$$

After transmission over an AWGN channel the received signal is

$$r(t) = s(t) + \eta(t) \quad (10.22)$$

where $\eta(t)$ is white noise with power density $R_\eta(f) = N_0/2$. The signal can then be represented in dimension i as

$$r_i = \int_{\mathbb{R}} r(t)\phi_i(t)dt = \int_{\mathbb{R}} (s(t) + \eta(t))\phi_i(t)dt = s_i + \eta_i \quad (10.23)$$

10.2. A note on the dimensionality

243

where $\eta_i = \int_{\mathbb{R}} \eta(t)\phi_i(t)dt$, and it has been used that

$$\begin{aligned} \int_{\mathbb{R}} s(t)\phi_i(t)dt &= \int_{\mathbb{R}} \sum_j s_j\phi_j(t)\phi_i(t)dt \\ &= \sum_j s_j \int_{\mathbb{R}} \phi_j(t)\phi_i(t)dt \\ &= \sum_j s_j\delta_{i-j} = s_i \end{aligned} \tag{10.24}$$

The noise parameter in the received dimension has the following mean and auto correlation

$$\begin{aligned} E[\eta_i] &= E\left[\int_{\mathbb{R}} \eta(t)\phi_i(t)dt\right] = \int_{\mathbb{R}} E[\eta(t)]\phi_i(t)dt = 0 \\ r_{\eta}(i, j) &= E[\eta_i\eta_j] = E\left[\int_{\mathbb{R}} \eta(t)\phi_i(t)dt \int_{\mathbb{R}} \eta(s)\phi_j(s)ds\right] \\ &= \int_{\mathbb{R}} \int_{\mathbb{R}} E[\eta(t)\eta(s)]\phi_i(t)\phi_j(s)dtds \\ &= \int_{\mathbb{R}} \frac{N_0}{2}\phi_i(t) \int_{\mathbb{R}} \delta(t-s)\phi_j(s)dsdt \\ &= \frac{N_0}{2} \int_{\mathbb{R}} \phi_i(t)\phi_j(t)dt = \frac{N_0}{2}\delta_{i-j} = \begin{cases} \frac{N_0}{2}, & i = j \\ 0, & i \neq j \end{cases} \end{aligned} \tag{10.25}$$

where it is used that $E[\eta(t)] = 0$ and $E[\eta(t)\eta(s)] = \frac{N_0}{2}\delta(t-s)$. This shows the noise component in each dimension is Gaussian with zero mean and variance $\frac{N_0}{2}$, i.e. $\eta_i \sim N(0, \sqrt{N_0/2})$. Hence, the N dimensions are equivalent to N transmissions over a Gaussian channel.

Denoting the total energy in an N -dimensional signal by E_s , the energy per dimension is $E_s^{(\text{dim})} = E_s/N$ and the signal to noise ratio

$$\text{SNR}_N^{(\text{dim})} = 2 \frac{E_s/N}{N_0} = \frac{2 E_s}{N N_0} \tag{10.26}$$

Signalling at the Nyquist sampling rate for a band limited signal with bandwidth W , the sampling rate is $F_s = 2W$. Thus, a vector with N samples will take the transmission time $T = \frac{N}{2W}$, i.e. $N = 2WT$. Hence, the SNR can be written as

$$\text{SNR}_N^{(\text{dim})} = \frac{E_s}{WTN_0} = \frac{P}{WN_0} \tag{10.27}$$

where in the second equality it is used that $E_s = TP$. Thus, the capacity per dimension is

$$C^{(\text{dim})} = \frac{1}{2} \log\left(1 + \frac{P}{WN_0}\right) \left[\frac{b}{\text{tr./dim}}\right] \tag{10.28}$$

and the capacity for the N dimensional signal construction

$$C^{(N)} = \frac{N}{2} \log\left(1 + \frac{P}{WN_0}\right) = WT \log\left(1 + \frac{P}{WN_0}\right) \quad [b/N \text{ dim tr.}] \quad (10.29)$$

Division by T gives the capacity in b/s as

$$C^{(N)} = W \log\left(1 + \frac{P}{WN_0}\right) \quad [b/s] \quad (10.30)$$

This means the capacity in bits per second for a band limited signal is independent of the dimensionality of the signal construction. Especially it is independent of the dimensionality of the signal constellation.

In the derivations it was seen that each amplitude in a signal constellation is equivalent to a sample in terms of the sampling theorem. In essence, N amplitudes gives N degrees of freedom, which can be translated to N samples. Each sample, or dimension, can transmit $\frac{1}{2} \log\left(1 + \frac{P}{WN_0}\right)$ bits per channel use. In this aspect one real amplitude in one dimension in the signal space is regarded as one sample. Hence, from an information theoretic view point there is no difference in transmitting N signals from a one-dimensional constellation during time T , or one signal from an N -dimensional constellation in the same time. They both represent an N dimensional signal space.

Even though the above derivations states that the dimensionality of the signal does not matter, one has to be a bit careful. The requirements in the derivations are that the basis functions are orthonormal and that the utilised bandwidth is unchanged. In the above description M-PAM signals are considered. An M-QAM constellation is essentially formed by using two orthogonal \sqrt{M} -PAM constellations, see Figure 10.7 which describes how two 4-PAM constellations are used to form a 16-QAM constellation. In general, such construction can be done using two real signals modulated in terms of a complex signal.

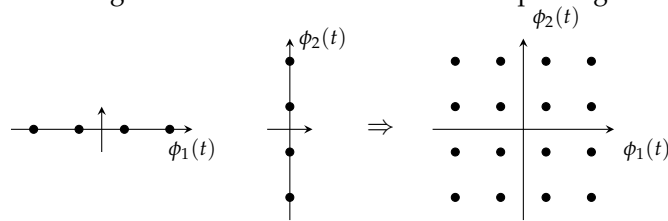


Figure 10.7: Two orthogonal 4-PAM considered as a 16-QAM constellation.

Consider a real base-band signal $s_b(t)$ with the positive bandwidth W , see Figure 10.8a. Since $s_b(t)$ is real its spectra is Hermitian symmetric. Denoting the positive frequency side of $S_b(f)$ by $S_+(f)$, the negative side is the complex conjugate mirror image $S_+^*(-f)$. A frequency shifted signal centered at carrier

10.2. A note on the dimensionality

frequency f_c is created from

$$s(t) = s_b(t)\cos 2\pi f_c t \tag{10.31}$$

Its Fourier transform is

$$S(f) = \frac{1}{2}S_b(f + f_c) + \frac{1}{2}S_b(f - f_c) \tag{10.32}$$

which is shown in Figure 10.8b. Since the inner half of the signal, i.e. for $f_c - W \leq |f| < f_c$ is a mirror image of the outer half, this can be filtered away without losing any information. The procedure is called single sideband modulation and is shown as the function $S_{SSB}(f)$ in Figure 10.8c. This means the effective bandwidth of both the baseband signal $s_b(t)$ and the frequency shifted version $s_{SSB}(t)$ is W . Hence, the capacity for the system is

$$C = \frac{1}{2} \log\left(1 + \frac{P}{N_0 W}\right) = \frac{1}{2} \log\left(1 + 2\frac{E_s}{N_0}\right) \text{ [b/transmission]} \tag{10.33}$$

or, by using the Nyquist sampling rate $F_s = 2W$,

$$C = W \log\left(1 + 2\frac{E_s}{N_0}\right) \text{ [b/s]} \tag{10.34}$$

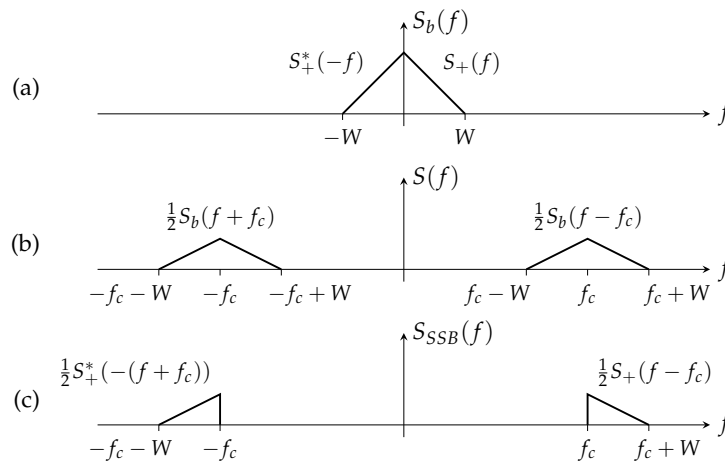


Figure 10.8: modulation of the signal $s_b(t)$ to a higher carrier frequency f_c using single sideband modulation.

Now, like in the QAM construction, consider two real signals in two dimensions. A natural way is to consider a complex signal with the two real baseband signals $s_{\mathcal{R}}(t)$ and $s_{\mathcal{I}}(t)$ with positive bandwidth W , written as

$$s_b(t) = s_{\mathcal{R}}(t) + js_{\mathcal{I}}(t) \tag{10.35}$$

Since the signal is complex there are no longer any symmetries in the frequency domain, and the complete bandwidth $-W \leq f \leq W$ is used for information. However, there is no way to transmit a complex signal directly since the signals are transmitted in a real world. Therefore the signal is shifted up in frequency using cosine for the real part and sine for the imaginary part as

$$s(t) = \text{Re}\{s_b(t)e^{j2\pi f_c t}\} = \frac{1}{2}s_{\mathcal{R}}(t) \cos 2\pi f_c t - \frac{1}{2}s_{\mathcal{I}}(t) \sin 2\pi f_c t \quad (10.36)$$

To view the signal in the frequency domain, use $\text{Re}\{x(t)\} = \frac{1}{2}(x(t) + x^*(t))$,

$$\begin{aligned} S(f) &= \mathcal{F}\{s(t)\} = \frac{1}{2}\mathcal{F}\{s_b(t)e^{j2\pi f_c t}\} + \frac{1}{2}\mathcal{F}\{s_b^*(t)e^{-j2\pi f_c t}\} \\ &= \frac{1}{2}S_b(f - f_c) + \frac{1}{2}S_b^*(-(f + f_c)) \end{aligned} \quad (10.37)$$

where the second equality follows from $\mathcal{F}\{x^*(t)\} = X^*(-f)$. The second term in (10.37), $\frac{1}{2}S_b^*(-(f + f_c))$ is a complex conjugated and mirrored version of $\frac{1}{2}S_b(f - f_c)$ centered around $-f_c$, meaning the negative frequency side of $S(f)$ is a Hermitian reflection of the positive frequency side, as it should for a real sequence. In this case the whole bandwidth $f_c - W \leq f \leq f_c + W$ contains information and the resulting bandwidth for the modulated signal is $W^{(2)} = 2W$. By assuming the power $P^{(2)}$ used over the signal, the resulting capacity is

$$C^{(2)} = \frac{1}{2} \log\left(1 + \frac{P^{(2)}}{N_0 W^{(2)}}\right) \quad [b/\text{tr}] \quad (10.38)$$

and, equivalently by using $F_s^{(2)} = 2W^{(2)}$

$$C^{(2)} = W^{(2)} \log\left(1 + \frac{P^{(2)}}{N_0 W^{(2)}}\right) \quad [b/s] \quad (10.39)$$

To compare the two signalling schemes, where one-dimensional or two-dimensional real signals are used, the constants in (10.39) need to be interpreted. Since the bandwidth is doubled in the second scheme, the power consumption will also double, $P^{(2)} = 2P$. Similarly, the energy used in the signalling will be divided over the two dimensions, and the energy in the second signal becomes $E_s^{(2)} = 2E_s$. Hence, the SNR for the second signalling can be expressed as

$$\text{SNR}^{(2)} = \frac{P^{(2)}}{N_0 W^{(2)}} = \frac{2P}{N_0 2W} = \frac{E_s^{(2)}}{N_0} \quad (10.40)$$

and,

$$C^{(2)} = 2W \log\left(1 + \frac{E_s^{(2)}}{N_0}\right) \quad \text{b/s} \quad (10.41)$$

This relation also reflect the relation between PAM and QAM signalling.

10.3 Shaping gain

The channel capacity in terms of bit rate in Figure 10.6, the gray line, is the maximum achievable bit rate on a Gaussian channel with the SNR measured in E_s/N_0 . To reach this limit the communication system must be optimised in all possible ways. One of many requirements is that the input signal must be chosen as a according to a continuous Gaussian distribution. In most communication systems the choice of signal is done according to uniform distribution over a discrete set of signals. In the figure this asymptotic loss is shown as the gap between the channel capacity and the dotted line. Since this reflects the gain that is possible to achieve by shaping the input distribution from uniform to Gaussian, it is called the *shaping gain* and often denoted γ_s . By viewing the total gain that is possible, viewed from the uncoded case, it can be split in two parts, the shaping gain γ_s and coding gain γ_c . Quite often it is easy to achieve a coding gain of a couple of dB by using some standard channel coding. But to achieve higher gains an alternative is to consider shaping of the constellation.

The ultimate shaping gain of 1.53 dB denoted in Figure 10.6 denotes the maximum shaping gain. To show this value consider the case when the signal to noise ratio, E_s/N_0 , becomes large. The interesting part of the plot is then the growth of the mutual information for M-PAM signalling before it flattens due to a finite number of signals. By letting the number of signal alternatives approaching infinity the distribution of X becomes the continuous rectangular distribution

$$f_u(x) = \frac{1}{2a}, \quad \text{where } -\frac{1}{a} \leq x \leq \frac{1}{a} \quad (10.42)$$

This should be compared to the case of a Gaussian distribution, $f_g(x)$. In this region of the plot, for high SNR the mutual information is dominated by the entropy of the input distribution. For the Gaussian case the average signal energy and the entropy is, see Appendix A,

$$P_g = E[X_g^2] = \sigma^2 \quad (10.43)$$

$$H(X_g) = \frac{1}{2} \log 2\pi e \sigma^2 = \frac{1}{2} \log 2\pi e P_g \quad (10.44)$$

For the uniform case the corresponding derivation gives

$$P_u = E[X_u^2] = \frac{a^2}{3} \quad (10.45)$$

$$H(X_u) = \log 2a = \frac{1}{2} \log 12P_u \quad (10.46)$$

For these two distributions to give the same entropy, the relation on input

power is

$$\gamma_s = \frac{P_u}{P_g} = \frac{\pi e}{6} \approx 1.62 = 1.53 \text{ dB} \tag{10.47}$$

The shaping from a uniform distribution of the signal alternative to a Gaussian gives the ultimately possible gain. In the next example it is shown that a fair amount of the gain can be reached just by considering a distribution that favours the low energy signals before the high energy. The mapping from a uniform to non-uniform distribution indicates that the shaping process can be seen as the dual of source coding, in the sense that perfect source coding gives a uniform distribution of the code symbols. One easy way to get unequal vectors from equally distributed bits, is to consider unequal lengths of the input vectors, and this mapping can be performed in a binary tree.

EXAMPLE 10.2 First, the unshaped system is defined as an 8-PAM system. Then the signal alternatives can be viewed as in Figure 10.9. If the signal alternatives are equally likely the energy derived as the second moment of the signal amplitudes is $E[X^2] = 21$ and for each signal three bits are transmitted.

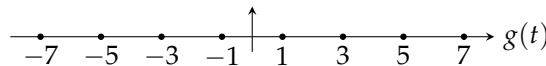


Figure 10.9: Signal alternatives in an 8-PAM constellation.

To find a constellation where the signals have lower average energy, the number of transmitted bits and the inter signal spacing should be unchanged. This corresponds to the obtained bit rate and the symbol error probability in the transmission, respectively. Instead the distribution of the signal alternatives should be chosen non-uniform. If the input sequence is considered as i.i.d. equiprobable bits, one way to alter the distribution is to have unequal length of the vectors mapping to the signal alternative. Here, these vectors are determined from the paths in a binary tree where there are no unused leaves. By choosing some vectors shorter than 3 and others longer, and by mapping high probable vectors to low energy signals, the total energy can be lowered. The tree in Figure 10.10 shows the mapping between signal alternatives s_i and the input vectors decided by the tree paths. Since the binary information is assumed to be equiprobable the probabilities for the nodes at each level is shown under the tree. The average length of the information vectors can then be determined by the path length lemma as

$$E[L] = 1 + 2\frac{1}{2} + 2\frac{1}{4} + 2\frac{1}{8} + 4\frac{1}{18} = 3 \tag{10.48}$$

In the tree there are 12 leaves corresponding to signal alternatives. Hence, the price paid is an expansion of the signal constellation, but the idea is to use the

10.3. Shaping gain

added high energy alternative with a low probability so in average there is a gain. The amplitudes and the corresponding probabilities of the signal alternatives are shown in Figure 10.11. The energy derived as the second moment is then

$$E[X_s^2] = 2\frac{1}{4} + 2\frac{1}{8}3^2 + 2\frac{1}{32}(5^2 + 7^2 + 9^2 + 11^2) = 20 \tag{10.49}$$

and the shaping gain is

$$\gamma_s = 10 \log_{10} \frac{21}{20} = 0.21\text{dB} \tag{10.50}$$

In Problem?? it is shown that the same construction when letting the tree grow even further can give an asymptotic shaping gain of $\gamma_s^{(\infty)} = 0.9177\text{dB}$.

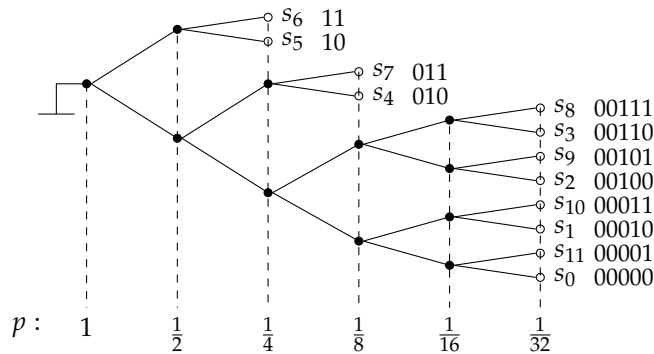


Figure 10.10: A binary tree for determining a shaping constellation.

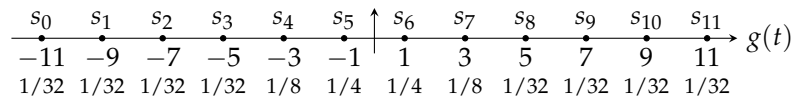


Figure 10.11: Signal alternatives and the probabilities in the shaped constellation.

In the example it was seen that the average energy can be decreased by shaping the probability distribution over the signal constellation. The shaping procedure can also be seen in another way. A vector of two symbols, each modulated by a 16-PAM signal constellation, result in a 256-QAM signal constellation, see upper left constellation of Figure 10.12. Since the QAM signal space has a square form the energy in the corner signal alternatives are rather high. If instead the 256 signal alternatives is chosen within a circle, the upper right

constellation in the figure, the total energy can be decreased. Assuming the distance between two signal points in the figure is 2 and that they are used with equal probability, the average energy derived as the second moment for the squared constellation is

$$E_{\text{QAM}} = 170 \tag{10.51}$$

Similarly, when the signal alternatives in the squared constellation are equally probable, the energy is

$$E_{\text{Sphere}} = 162.75 \tag{10.52}$$

The resulting shaping gain is

$$\gamma_s = 10 \log_{10} \frac{E_{\text{QAM}}}{E_{\text{Sphere}}} = 0.189\text{dB} \tag{10.53}$$

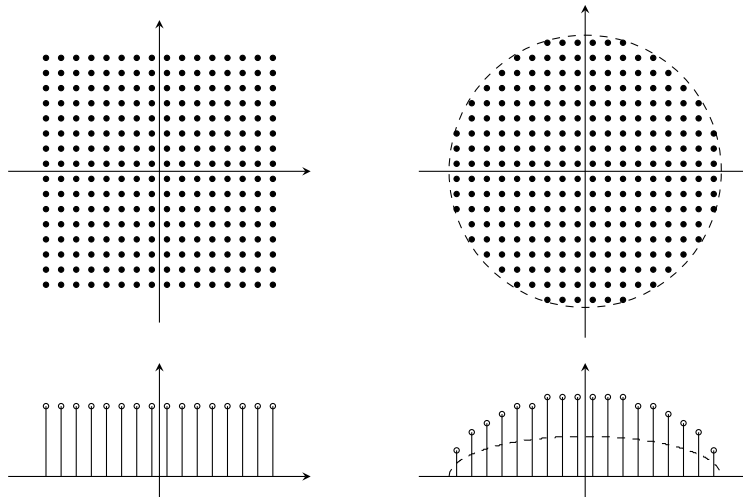


Figure 10.12: The signal alternatives in the two dimensional constellations 256-QAM and a spherical shaped version. Below are the distributions projected to one dimension. The dashed distribution is the projection of the continuous spherical constellation.

In Figure 10.12 it is assumed that the signal alternatives are equally likely in both the squared and the spherical case. The distributions below the signal constellations are the probability functions projected in one dimension. Clearly for the square case there are 16 equally likely alternatives. In the spherical case there are 18 alternatives where the low energy alternatives has highest probability. Hence, by choosing a spherical constellation in two dimensions the distribution is shaped when projected into one dimension.

10.3. Shaping gain

251

To find the maximum shaping, a cubic constellation in N dimensions is compared with a spherical constellation in N dimensions. When the number of signal alternatives grows the discrete constellations can be replaced with continuous distributions without any essential loss of accuracy. Then the second moment of a uniform distribution over an N dimensional cube should be compared with the second moment of a uniform distribution over an N dimensional sphere. To compare the two distributions they should have the same volume, and therefore they are normalised to unit volume.

Starting with the cubic constellation, the volume of an N dimensional cube with side A is

$$V_{\square} = \int_{\square} d\mathbf{x} = \int_{-A/2}^{A/2} \cdots \int_{-A/2}^{A/2} dx_1 \dots dx_N = A^N \quad (10.54)$$

where $\mathbf{x} = (x_1, \dots, x_N)$ is an N dimensional vector. Normalising to a unit volume cube gives that $A = 1$. Since the N -cube is the boundary for a uniform distribution, the probability function is $f_{\square}(\mathbf{x}) = 1/V_{\square} = 1$. Hence, the second moment, or the energy, for the cubic constellation in N dimensions can be derived as

$$\begin{aligned} E_{\square}^{(N)} &= \int_{\square} |\mathbf{x}|^2 d\mathbf{x} = \int_{-1/2}^{1/2} \cdots \int_{-1/2}^{1/2} (x_1^2 + \cdots + x_N^2) dx_1 \dots dx_N \\ &= N \int_{-1/2}^{1/2} x^2 dx = N \left[\frac{x^3}{3} \right]_{-1/2}^{1/2} = N \frac{1}{12} \end{aligned} \quad (10.55)$$

To do similar derivations for the spheric constellation in N dimensions a useful integral relation from [24], formula 4.642, is noted

$$\int_{|\mathbf{x}|^2 \leq R^2} f(|\mathbf{x}|) d\mathbf{x} = \frac{2\pi^{N/2}}{\Gamma(\frac{N}{2})} \int_0^R x^{N-1} f(x) dx \quad (10.56)$$

where $\Gamma(n)$ is the Gamma function, see Section A.2. By letting $f(x) = 1$, the volume of an N -dimensional sphere is

$$\begin{aligned} V_{\circlearrowleft} &= \int_{\circlearrowleft} d\mathbf{x} = \int_{|\mathbf{x}|^2 \leq R^2} d\mathbf{x} = \frac{2\pi^{N/2}}{\Gamma(\frac{N}{2})} \int_0^R x^{N-1} dx \\ &= \frac{2\pi^{N/2}}{\Gamma(\frac{N}{2})} \left[\frac{x^N}{N} \right]_0^R = \frac{2\pi^{N/2}}{\Gamma(\frac{N}{2})} \frac{R^N}{N} \end{aligned} \quad (10.57)$$

Setting $V_{\circlearrowleft} = 1$ yields the radius

$$R = \frac{1}{\sqrt{\pi}} \left(\frac{N}{2} \Gamma\left(\frac{N}{2}\right) \right)^{1/N} = \frac{1}{\sqrt{\pi}} \left(\Gamma\left(\frac{N}{2} + 1\right) \right)^{1/N} \quad (10.58)$$

and then the normalised energy can be derived as

$$\begin{aligned}
 E_{\bigcirc}^{(N)} &= \int_{\bigcirc} |\mathbf{x}|^2 d\mathbf{x} = \frac{2\pi^{N/2}}{\Gamma(\frac{N}{2})} \int_0^R x^{N-1} x^2 dx \\
 &= \frac{2\pi^{N/2}}{\Gamma(\frac{N}{2})} \frac{R^{N+2}}{N+2} = \frac{2\pi^{N/2} R^N}{\underbrace{\Gamma(\frac{N}{2}) N}_{V_{\bigcirc}=1}} \frac{N}{N+2} R^2 \\
 &= \frac{N}{(N+2)\pi} \left(\Gamma\left(\frac{N}{2} + 1\right)\right)^{2/N} \tag{10.59}
 \end{aligned}$$

The shaping gain for the N -dimensional case when comparing the cubic and the spherical constellations is

$$\gamma_s^{(N)} = \frac{E_{\square}^{(N)}}{E_{\bigcirc}^{(N)}} = \frac{\pi(N+2)}{12\Gamma\left(\frac{N}{2} + 1\right)^{2/N}} \tag{10.60}$$

The Gamma-function generalises the factorial function to positive real values with a smooth curve where $n! = \Gamma(n+1)$, for integer n . Therefore it is reasonable to use Stirling's approximation to get

$$\Gamma(n+1) = n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \tag{10.61}$$

Hence, for large N ,

$$\begin{aligned}
 \gamma_s &\approx \frac{\pi(N+2)}{12 \left((2\pi \frac{N}{2})^{1/2} (\frac{N}{2e})^{N/2} \right)^{2/N}} \\
 &= \frac{\pi e}{6} \frac{N+2}{N} \left(\frac{1}{\pi N}\right)^{1/N} \rightarrow \frac{\pi e}{6}, \quad N \rightarrow \infty \tag{10.62}
 \end{aligned}$$

which is the same ultimate shaping gain as when comparing uniform and Gaussian distributions for the input symbols. Actually, as will be seen in Problem ??, the projection from a uniform distribution over a multi-dimensional sphere to one dimension will be a Gaussian distribution when the dimensionality grows to infinity. Therefore comparing the shaping gain between multidimensional cubic and spherical uniform distributions is the same as comparing the one-dimensional uniform and Gaussian distributions.

10.4 SNR gap

When describing the capacity formula for discrete input constellations like M-PAM, it is also natural to consider the *SNR gap*. This is a measure of how far

10.4. SNR gap

from the capacity limit a system is working for a specific achieved probability of error. Here PAM signalling is considered to derive the achieved bit error rate. Then the SNR gap describes the possible gain in SNR by approaching the capacity.

Previously, the signal constellation for 2, 4 and 8-PAM has been considered, see Figure 10.1 and 10.9. In general, for an M-PAM constellation the signal amplitudes are determined by

$$A_i = M - 1 - 2i, \quad i = 0, 1, \dots, M - 1 \quad (10.63)$$

Then the valid amplitudes will be as described in Figure 10.13. The pulse shape is determined by the function $g(t) = \sqrt{E_g}\phi(t)$ where $\phi(t)$ has unit energy.

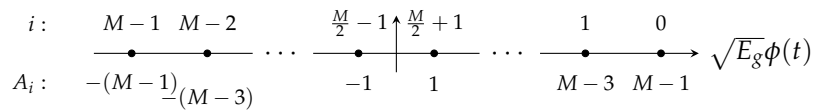


Figure 10.13: Signal alternatives in an M-PAM constellation.

The received signal, distorted by AWGN, is given as $y = A_i + \eta$, where $\eta \sim N(0, \sqrt{N_0}/2)$. An ML decoding rule chooses the signal amplitude closest to the received signal in terms of Euclidian distance in Figure 10.13. There will be a decoding error in the case when the received signal is not closest to the transmitted signal alternative. For the $M - 2$ inner signal alternatives, $i \in \{1, \dots, M - 2\}$, this will happen when the noise component η is either larger than $\sqrt{E_g}$ or smaller than $-\sqrt{E_g}$. In both cases the probability is $P(\eta > \sqrt{E_g})$. For the outer signal alternatives, $i \in \{0, M - 1\}$, it will only be error in one of the cases. That means the error probability conditioned on the signal alternative is

$$P_{e|i} = \begin{cases} 2P(\eta > \sqrt{E_g}), & i = 1, \dots, M - 2 \\ P(\eta > \sqrt{E_g}), & i = 0, M - 1 \end{cases} \quad (10.64)$$

With equally likely signal alternatives the error probability is

$$\begin{aligned} P_e &= \sum_i \frac{1}{M} P_{e|i} = \frac{1}{M} \left((M - 2)2P(\eta > \sqrt{E_g}) + 2P(\eta > \sqrt{E_g}) \right) \\ &= 2 \frac{M - 1}{M} P(\eta > \sqrt{E_g}) = 2 \left(1 - \frac{1}{M} \right) P(\eta > \sqrt{E_g}) \end{aligned} \quad (10.65)$$

The above probability is given in terms of the energy in the pulse shape E_g . the energy in signal alternative i is $A_i^2 E_g$, and hence, the average signal energy is given by

$$E_s = \frac{1}{M} \sum_{i=0}^{M-1} A_i^2 E_g = \frac{E_g}{M} \sum_{i=0}^{M-1} (M - 1 - 2i)^2 = \frac{E_g}{3} (M^2 - 1) \quad (10.66)$$

or, equivalently,

$$E_g = \frac{3E_s}{M^2 - 1} \quad (10.67)$$

Then, together with the noise variance of $N_0/2$, the signal error probability can be expressed as

$$\begin{aligned} P_e &= 2\left(1 - \frac{1}{M}\right)P\left(\eta > \sqrt{\frac{3E_s}{M^2-1}}\right) \\ &= 2\left(1 - \frac{1}{M}\right)Q\left(\sqrt{\frac{3E_s}{(M^2-1)N_0/2}}\right) \\ &= 2\left(1 - \frac{1}{M}\right)Q\left(\sqrt{\frac{3}{M^2-1}2\frac{E_s}{N_0}}\right) \\ &= 2\left(1 - \frac{1}{M}\right)Q\left(\sqrt{\frac{3}{M^2-1}\frac{P}{WN_0}}\right) \end{aligned} \quad (10.68)$$

When transmitting binary vectors the number of signal alternatives should be a power of two, $M = 2^k$, where k is the number of transmitted bits per channel use. To have reliable communication this number should be less than the capacity, in bits per transmission,

$$k \leq C = \frac{1}{2} \log(1 + \text{SNR}) \quad (10.69)$$

By rearranging the relation between the capacity and the transmitted bits above, it is seen that

$$\frac{\text{SNR}}{2^{2k} - 1} \geq 1 \quad (10.70)$$

Therefore, it is reasonable to define a normalised SNR as

$$\text{SNR}_{\text{norm}} = \frac{\text{SNR}}{2^{2k} - 1} \quad (10.71)$$

where the signal to noise ratio is

$$\text{SNR} = \frac{P}{WN_0} = 2\frac{E_s}{N_0} = 2k\frac{E_b}{N_0} \quad (10.72)$$

As $k = C$ the normalised SNR is one since $C = \frac{1}{2} \log(1 + \text{SNR})$ gives $\frac{\text{SNR}}{2^{2C}-1} = 1$. Thus,

$$\text{SNR}_{\text{norm}} \begin{cases} = 0\text{dB}, & k = C \\ > 0\text{dB}, & k < C \end{cases} \quad (10.73)$$

which means the normalised SNR can be seen as a measure of how far from the capacity a system works.

10.4. SNR gap

255

Since $M = 2^k$ the normalised SNR can be written as $\text{SNR}_{\text{norm}} = \frac{\text{SNR}}{M^2 - 1}$, and the error probability for the M-PAM constellation becomes

$$P_e = 2 \left(1 - \frac{1}{M}\right) Q(\sqrt{3 \cdot \text{SNR}_{\text{norm}}}) \tag{10.74}$$

For large M it is simplified to

$$P_e = 2Q(\sqrt{3 \cdot \text{SNR}_{\text{norm}}}) \tag{10.75}$$

In Figure 10.14 the error probability is plotted as a function of the normalised SNR for 2-PAM, 4-PAM, 8-PAM and M-PAM, where M is large. At an error probability of 10^{-6} the normalised SNR is close to 9 dB for large M . For a 2-PAM system it is for the same error probability 8.8 dB. The conclusion from this is that a PAM system working at an error probability of 10^{-6} has a gap to the capacity limit of 9 dB.

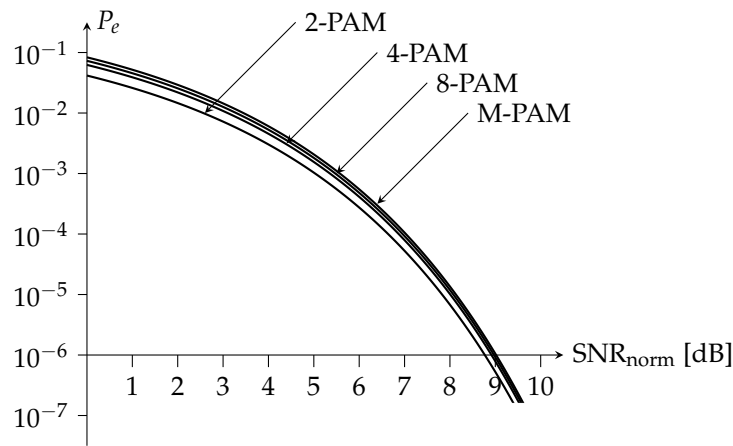


Figure 10.14: Symbol error probability for M-PAM signals as a function of the normalised SNR.

Quite often the SNR gap is used when estimating the bit rate achieved by a PAM (or QAM) system. Then it is viewed from another angle. Starting with (10.75) the symbol error probability for large M , the normalised SNR can be written as

$$\text{SNR}_{\text{norm}} = \frac{1}{3} \left(Q^{-1}(P_e/2)\right)^2 \tag{10.76}$$

Since the normalised SNR is $\text{SNR}_{\text{norm}} = \frac{\text{SNR}}{2^{2k} - 1}$, this gives

$$2^{2k} = 1 + \frac{\text{SNR}}{\frac{1}{3} \left(Q^{-1}(P_e/2)\right)^2} \tag{10.77}$$

or, equivalently, the number of bits per transmission

$$k = \frac{1}{2} \log \left(1 + \frac{\text{SNR}}{\frac{1}{3}(Q^{-1}(P_e/2))^2} \right) = \frac{1}{2} \log \left(1 + \frac{\text{SNR}}{\Gamma} \right) \quad (10.78)$$

where $\Gamma = \frac{1}{3}(Q^{-1}(P_e/2))^2$ is the same SNR gap for PAM (or QAM) constellations as derived earlier from Figure 10.14, in terms of normalised SNR. Going back to Figure 10.6, the circles on the curves for the mutual information for the M-PAM systems correspond to the SNR where $P_e = 10^{-6}$. For high SNR the horizontal difference in SNR between the capacity and the circular mark is the SNR gap Γ . If the required error probability is decreased, the uncoded M-PAM system can tolerate a lower SNR and the gap is decreased. Similarly, if the required error probability is increased the SNR for the M-PAM system is moved to the right and the gap is increased. It is here worth noticing that the capacity, as well as the mutual information plots in the figure, are bounds for when it is possible to achieve arbitrarily low error probability, while the circular point denote the uncoded system when the error is 10^{-6} . To maintain the same error probability for a lower SNR some coding and/or shaping is required. The capacity limit is in this aspect the limit of the lowest possible SNR for which it is possible to transmit this number of bits.

In Figure 10.15 the SNR gap Γ is plotted as a function of the symbol error rate P_e . For $P_e = 10^{-6}$ it becomes $\Gamma \approx 9$ dB, which is often assumed in bit rate estimations. For the bit rate in bits per seconds, the Nyquist sampling rate of $2W$ can be assumed to get

$$R_b = W \log \left(1 + \frac{\text{SNR}}{\Gamma} \right) = W \log \left(1 + \frac{P}{\Gamma W N_0} \right) \quad (10.79)$$

EXAMPLE 10.3 Consider an LTE like communication system using the bandwidth 20 MHz. Assume that the received signal level is -70 dBm over the complete band. That gives the signal power $P = 10^{-70/10}$ mW. While this level can be regarded as pessimistic, it can be compensated by an optimistic view of the noise. As soon as an electrical current flows through a conductor the thermal noise is added. So just by receiving the signal in the antenna a noise level of -174 dBm/Hz is added to the signal. This means that $N_0 = 10^{-174/10}$ mW/Hz. The signal to noise ratio over the bandwidth is then

$$\text{SNR} = \frac{P}{N_0 W} = \frac{10^{-7}}{10^{-17.4} \cdot 20 \cdot 10^6} = 1.26 \cdot 10^3 = 31 \text{ dB} \quad (10.80)$$

The capacity for this system is then

$$C = W \log(1 + \text{SNR}) = 206 \text{ Mb/s} \quad (10.81)$$

10.4. SNR gap

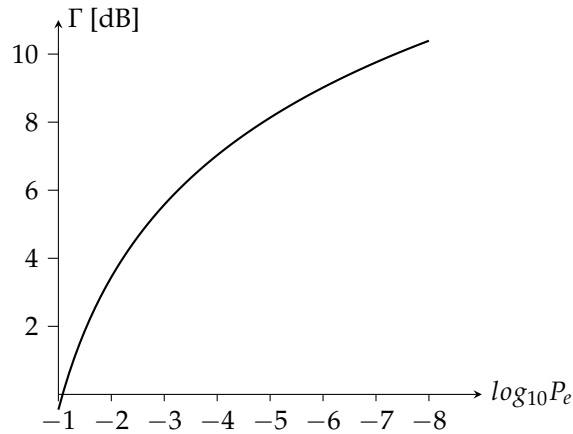


Figure 10.15: SNR gap Γ as a function of the symbol error P_e .

To estimate the achieved bit rate a closer look on the system is required. First there is an error correcting code. In the LTE system there are different strengths of the coding that can be chosen. Here it is assumed a coding gain of $\gamma_c = 4$ dB, which is relatively high. Further, assume the system is working at an error probability of $P_e = 10^{-6}$, which gives an SNR gap of $\Gamma = 9$ dB. Even though it is not included in the LTE system as standardised today a shaping gain is assumed of $\gamma_s = 0.5$ dB. This means the effective SNR of the transmission is

$$\text{SNR}_{\text{eff}} = \text{SNR} - \Gamma + \gamma_c + \gamma_s = 31 - 9 + 4 + 0.5 = 26.5 \text{ dB} = 447 \quad (10.82)$$

That means an estimation of the bit rate in the system can be obtained as

$$R_b = W \log(1 + \text{SNR}_{\text{eff}}) = 122 \text{ Mb/s} \quad (10.83)$$

In reality there are some more things to consider in the LTE system and the true bit rate for the 20 MHz band is about 100 Mb/s.
