
Chapter 5 Coded-Modulation for Band-Limited AWGN Channels

We now introduce the bandwidth-efficient coded-modulation techniques for ideal AWGN channels. The idea of combined coding and modulation design was first suggested by J. L. Massey in 1974, and then realized with stunning results by Ungerboeck and Imai. The common core is to optimize the code in Euclidean space.

On band-limited channels, nonbinary signal alphabets such as M -PAM must be used. The M -ary signaling and the potential coding gain in the bandwidth-limited regime have been discussed in Chapter 2 from the information-theoretic point of view.

5.1 编码调制的基本原理

Traditionally, coding and modulation have been considered as two separate parts of a digital communication system. At the transmitter, an error-correcting encoder is followed by a simple modulator; at the receiver, the received waveform is first demodulated, and then the error correction code is decoded. In this scenario, the modulator and demodulator are usually devised to convert a waveform channel into a discrete channel, and the error correction encoder/decoder are designed, based on maximizing the minimum Hamming distance, to correct the errors that occurred in the discrete channel. Higher improvement in performance is achieved by lowering the code rate at a cost of bandwidth expansion.

最近, 随着数据速率的日益提高, 要求通信系统具有较高的频谱利用率。为了在提高系统功率效率的同时, 不宽展系统所占用的带宽, 人们提出了编码调制技术。With coded modulation schemes, significant coding gains (so the BER performance improvement) can be achieved without increasing bandwidth (or sacrificing bandwidth efficiency).

编码调制遵循下面两个基本原理:

- 通过扩展信号星座 (即增加调制信号集中的信号个数) 而不是通过增加系统的带宽来提供编码所要求的信号冗余。

Example 5.1: Consider the situation where a stream of data is to be transmitted with throughput of 2 bits/s/Hz over an AWGN channel. One possible solution is to use an uncoded QPSK system. As a coded solution, we may employ a rate-2/3 convolutional code with an 8-PSK signal set. Note that this coded 8-PSK scheme yields the same throughput as uncoded QPSK, i.e., 2 bits/s/Hz; moreover, both schemes require the same bandwidth.

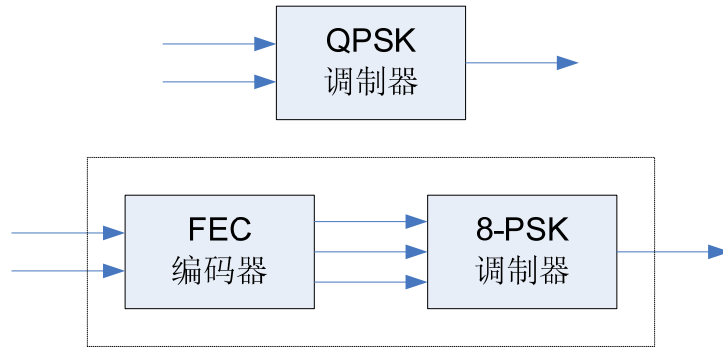


Figure 5.1.1

从第二章中的调制信号星座的容量分析可知，信号星座点个数增加一倍所提供的冗余已足够实现在不增加系统带宽的条件下，逼近容量限的性能。再进一步扩展星座，所得到的性能增益将很少。因此，在通常的编码调制系统中采用的是码率为 $k/(k+1)$ 的信道码。

■ 将编码与调制作为一个整体进行联合优化设计。

因为 Although the expansion of a signal set (e.g., from QPSK to 8-PSK) provides the redundancy required for coding, it shrinks the distance between the signal points if the average signal energy is kept constant. This reduction in distance should be compensated by coding advantage if the coded scheme is to provide a benefit.

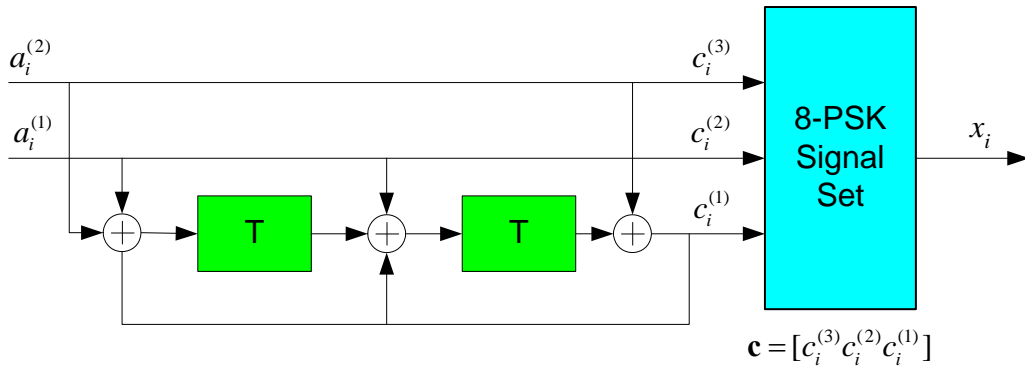
如果是按照传统方法，简单地在一个纠错编码器后级联一个M元调制器，而纠错编码器是基于汉明距离准则进行设计，则所得到的结果往往会令人失望。

另外，The use of hard-decision demodulation prior to the decoding in a coded scheme causes a loss of SNR. To avoid such a hard-decision loss, it is necessary to employ soft-output detector. TCM integrates demodulation and decoding in a single step and decoder operates directly on the soft-output samples of the channel. The decision rule of the optimum decoder depend on the Euclidean distance.

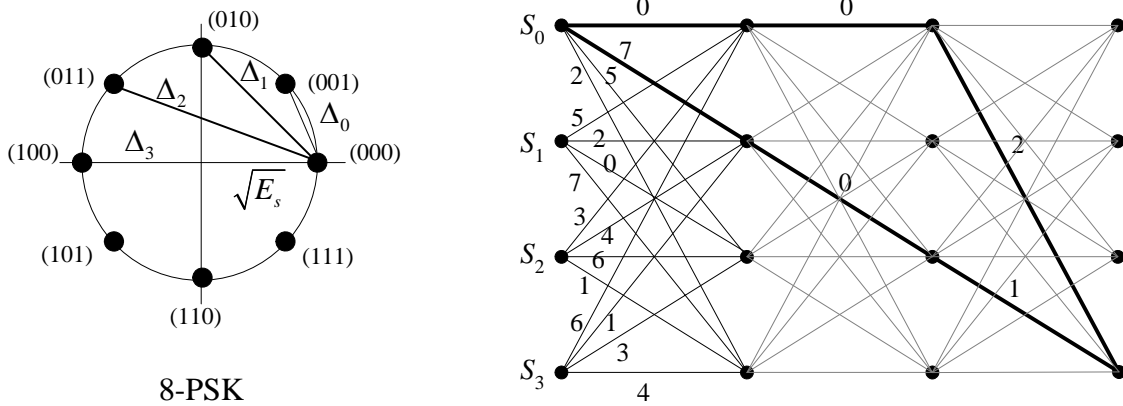
Example 5.2: (对于上例中的)

For the coded 8-PSK scheme above, if we choose the rate-2/3 convolutional code of Fig.5.1.2(a) which is designed based on maximizing free Hamming distance, and the mapping of 3 output bits of the convolutional encoder to the 8-PSK signal points is done as shown in Fig.5.1.2(b), then we can find that the minimum Euclidean distance between a pair of paths forming an error event (which is sometimes called *free Euclidean distance*) is

$$\begin{aligned}
 d_{free}^2 &= \min_{\{x_i\} \neq \{x'_i\}} \sum_i d_E^2(x_i, x'_i) \\
 &= d_E^2(x_0, x_7) + d_E^2(x_0, x_0) + d_E^2(x_2, x_1) \\
 &= \Delta_0^2 + 0 + \Delta_0^2 = 1.172E_s
 \end{aligned}$$



(a) Encoder structure



(b) Signal mapping rule and the trellis diagram of the coded scheme

Fig. 5.1.2 A rate-2/3 convolutional coded 8-PSK scheme

To compare the coded and uncoded schemes it is common to use the *coding gain* parameter, which is defined as the difference in SNR for an objective target bit error rate between a coded system and an uncoded system.

$$\text{coding gain} \triangleq \text{SNR} |_{\text{uncoded}} - \text{SNR} |_{\text{coded}}$$

At high SNR, this gain is termed the *asymptotic coding gain* (ACG) and is expressed as

$$\gamma = 10 \log_{10} \frac{(d_{\text{free}}^2 / E_s)_{\text{coded}}}{(d_{E,\text{min}}^2 / E_s)_{\text{uncoded}}} \quad \text{dB}$$

For the coded scheme above, $\gamma = 10 \log_{10} \frac{1.172}{2} = -2.3 \text{ dB}$. This result shows the performance degradation of the coded scheme (optimized based on the free Hamming distance) compared to the uncoded one.

Massey pointed out that it was necessary to integrate the design of encoder and modulator, and to treat the code and modulation scheme as an entirety, as shown in Fig. 5.1. Thus, 系统整体方案就应该基于 *maximizing the minimum Euclidean distance between coded signal*

*sequences rather than Hamming distance*来设计. More recently, it has been recognized that the design of coded modulation schemes for the AWGN channel is a problem that is best viewed in the geometric signal-space context. 所以编码调制也称为信号空间编码。

In TCM schemes, the code and an expanded signal set are jointly designed as a physical unit. The design criterion is to maximize the free Euclidean distance between coded signal sequences *rather than Hamming distance*. The resulting code can provide a significant coding gain and the loss from the expansion of the signal set can be overcome.

Example 5.3:

我们从编码调制的角度，考虑图 5.1.1 中的编码器与调制器的联合设计。As an alternative coded scheme, we may use the 8-PSK TCM scheme shown in Fig. 5.1.3, which was introduced by Ungerboeck. We will see in Section 5.5 that this TCM scheme can provide an asymptotic coding gain of $\gamma = 3$ (dB).

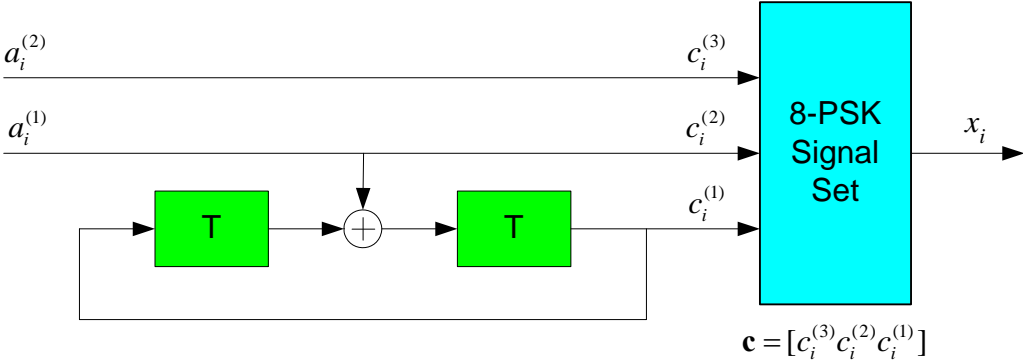


Fig. 5.1.3 The 4-state TCM encoder for 8-PSK

The performances of various TCM schemes are shown in Fig. 5.1.4. It is seen from Fig. 5.1.4 that the improvement of coding is evident. Note that the coding schemes shown in Figure 5.1.4 achieves the coding gains without requiring more bandwidth than the uncoded QPSK system.

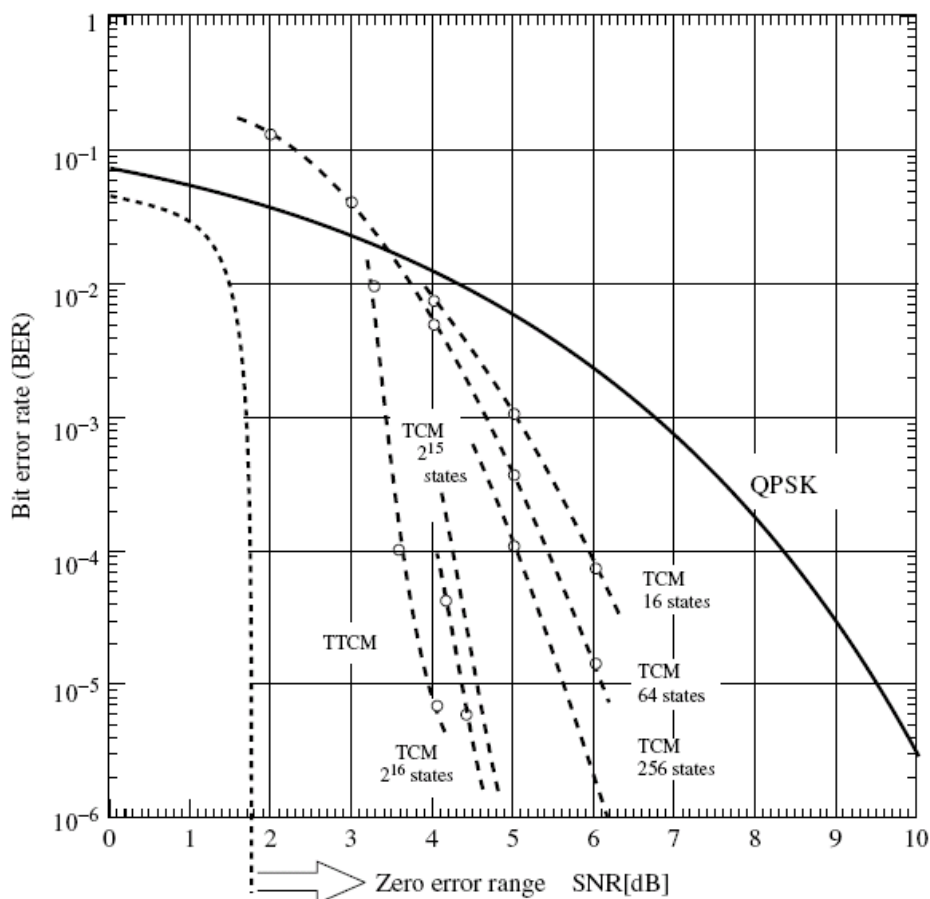


Figure 5.1.4: Bit error probability of Quadrature Phase-Shift Keying (QPSK) and selected 8-PSK trellis-coded modulation (TCM), trellis-turbo coded (TTCM), and block turbo coded (BTC) systems as a function of the normalized signal-to-noise ratio.

5.1.1 两种基本实现方法

Similar to the case of binary codes, we introduce interdependences between consecutive signal points in order to increase the distance between the closest pair of sequences of signal points. A perspective from *signal-space coding* may provide more insight into coded modulation schemes. In order to obtain large coding gain, the codes should be designed in a subspace of signal space with high dimensionality, where a larger minimum distance in relation to signal power can be obtained. The dimensionality $2BT_0$ can be increased for fixed bandwidth B by increasing the time interval T_0 , making it multiple symbol intervals.

For moderate coding gain at moderate complexity, Two basic ways to generate modulation (or signal-space) codes in conjunction with passband QAM modulation are as follows:

- 直接来自于几何考虑: A sequence of $N/2$ two-dimensional transmitted symbols can be considered as a single point in an N -dimensional constellation. Each element of the constellation alphabet (called a codeword) is a vector in \mathbb{R}^N (or $\mathbb{C}^{N/2}$). A set of K input bits are used to select one of 2^K codewords in the multidimensional constellation. A typical example of the multidimensional constellation is the lattice code. 它类似于二进制编码

中的分组码。

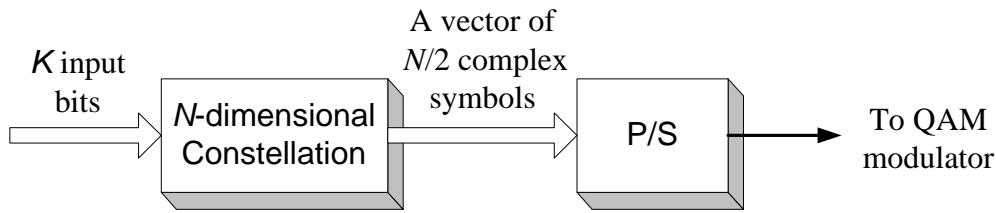


Fig. 5.1.5

- Another way is to extend the dimensionality of the transmitted signal by basing it on a finite-state machine (FSM). The extra bits produced by the FSM implies an inherent increase in the number of points in the constellation.

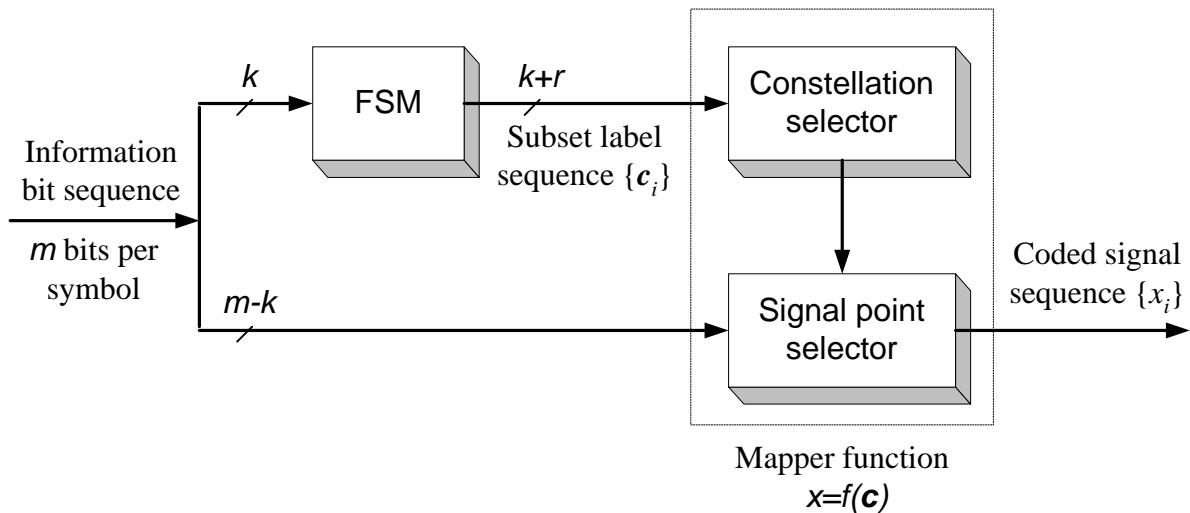


Fig. 5.1.6

5.1.2 Overview of Coded Modulation Techniques

The existing coded modulation schemes for band-limited AWGN channels can be broadly classified into four categories:

- 1) Lattice codes
- 2) Trellis/TCM codes (trellis-coded modulation)

TCM was proposed by Ungerboeck in 1982, in which a convolutional code is usually used as the underlying FSM. The term *trellis-coded modulation* originates from the fact that these coded sequences consist of modulated symbols rather than binary digits. In other words, in TCM schemes, the trellis branches are labeled with redundant nonbinary modulated symbols rather than with binary coded symbols. TCM codes can be decoded by a maximum-likelihood decoder using Viterbi algorithm.

Trellis codes are to lattices as binary convolutional codes are to block codes.

- 3) Turbo-TCM
- 4) Multilevel codes (also known as BCM)

Multilevel codes was proposed by H. Imai in 1977. The underlying strategy is to protect

each label bit of the signal point by an individual binary code, so multiple encoders (at different levels) are employed. At the receiver, the received sequence of signal points are usually decoded by a multistage decoder.

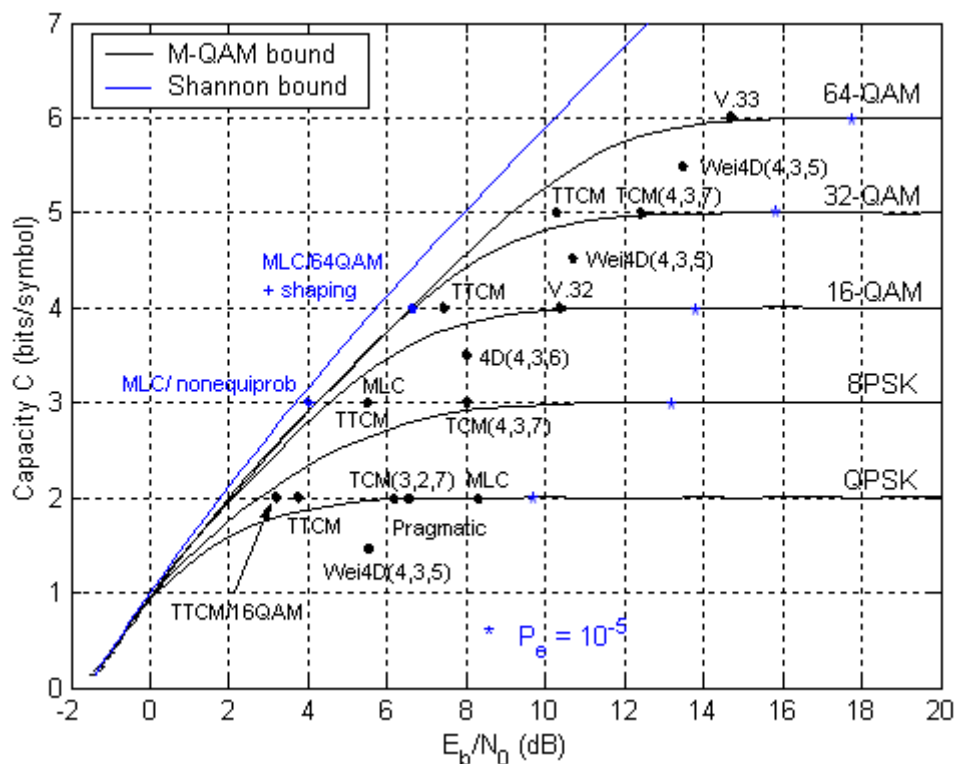
In multilevel coding (MLC) schemes, any code, e.g., block codes, convolutional codes, or concatenated codes, can be used as component codes. Since in early MLC schemes the FEC codes used were usually block codes, the MLC scheme is also referred to as block coded modulation (BCM).

5) Bit-interleaved coded modulation (BICM) (with iterative decoding)

BICM was first proposed by Zehavi in 1992 for coding for fading channels, in which the output stream of a binary encoder is bit-interleaved and then mapped to an M -ary constellation. Its basic idea is to increase the code diversity. (For Rayleigh fading channels, the code performance depends strongly its minimum Hamming distance rather than the minimum Euclidean distance.) The information-theoretic aspects of BICM have been analyzed by Caire.

Recently, it has been recognized that the BICM based on turbo-like codes and iterative decoding provides an effective realization method for Gallager's coding theorem (proposed in 1968 for discrete memoryless channels). With this scheme, very good performance can be achieved on both AWGN and fading channels.

Fig. 5.1.7 depicts the performances of some of typical coded modulation schemes.



constellation designs.

- The first idea is to change the relative spacing of points in the constellation. The hexagonal constellation leads to the reduced variance with the same minimum distance. (Alternatively, we could keep the variance constant, in which case the hexagonal constellation would have a larger minimum distance than the squared constellation.) This decrease in power for the same minimum distance or increase in minimum distance for the same power through changing the relative spacing of the points is called *coding gain*.
- The 2nd approach is to change the shape or outline of the constellation without changing the relative positioning of points. The circular constellation will have a lower variance than the squared constellation. On the same grounds, a circular constellation will have the lowest variance of any shaping region for a square grid of points. The resulting reduction in power is called *shaping gain*.

Significantly, shaping the constellation changes the marginal density of the data symbol. This is illustrated in Fig. 5.9.

Coding and shaping gain can be combined, e.g., by changing the points in the circularly shaped constellation to a hexagonal grid while retaining the circular shaping. Usually, channel coding deals with the internal arrangement of the points, whereas shaping treats regions.

- The 3rd idea is to employ multidimensional signal constellation. A data symbol drawn from an N -dim constellation is transmitted once every $N/2$ signaling interval. When we design a 2D constellation, and choose the $N/2$ successive symbols to be an arbitrary sequence of 2D symbols drawn from that constellation, the resulting N -dim constellation is the $N/2$ -fold Cartesian product of 2D constellations. From Chapter 2, we know that the performance of this N -dim constellation is the same as the underlying 2D constellation. An alternative is to design an N -dim signal constellation that is not constrained to have this Cartesian product structure. When $N > 2$, it is called a multidimensional signal constellation.

Greater shaping and coding gains can be achieved with a multidimensional signal constellation than with a 2D constellation. However, multidimensional constellations suffer from a complexity that increases exponentially with dimensionality.

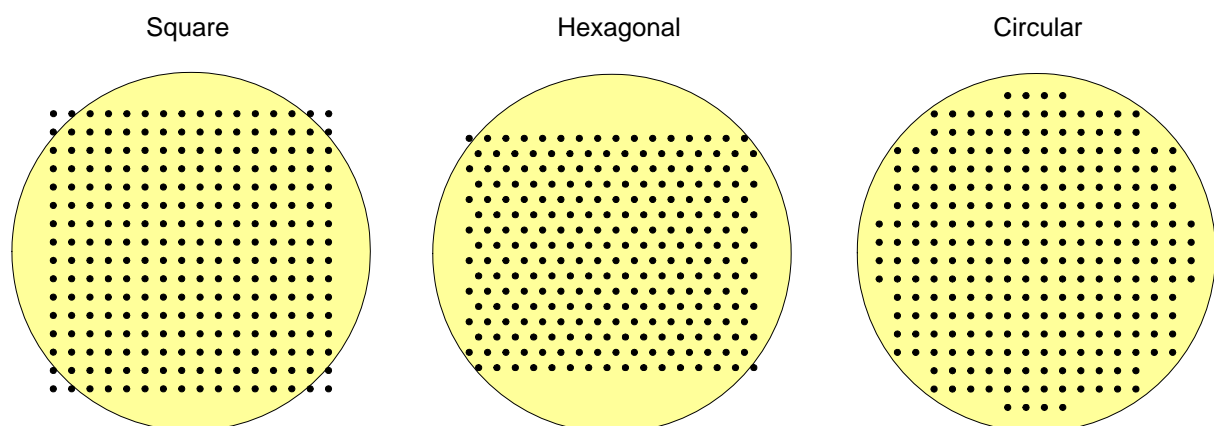


Fig.5.8 Three 2D constellations with the same minimum distance and 256 points.

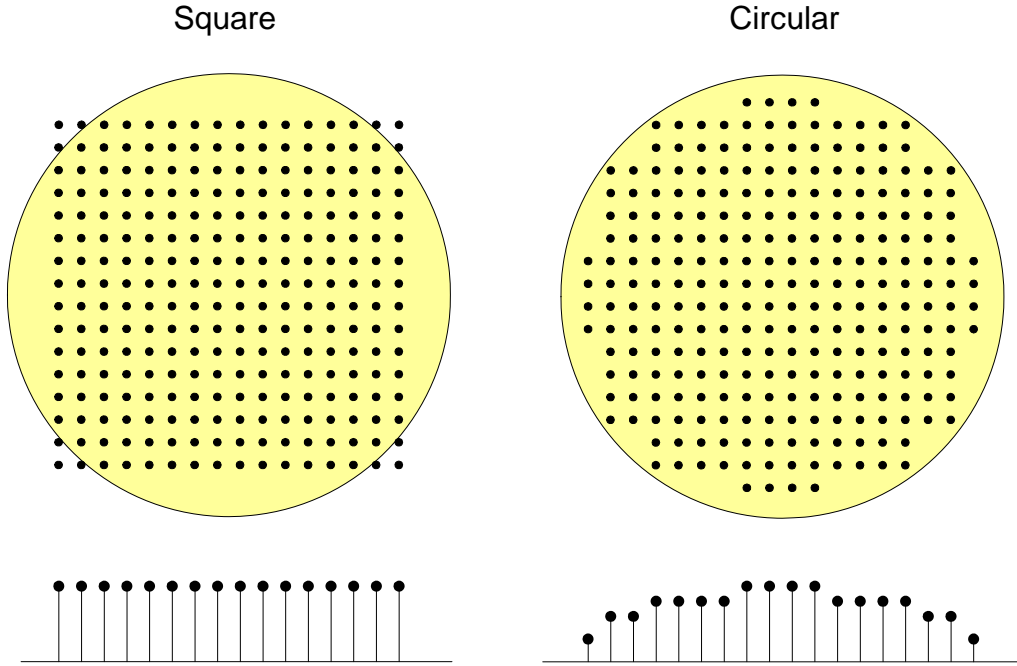


Fig. 5.9 The one-dimensional probability distribution for an unshaped and shaped 2D constellation.

5.2.1 Ultimate Shaping Gain

Here, we just provide a preliminary discussion on the maximum shaping gain; we will discuss this in detail in the next section. Without loss of generality, the derivations are done for one-dimensional constellations.

The baseline system is assumed to use a uniform distribution of signal points. The shaped system should exhibit a (discrete) Gaussian distribution. In order to transmit at the same rate, both distributions have to have the same entropy.

When considering constellations with a large number of signal points, it is convenient to approximate the distribution by a continuous probability density function (pdf). Hence we compare a continuous uniform pdf with a Gaussian one.

Let E_u be the average energy of the reference system. Then the differential entropy of its transmitted symbols x is given by

$$h(X) = \frac{1}{2} \log_2(12\sigma_x^2) = \frac{1}{2} \log_2(12E_u)$$

If $x \sim \mathcal{N}(0, E_g)$ (Gaussian distributed with average energy E_g), its entropy is equal to

$$h(X) = \frac{1}{2} \log_2(2\pi e E_g)$$

Since the above entropies should be equal, we have

$$\gamma_{s,\infty} \equiv \frac{E_u}{E_g} = \frac{\pi e}{6} \approx 1.53 \text{ dB}$$

The quantity $\gamma_{s,\infty}$ is called the *ultimate shaping gain*, which is achieved for a continuous Gaussian pdf.

5.3 Lattice Constellations

It is known from Shannon's capacity theorem that an optimal block code for a bandwidth-limited AWGN channel consists of a dense packing of code points within a sphere in a high-dimensional Euclidean space. Most of the densest known packings are lattices.

An n -dimensional (n -D) lattice Λ is an infinite discrete set of points (vectors, n -tuples) in the real Euclidean n -space \mathbb{R}^n that has the group property.

Example 5.3.1: The set of all integers, $\mathbb{Z} = \{0, \pm 1, \pm 2, \dots\}$, is a one-dimensional lattice, since \mathbb{Z} is a discrete subgroup of \mathbb{R} . The set \mathbb{Z}^2 of all integer-valued two-tuples (n_1, n_2) with $n_i \in \mathbb{Z}$ is a 2-dimensional lattice. More generally, the set \mathbb{Z}^n of all integer-valued n -tuples is an n -D lattice.

The lattice $R\mathbb{Z}^2$, where $R = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$, is obtained by rotating \mathbb{Z}^2 by $\pi/4$ and scaling by $\sqrt{2}$. Clearly, $R^2\mathbb{Z}^2 = 2\mathbb{Z}^2$.

Definition 1: Let $\{\mathbf{g}_j, 1 \leq j \leq m\}$ be a set of linearly independent vectors in \mathbb{R}^n (so that $m \leq n$). The set of points

$$\Lambda = \left\{ \mathbf{x} = \sum_{j=1}^m a_j \mathbf{g}_j \mid a_j \in \mathbb{Z} \right\} \quad (5.1)$$

is called an m -dimensional *lattice*, and $\{\mathbf{g}_j, 1 \leq j \leq m\}$ is called a *basis* of the lattice.

That is, Λ 是基向量的整数线性组合。The matrix with \mathbf{g}_j as rows

$$G = \begin{bmatrix} \mathbf{g}_1 \\ \mathbf{g}_2 \\ \vdots \\ \mathbf{g}_m \end{bmatrix}$$

is called a *generator matrix* for the lattice. 在后续讨论中, we will deal with full-rank lattices, i.e., $m=n$. So a general n -D lattice that spans \mathbb{R}^n may be expressed as

$$\Lambda = \left\{ \mathbf{x} = \mathbf{a}G \mid \mathbf{a} \in \mathbb{Z}^n \right\} \quad (5.1)$$

例如, the lattice \mathbb{Z}^2 has the generator $G = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$.

- A *coset* of a lattice Λ , denoted by $\Lambda + \mathbf{x}$, is a set of all points obtained by adding a fixed point \mathbf{x} to all lattice points $\mathbf{a} \in \Lambda$. Geometrically, the coset $\Lambda + \mathbf{x}$ is a translate of Λ by \mathbf{x} . If

$\mathbf{x} \in \Lambda$, by the group property, then $\Lambda + \mathbf{x} = \Lambda$. This implies that a lattice is “geometrically uniform;” every point of the lattice has the same number of neighbors at each distance, and all decision regions of a minimum distance decoder (“Voronoi regions”) are congruent and form a tessellation of \mathbb{R}^n .

- A *sublattice* Λ' of a given lattice Λ is a subset of the points in Λ that is itself a lattice. The set of all cosets of a sublattice is denoted by Λ/Λ' and is called a partition of Λ . In other words,

$$\Lambda = \Lambda' \cup [\Lambda/\Lambda'] = \Lambda' \cup (\Lambda' + \mathbf{x}) \quad (5.2)$$

where \mathbf{x} is chosen such that $(\Lambda' + \mathbf{x}) \in \Lambda$. (There are q cosets in a q -ary partition.) For

example, $\mathbb{Z}^2 = R\mathbb{Z}^2 + \{(0,0), (0,1)\}$.

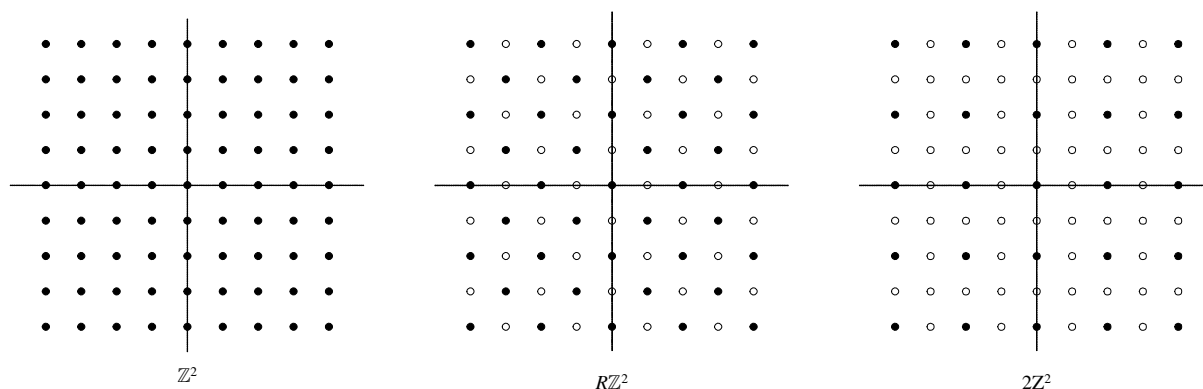


Fig. 5.10 Illustration of the binary partition chain $\mathbb{Z}^2 / R\mathbb{Z}^2 / 2\mathbb{Z}^2$

- The nearest neighbor quantizer $Q_\Lambda(\cdot)$ is defined by

$$Q_\Lambda(\mathbf{y}) = \mathbf{x} \in \Lambda, \quad \text{if } \|\mathbf{y} - \mathbf{x}\| \leq \|\mathbf{y} - \mathbf{x}'\|, \quad \forall \mathbf{x}' \in \Lambda$$

The *fundamental Voronoi region* of Λ is the set of points in \mathbb{R}^n closest to the zero codeword; i.e.,

$$\mathcal{V}_0 = \{\mathbf{y} \in \mathbb{R}^n \mid Q_\Lambda(\mathbf{y}) = \mathbf{0}\}$$

The *Voronoi region* associated with $\mathbf{x} \in \Lambda$ is the set of points \mathbf{y} such that $Q_\Lambda(\mathbf{y}) = \mathbf{x}$, and is given by a shift of \mathcal{V}_0 by \mathbf{x} . Note that other fundamental regions exist.

- A *fundamental parallelotope* of the lattice is the parallelotope (超平行体) that consists of the points $\{\mathbf{a}G \mid \mathbf{a} \in [0,1)^n\}$.

A fundamental parallelotope is an example of a fundamental region for the lattice; i.e., a building block which when repeated many times fills the whole space with just one lattice

point in each copy. Fig. 5.11 shows that around each lattice point is a region known as the fundamental parallelotope (用阴影表示).

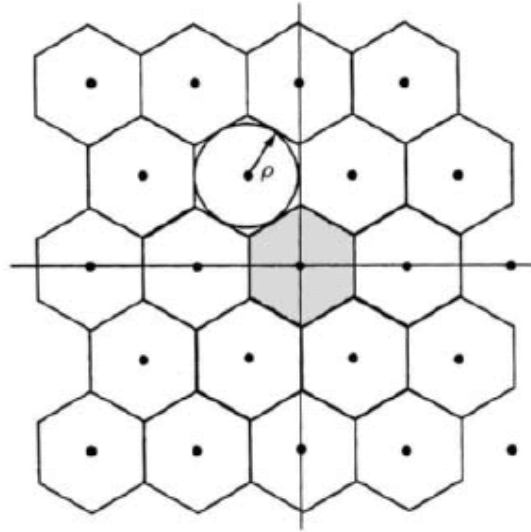


Fig. 5.11 The fundamental parallelotopes around the lattice points.

The key geometrical parameters of a lattice Λ are:

- (a) the *minimum squared Euclidean distance* $d_{\min}^2(\Lambda)$ between lattice points;
- (b) the *kissing number* $K_{\min}(\Lambda)$, which is the number of nearest neighbors to any lattice point;
- (c) the *fundamental volume* $V(\Lambda)$, which is the volume of the n -space corresponding to each lattice point. As indicated in Fig. 5.11, this volume is the volume of the fundamental region. Let $A = GG^T$. It can be shown that $V(\Lambda) \equiv |\det(A)|^{1/2} = \det(G)$ for any generator matrix G of Λ .

These parameters will directly affect the performance of a lattice constellation (lattice code). A normalized density parameter

$$\gamma_c(\Lambda) = \frac{d_{\min}^2(\Lambda)}{V(\Lambda)^{2/n}}$$

will be identified as the nominal coding gain.

Definition 2: A lattice constellation

$$C(\Lambda, \mathcal{R}) = (\Lambda + \mathbf{t}) \cap \mathcal{R}$$

is the finite set of points in a lattice translate $\Lambda + \mathbf{t}$ that lie within a compact bounding region \mathcal{R} of n -space.

(Note: A lattice is constrained to have a point at zero. The translate vector \mathbf{t} make the resulting constellation has no point at zero. The intersection of Λ with the region \mathcal{R} results in a finite

number of points in the constellation.)

Example 5.3.2: An M -PAM constellation $\{\pm 1, \pm 3, \dots, \pm(M-1)\}$ is a one-dimensional lattice constellation $C(2\mathbb{Z}, \mathcal{R})$ with $\Lambda+t=2\mathbb{Z}+1$ and $\mathcal{R}=[-M, M]$.

Some 2D lattice constellations are shown in Fig. 12.

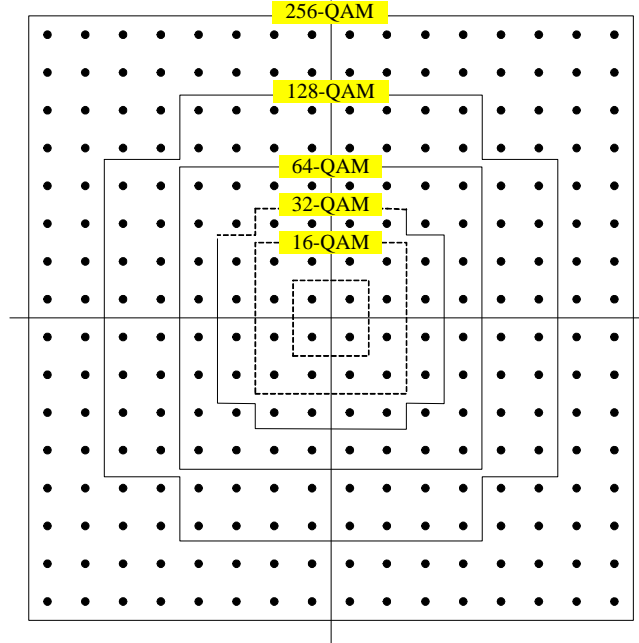


Fig.12 Two-dimensional constellations based on the integer lattice \mathbb{Z}^2 .

The key geometric properties of the region \mathcal{R} are

(a) its volume $V(\mathcal{R}) = \int_{\mathcal{R}} d\mathbf{x}$;

(b) the average energy $P(\mathcal{R})$ per dimension of a uniform probability density function over \mathcal{R} :

$$P(\mathcal{R}) = \int_{\mathcal{R}} \frac{\|\mathbf{x}\|^2}{n} \frac{d\mathbf{x}}{V(\mathcal{R})} \quad (5.3)$$

For performance analysis of large lattice constellations, we use the following approximations (Forney call this the *continuous approximation*).

■ *The continuous approximation:*

(a) The *size* of the constellation $C(\Lambda, \mathcal{R})$ (i.e., the number of signal points in $C(\Lambda, \mathcal{R})$) is well approximated by $V(\mathcal{R})/V(\Lambda)$.

(b) When the number of points in $C(\Lambda, \mathcal{R})$ is large, a uniform discrete distribution of the points over $C(\Lambda, \mathcal{R})$ is well approximated by a uniform continuous distribution over

the region \mathcal{R} . Thus, the *average energy per dimension* of a signal constellation is

$$P(C(\Lambda, \mathcal{R})) \approx P(\mathcal{R})$$

(c) The *average number of nearest neighbors* to any point in $C(\Lambda, \mathcal{R})$ is $\approx K_{\min}(\Lambda)$.

Example 5.3.3: For $\mathcal{R}=[-M, M]$, the parameters are $V(\mathcal{R}) = 2M$, $P(\mathcal{R}) = M^2/3$.

When $n=2$, and the shaping gain is an $2r \times 2r$ square, we have $V(\mathcal{R})=(2r)^2$, and

$$= \frac{1}{2} \times \frac{1}{4r^2} \int_{-r}^r \int_{-r}^r (x_1^2 + x_2^2) dx_1 dx_2 = \frac{r^2}{3}. \text{ When the shaping region } \mathcal{R} \text{ is a circle with radius } r,$$

$V(\mathcal{R}) = \pi r^2$. If $\mathbf{X}=(X_1, X_2)$ is uniformly distributed over the circle, then $P(\mathcal{R}) = E[\|\mathbf{X}\|^2]/2$

$$= E[X_1^2], \text{ where } E[X_1^2] = \frac{1}{\pi r^2} \int_{-r}^r x_1^2 \int_{-\sqrt{r^2-x_1^2}}^{\sqrt{r^2-x_1^2}} dx_2 dx_1 = \frac{r^2}{4}.$$

5.3.1 Shaping and Coding Gain

Under the continuous approximation, the coding gain for a lattice code is separable into two parts:

- 1) The *fundamental coding gain* $\gamma_c(\Lambda)$, which depends on the relative spacing of points in Λ , but is independent of \mathcal{R} .
- 2) The *shaping gain* $\gamma_s(\mathcal{R})$, which is determined by the choice of the signal constellation bounding region \mathcal{R} .

Consider the probability of decoding error for a lattice constellation Λ used over an AWGN channel. The union bound estimate (UBE) is

$$P(e) \approx K_{\min}(\Lambda) Q\left(\sqrt{\frac{d_{\min}^2(\Lambda)}{4\sigma^2}}\right)$$

where σ^2 is the noise power per dimension.

$$\text{Since } \frac{d_{\min}^2(\Lambda)}{4\sigma^2} = \frac{d_{\min}^2(\Lambda)(2^\rho - 1)}{4P(C(\Lambda, \mathcal{R}))} \cdot \frac{P(C(\Lambda, \mathcal{R}))}{\sigma^2(2^\rho - 1)} = \gamma \cdot SNR_{\text{norm}}, \text{ where } \gamma = \frac{d_{\min}^2(\Lambda)(2^\rho - 1)}{4P(C(\Lambda, \mathcal{R}))}$$

is a parameter of the constellation, the UBE can be written as

$$P(e) \approx K_{\min}(\Lambda) Q\left(\sqrt{\gamma \cdot SNR_{\text{norm}}}\right)$$

Using the continuous approximation, we have

$$\rho = \frac{2}{n} \log_2 |C(\Lambda, \mathcal{R})| \approx \frac{2}{n} \log_2 \frac{V(\mathcal{R})}{V(\Lambda)}$$

$$SNR_{\text{norm}} = \frac{P(C(\Lambda, \mathcal{R}))}{\sigma^2(2^\rho - 1)} \approx \frac{V(\Lambda)^{2/n}}{V(\mathcal{R})^{2/n}} \cdot \frac{P(\mathcal{R})}{\sigma^2}$$

$$\gamma \approx \frac{d_{\min}^2(\Lambda) \left(\frac{V(\mathcal{R})}{V(\Lambda)} \right)^{2/n}}{4P(\mathcal{R})} = 3\gamma_c(\Lambda)\gamma_s(\mathcal{R})$$

where

$$\gamma_c(\Lambda) = \frac{d_{\min}^2(\Lambda)}{V(\Lambda)^{2/n}} \quad (5.4)$$

is defined as the *nominal (or asymptotic) coding gain* of Λ , and

$$\gamma_s(\mathcal{R}) = \frac{V(\mathcal{R})^{2/n}}{12P(\mathcal{R})} \quad (5.5)$$

is defined as the *shaping gain* of \mathcal{R} .

The nominal coding gain $\gamma_c(\Lambda)$ measures the increase in density of Λ over the baseline integer lattice \mathbb{Z}^n . The shaping gain $\gamma_s(\mathcal{R})$ measures the decrease in average energy of \mathcal{R} relative to an n -cube $[-b, b]^n$. It can be shown that, given any lattice constellation C , the nominal coding and shaping gains of any K -fold Cartesian product constellation C^K is the same as those of C .

Example 5.3.4: For the square constellation, $n=2$, shaping region \mathcal{R} is a $2r \times 2r$ square. Then

$$\gamma_s(\mathcal{R}) = \frac{V(\mathcal{R})}{12P(\mathcal{R})} = \frac{4r^2}{12 \times r^2/3} = 1. \text{ Since } V(\Lambda) = d_{\min}^2(\Lambda), \text{ we have } \gamma_c(\Lambda)=1, \text{ and so } \gamma=3.$$

From the discussion above, the probability of block decoding error per two dimensions is

$$P_s(e) \approx K_s(\Lambda) Q\left(\sqrt{\gamma_c(\Lambda)\gamma_s(\mathcal{R}) \cdot 3SNR_{\text{norm}}}\right) \quad (5.6)$$

where $K_s(\Lambda) = 2K_{\min}(\Lambda)/n$ is the normalized error coefficient per two dimensions.

Note that the nominal coding gain is based solely on the argument of the $Q(\cdot)$ function in the UBE, which becomes infinite for dense n -dimensional lattices as $n \rightarrow \infty$. On the other hand, as before, the effective coding gain is limited by the number of nearest neighbors; i.e., the error coefficient $K_{\min}(\Lambda)$, which becomes very large for high-dimensional dense lattices. In fact, the Shannon limit shows that no lattice can have a combined effective coding gain and shaping gain greater than 9dB at $P(e)=10e-6$. This limits the maximum possible effective

coding gain to 7.5dB, because shaping gain can contribute up to 1.53dB (which will be discussed in the next subsection).

5.3.2 Maximum Shaping Gain

The maximum shaping gain is achieved by a spherical multidimensional constellation. Let $S_n(r)$ be an n -sphere of radius r , and let \mathbf{x} be an n -vector with real-valued components. Then

$$S_n(r) = \{\mathbf{x} : \|\mathbf{x}\| \leq r\} = \left\{ \mathbf{x} \left| \sum_{i=1}^n x_i^2 \leq r^2 \right. \right\}$$

and the volume is

$$V[S_n(r)] = \int_{S_n(r)} d\mathbf{x}.$$

Changing the variable of integration to $r' = \mathbf{x}/r$, this volume can be expressed as

$$V[S_n(r)] = V[S_n(1)]r^n.$$

Suppose that \mathbf{X} is a random vector uniformly distributed over $S_n(r)$. The pdf of \mathbf{X} is $p_{\mathbf{x}}(\mathbf{x}) = 1/V[S_n(r)]$, for $\mathbf{x} \in S_n(r)$ and zero elsewhere. The marginal density of one component of \mathbf{X} , say X_1 , is

$$p_{X_1}(x_1) = \frac{V[S_{n-1}(\sqrt{r^2 - x_1^2})]}{V[S_n(r)]} = \frac{V[S_{n-1}(1)]}{V[S_n(1)]} \cdot \frac{1}{r} \cdot \left[1 - \left(\frac{x_1}{r} \right)^2 \right]^{(n-1)/2}$$

Since $\int p_{X_1}(x_1) dx_1 = 1$, it follows that

$$\frac{V[S_{n-1}(1)]}{V[S_n(1)]} = \int_{-1}^1 (1 - \beta^2)^{\frac{n-1}{2}} d\beta$$

where $\beta = x_1/r$. For even n , we have

$$V[S_2(1)] = \pi, \quad \frac{V[S_n(1)]}{V[S_{n-2}(1)]} = \frac{2\pi}{n}, \quad V[S_n(1)] = \frac{\pi^{n/2}}{(n/2)!}$$

The average energy per dimension can be determined by

$$P[S_n(r)] = E[\|\mathbf{X}\|^2]/n = E[X_1^2] = \frac{r^2 V[S_{n-1}(1)]}{V[S_n(1)]} \int_{-1}^1 \beta^2 (1 - \beta^2)^{\frac{n-1}{2}} d\beta$$

For integer n ,

$$P[S_n(r)] = \frac{r^2}{n+2}.$$

Thus the shaping gain of an n -sphere is

$$\gamma_{S_n(r)} = \frac{\pi(n+2)}{12[(n/2)!]^{2/n}} \quad (5.7)$$

By Stirling's approximation, $m! \approx (m/e)^m$ as $m \rightarrow \infty$, we obtain that

$$\gamma_{S_n(r)} \rightarrow \frac{\pi e}{6} \text{ (1.53dB) as } n \rightarrow \infty$$

which is called the *ultimate shaping gain*.

Fig. 13 shows the shaping gain of an n -sphere for dimensions $n \leq 512$. Note that the shaping gain of a 16-sphere is about 1dB.

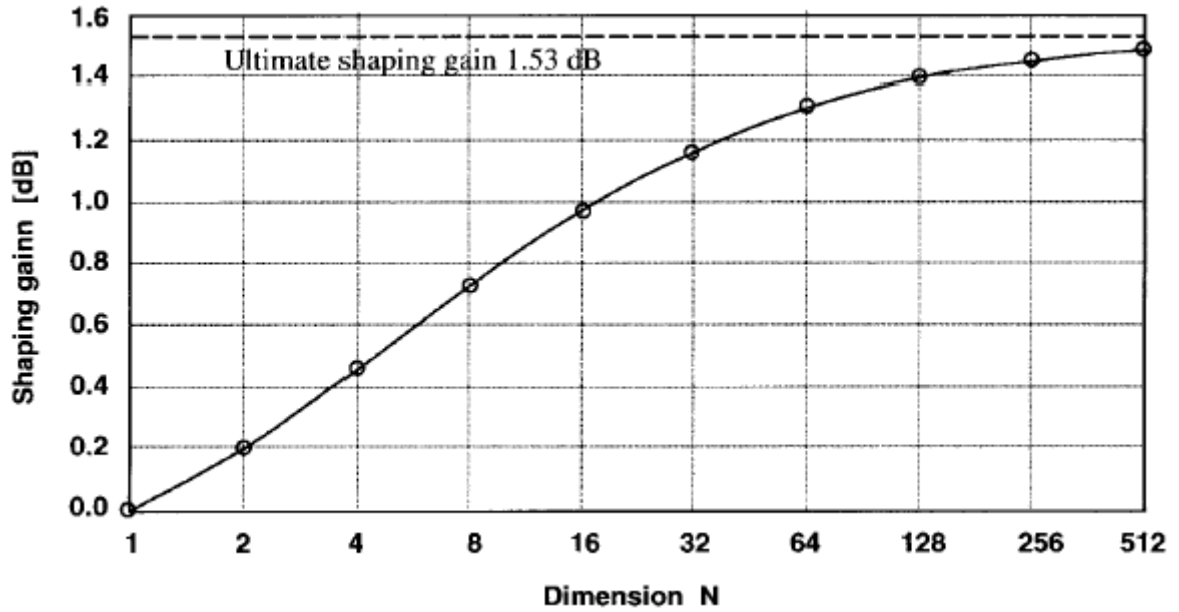


Fig. 13 Shaping gain of n -spheres

5.3.3 Marginal Density for Spherical Shaping

We now show that the projection of a uniform probability distribution over an n -sphere onto one or two dimensions is a nonuniform probability distribution that approaches a Gaussian distribution as $n \rightarrow \infty$.

By defining a normalized unit-variance random variable $Y_1 = X_1 \sqrt{n+2}/r$, we can see that the density of Y_1 is given by

$$p_{Y_1}(y_1) = A \cdot \left[1 - \left(\frac{y_1}{\sqrt{n+2}} \right)^2 \right]^{(n-1)/2}$$

where A is a constant. As $n \rightarrow \infty$,

$$p_{Y_1}(y_1) \rightarrow A \cdot \left[1 - \frac{y_1^2}{n} \right]^{n/2} \rightarrow A e^{-\frac{y_1^2}{2}}$$

which implies that X_1 approaches Gaussian. Furthermore, it can be shown that a

K -dimensional marginal density of an n -dimensional spherically uniform random vector approaches a joint Gaussian density with independent components for fixed K as $n \rightarrow \infty$. Fig.14 depicts the 1D marginal density for some even values of n .

5.3.4 Shaping Techniques

In principle, with spherical shaping, the lower-dimensional constellation will become arbitrarily large, even with fixed average power. However, in practice, the lower-dimensional constellation is constrained by design to a certain region \mathcal{R} to limit “shaping constellation expansion.” Thus, the resulting lower-dimensional probability distribution approaches a truncated Gaussian distribution within the region \mathcal{R} .

With large constellations, shaping can be implemented almost independently of coding by operations on the MSB of M -PAM or $(M \times M)$ -QAM constellation labels, which affect the gross shape of the n -D constellation. In contrast, coding affects the LSB and determines fine structure.

Two commonly used schemes are *trellis shaping* and *shell mapping*, both of which can easily obtain shaping gains of 1dB or more while limiting 2D shaping constellation expansion to a factor 1.5 or less. For details, refer to [Forney92] and [Khandani93].

The V.34 modem uses 16D shell mapping and obtains shaping gains of the order of 0.8dB with 2D shaping constellation expansion limited to 25% [Forney96].

5.3.5 Important Lattices

Table 5.1 identifies some interesting higher dimensional lattices and their parameters. All of these are sublattices of \mathbf{Z}^N and can thus be based on a processor designed for ordinary rectangular N -D QAM

Table 5.1 Important coding lattices and their parameters [3]

Lattice	Common name	Kissing number	Fundamental volume	Lattice coding gain
\mathbf{Z}^N	Cubic	$2N$	1	1
\mathbf{D}_N	Checkerboard	$2N(N-1)$	$2^{1-N/2}$	$2^{1-2/N}$
\mathbf{A}_2	2D Sphere packing (hexagonal)	6	$\sqrt{3}/2$	1.155
\mathbf{E}_8	Gosset	240	1/16	2
Λ_{16}	Barnes-Wall	4320	2.33×10^{-4}	2.829
Λ_{24}	Leetch	196560	5.96×10^{-8}	4
\mathbf{D}_{32}	Barnes-Wall	208320		4

Note: 1) $\mathbf{D}_2 = \mathbf{RZ}^2$, where $\mathbf{R} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}$ is a 2-D rotation operation. A lattice with generator

\mathbf{RG} is a rotation CW by $\pi/4$ of the lattice with \mathbf{G} , combined with an expansion by $\sqrt{2}$.

2) The checkerboard lattice \mathbf{D}_N in N dimensions can be defined as those points in \mathbf{Z}^N for

which $\sum_{j=1}^N i_j$ is even.

5.4 Trellis-Coded Modulation (TCM)

Lattice codes are useful for a modest number of dimensions. For large dimensionality, their implementation complexity becomes excessive. The approach using an FSM at the transmitter (as illustrated in Fig. 5.1.6) provides an alternative one with lower complexity. Significant coding gains can also be achieved by using this method (possibly in conjunction with a multidimensional constellation). TCM codes (which also are simply called trellis codes) are a class of Euclidean space codes based on an FSM. Its encoder structure is similar to Fig. 5.1.6 with a binary convolutional encoder as the FSM.

TCM was proposed by Ungerboeck in 1982 []. The essential new concept of TCM was to use signal-set expansion to provide redundancy for coding, and to design coding and signal-mapping functions jointly so as to maximize directly the “free distance” between coded signal sequences. The term *trellis-coded modulation* originates from the fact that these coded sequences consist of modulated symbols rather than binary digits. In other words, in TCM schemes, the trellis branches are labeled with redundant nonbinary modulated symbols rather than with binary coded symbols.

5.4.1 Notation and Definitions

For a trellis code \mathcal{C} (of length n), the minimum squared Euclidean distance between two different sequences of signal points is referred to as its *free squared Euclidean distance*; i.e.,

$$d_{\text{free}}^2(\mathcal{C}) = \min_{m \neq m'} \|\mathbf{x}_m - \mathbf{x}_{m'}\|^2, \quad \mathbf{x}_m, \mathbf{x}_{m'} \in \mathcal{C} \subseteq \mathbb{R}^n \text{ (or } \mathbb{C}^n)$$

The *asymptotic coding gain* (including shaping gain) is defined to be

$$\gamma \equiv 10 \log_{10} \left(\frac{d_{\text{free}}^2(\mathcal{C}) / E_s}{d_{E,\text{min}}^{2(u)} / E_s^{(u)}} \right) \text{ dB} \quad (5.1)$$

where $d_{E,\text{min}}^{2(u)}$ denote the minimum squared Euclidean distance between signal points in the uncoded scheme, and E_s and $E_s^{(u)}$ denote the average signal energies of the coded and uncoded schemes, respectively.

Eq. (5.1) can be rewritten as

$$\gamma = 10 \log_{10} \left(\frac{E_s^{(u)}}{E_s} \cdot \frac{d_{\text{free}}^2(C)}{d_{E,\min}^{2(u)}} \right) = 10 \log_{10} (\gamma_c^{-1} \cdot \gamma_d)$$

where $\gamma_c = E_s / E_s^{(u)}$ is the constellation expansion factor, and γ_d is the distance gain factor.

Ungerboeck introduced “mapping by set partitioning” to break down a signal constellation into a suitable number of subsets such that in each subset the signal points are farther apart. Assume that we partition the original constellation \mathcal{A} into L subsets, $\{\mathcal{A}_0, \mathcal{A}_1, \dots, \mathcal{A}_{L-1}\}$, of signal points. Then the intra-subset minimum squared Euclidean distance is given by

$$d_{\min}^2(\mathcal{A}_l) \equiv \min_{s, s' \in \mathcal{A}_l} \{d_E^2(s, s')\}, \quad 0 \leq l \leq L-1$$

where $d_E^2(s, s')$ is the squared Euclidean distance between the signal points s and s' within the subset \mathcal{A}_l . The smallest squared Euclidean distance between two different sequences of signal points for which the subset label sequences (i.e., FSM output) are the same is equal to the minimum one of all $d_{\min}^2(\mathcal{A}_l)$'s. Let $d_{\min}^2(\mathcal{S})$ denote this squared Euclidean distance; that is

$$d_{\min}^2(\mathcal{S}) \equiv \min_{0 \leq l \leq L-1} \{d_{\min}^2(\mathcal{A}_l)\}.$$

Let

$$d_{\min}^2(\mathcal{S}, \mathcal{S}') \equiv \min_{\substack{s \in \mathcal{S} \\ s' \in \mathcal{S}'}} \{d_E^2(s, s')\}$$

denote the smallest squared Euclidean distance between a signal point in subset \mathcal{S} and a signal point in subset \mathcal{S}' , where $\mathcal{S}, \mathcal{S}' \in \{\mathcal{A}_0, \mathcal{A}_1, \dots, \mathcal{A}_{L-1}\}$. Denote by $d_{\text{free}}^2(\mathcal{S})$ the minimum squared Euclidean distance between any two different sequences of subsets (i.e., sequences whose elements are the subsets); then

$$d_{\text{free}}^2(\mathcal{S}) \equiv \min_{\mathcal{S} \neq \mathcal{S}'} \left\{ \sum_{i=1}^{\infty} d_{\min}^2(\mathcal{S}_i, \mathcal{S}'_i) \right\}$$

where $\mathcal{S} = (\mathcal{S}_0, \mathcal{S}_1, \dots)$ and $\mathcal{S}' = (\mathcal{S}'_0, \mathcal{S}'_1, \dots)$ are two infinite code sequence stemming from the same state and whose elements are the subsets $\{\mathcal{A}_0, \mathcal{A}_1, \dots, \mathcal{A}_{L-1}\}$.

Thus, we have

$$d_{\text{free}}^2(C) = \min \{d_{\text{free}}^2(\mathcal{S}), d_{\min}^2(\mathcal{S})\} \quad (5.2)$$

If a lattice constellation is used in TCM codes, we will refer to the resulting code \mathcal{C} as the lattice-type trellis code. Also, we will use the notation $d_{\text{free}}^2(\Lambda')$ and $d_{\text{min}}^2(\Lambda')$ instead of $d_{\text{free}}^2(\mathcal{S})$ and $d_{\text{min}}^2(\mathcal{S})$, respectively.

5.4.2 Design of TCM Schemes

From (5.2), we can see that maximizing the minimum Euclidean distance between sequences of signal points requires to maximize both intra-subset and inter-sequence-of-subsets Euclidean distances. Usually, this is done by following Ungerboeck's principle of *mapping by set partitioning*.

5.4.2.1 Mapping by Set Partitioning

The mapping by set partitioning rule is based on successive partitioning of the original constellation into subsets with increased intrasubset minimum squared distances. Usually, set partitioning of an N -D constellation \mathcal{A} into 2^n subsets (of equal size) is performed by n partitioning steps, with each two-way selection being identified by a label bit $c^{(j)} \in \{0,1\}, 1 \leq j \leq n$; i.e.,

$$A \xrightarrow{c^{(1)}} B(c^{(1)}) \xrightarrow{c^{(2)}} C(c^{(2)}c^{(1)}) \xrightarrow{c^{(3)}} \dots \xrightarrow{c^{(n)}} S(c^{(n)} \dots c^{(1)})$$

$$\text{with } \Delta_0^2 \leq \Delta_1^2 \leq \Delta_2^2 \leq \dots \leq \Delta_n^2$$

where $\Delta_j^2, j=0,1,\dots,n$ denote the intrasubset minimum squared distances of the j th level subsets A, B, \dots , of the partitioning chain. The above set partitioning is continued until the minimum squared distances between signal points at the final level of partitioning, Δ_n^2 , is at least as large as the desired minimum squared Euclidean distance $d_{\text{free}}^2(\mathcal{C})$ of the TCM code to be designed. The binary n -tuples $\mathbf{c} = [c^{(n)} \dots c^{(1)}]$ are called the subset labels of the final subsets S .

As an example, Fig. 5.4.2 shows the set partition for 8-PSK.

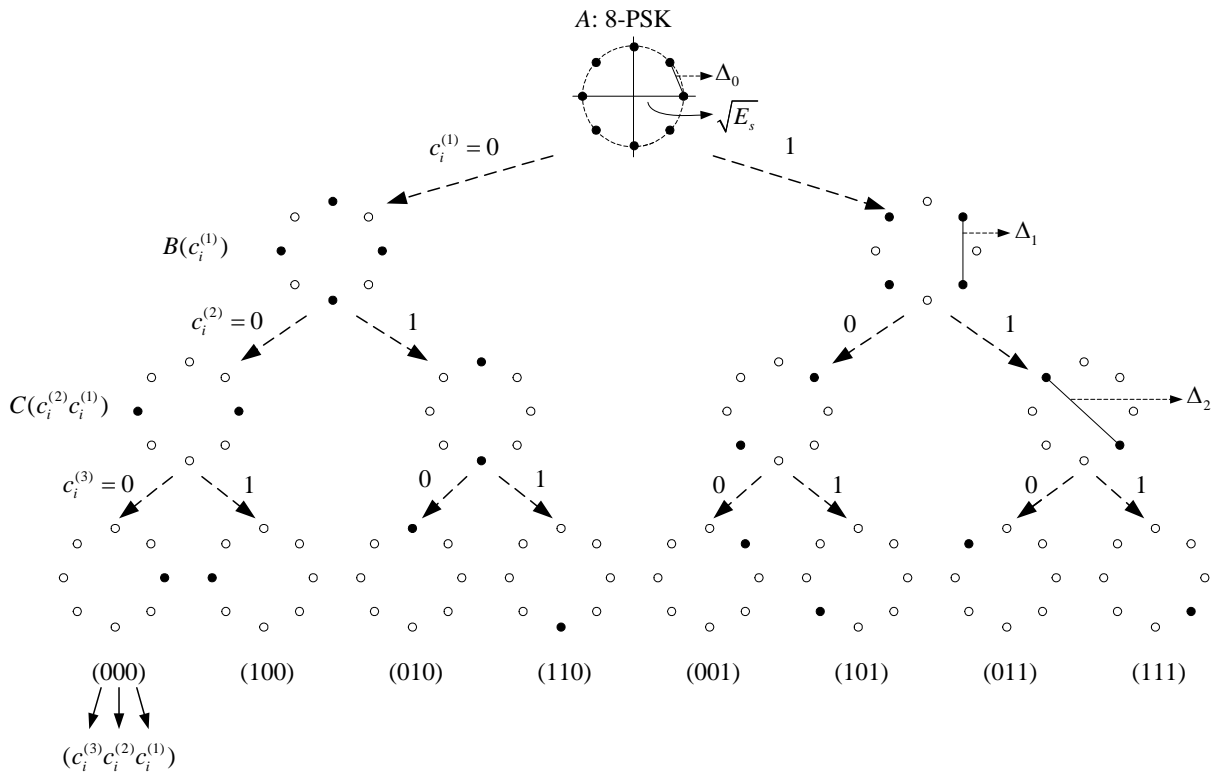


Figure 5.4.2 Set partitioning of an 8-PSK constellation

5.4.2.2 Structure of a TCM Encoder

With the above set partitioning, the intrasubset minimum Euclidean distance is maximized. In the following we will discuss the assignment of signal subsets to trellis branches such that the minimum Euclidean distance between any two different sequences of subsets is maximized. Before this, we first describe the principle of a TCM encoder in more detail.

A general encoder for a trellis code \mathcal{C} is depicted in Fig.5.4.3. At each time instant i , an input symbol consisting of m information bits enters the TCM encoder. Based on the sufficient intrasubset Euclidean distance at a certain level of the set-partitioning chain, we encode k out of m input bits and leave $(m-k)$ bits uncoded. The k input bits are encoded by a rate- $k/(k+r)$ binary convolutional encoder with 2^v states into a subset label \mathbf{c}_i consisting of $k+r$ coded bits. The label \mathbf{c}_i is then used to select signal subset $S(\mathbf{c}_i)$. The remaining $m-k$ input bits are used to select one signal x_i from 2^{m-k} signals in the subset $S(\mathbf{c}_i)$. (If there is any shaping, it is done at this level.) The size of constellation \mathcal{A} is therefore 2^{m+r} .

In practice, the rate of the convolutional code is always chosen to be $k/(k+1)$, i.e., only one redundancy bit per N -D symbol (assuming that an N -D constellation is used). So the coding constellation expansion factor is 2 per N -dimension. The least significant label bit $c_i^{(1)}$ is then the sole parity-check bit.

(Note: Suppose that a constellation of size M is used for an uncoded system. From the analysis of constrained-capacity for AWGN channels with discrete-input, we can see that most of coding gain can be achieved by using a constellation of size $2M$ plus a channel coding

algorithm. It is not greatly advantageous to employ a constellation of size greater than $2M$. This implies one redundancy bit is sufficient.)

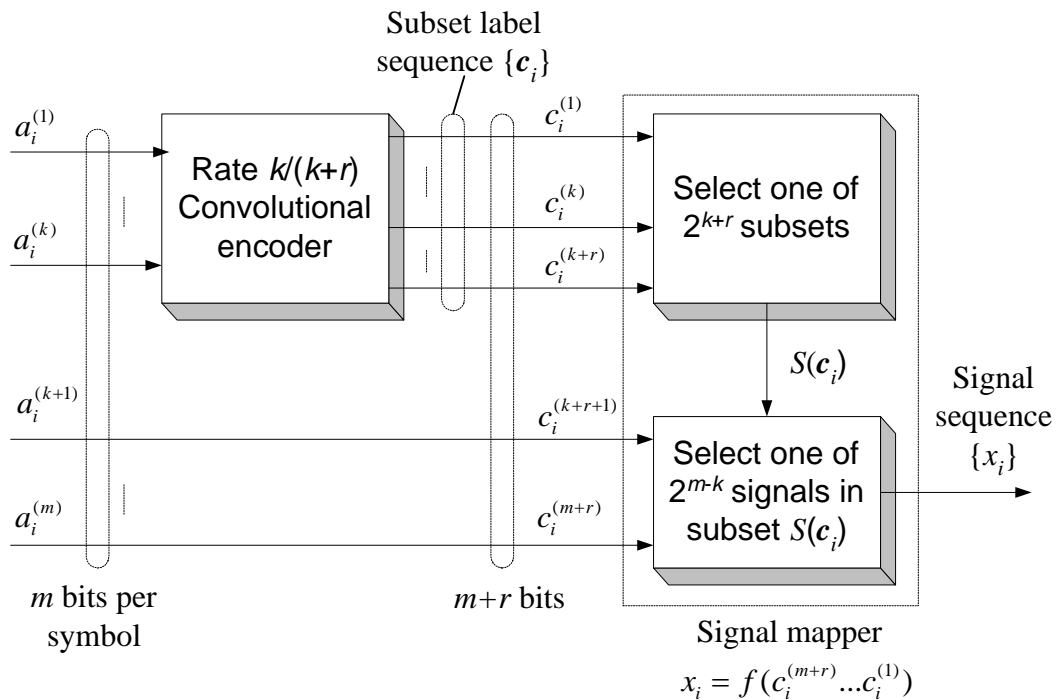


Figure 5.4.3 General structure of a TCM encoder (记号 a_i 改为 u_i)

In this way, all branches (state transitions) in the trellis diagram for a TCM code may be labeled with modulated symbols. The number of transitions between two consecutive states depends on the number of uncoded bits. For a TCM encoder with $m > k$, there are 2^{m-k} parallel transitions between states, which are associated with the 2^{m-k} signals of the subsets at the final level of partitioning. The minimum squared Euclidean distance between parallel transitions is $d_{\min}^2(\mathcal{S}) = \Delta_n^2$.

Fig. 5.4.4 depicts the encoder of a rate-2/3 8-PSK TCM code.

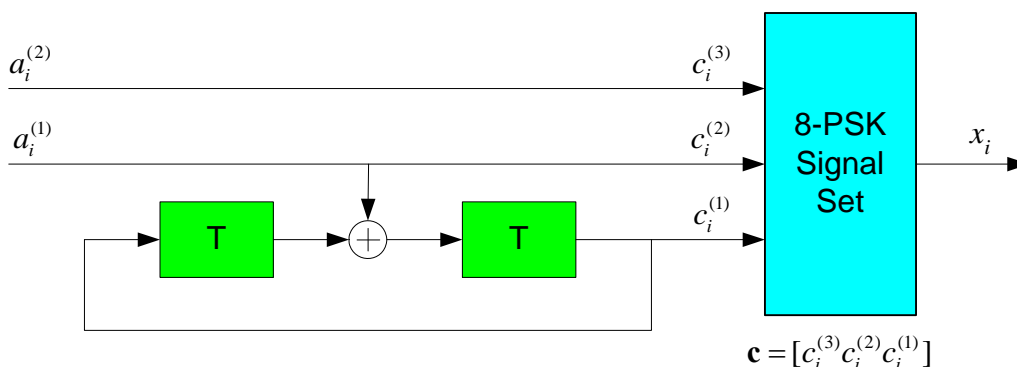


Figure 5.4.4 The 4-state TCM encoder for 8-PSK

- 关于采用 N -D constellation 的进一步讨论
如果使用 N 维信号星座，则每星座点是一个 N 维实向量（称为 N 维调制符号）。传

输时，我们将其作为包含 $L=N/2$ 个 QAM 符号的复符号序列进行信号调制及在信道上传输，如图 5.4.5 所示。这样，对于 TCM 编码器来说，每一个时钟节拍产生 L 个已调(QAM)符号。令 T_s 为发送一个 QAM 符号的时间间隔，则卷积编码器的工作时钟周期是 $T=L \cdot T_s$ ；在每一个工作节拍， m 个信息比特输入 TCM 编码器，产生 L 个发送符号。因此，系统的谱效率是 $\rho = \frac{m}{L} = \frac{2m}{N}$ bits per channel use. 图 5.4.5 给出了一个 N 维 TCM 编码器的示意图。

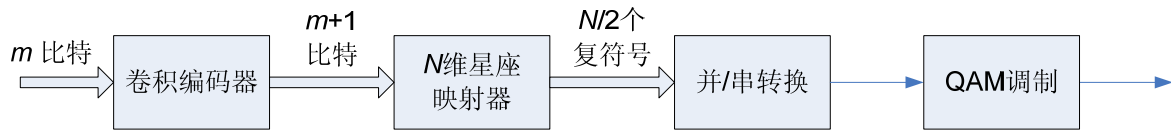


Figure 5.4.5

5.4.2.3 Ungerboeck TCM Design Rules

We now proceed to consider the assignment of signal subsets to trellis transitions. The design criterion for TCM codes is based on optimizing the Euclidean distance spectrum. Specifically, the convolutional code and the labeling of the subsets are chosen primarily to maximize the minimum squared Euclidean distance $d_{\text{free}}^2(\mathcal{S})$ between sequences of signal points in any possible encoded subset sequence, and secondly to minimize the maximum possible number $K_{\min}(C)$ of nearest-neighbor sequences.

The following heuristic design rules were introduced by Ungerboeck to maximize the free squared Euclidean distance $d_{\text{free}}^2(C)$ of the TCM code C .

- **Rule 1:** Parallel transitions, if present, are associated with signals of the subsets at the final level of partitioning chain; i.e., signals within $S(c_i)$. These signals have minimum squared Euclidean distance $d_{\min}^2(\mathcal{S}) = \Delta_n^2$.
- **Rule 2:** The transitions originating from or merging into one state are associated with signals of subsets at the first level of set partitioning chain; i.e., $B(c^{(1)})$.
- **Rule 3:** All signals are used with equal frequency in the trellis diagram.

As an example, consider the 4-state 8-PSK TCM in Fig. 5.4.4, and the set partitioning in Fig. 5.4.2. The corresponding trellis diagram is shown in Fig. 5.4.6, where the signal points are represented by the decimal expressions of their binary labels; i.e., the signals with the

labels $[c_i^{(3)}c_i^{(2)}c_i^{(1)}]$ are denoted by $b = \sum_{j=1}^3 2^{j-1} c_i^{(j)}$. The parallel transitions are associated with

signals from one of the four subsets, $C(00)$, $C(01)$, $C(10)$, $C(11)$, with minimum squared Euclidean distance $\Delta_2^2 = 4E_s$. The signals associated with transitions diverging from or

merging into one state are from $B(0)$ and $B(1)$, with minimum squared Euclidean distance $\Delta_1^2 = 2E_s$.

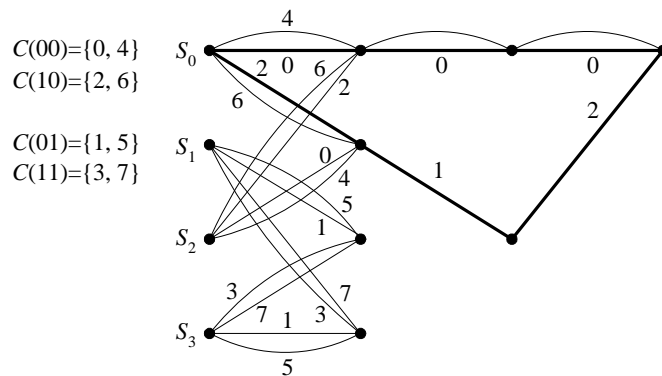


Fig. 5.4.6 Trellis diagram of the TCM scheme shown in Fig. 5.4.4.

The error event corresponding to $d_{\text{free}}^2(\mathcal{S})$ is also shown in Fig. 5.4.6 by bold lines. We can see that

$$d_{\text{free}}^2(\mathcal{S}) = d_E^2(0,2) + d_E^2(0,1) + d_E^2(0,2) = 4.586E_s$$

Thus, the free squared Euclidean distance of this TCM code is

$$d_{\text{free}}^2(C) = \min\{d_{\text{free}}^2(\mathcal{S}), \Delta_2^2\} = 4E_s.$$

Compared with an uncoded QPSK scheme with the minimum squared Euclidean distance $2E_s$ between signal points, this TCM scheme can provide an asymptotic coding gain of

$$\gamma = 10 \log \frac{4}{2} = 3 \text{ (dB)}$$

5.4.2.4 Design Examples

Example 1 (8-state 8-PSK code): From the above example, we can see that the free distance of the designed TCM code is limited to Δ_2^2 by the parallel transitions. As an alternative example, we consider the design of an 8-state TCM scheme for 8-PSK with throughput of 2 bits/symbol. For obvious reasons, the parallel transitions should be avoided in the trellis diagram of this code. A suitable trellis diagram is shown in Fig. 5.4.7. The trellis branches are associated with a partitioned 8-PSK signal set (see Fig. 5.10) according to Ungerboeck's design rules. Observe that the signals from subset $B(0)$ and $B(1)$, at the first level of set partitioning, are assigned to the transitions originating from the states with even and odd subscripts, respectively, in Fig. 5.4.7. This assignment guarantees that the signals merging into one state are selected from either $B(0)$ or $B(1)$.

The squared free Euclidean distance of this scheme is given by

$$d_{\text{free}}^2(C) = d_{\text{free}}^2(\mathcal{S}) = \Delta_1^2 + \Delta_0^2 + \Delta_1^2 = 4.586E_s$$

which, when compared with an uncoded QPSK with the minimum squared Euclidean distance

$2E_s$ between signal points, provides

$$\gamma = 10 \log \frac{4.586}{2} = 3.6 \text{ (dB)}$$

asymptotic coding gain. Thus, the 8-state 8-PSK TCM code can provide 0.6dB more asymptotic coding gain than 4-state 8-PSK code at the cost of increased complexity.

Two equivalent realizations of the encoder of this TCM scheme are shown in Fig. 5.4.8.

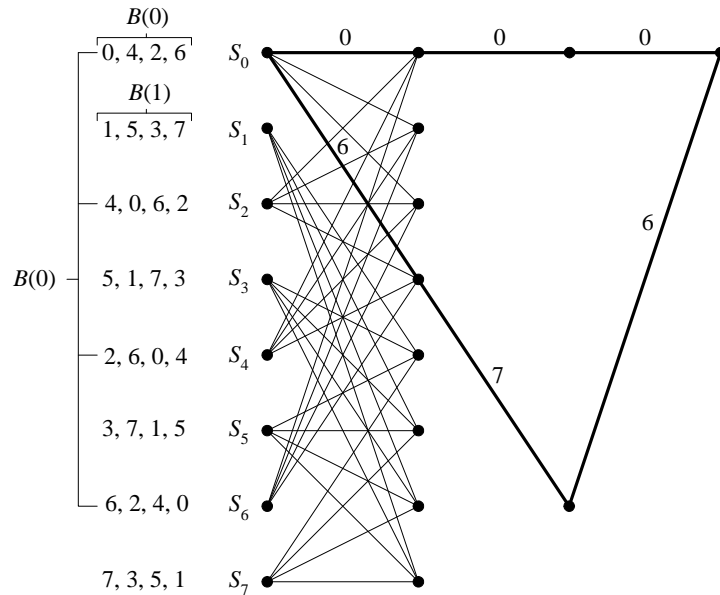
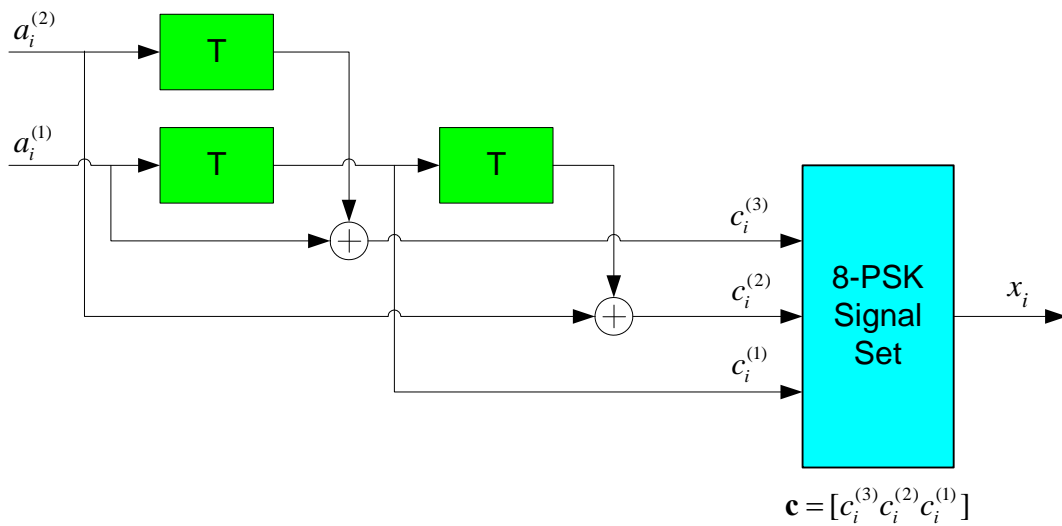
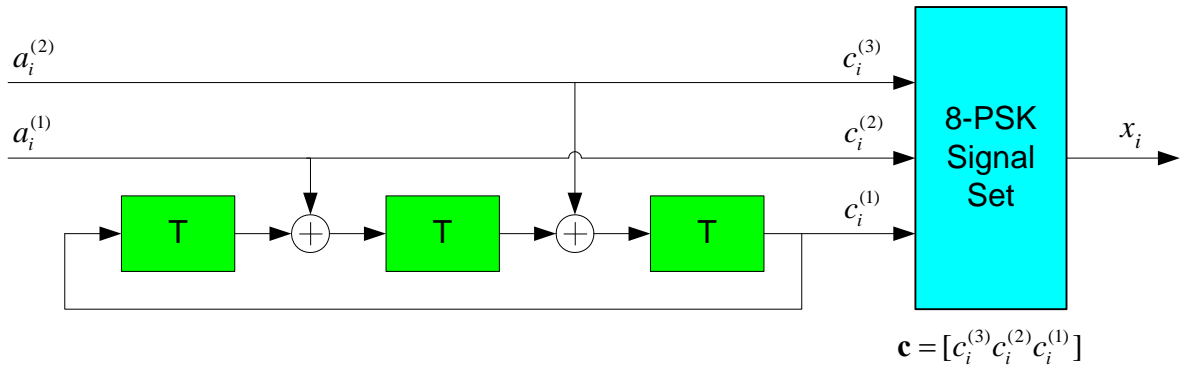


Fig. 5.4.7 Trellis diagram of the 8-state TCM scheme for 8-PSK with throughput of 2 bits/signal.



(a)



(b)

Fig. 5.4.8 Two equivalent realizations of the 8-state TCM scheme for 8-PSK. (a) Feedforward encoder. b) Systematic encoder with feedback.

Example 2 (8-state TCM code for 16-QAM): A 16-QAM signal constellation is shown in Fig. 5.4.9, with the squared Euclidean distance Δ_0^2 between adjacent signal points. For the 8-state 16-QAM trellis code with throughput of 3 bits/symbol, we choose $k=2$ and use a convolutional code of rate $2/3$. The encoder structure is shown in Fig.5.4.10, and the set partitioning is illustrated in Fig.5.4.11. The corresponding trellis diagram is shown in Fig. 5.4.12. The squared free Euclidean distance of this code is given by

$$d_{\text{free}}^2(C) = \min\{\Delta_1^2 + \Delta_0^2 + \Delta_1^2, \Delta_3^2\} = \min\{5\Delta_0^2, 8\Delta_0^2\} = 2E_s$$

Clearly, it is limited by the minimum squared Euclidean distance between sequences of subsets (rather than the parallel paths). The asymptotic coding gain (compared to an uncoded 8-PSK) is equal to

$$\gamma = 10 \log \frac{2}{2 - \sqrt{2}} = 5.3 \text{ (dB)}$$

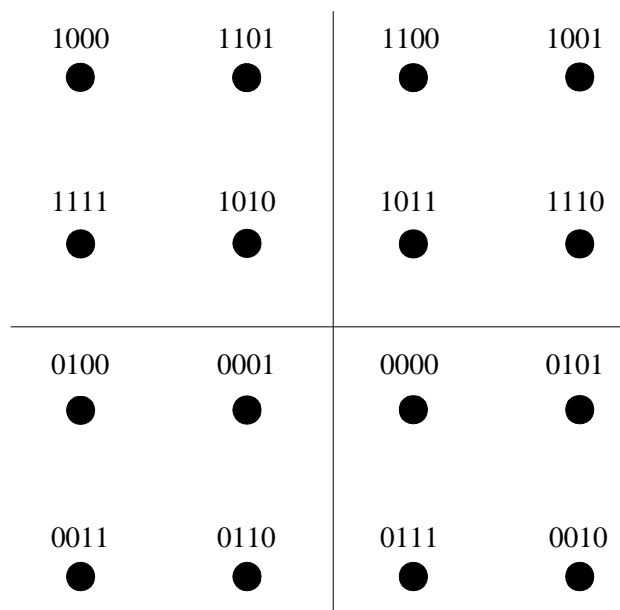
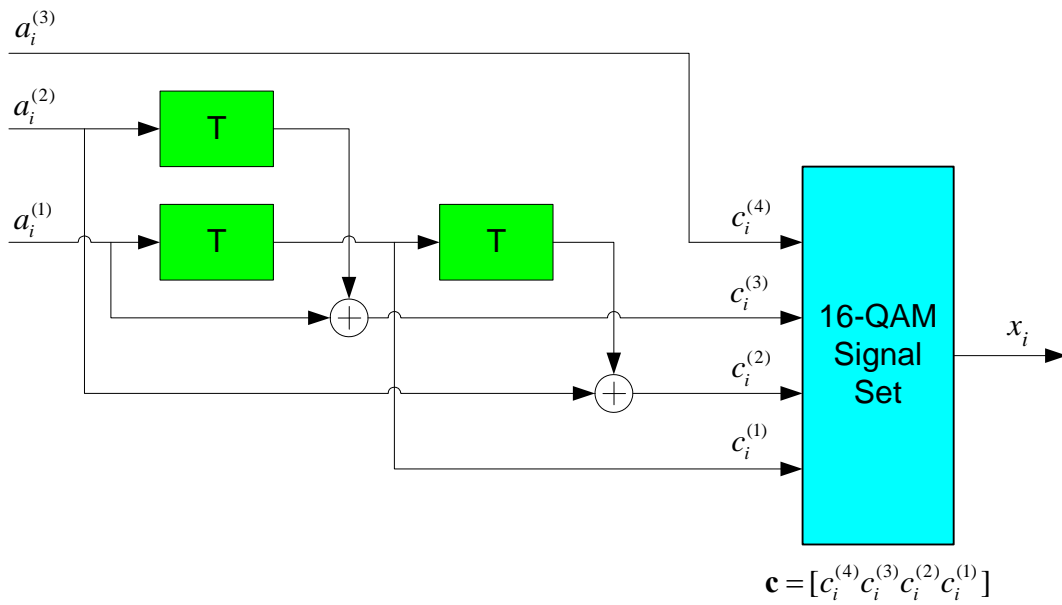
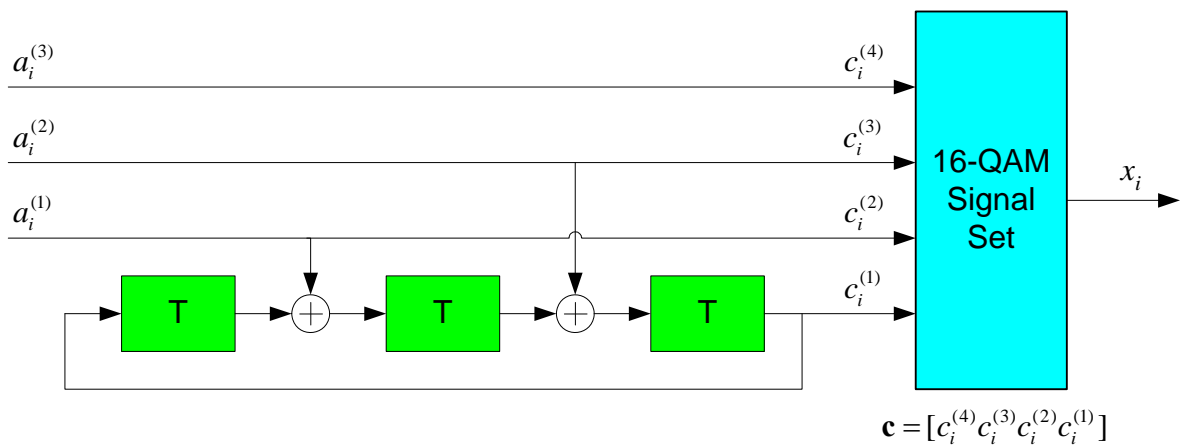


Fig. 5.4.9 A 16-QAM constellation



(a) Feedforward encoder (in controller form)



(b) Systematic encoder with feedback in observer form

Fig. 5.4.10 Two equivalent realizations of the encoder of the 8-state 16-QAM TCM code.

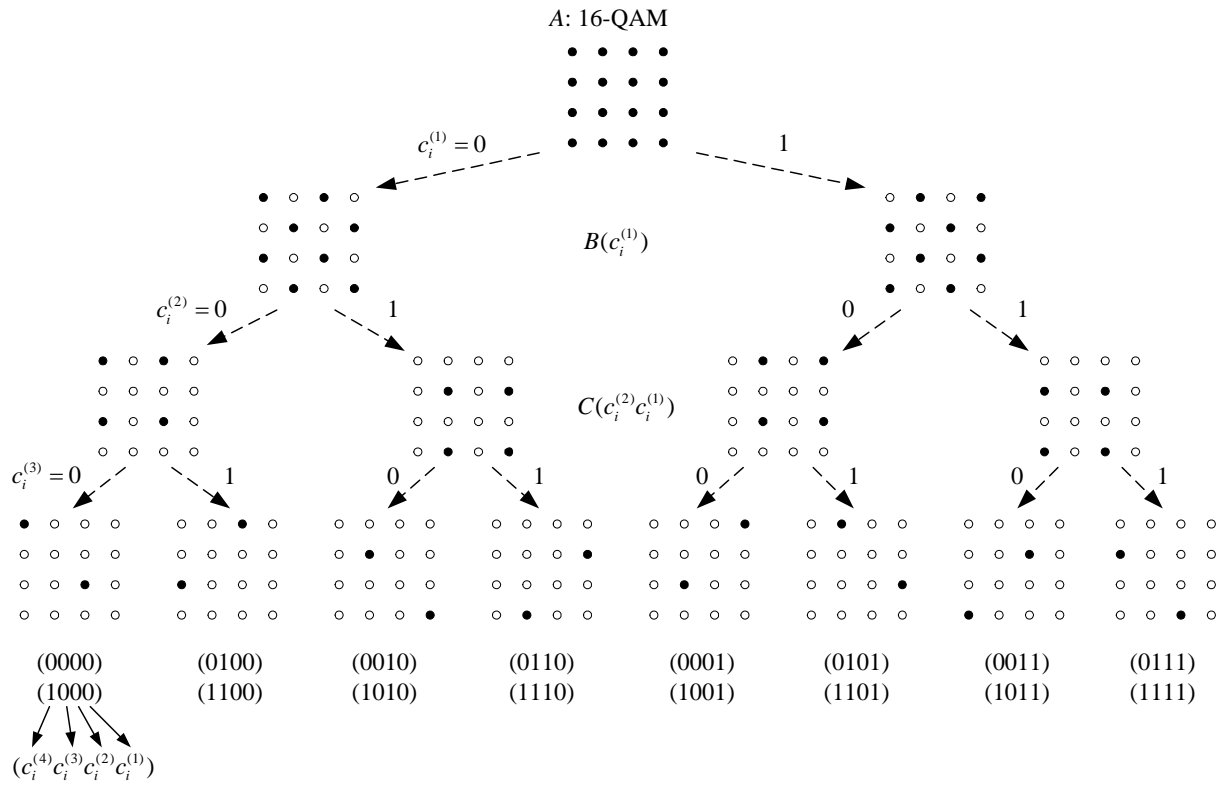


Fig. 5.4.11 Set partitioning of a 16-QAM constellation

In [Unger87], Ungerboeck provided some TCM codes with the number of states less than 256, which are results of a code-search program.

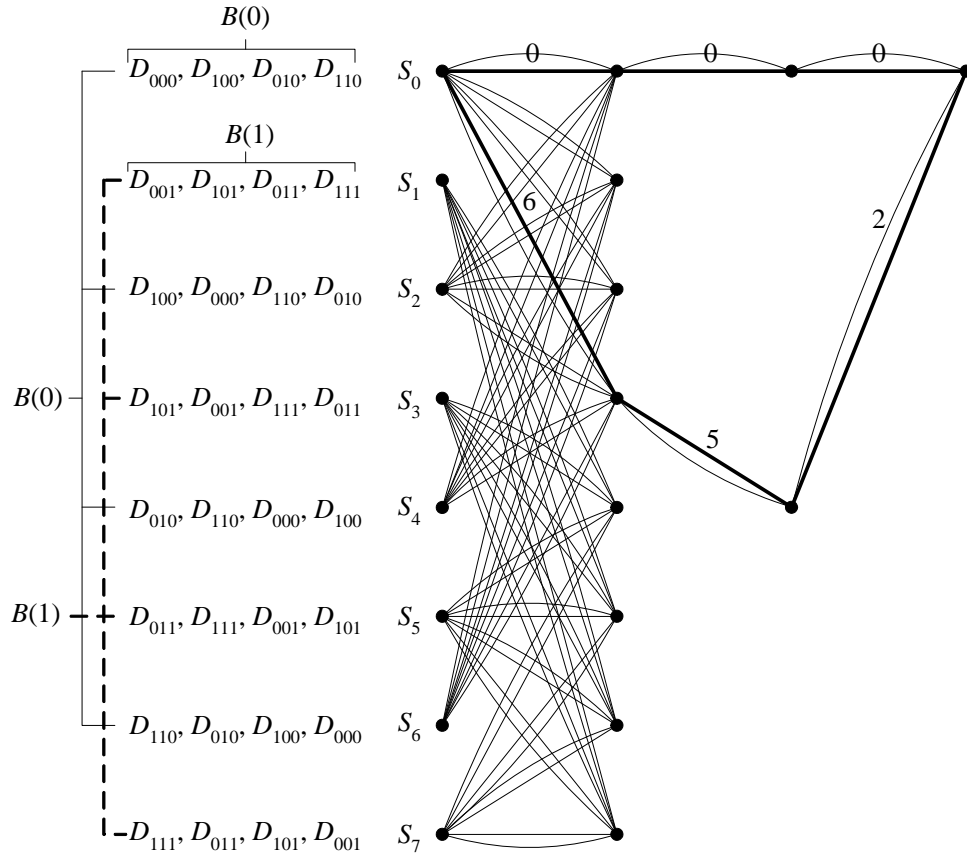


Fig. 5.4.12 Trellis diagram of the 8-state TCM scheme for 16-QAM

■ TCM Codes for Higher-Order Constellations

The above discussed method can be directly extended to higher-order constellations, such as 64-QAM. We now consider a low-complexity alternative.

如前所述，对于高阶调制（如16-QAM以上），由于工作信噪比较高，在子集划分链的某一层上，信号子集内欧氏距离已足够大，使得子集内信号点在没有编码保护的条件下，它们之间的错误概率已足够小（或满足目标错误概率），这样，我们就不需要再往下进行子集划分，从而对剩余的输入信息比特就不进行编码。那么如何确定子集划分到何时为止呢？前面给出了依据 $d_{\min}^2(\mathcal{S}) \geq d_{\text{free}}^2(\mathcal{C})$ 来确定划分层数的方法，下面介绍一种更为简单的近似判断方法，它是一种启发式规则[Robert98]。

该规则基于这样的经验：在AWGN信道中，TCM方案的错误率要达到 $\text{BER} \leq 10^{-5}$ ，至少需要 E_b/N_0 比相应谱效率的Shannon容量限高4dB左右，即

$$\left. \frac{E_b}{N_0} \right|_{\text{operating}} \geq \left. \frac{E_b}{N_0} \right|_{\text{capacity}} + 4\text{dB}$$

假定原始信号星座为 \mathcal{A} ，其信号平均能量为 E_s 。现在对星座 \mathcal{A} 进行子集划分。对每一层

子集划分，我们依据下面的公式近似计算一个子集内比特错误概率

$$P_b(\Delta_j) \approx Q\left(\sqrt{\frac{d_{E,\min}^2}{2N_0}}\right) = Q\left(\sqrt{\frac{\Delta_j^2}{2E_s} \cdot \left(\frac{E_s}{N_0}\right)}\right) = Q\left(\sqrt{\frac{\Delta_j^2 \rho}{2E_s} \cdot \left(\frac{E_b}{N_0}\right)_{\text{operating}}}\right) \quad (5.3)$$

其中 Δ_j^2 是信号星座 \mathcal{A} 的子集内最小欧氏距离， ρ 是频谱效率。这里，我们没有计入上一层子集对下一层子集的判决错误传播。

由式(5.3)即可决定 m 个输入比特中需要编码保护的比特数。假定在第 k 层子集划分中， $P_b(\Delta_k) < P_{b,\text{target}}$ ，则我们只需要对 k 个比特进行卷积编码，留下 $m-k$ 个 leave ($n - \tilde{n}$) bits uncoded for each information symbol. The coded bits are used to define the signal subsets and the uncoded bits are used to select signal points from a subset.

Example: Consider a CT-TCM code for 64-QAM with 5 bits/symbol. In this case, the channel capacity is $E_s/N_0 = 16.2\text{dB}$. According to [Robert98], the intra-subset error probability is $P_b \leq 10^{-6}$ after three levels of Ungerboeck-type set-partition. As a result, we adopt $\tilde{n} = 3$ and employ the previous design for 16-QAM. The constellation used is illustrated in Fig. 5.4.14. The performance of this code is shown in Fig. 5.4.15.

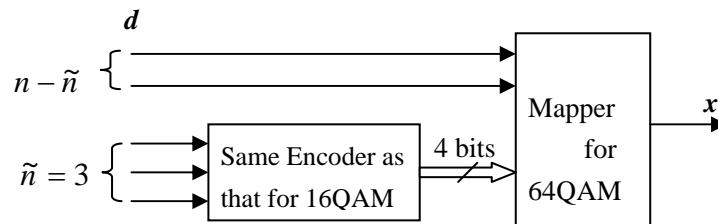


Fig. 5.4.13 Structure of the component encoder for 64-QAM

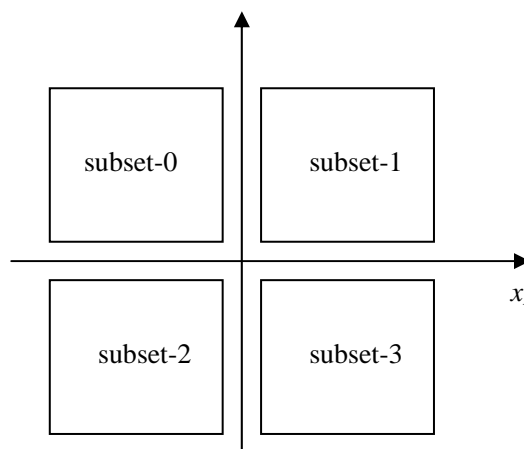


Fig. 5.4.14. A 64-QAM constellation. Each subset shown is a shifted version of the 16-QAM constellation (see Fig. 5.4.9). For each input symbol of 5 bits, the first $\tilde{n} = 3$ bits are used to select a point within a subset and the remaining $n - \tilde{n} = 2$ bits to select an appropriate subset.

Part. level	0	1	2	3	4	5	6
Δ_i	0.095	0.19	0.38	0.76	1.52	3.05	∞
$P_b(\Delta_i)$	0.06	0.013	$8 \cdot 10^{-4}$	$4 \cdot 10^{-6}$	$1.3 \cdot 10^{-10}$	$1.8 \cdot 10^{-19}$	-

5.4.2.5 Best TCM Schemes and the Realizations of TCM Encoders

TCM schemes with a small number of states can be constructed using the heuristic design rules aforementioned. To construct a larger TCM scheme, however, a computer code-search program is necessary. Table 5.2 provides the results of such a computer search [Unger82], where the TCM schemes are denoted by the coefficients of the parity-check polynomials in octal form. For example, for the 8-state 8-PSK TCM scheme,

$$\mathbf{H}(D) = [H_1(D) \ H_2(D) \ H_3(D)] \text{ and}$$

$$H_1(D) = D^2 \Rightarrow h_1 = (000 \ 100)_2 = (04)_8$$

$$H_2(D) = D \Rightarrow h_2 = (000 \ 010)_2 = (02)_8$$

$$H_3(D) = D^3 + 1 \Rightarrow h_3 = (001 \ 001)_2 = (11)_8$$

The sets $\{g_1, g_2, g_3\}$ are coefficients of the generator polynomials, which are related to $\{h_1, h_2, h_3\}$ by

$$\mathbf{G}(D) \cdot \mathbf{H}^T(D) = \mathbf{0}$$

where $\mathbf{G}(D) = [G_1(D) \ G_2(D) \ G_3(D)]$. The above relation yields that the sets g_1, g_2, g_3 are the respective h_l taken in forward order; i.e., if $h_l = (h_l^{(m)}, h_l^{(m-1)}, \dots, h_l^{(1)}, h_l^{(0)})$, then $g_l = (h_l^{(0)}, h_l^{(1)}, \dots, h_l^{(m-1)}, h_l^{(m)})$. For example,

$$h_1 = (000 \ 100)_2 \Rightarrow g_1 = (001 \ 000)_2 = (10)_8$$

$$h_2 = (000 \ 010)_2 \Rightarrow g_2 = (010 \ 000)_2 = (20)_8$$

$$h_3 = (001 \ 001)_2 \Rightarrow g_3 = (100 \ 100)_2 = (44)_8$$

Traditionally, there are two realizations for the encoder of TCM codes: in either observer canonical form with $\{h_l\}$ denoting a tap set, or controller canonical form with $\{g_l\}$ denoting such a tap set. In Table 5.2, 5.3 and 5.4, we give both these tap sets, and the coefficients h_1, h_2, h_3 are listed (in reverse order, 即连接多项式系数的低位在右, 高位在左) as right-justified octals; but the sets g_1, g_2, g_3 are listed (in forward order) as left-justified octals.

Table 5.2 Trellis-coded 8-PSK schemes []. Observer form taps \mathbf{h} shown as right-justified octals; controller form taps \mathbf{g} shown as left-justified octals. # indicates that free Euclidean distance is determined by parallel paths.

No. of states	No. of coded	h_0	h_1	h_2	g_1	g_2	g_3	Normalized free	Asy. Coding
---------------	--------------	-------	-------	-------	-------	-------	-------	-----------------	-------------

bits, k								distance	gain γ
								d_f^2 / E_s	(dB)
4	1	5	2	-	5	2	-	4.000#	3.01
8	2	11	02	04	44	20	10	4.586	3.60
16	2	23	04	16	62	10	34	5.172	4.13
32	2	45	16	34				5.758	4.59
64	2	103	030	066				6.343	5.01
128	2	277	054	122				6.586	5.17
256	2	435	072	130				7.515	5.75

Table 5.3 Trellis-coded 16-PSK schemes [1].

No. of states	No. of coded bits, k	h_1	h_2	h_3	g_1	g_2	g_3	Normalized free distance	Asy. Coding gain γ
								d_f^2 / E_s	(dB)
4	1	5	2	-	5	2	-	1.324	3.54
8	1	13	04	-	64	10	-	1.476	4.01
16	1	23	04	-	62	10	-	1.628	4.44
32	1	45	10	-	51	04	-	1.910	5.13
64	1	103	024	-	604	120	-	2.000#	5.33
128	1	024	203	-	120	604	-	2.000#	5.33
256	2	427	176	374	721	374	176	2.085	5.51

Table 5.4 Trellis-coded M -QAM schemes

states	coded bits, k	h_1	h_2	h_3	g_1	g_2	g_3	Asy. coding gain $[\gamma]$ (dB)				
								16Q	32CR	64Q	128CR	256Q
4	1	5	2	-	5	2	-	[4.4]	[3.0]	[2.8]	[3.1]	[2.9]
8	2	11	02	04	44	20	10	[5.3]	[4.0]	[3.8]	[4.1]	[3.8]
16	2	23	04	16	62	10	34	[6.1]	[4.8]	[4.6]	[4.9]	[4.6]
32	2	41	06	10	41	30	04	[6.1]	[4.8]	[4.6]	[4.9]	[4.6]
64	2	101	016	064	404	340	130	[6.8]	[5.4]	[5.2]	[5.5]	[5.3]
128#	2	203	014	042	606	140	210	[7.4]	[6.0]	[5.8]	[6.1]	[5.9]
256#	2	401	056	304	401	350	106	[7.4]	[6.0]	[5.8]	[6.1]	[5.9]
512#	2	1001	0346	0510	4004	3160	0450	[7.4]	[6.0]	[5.8]	[6.1]	[5.9]

In Tables 5.2, 5.3 and 5.4, only the convolutional codes (i.e., subset selector) of rates of 1/2 and 2/3 are listed. For the rectangular-grid QAM family, the square free distances of the codes for 16-QAM, ..., 256-QAM are shown against the same subset selector. This is due to the fact that all these constellations are the subsets of the same rectangular point pattern; they all have the same worst-case local neighbor structure after 1 or 2 splits. The generic encoder structure is shown in Fig. 5.4.16.

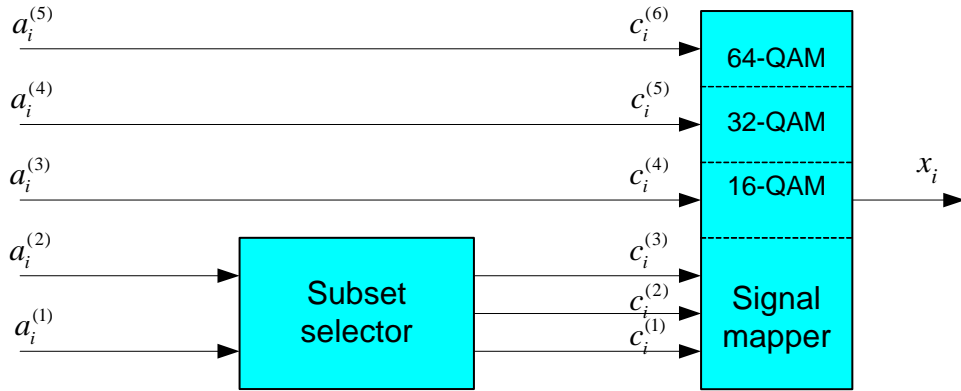
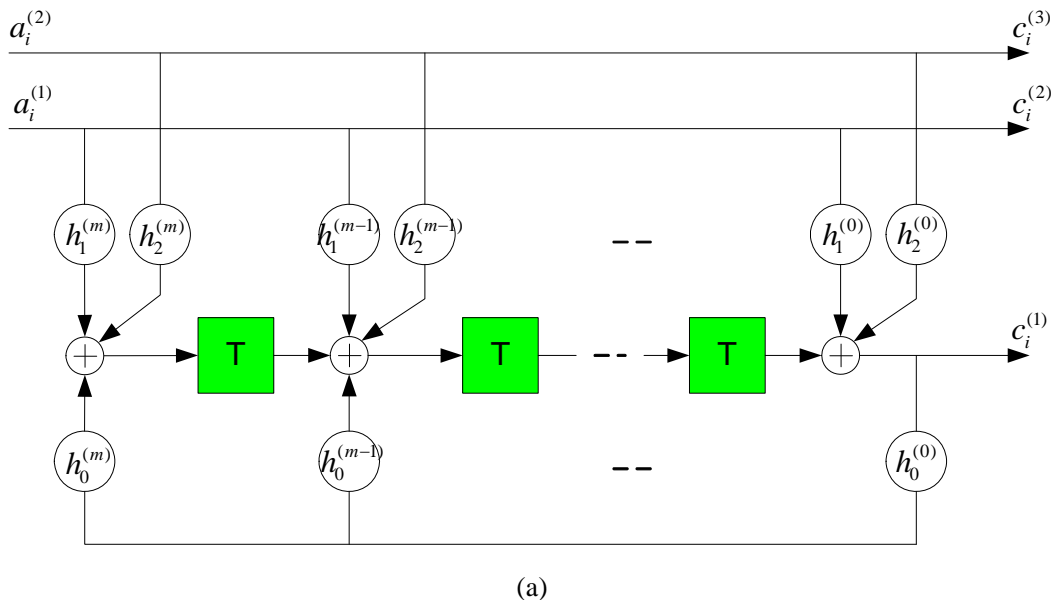


Figure 5.4.16 Generic encoder for QAM signal constellations.

■ Realizations of TCM encoders

The encoder of Ungerboeck codes may be realized in either observer canonical form or controller canonical form. In ordinary convolutional error correction, a (nonsystematic) feedforward convolutional encoder always exists that creates the same codeword set as a systematic feedback encoder. In TCM, these words map to subset sequences, resulting in the same subset sequences set, with the same inter-subset minimum distance. However, note that the mapping between information sequences and codewords is different in the two cases. In the following we show how to obtain an encoder realization if the coefficients of the parity-check polynomials are given.

The observer-form encoder may be obtained from the coefficients $\{h_1, h_2, h_3\}$ of the parity-check polynomials, as shown in Fig.5.4.17. It is straightforward to obtain the controller-form encoder from the generator polynomial matrix.



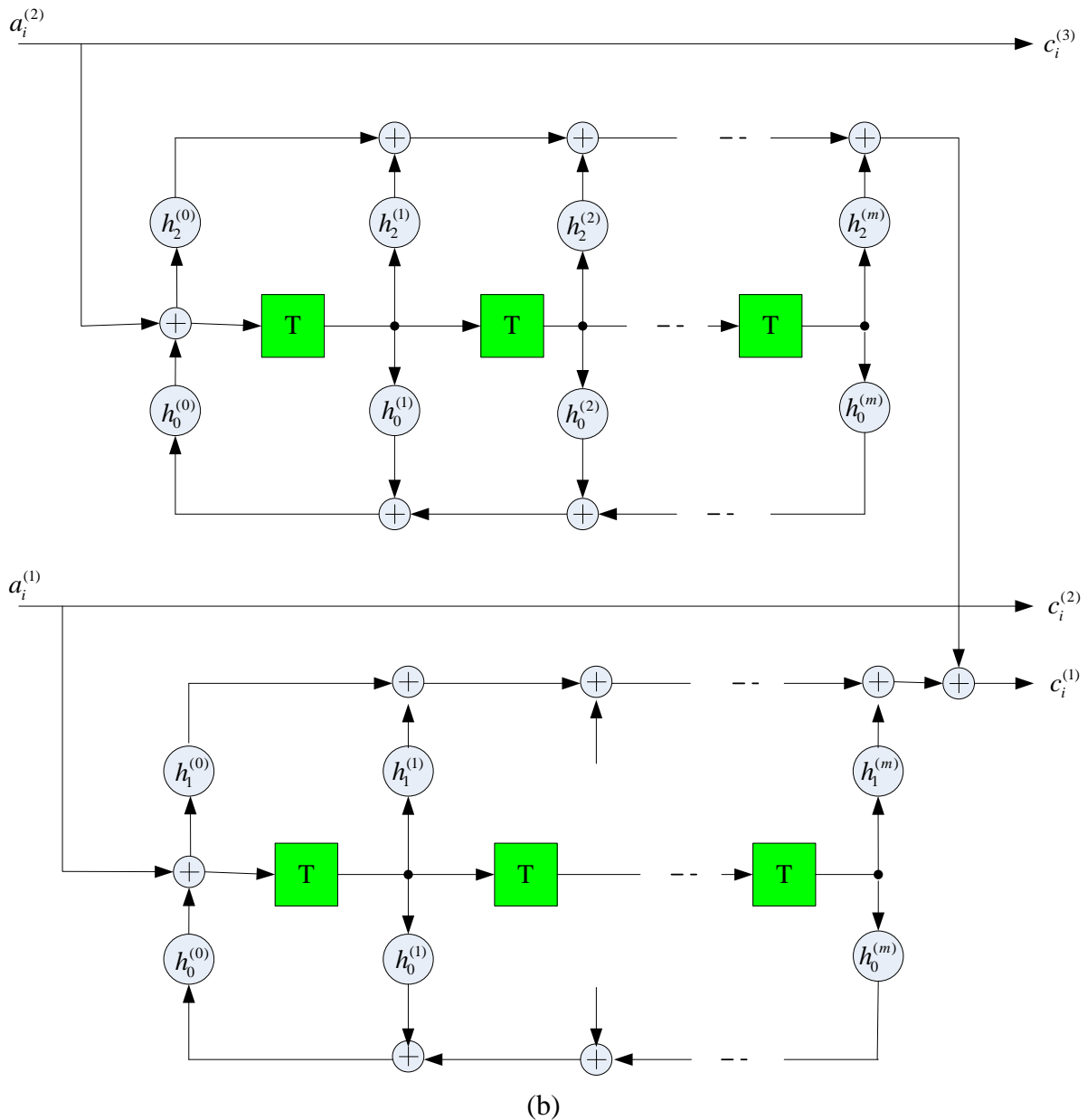


Fig.5.4.17 Realizations of a rate-2/3 systematic feedback convolutional encoder. (a) Observer canonical form. (b) Controller canonical form, the sets g_0, g_1, g_2 are the respective h_l in forward order.

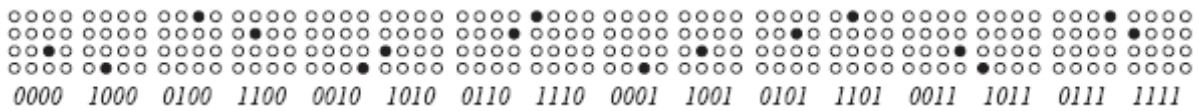
For observer canonical form encoder realizations, the lowest-degree (constant) terms in the generator (parity-check) polynomials represent the connections at the right ends of the shift registers, whereas the highest-degree terms represent the connections at the left ends of the shift registers. It is important to note that this is exactly the opposite of the correspondence between polynomial coefficients and delay elements in the case of a controller canonical form realizations. For this reason, when an encoder is realized in observer canonical form, it is common to write the generator (parity-check) polynomials in the opposite of the usual order; i.e., from highest degree to lowest degree. (See $\mathbf{H(D)}$ on last page)

下表给出了采用 rate-3/4 的卷积码、16-QAM 调制的 TCM 码的有关参数[Schlegel04],

其中 A_{free} is the average number of paths at distance d_{free}^2 , B_{free} is the average number of bit errors on those paths. 跟表 5.4 中的码相比（在下表中用”Ung”指示原来表 5.4 中的码），新的 TCM 码改进了错误系数（即减小了 A_{free} 和 B_{free} ）。

	2^v	$h^{(0)}$	$h^{(1)}$	$h^{(2)}$	$h^{(3)}$	d_{free}^2	$A_{d_{\text{free}}}$	$B_{d_{\text{free}}}$
Ung 8	8	11	2	4	0	5.0	3.656	18.313
	8	13	4	2	6	5.0	3.656	12.344
Ung 16	16	23	4	16	0	6.0	9.156	53.5
	16	25	12	6	14	6.0	9.156	37.594
Ung 32	32	41	6	10	0	6.0	2.641	16.063
	32	47	22	16	34	6.0	2	6
Ung 64	64	101	16	64	0	7.0	8.422	55.688
	64	117	26	74	52	7.0	5.078	21.688
Ung 128	128	203	14	42	0	8.0	36.16	277.367
	128	313	176	154	22	8.0	20.328	100.031
Ung 256	256	401	56	304	0	8.0	7.613	51.953
	256	417	266	40	226	8.0	3.273	16.391

Nonstandard mapping used for the improved 16-QAM codes above:



5.4.3 Decoding of TCM codes

A TCM code is described by using a trellis diagram, so it may be maximum-likelihood decoded using Viterbi algorithm. The metric is defined as the Euclidean distance between the coded sequence of signal points and the received sequence.

Specifically, given a received symbol $y_i \in \mathbb{R}^n$, the receiver first finds the closest signal point $\hat{x}_i(\mathbf{c})$ in the subset $S(\mathbf{c})$, and stores its metrics as the representative of parallel branches.

This is called *subset decoding*. A VA decoder then finds the code sequence $\{\mathbf{c}_i\}$ for which the signals chosen in the subsets are closest to the entire received sequence $\{y_i\}$. The decoding complexity is dominated by the complexity of the VA decoder.

5.5 Performance Evaluation of TCM

The performance evaluation of TCM schemes over AWGN channels involves the computation of several important parameters:

- Minimum Euclidean distance d_{free} ,
- Number of nearest neighbors,
- Distance spectrum,
- Error event probability, and

Bit error probability

The complexity of the algorithms for computing these parameters depends on the degree of symmetry of a code.

5.5.1 Symmetry Properties of TCM Schemes

Definition: An isometry T of \mathbb{R}^N is a transform $T: \mathbb{R}^N \rightarrow \mathbb{R}^N$ ($x \rightarrow T(x)$) that preserves Euclidean distances,

$$\|T(x) - T(y)\|^2 = \|x - y\|^2, \quad x, y \in \mathbb{R}^N$$

Translations, rotations and reflections are typical examples of isometries. In fact every isometry in \mathbb{R}^N can be decomposed into a sequence of translations, rotations, and reflections.

Definition: Two signal sets \mathcal{S} and \mathcal{S}' are said to be (geometrically) congruent if there exists an isometry T such that $T(\mathcal{S}) = \mathcal{S}'$. Furthermore, an isometry T with the property that $T(\mathcal{S}) = \mathcal{S}$ is called a *symmetry* of \mathcal{S} .

For example, for a QPSK signal set, the symmetries are the rotations of $0, \pi/2, \pi, 3\pi/2$, the two reflections by the main axes, etc.

Definition: A signal set \mathcal{S} is *geometrically uniform* if, for any two points $x_i, x_j \in \mathcal{S}$, there exists an isometry $T_{i \rightarrow j}$ that transforms x_i to x_j while \mathcal{S} invariant; i.e.,

$$\begin{aligned} T_{i \rightarrow j}(x_i) &= x_j \\ T_{i \rightarrow j}(\mathcal{S}) &= \mathcal{S} \end{aligned}$$

Hence, a signal set is geometrically uniform if we can take an arbitrary signal point and, through application of isometries $T_{i \rightarrow j}$, generate the entire signal constellation \mathcal{S} . Roughly speaking, a geometrically uniform signal set looks identical from every signal point, which includes M -PSK constellations and all the lattice-based constellations (if we disregard the restriction to a finite number of points) as typical examples.

If the squared Euclidean distance between two signal points depends only on the difference between their labels; i.e.,

$$d_E^2 = \|f(\mathbf{c}_i) - f(\mathbf{c}_j)\|^2 = \|f(\mathbf{c}_i \oplus \mathbf{c}_j) - f(\mathbf{0})\|^2 = \|f(\mathbf{e}_i) - f(\mathbf{0})\|^2$$

then the signal $f(\mathbf{0})$ can always be used as a reference signal. Signal mappings with this property are called *regular*.

■ Geometrically Uniform Codes

Denote by B_0 the set of 2^m vectors \mathbf{c} corresponding to the all-zero state of the

convolutional encoder. Every state in the trellis diagram has either B_0 or the coset $B_1 = B_0 + \mathbf{c}'$ associated with it. A sufficient condition for uniformity is that the following one-to-one correspondence

$$f(\mathbf{c}) \rightarrow f(\mathbf{c} \oplus \mathbf{c}'), \quad \mathbf{c} \in B_0 \quad (5.4)$$

be an isometry. Thus

$$\|f(\mathbf{c}) - f(\mathbf{c} \oplus \mathbf{e})\|^2 = \|f(\mathbf{c} \oplus \mathbf{c}') - f(\mathbf{c} \oplus \mathbf{c}' \oplus \mathbf{e})\|^2, \quad \text{for any } \mathbf{c} \in B_0 \text{ and } \mathbf{e}$$

Eq. (5.4) defines a correspondence between the two subsets of signals that are obtained in the first set-partition level. Thus, the condition for uniformity implies that these two subsets must be related by an isometry (e.g., a rotation or a reflection).

Next we generalize the above discussion by considering a signal set which is partitioned at the first level of the partitioning chain by the value of k th component, $c_i^{(k)}$, into two subsets $B_0^{(k)}$ and $B_1^{(k)}$.

Definition: A TCM scheme is called k -isometric if

- a) the subsets $B_0^{(k)}$ and $B_1^{(k)}$ are related by an isometry, and
- b) the labels of two corresponding signals $f(\mathbf{c})$ and $f(\tilde{\mathbf{c}})$ in the isometry, i.e., \mathbf{c} and $\tilde{\mathbf{c}}$, differ only in their k th component.

A k -isometric TCM code having the following properties is called *uniform*:

- a) The coded m -tuples formed from the label \mathbf{c}_i by deleting the k th component, i.e., $\mathbf{c}_i^* = (c_i^{(m+1)}, c_i^{(m)}, \dots, c_i^{(k+1)}, c_i^{(k-1)}, \dots, c_i^{(1)})$, are independent sequences.
- b) The sequence of $(m+1)$ -tuples \mathbf{c}_i are uniquely determined from the sequence of m -tuples \mathbf{c}_i^* .

In [Forney91], Forney defined geometrically uniform trellis codes. It follows that a geometrically uniform trellis code is regular. That is, the distance between any two sequences $d_E^2(f(\mathbf{c}), f(\mathbf{c}')) = d_E^2(f(\mathbf{c} \oplus \mathbf{c}'), f(\mathbf{0}))$, where \mathbf{c} and \mathbf{c}' are label sequences.

Theorem: The set of codewords of a geometrically uniform trellis code \mathcal{C} is geometrically uniform; that is, for any two sequences $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{C}$, there exists an isometry $T_{i \rightarrow j}$ that transforms \mathbf{x}_i to \mathbf{x}_j while \mathcal{C} invariant.

Thus, all code sequences in a geometrically uniform trellis code have the same error performance over AWGN channels. Some examples of geometrically uniform trellis codes include:

- 1) All convolutional codes are geometrically uniform.
- 2) Trellis codes using rate-1/2 convolutional codes and QPSK signal set with Gray mapping are geometrically uniform.

5.5.2 Error Probability of TCM Schemes

TCM 码的性能分析方法类似于卷积码，我们用 union bound 来表示有关错误概率。但是，由于 TCM 码一般不是线性的，the problem is complicated.

Consider the Ungerboeck model for rate- $m/(m+1)$ TCM schemes. The m input bits $\mathbf{a}_i = (a_i^{(m)}, \dots, a_i^{(2)}, a_i^{(1)})$ are first encoded convolutionally into a sequence of signal labels $\mathbf{c}_i = (c_i^{(m+1)}, \dots, c_i^{(2)}, c_i^{(1)})$, which may be linear operation. The label \mathbf{c}_i is then mapped to a $2^{(m+1)}$ -ary constellation, yielding the transmitted signal $x_i = f(\mathbf{c}_i)$, where $f()$ denotes the mapping function. This operation is generally nonlinear, thus making the overall TCM code nonlinear. In the following we will derive the upper bound on error probability of a TCM code.

■ Error event probability

Refer to Fig. 5.4.20. Denote by \mathbf{x} the correct path through the trellis. Suppose that decoding is correct up to a trellis node at time j . Let $\mathbf{x}_{i,j}$ be the i th incorrect path that diverges from \mathbf{x} at time j , then remerges later. Let e_{ij} be the event that $\mathbf{x}_{i,j}$ is chosen by the decoder. Then the event-error probability when \mathbf{x} is the correct path is

$$P(e | \mathbf{x}) = P\left(\bigcup_j \bigcup_i e_{ij} | \mathbf{x}\right)$$

The average probability of error is then given by

$$P(e) = \sum_{\mathbf{x}} P(\mathbf{x}) P(e | \mathbf{x}) = \sum_{\mathbf{x}} P(\mathbf{x}) P\left(\bigcup_j \bigcup_i e_{ij} | \mathbf{x}\right)$$

where $P(\mathbf{x})$ is the probability of transmitting \mathbf{x} . Applying the union bound, we have

$$P(e) \leq \sum_{\mathbf{x}} P(\mathbf{x}) \sum_j P\left(\bigcup_i e_{ij} | \mathbf{x}\right)$$

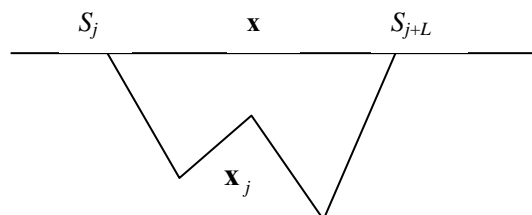


Fig. 5.4.20. An error event of length L

Let n denote the trellis length. Then the rate at which error events occur is given by

$$P_e = \lim_{n \rightarrow \infty} \frac{1}{n} P(e)$$

Averaged over an infinite trellis, every node has the same characteristics, so the dependence on an individual node j can be removed to write

$$P_e \leq \sum_{\mathbf{x}} P(\mathbf{x}) P\left(\bigcup_i e_i | \mathbf{x}\right) \quad (5.5)$$

where e_i is the event that an error starts at an arbitrary but fixed time unit, say j . Equation (5.5) can also be interpreted as the *first event error probability*. Applying the union bound again yields

$$P_e \leq \sum_{\mathbf{x}} P(\mathbf{x}) \sum_{e_i} P(e_i | \mathbf{x}) \quad (5.6)$$

The probability $P(e_i | \mathbf{x})$ is the pairwise error probability, which we denote by $P_2(\mathbf{x} \rightarrow e_i)$.

For the AWGN channel, we have

$$P_2(\mathbf{x} \rightarrow e_i) = Q\left(\sqrt{\frac{d_E^2(\mathbf{x}, e_i)}{2N_0}}\right)$$

where $d_E^2(\mathbf{x}, e_i)$ is the squared Euclidean distance between the signals on the error path e_i and the signals on the correct path \mathbf{x} . Thus the upper bound becomes

$$P_e \leq \sum_{\mathbf{x}} P(\mathbf{x}) \sum_{e_i | \mathbf{x}} Q\left(\sqrt{\frac{d_E^2(\mathbf{x}, e_i)}{2N_0}}\right) \quad (5.7)$$

By rearrange the sum in (5.7), we have

$$P_e \leq \sum_{d_i} A_{d_i} Q\left(\sqrt{\frac{d_i^2 E_s}{2N_0}}\right) \quad (5.8)$$

where A_{d_i} is the average number of paths \mathbf{x}' that are at distance d_i from \mathbf{x} , and the sum is over all the distances; 并且 d_i 是在单位平均能量信号集的假设下计算的。The set of pairs (d_i, A_{d_i}) is known as the (平均) distance spectrum of the code. The smallest distance d_i is the free distance of the code.

Using $Q(\sqrt{x+y}) \leq Q(\sqrt{x})e^{-y/2}$, we obtain

$$Q\left(\sqrt{\frac{d_{free}^2 + (d_E^2 - d_{free}^2)}{2N_0}}\right) \leq Q\left(\sqrt{\frac{d_{free}^2}{2N_0}}\right) \exp\left(\frac{d_{free}^2 - d_E^2}{4N_0}\right)$$

$$= Q\left(\sqrt{\frac{d_{free}^2}{2N_0}}\right) \exp\left(\frac{d_{free}^2}{4N_0}\right) \exp\left(-\frac{d_E^2}{4N_0}\right)$$

Substituting into (5.8), we have the upper bound on event-error probability

$$P(e) \leq Q\left(\sqrt{\frac{d_{free}^2 E_s}{2N_0}}\right) \exp\left(\frac{d_{free}^2 E_s}{4N_0}\right) \sum_{d_i} A_{d_i} \exp\left(-\frac{d_i^2 E_s}{4N_0}\right)$$

where d_i and d_{free} 是在单位平均能量信号集的假设下计算的。 (d_i, A_{d_i}) 可以通过 TCM 编码器的平均重量枚举函数来获得[Rouanne-Costello89].

■ Bit Error Probability

Each event-error causes a certain number of bit errors in the decoded information bits. Let B_{d_i} denote the average number of bit errors on error paths with distance d_i . Since the trellis code encodes m bits per symbol, the average BER is upper-bounded by

$$P_b \leq \sum_{d_i} \frac{1}{m} B_{d_i} Q\left(\sqrt{\frac{d_i^2 E_s}{2N_0}}\right)$$

Note: If a trellis code is regular (or geometrically uniform), the averaging over \mathbf{x} in (5.7) is not necessary and any code sequence may be taken as reference sequence.

图 5.4.21 是使用 8-PSK 信号集的 TCM 码的性能比较。

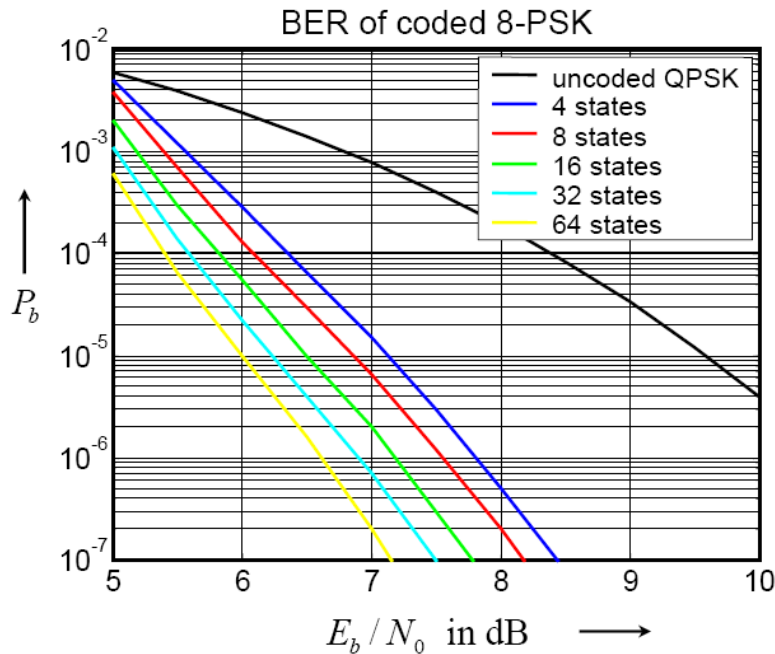


图 5.4.21

5.6 Lattice-Type Trellis Codes

Consider a chain of K binary lattice partition,

$$\Lambda_1/\Lambda_2/\dots/\Lambda_K$$

There are 2^K cosets of Λ_K whose union makes up the original lattice Λ_1 . Each such coset can be identified by a binary K -tuple $c = (c^{(1)}, \dots, c^{(K)})$; i.e., each coset is given by $(\Lambda_K + \mathbf{t}(c))$, where

$$\mathbf{t}(c) = \sum_{i=1}^K c^{(i)} \mathbf{t}^{(i)}$$

Here, $\mathbf{t}^{(i)}$ is an element of Λ_i but not of Λ_{i+1} . The K -tuples c will be called the coset labels of the final cosets $(\Lambda_K + \mathbf{t}(c))$. As an example, Fig.5.4.22 shows the partition $\mathbb{Z}^2 / 2\mathbb{Z}^2$, where $\mathbf{t}^{(0)} = (0, 1)$ and $\mathbf{t}^{(2)} = (1, 1)$.

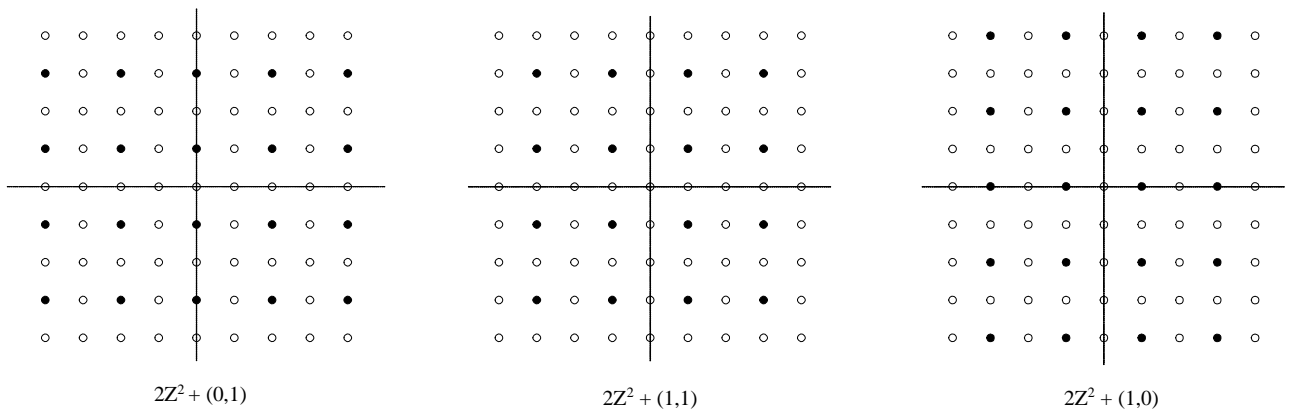


Fig.5.4.22 The four cosets of $2\mathbb{Z}^2$ in the partition $\mathbb{Z}^2 / 2\mathbb{Z}^2$.

With the above lattice partitions, we can perform partitioning of lattice constellations as follows. The n -D lattice constellation $C(\Lambda, \mathcal{R})$ is partitioned into 2^K subsets of equal size that are consistent with the lattice partition Λ/Λ' with $|\Lambda/\Lambda'|=2^K$. The 2^K subsets of points of $C(\Lambda, \mathcal{R})$ are then form sublattice constellations of the form $C(\Lambda', \mathcal{R})$. The region \mathcal{R} must be chosen so that there are an equal number of points in each subset. The sublattice Λ' is usually chosen to be as dense as possible.

With the lattice formulation, we can describe trellis encoders in the general form of Fig. 5.4.23.

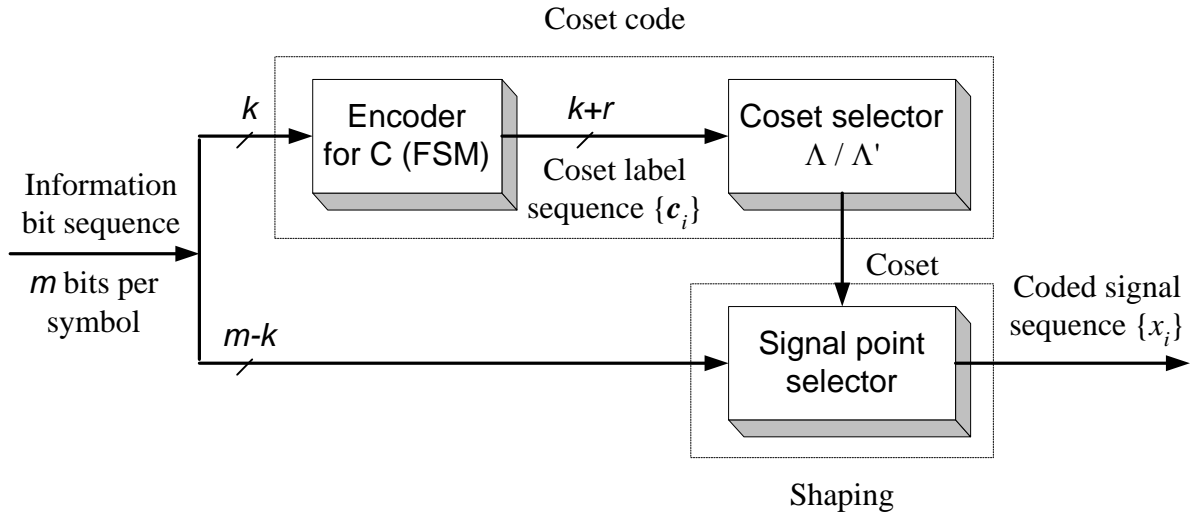


Fig. 5.4.23 Generic trellis encoder block diagram using lattice notation. The signal point selector will use some shaping region \mathcal{R} .

Definition 1: A coset code $\mathcal{C}(\Lambda/\Lambda', C)$ is the set of all signal points that lie within a sequence of cosets of Λ' that could be specified by a sequence of coded bits from C .

Coset codes allow us to separate the coding gain due to the lattice, the shaping gain, and additional coding gain due to C . Due to the code C , only particular sequences of cosets are to be delivered to the signal point selector. The encoder for the code C should be designed to maximize the minimum distance between sequences of cosets.

If C is a block code, then we call the coset code $\mathcal{C}(\Lambda/\Lambda', C)$ a *lattice code*.

Definition 2: A *lattice-type trellis code* is a coset code $\mathcal{C}(\Lambda/\Lambda', C)$, where C is a rate- $k/(k+r)$ convolutional code.

With this description, we can see that Ungerboeck's original one- and two-dimensional PAM codes are based on the four-way partition $\mathbb{Z}/4\mathbb{Z}$ and the eight-way partition $\mathbb{Z}^2/2R\mathbb{Z}^2$, respectively.

Define the redundancy $r(C)$ of the convolutional code C to be the number of redundant bits generated by convolutional encoder per N dimensions. The normalized redundancy per two dimensions is

$$v(C) = \frac{r(C)}{N/2}$$

The coding constellation expansion causes that the transmitted power is increased by approximately $2^{v(C)}$. Note that with the typical $r(C)=1$, the transmitted power is increased by

$2^{2/N}$. On the other hand, the convolutional encoder has increased the minimum distance by a factor of $d_{\text{free}}^2(\mathcal{C})/d_{\text{min}}^2(\Lambda)$. Therefore, the coding gain due to the convolutional code C is

$$\gamma_c = \frac{d_{\text{free}}^2(\mathcal{C})}{d_{\text{min}}^2(\Lambda)2^{v(\mathcal{C})}}$$

References

- [1] G. D. Forney, Jr. and G. Ungerboeck, "Modulation and coding for linear Gaussian channels," *IEEE Trans. Inform. Theory*, vol.44, no.6, pp. 2384-2415, Oct. 1998.
- [2] G. D. Forney, *Principles of Digital Communication (II)*. Course notes. MIT, 2005.
- [3] J. H. Conway and N. J. A. Sloane, *Sphere Packings, Lattices, and Groups*, 2nd ed. New York: Springer-Verlag, 1993.
- [4] E. Biglierli, D. Divsalar, P. J. McLane, and M. K. Simon, *Introduction to Trellis-Coded Modulation with Applications*. New York: MacMillan, 1991.
- [5] S. H. Jamali and T. Le-Ngoc, *Coded-Modulation Techniques for Fading Channels*. Kluwer Academic Publishers, 1994.
- [6] C. Schlegel and L. Perez, *Trellis and Turbo Coding*. IEEE Press, 2004.
- [7] G. Ungerboeck, "Channel coding with multilevel/phase signaling," *IEEE Trans. Inform. Theory*, vol.25, pp.55-67, Jan. 1982.
- [8] G. Ungerboeck, "Trellis-coded modulation with redundant signal sets," Parts I-II, *IEEE Commun. Mag.*, vol.25, pp.5-21, Feb. 1987.
- [9] G. D. Forney, Jr., R. G. Gallager, G. R. Lang, F. M. Longstaff, and S. U. Qureshi, "Efficient modulation for band-limited channels," *IEEE J-SAC*, vol.2, no.5, pp.632-647, Sept. 1984.