# Renaissance Lackings

symmetry & structure in a multidimensional world

by Thomas Bewley, Laul Belitz, & Joseph Cessna





Renaissance Packings: symmetry & structure in a multidimensional world, by Thomas Bewley, Paul Belitz, & Joseph Cessna, published by Renaissance Press.

First edition (widescreen version, English language), 2011. ISBN # ?. Library of Congress Control Number: ?



Published by:

Renaissance Press, 3368 Governor Drive F #244, San Diego CA 92122 USA Online at: http://renaissance-press.com



This book is licensed by T Bewley, P Belitz, & J Cessna under the Creative Commons Attribution-Noncommercial-No Derivative Works 3.0 Unported License (see http://creativecommons.org/licenses/by-nc-nd/3.0/).



All numerical codes included in the Renaissance Packings Codebase, including all codes listed in this book as well as the related free software projects described herein, are copyright (C) 2011 by T Bewley, P Belitz, & J Cessna, and are available online at http://renaissance-packings.com. The Renaissance Packings Codebase is free software: you can redistribute it and/or modify any and all codes included therein

under the terms of the GNU General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version. The Renaissance Packings Codebase is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License (http://www.gnu.org/licenses/gpl-3.0.html) for more details.

# **Summary**

Part II. Derivative-free optimization

**Kriging interpolation** 

Part I. Fundamental concepts & constructions	1
1 Historical retrospective	3
2 Dense lattice packings for $n < 24$	15

rart 1. Fundamental concepts & constructions			
Historical retrospective	3		
<b>Dense lattice packings for</b> $n \le 24$	15		
Characteristics of exemplary lattice and nonlattice packings and nets	33		
	Historical retrospective  Dense lattice packings for $n \le 24$		

3	Characteristics of exemplary lattice and nonlattice packings and nets	3.
4	<b>Rare nonlattice packings &amp; nets for</b> $n \le 8$	4

,	Characteristics of exemplary fattice and nomattice packings and nets	•
1	Rare nonlattice packings & nets for $n \le 8$	4
5	Coding theory	6

1	Rare nonlattice packings & nets for $n \le 8$	4
5	Coding theory	6
	Funther connections between lettics theory and coding theory	0

5	Coding theory	61
5	Further connections between lattice theory and coding theory	99

	country	
5	Further connections between lattice theory and coding theory	99

Extending lattice theory for coordinated derivative-free optimization

108

109

129

9 Global optimization leveraging Kriging-based interpolation	139			
10 Lattice-based derivative-free optimization via global surrogates (LABDOGS)				
11 Optimization via incomplete function evaluations (αDOGS)	161			
12 Optimization in the presence of complex constraint boundaries (latticeMADS)	163			
Part III. Structured computational interconnects	164			
13 Joe chap 1	165			
14 Joe chap 2	167			
15 Joe chap 3	169			
16 Joe chap 4	171			

## **Preface**

The field of n-dimensional sphere packings is elegant and mature in its mathematical development and characterization. However, it is still relatively limited in its practical applications, especially for n > 3. The present text intends to open up two broad new areas for profitable application of this powerful body of mathematical literature in science and engineering. Towards this end, Part I reviews the essential results available in this field (reconciling the theoretical literature for dense and rare sphere packings, which are today largely disjoint), catalogs the key properties of the principle dense and rare sphere packings and corresponding nets available up to n = 24 (including hundreds of values not previously known), and extends the study of regular rare sphere packings and nets to n > 3 dimensions (an area which up to now has been largely unexplored).

Part II then builds from the presentation in Part I to develop three new algorithms (LABDOGS,  $\alpha$ DOGS, and latticeMADS) for performing efficient derivative-free optimization in non-differentiable problems with expensive function evaluations, leveraging the lattices derivd from dense sphere packings as an alternative to Cartesian grids to coordinate the search. We pay particular attention to the improved uniformity and nearest-neighbor configuration of the lattices used over their Cartesian alternatives, and the improvements in

efficiency of optimization algorithms coordinated by such lattices that follow as a direct consequence.

Finally, Part III builds from the presentation in Part I to develop new interconnect strategies for switchless multiprocessor computer systems, leveraging the nets derived from rare sphere packings as alternatives to Cartesian grids to establish structured, fast, and inexpensive interconnects. We pay particular attention to the improved coordination sequences facilitated by such nets over their Cartesian alternatives, and the improvements in the rate of spread of information across the computer system that follow as a direct consequence.

In the applications discussed in Parts II and III, Cartesian grids are used as the default choice today in almost all related realizations. A primary goal of this text is to subvert this dominant Cartesian paradigm and to establish, via the examples we have chosen to highlight, that significant performance gains may be realized in practical engineering applications by leveraging *n*-dimensional sphere packings appropriately.

### A gentle introduction to sphere packing theory

An *n*-dimensional infinite *sphere packing* is simply an array of nodal points in  $\mathbb{R}^n$  obtained via the packing of identical *n*-dimensional spheres. By *packing*, we mean an equilibrium configuration of spheres, each with at least 2 nearest neighbors, against which a repellant force is applied. Many packings investigated in the literature are *stable* packings, meaning that there is a restoring force associated with any small movement of any node of the packing; this requires each sphere in the (*n*-dimensional) packing to have at least n+1 neighbors. Unstable packings with lower nearest-neighbor counts are also of interest. By replacing each sphere in an *n*-dimensional packing with a nodal point (representing, e.g., a computer), and connecting those nodal points which are nearest neighbors, a *net* (a.k.a. *interconnect* or *contact graph*) is formed <sup>1</sup>.

<sup>&</sup>lt;sup>1</sup>As mentioned in the second-to-last paragraph of §2.3, it is natural with certain sphere packings (for example,  $D_n^*$ ,  $A_n^r$ , and the packings associated with the  $T_n^{90}$  and  $T_n^{60}$  nets) to define nets which are *not* contact graphs of the corresponding sphere packings by

which looks identical upon shifting any nodal point to the origin); this shift invariance generally makes lattice packings simpler to describe and enumerate than their nonlattice alternatives. Note that there are many regular<sup>2</sup> sphere packings which are *not* shift invariant [the nonlattice packings corresponding to the honeycomb net in 2D and the diamond and quartz nets in 3D are some well-known examples]. We will focus our attention

An *n*-dimensional real *lattice* (a.k.a. *lattice packing*) is a sphere packing which is shift invariant (that is,

lar<sup>2</sup> sphere packings which are *not* shift invariant [the nonlattice packings corresponding to the honeycomb net in 2D and the diamond and quartz nets in 3D are some well-known examples]. We will focus our attention in this text on those packings and nets which are at least *uninodal* (that is, which look identical upon shifting any nodal point to the origin and rotating and reflecting appropriately). For *dense* sphere packings, from a practical perspective, lattice packings are essentially<sup>3</sup> as good a choice as their more cumbersome nonlattice alternatives for  $n \le 24$  in terms of the four metrics defined below (that is, for maximizing packing density and kissing number and minimizing covering thickness and quantization error). However, the best *rare* sphere packings (with small kissing number) are all nonlattice packings.

As illustrated in Table P.1 and Figure P.1, we may introduce the subject of *n*-dimensional sphere packings by focusing our attention first on the n = 2 case: specifically, on the *triangular*<sup>4</sup> lattice  $(A_2)$ , the *square* lattice  $(\mathbb{Z}^2)$ , and the *honeycomb* nonlattice packing  $(A_2^+)$ . The characteristics of such sphere packings may be quantified by the following measures:

• The packing radius (a.k.a. error-correction radius) of a packing,  $\rho$ , is the maximal radius of the spheres in a set of identical nonoverlapping spheres centered at each nodal point.

connecting non-nearest-neighbor points.

<sup>&</sup>lt;sup>2</sup>The regularity of a nonlattice packing is quantified precisely in §4.1.

<sup>&</sup>lt;sup>3</sup>For n = 10, 11, 13, 18, 20, and 22, there exist nonlattice packings (denoted  $P_{10c}$ ,  $P_{11a}$ ,  $P_{13a}$ ,  $\mathcal{B}_{18}^*$ ,  $\mathcal{B}_{20}^*$ ,  $\mathcal{A}_{22}^*$ ) that are 8.3%, 9.6%,

<sup>9.6%, 4.0%, 5.2%,</sup> and 15.2% denser then the corresponding best known lattice packings (Conway & Sloane 1998, p. xix); to put this into perspective, the density of  $\Lambda_{22}$  is over  $10^6$  times the density of  $\mathbb{Z}^{22}$ .

Anote that many in this field refer to the  $A_2$  lattice (Figure P.1a,b) as "hexagonal". We prefer the unambiguous name "triangular" to avoid confusion with the honeycomb nonlattice packing (Figure P.1e,f).

n	packing	name	Δ	Θ	G	τ	$td_{10}$
	$A_2$	triangular	0.9069	1.2092	0.08019	6	331
2	$\mathbb{Z}^2$	square	0.7854	1.5708	0.08333	4	221
	$A_2^+$	honeycomb	0.6046	2.4184	0.09623	3	166
	$E_8$	Gosset	0.2537	4.059	0.07168	240	1,006,201,681
	$\mathbb{Z}^8$	Cartesian	0.01585	64.94	0.08333	16	1,256,465
8	8 (unstable)	V <sub>8</sub> <sup>90</sup>	5.590e-4	49.89	0.09206	4	37,009
		Y <sub>8</sub> <sup>90</sup>	2.327e-4	87.31	0.09266	3	2290
2.4	$\Lambda_{24}$	Leech	0.001930	7.904	0.06577	196560	> 10 <sup>15</sup>
24	$\mathbb{Z}^{24}$	Cartesian	1.150e-10	4,200,263	0.08333	48	24,680,949,041

Table P.1. Characteristics of selected lattice and uninodal nonlattice packings and nets. Note that n=24 is a natural stopping point in this study. It is special because it is the only integer n>1 that satisfies the equation  $1^2+2^2+\ldots+n^2=m^2$  where m is itself an integer; as a consequence, a particularly uniform lattice with a large number of symmetries is available in this dimension.

- The packing density of a packing,  $\Delta$ , is the fraction of the volume of the domain included within a set of identical non-overlapping spheres of radius  $\rho$  centered at each nodal point on the packing. Packings that maximize this metric are referred to as *close-packed*.
- The covering radius of a packing, R, is the maximum distance between any point in the domain and its

nearest nodal point on the packing. The *deep holes* of a packing are those points which are at a distance R from all of their nearest neighbors. Typical vectors from a nodal point to the nearest deep holes in a lattice packing are often denoted [1], [2], etc.

- The *covering thickness* of a packing,  $\Theta$ , is the number of spheres of radius R centered at each nodal point containing an arbitrary point in the domain, averaged over the domain.
- The *Voronoï cell* of a nodal point in a packing,  $\Omega(P_i)$ , consists of all points in the domain that are at least as close to the nodal point  $P_i$  as they are to any other nodal point  $P_j$ .
- The mean squared quantization error per dimension of a lattice or uninodal nonlattice packing, G, is the average mean square distance of any point in the domain to its nearest nodal point, normalized by n times the appropriate power of the volume, V, of the Voronoï cell. Shifting the origin to be at the centroid of a Voronoï cell  $\Omega(P_i)$ , it is given by

$$G = \frac{S}{nV^{\frac{n+2}{n}}} \quad \text{where} \quad S = \int_{\Omega(P_i)} |\mathbf{x}|^2 d\mathbf{x}, \quad V = \int_{\Omega(P_i)} d\mathbf{x}. \tag{1}$$

- The kissing number (a.k.a. error coefficient) of a lattice or uninodal nonlattice packing,  $\tau$ , is the number of nearest neighbors to any given nodal point in the packing. That is, it is the number of spheres of radius  $\rho$  centered at the nodal points of the packing that touch, or "kiss", the sphere of radius  $\rho$  at the origin.
- The coordination number of a net (derived from a sphere packing, as discussed previously) is the first number of the net's coordination sequence, the k'th element of which is given by  $td_k td_{k-1}$ , where  $td_k$ , which quantifies the net's local topological density, is the total number of nodes reached via k hops or less from the origin in the net<sup>5</sup>.

 $<sup>^{5}</sup>$ In most cases, the natural net to form from a sphere packing is the contact graph; in such cases, the kissing number,  $\tau$ , and the

with Cartesian approaches for such applications. First, the discretization of space is significantly less uniform when using the Cartesian grid as opposed to the available alternatives, as measured by the packing density  $\Delta$ , the covering thickness  $\Theta$ , and the mean-squared quantization error per dimension, G (see Table P.1). Second, the configuration of nearest-neighbor gridpoints is significantly more limited when using the Cartesian grid,

Certain applications, such as those explored in Part II, require dense lattices. There are two key drawbacks

- as measured by the kissing number  $\tau$ , which is an indicator of the degree of flexibility available when selecting from nearest-neighbor points. As seen by comparing the n=2, n=8, and n=24 cases in Table P.1, these drawbacks become increasingly substantial as the dimension n is increased; by the dimension n=24, the Cartesian grid has
- a factor of 0.001930/1.1501e − 10 ≈ 17,000,000 worse (lower) packing density,
   a factor of 4,200,263/7.9035 ≈ 530,000 worse (higher) covering thickness,
- a factor of  $0.08333/0.0658 \approx 1.27$  worse (higher) mean-squared quantization error, and
- a factor of  $196560/48 \approx 4100$  worse (lower) kissing number

this distinction.

tions requiring dense (that is, uniform) lattices with n > 3 is simply untenable.

Other applications, such as those explored in Part III, require regular nets which, with low coordination number, connect to a large number of nodes with each successive hop from the origin, as quantified by the

than the densest available alternative lattice. Thus, the selection of the Cartesian grid, by default, for applica-

number, connect to a large number of nodes with each successive hop from the origin, as quantified by the net's coordination sequence. As mentioned previously, a useful measure of a net's topological density is given, e.g., by  $td_{10}$ , which is the number of distinct nodes within 10 hops of the origin. Note that the coordination

coordination number are equal. As mentioned previously, it is natural with certain sphere packings to define nets which are *not* contact graphs by connecting non-nearest-neighbor points; in such cases, the kissing number (a property of the sphere packing) and the coordination number (as defined here, a property of a corresponding net) are, in general, *not* equal. We find this clear semantical distinction to be useful to prevent confusion between these two distinct concepts; note that some authors (e.g., Conway & Sloane 1998) do not make

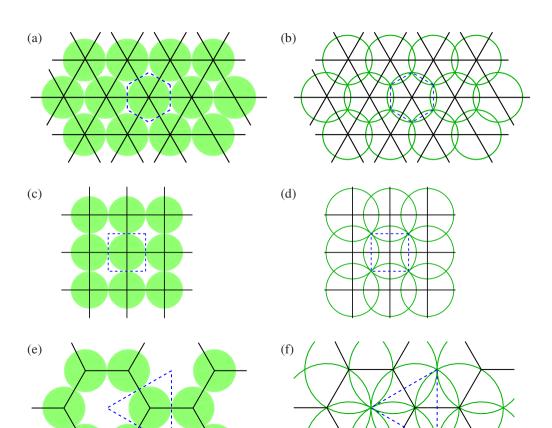
constructions introduced in §4 are as small as 3 or 4, while the topological density increases rapidly as n is increased (compare, e.g., the values of  $td_{10}$  for  $A_2^+$  and  $\mathbb{Z}^2$ , with  $\tau = 3$  and  $\tau = 4$  respectively, to those for

number of the *n*-dimensional Cartesian grid is 2n; the coordination number of the alternative n-dimensional

 $V_9^{90}$  and  $V_9^{90}$  in Table P.1); it is thus seen that, for applications requiring graphs with low coordination number and high topological density, the selection of the Cartesian grid, by default, is also untenable. We are thus motivated to make the fundamental results of both dense and rare n-dimensional sphere packing theory more broadly accessible to the science and engineering community, and to illustrate how this powerful body of theory may be put to use in important new applications of practical relevance. Towards

this end, Part I succinctly reviews and extends several significant results in this mature and sophisticated field, inter-relating the literature on dense and rare packings, which is today largely disjoint. These results are leveraged heavily in the applications described in Parts II and III. We note that, beyond providing an up-todate and synthetic review of this otherwise difficult subject in a (hopefully) accessible language, a significant number of new computations, constructions, algorithms, metrics, and codes are also reported in Part I [the

reader is referred specifically to §3, §4.4.1 through §4.4.7, §4.5, and §6.1.5].



# Part I

Fundamental	concepts & constructions,
fro	m dense to rare

Characteristics of exemplary lattice and nonlattice packings and nets

Further connections between lattice theory and coding theory

15

33

41

61

99

Historical retrospective

**Coding theory** 

Dense lattice packings for  $n \le 24$ 

**Rare nonlattice packings & nets for**  $n \le 8$ 

## Chapter 1

## Historical retrospective

(	1	16	11	ŀe	n	ts

1.1	Finite packings	4
1.2	Infinite packings	11

The mathematical characterization of sphere packings has a long and rich history. Some recent articles and popular books recount this history in detail, including Zong (1999), Szpiro (2003), Hales (2006), and Aste & Weaire (2008). The purpose of the present Part I is not to repeat these historical retrospectives, which these

sources do quite adequately, but to characterize, catalog, and extend the infinite packings available today to facilitate their practical application in new fields. Nonetheless, we would remiss if we didn't at least provide a brief historical context to this field, which we attempt in this short chapter.

#### 1.1 Finite packings

**Mystic marbles**. We begin by defining, for  $m \ge 1$ , a notation to build from:

$$T_{0,m} \triangleq 1$$
,  $T_{1,m} \triangleq \sum_{k=1}^{m} T_{0,k} = m$  (the positive integers).

In the sixth century BC, Pythagoras and his secret society of numerologists, the Pythagoreans, discovered geometrically (see Figure 1.1, and pp. 43-50 of Heath 1931) the formula for the number of marbles placed in a (2D) triangle (that is, the "triangular numbers"):

$$T_{2,m} \triangleq \sum_{k=1}^{m} T_{1,k} = m(m+1)/2.$$

**Stacked spheres**. The earliest known mathematical work to discuss the (3D) stacking of objects is a Sanskrit document *The Aryabhatiya of Aryabhata* (499 AD; see Clark 1930, p. 37), which states:

"In the case of an *upaciti* [lit., 'pile'] which has ... the product of three terms, having the number of terms for the first term and one as the common difference, divided by six, is the *citighana* [lit., 'cubic contents of the pile']. Or, the cube of the number of terms plus one, minus the cube root of this cube, divided by six."

#### 1.1. FINITE PACKINGS

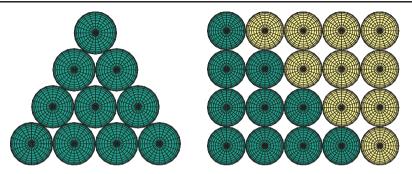


Figure 1.1: (left) Ten marbles placed in a triangle [referred to by the Pythagoreans as a  $\tau \epsilon \tau \rho \alpha \kappa \tau \psi \zeta$ , and upon which they placed a particular mystic significance], and (right) the Pythagoreans' placement of two triangular groups of marbles into an "oblong"  $m \times (m+1)$  rectangle, from which the formula for  $T_{2,m}$  follows immediately.

Thus, Aryabhata establishes, in words, two equivalent expressions for the number of objects ("cubic contents") in a (3D) triangular-based pyramid ("pile") with *m* objects on each edge:

$$T_{3,m} = \frac{m(m+1)(m+2)}{3!} = \frac{(m+1)^3 - (m+1)}{6};$$

note also that  $T_{3,m} \triangleq \sum_{k=1}^{m} T_{2,k}$ .

Thomas Harriot was apparently the first to frame the problem of sphere packing mathematically in modern times (see, e.g., the biography of Harriot by Rukeyser 1972). At the request of Sir Walter Raleigh, for whom

Harriot served, among other capacities, as an instructor of astronomical navigational and on various problems related to gunnery, Harriot (on December 12, 1591) computed, but did not publish, the number of cannonballs in a pile with a triangular, square  $[m \times m]$ , and rectangular  $[m \times (m+1)$ , a.k.a. "oblong"] base, as illustrated in Figure 1.2, obtaining  $T_{3,m}$ ,  $S_m$ , and  $R_m$  respectively, where

$$S_m = \sum_{k=1}^m k^2 = \frac{m(m+1)(2m+1)}{6}, \quad R_m = \sum_{k=1}^m k(k+1) = S_m + T_{2,m} = \frac{m(m+1)(2m+4)}{6}.$$

In 1614, Harriot wrote *De Numeris Triangularibus Et inde De Progressionibus Artithmeticis: Magisteria magna (On triangular numbers and thence on arithmetic progressions: the great doctrine)*<sup>1</sup>. Looking closely at the triangular table of binomial coefficients<sup>2</sup> on pp. 1-3 (folios 108-110) of this remarkable document, it is seen that Harriot understood the *geometric* relationship between the positive integers  $T_{1,m}$ , the "triangular numbers"  $T_{2,m}$  [that is, the number of spheres in a (2D) triangle with m spheres on each edge], the "pyramidal numbers"  $T_{3,m}$  [that is, the number of spheres in a (3D) trianglar-based pyramid with m spheres on each edge], and the next logical steps in this arithmetic progression, given by:

$$T_{4,m} \triangleq \sum_{k=1}^{m} T_{3,k} = \frac{m(m+1)(m+2)(m+3)}{4!}, \ T_{5,m} \triangleq \sum_{k=1}^{m} T_{4,k} = \frac{m(m+1)(m+2)(m+3)(m+4)}{5!},$$

etc. In particular, Harriot noticed that the (n+1)'th element of the (n+m)'th row of this triangular table is  $T_{n,m}$ . Accordingly, we may think of  $T_{n,m}$  as the number of spheres in an "n-dimensional pyramid" with m spheres on each edge, with  $T_{n,2}$  representing n+1 spheres configured at the corners of a regular n-dimensional

<sup>&</sup>lt;sup>1</sup>Harriot (1614) passed through several hands before finally being published in 2009, almost 4 centuries later.

<sup>&</sup>lt;sup>2</sup>This famous triangular table of binomial coefficients is incorrectly attributed by many in the west to Blaise Pascal (b. 1623), though it dates back to several earlier sources, the earliest being Pingala's Sanskrit work *Chandas Shastra*, written in the fifth century BC.

#### 1.1. FINITE PACKINGS

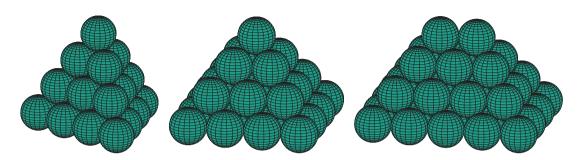


Figure 1.2: Pyramidal stacks of spheres with triangular, square, and "oblong" (rectangular) bases. All three stacks are subsets of the face-centered cubic lattice, discussed further in §2.3.

simplex. It is thus natural to credit Harriot (1614) with the first important steps towards the discovery of laminated lattices, discussed further in §2.4 and §2.6.

Harriot also introduced the packing problem to Johannes Kepler, ultimately leading Kepler (1611), in another remarkable document *Strena seu de nive sexangula* (*The six-cornered snowflake*), which also hypothesized about a related atomistic physical basis for hexagonal symmetry in crystal structures of water, to conjecture that

"The (cubic or hexagonal close) packing is the tightest possible, such that in no other arrangement can more spheres be packed into the same container."

Kepler's conjecture is patently false if considered in a finite container of a specified shape. For instance, a  $2d \times 2d \times 2d$  cubic container can fit 8 spheres of diameter d if arranged in Cartesian configuration, but can

only fit 5 spheres if arranged in a "close-packed" configuration<sup>3</sup>. It is presumed that Kepler in fact recognized this, and thus Kepler's conjecture is commonly understood as a conjecture regarding the densest packing possible in the limit that the size of the container is taken to infinity (for further discussion, see §1.2).

**Permuted planets**. Note in Figure 1.2 that any sphere (referred to as a "sun") on the interior of the piles has 12 nearest neighbors (referred to as its "planets"). Considering this sun and its 12 planets in isolation, there is in fact adequate room to permute the planets to different positions while keeping them in contact with the sun, something like a 12-cornered Rubik's cube with spherical pieces (see Figure 1.3). Due to the extra space available in this configuration, it is unclear upon first inspection whether or not there is sufficient room to fit a 13'th planet in to touch the sun while keeping all of the other 12 planets in contact with it. In 1694, Isaac Newton conjectured this could not be done, in a famous disagreement with David Gregory, who thought it could. Newton turned out to be right, with a complete proof first given in Schütte & van der Waerden (1953), and a substantially simplified proof given in Leech (1956).

Cartoned cans. Moving from 16th-century stacks of cannonballs to 21st-century commerce, the question of dense finite packings of circles and spheres finds practical relevance in a variety of packaging problems. For example, to form a rectangular cardboard carton for 12 fl oz soda cans, 164 cm<sup>2</sup> of cardboard per can is needed if 18 cans are placed in a cartesian configuration with 3 rows of 6 cans per row, whereas 3.3% less cardboard per can is needed if 18 cans are placed in a triangular configuration (within a rectangular box) with 5 rows of {4,3,4,3,4} cans per row. If an eye-catching (stackable, strong, "green"...) hexagonal cardboard carton for the soda cans is used, with 19 cans (described in marketing terms as "18 plus 1 free") again placed in a triangular configuration, 17.7% less cardboard per can is required.

<sup>&</sup>lt;sup>3</sup>For larger containers, the arrangements which pack in the greatest number of spheres (or other objects) must in general be found numerically (see Gensane 2004, Schürmann 2006, and Friedman 2009).

#### 1.1. FINITE PACKINGS

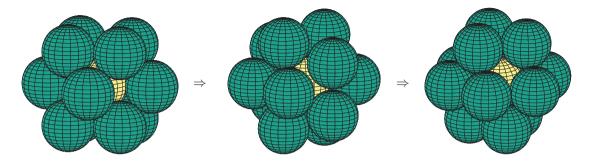


Figure 1.3: Illustration of the 13 spheres (a.k.a. Newton-Gregory) problem and planetary permutations. Configuration (a) is 13 of the spheres taken from the second, third, and fourth layers of the stack in the orientation shown in Figure 1.2b, whereas configuration (c) is 13 of the spheres taken from the third, fourth, and fifth layers of the stack in the orientation shown in Figure 1.2a [extended by one additional layer]. In both configurations, the 12 "planets" (positioned around the central "sun") are centered at the vertices of a cuboctahedron. The planets can be permuted by "pinching" together two of the four planets on the corners of each square face, in an alternating fashion, to form a symmetric icosahedral configuration with significant space between each pair of planets [configuration (b)], then "pushing" apart pairs of planets in an analogous fashion to form a different cuboctahedron. Alternatively, starting from configuration (b), identifying any pair of opposite planets as "poles", and slightly shifting the five planets in each of the "tropics" as close as possible to their nearest respective poles, the resulting northern and southern groupings of planets can be rotated in relation to each other along the equator. Repeated application of these two fundamental motions can be used to permute the planets arbitrarily.

Catastrophic sausages. Two new questions arise when one "shrink-wraps" a number (m) of n-dimensional spheres (resulting in a convex, fitted container), namely: what configuration of the spheres minimizes the surface area of the resulting container, and what configuration minimizes the volume of the resulting container? Both questions remain open, and are reviewed in Zong (1999). Regarding the minimim surface area question, it was conjectured by Croft, Falconer, & Guy (1991) that the minimum surface area, for n > 2 and large m, is achieved with a roughly spherical arrangement. In contrast, regarding the minimim volume question, it was conjectured by L. Fejes Tóth (1975) that the minimum volume, for  $n \ge 5$  and any m, is achieved by placing the spheres in a line, leading to a shrink-wrapped container in the shape of a "sausage". For n = 3, it has been shown that a roughly spherical arrangement minimizes the volume for m = 56, m = 59 to 62, and  $m \ge 65$ , and it is conjectured that a sausage configuration minimizes the volume for all other m (see Gandini & Willis 1992); for n = 4, there appears to be a similar "catastrophe" in the volume-minimizing solution, from a sausage configuration to a roughly spherical configuration, as m is increased beyond a critical value (Willis 1983 conjectures this critical value to be  $m \approx 75000$ , whereas Gandini & Zucco 1992 conjectures it to be m = 375769).

Concealed origins. Finally, L. Fejes Tóth (1959) presents a curious set of questions that arise when considering the blocking of light with a finite number of opaque unit spheres packed around the origin. The first such question, known as Hornich's Problem, seeks the smallest number of opaque unit spheres that completely conceal light rays emanating from a point source at the center of a transparent unit sphere at the origin. A related question, known as L. Fejes Tóth's Problem, seeks the smallest number of opaque spheres that completely conceal light rays emanating from the surface of a unit sphere at the origin (e.g., in Figure 1.3, adding additional outer planets to completely conceal the view of the sun from all angles). In 2D, the (trivial) answer to both problems is 6, via the triangular packing indicated in Figure P.1a. In higher dimensions, both

questions remain open, and the answer differs depending on whether or not the sphere centers are restricted to

#### 1.2. INFINITE PACKINGS

the nodal points of a lattice. For the L. Fejes Tóth's Problem, for  $n \ge 3$ , the answer is unbounded if restricted to lattice points, and bounded if not. For Hornich's Problem, the answer is bounded in both cases, with the number of spheres, h, required in the 3D case, when not restricted to lattice points, being somewhere in the range  $30 \le h \le 42$ . Zong (1999) derives several of the known bounds available in both problems.

#### 1.2 Infinite packings

In the last 300 years, *many* different constructions of infinite lattice and nonlattice packings have been proposed in each dimension. These packings each have different packing density, covering thickness, mean-squared quantization error, and kissing number, and their corresponding nets each have different topological density; knowledge of these properties is essential when selecting a packing or net for any given application. We have thus attempted to catalog these constructions and their properties thoroughly in this review (see §3).

 $E_6$ ,  $E_7$ ,  $E_8$ , and  $\Lambda_{24}$  constructions given in §2 have been proven to be of maximum density, in Lagrange (1773) for n=2, Gauss (1831) for n=3, Korkine & Zolotareff (1873, 1877) for n=4 and 5, Blichfeldt (1935) for n=6 through 8, and Cohn & Kumar (2009) for n=24. There are no such proofs of optimality for other values of n, though the lattices  $\Lambda_n$  and  $K_n$  introduced in §2.6 are likely candidates in the range 9 < n < 23.

In the characterization of density, amongst all *lattice* packings of a given dimension, the  $A_2$ ,  $A_3$ ,  $D_4$ ,  $D_5$ ,

Remarkably, if one considers both lattice and nonlattice packings, proof of which packing is of maximum density in a given dimension is still open for n > 3. It was established in Thue (1892) that  $A_2$  has the maximum density amongst all lattice and nonlattice packings for n = 2. Considerable attention has been focused over the centuries on the corresponding question for  $A_3$  in dimension n = 3, that is, on Kepler's conjecture (posed in 1611) in the limit that the container size is taken to infinity. Indeed, David Hilbert, in his celebrated list of 23 significant open problems in mathematics in 1900, included a generalization of Kepler's conjecture as part of his 18th problem (see, e.g., Milnor 1976).

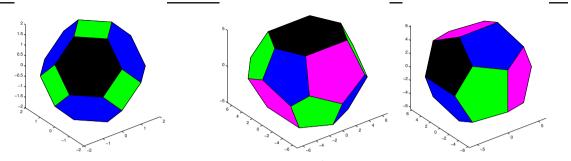


Figure 1.4: (a) A regular truncated octahedron, used to tile  $\mathbb{R}^3$  in Kelvin's conjecture; (b) an irregular tetra-kaidecahedron and dodecahedron, used to tile  $\mathbb{R}^3$  in the Weaire-Phelan structure.

Note that it is not at all obvious that an infinite packing as regular as  $A_3$  would necessarily be the packing that maximizes density. Indeed, as mentioned in footnote 3 on page vii, nonlattice packings are known in dimensions n = 10, 11, 13, 18, 20, and 22 that are each slightly denser than the densest known lattice packings in these dimensions.

In three dimensions, physiologist Stephen Hales (1727), in his groundbreaking work *Vegetable Staticks*, reported a curious experiment:

"I compressed several fresh parcels of Pease in the same Pot,... by the great incumbent of weight, pressed into the interstices of the Pease, which they adequately filled up, being therefore formed into pretty regular dodecahedrons."

This report implied that many of the dilated peas in this experiment had 12 nearest neighbors and/or pentagonal faces. However, the "pretty regular" qualification left a certain ambiguity, and this experiment left

#### 1.2. INFINITE PACKINGS

mathematicians puzzled, as it is patently impossible to tile  $\mathbb{R}^3$  with regular dodecahedra. Kelvin (1887) formalized the question inherent in Hales' dilated pea experiment by asking how R<sup>3</sup> could be divided into regions of equal volume while minimizing the partitional area. He conjectured the answer to be a regular tiling of  $\mathbb{R}^3$  with truncated octahedra, which are in fact the Voronoï cells of the  $A_3^*$  lattice (see §4.4.3). [Note that the Voronoï cell of the  $A_3$  lattice is the (face-transitive) rhombic dodecahedron, which is dual to the cuboctahedron illustrated in Figures 1.3a,c and tiles  $\mathbb{R}^3$  with slightly greater partitional area than does the tiling with truncated octahedra.] Kelvin's conjecture stood for over 100 years, until Weaire & Phelan (1994) discovered a tiling of R<sup>3</sup> based on irregular tetrakaidecahedra (with 2 hexagonal faces and 12 pentagonal faces) and irregular dodecahedra (with 12 pentagonal faces); this tiling has 0.3% less partitional area than the much more regular tiling with truncated octahedra considered by Kelvin (see Figure 1.4). In hindsight, it is quite possible that Hales might have in fact stumbled upon the Weaire-Phelan structure in his cooking pot (in 1727!) and, seeing all of those pentagonal faces and 12-sided (as well as 14-sided) dilated peas, asserted that what he was looking at was a culinary approximation to a tiling of  $\mathbb{R}^3$  with regular dodecahedra, even though such a tiling is impossible.

Returning to Kepler's conjecture, in 1998, Thomas Hales (no relation to Stephen) announced a long-sought-after proof, in a remarkably difficult analysis making extensive use of computer calculations. This proof was spread over a sequence of papers published in the years that followed (see Hales 2005). An extensive discussion of this proof, which is still under mathematical scrutiny, is given in Szpiro (2003). Inspiration for this proof was based, in part, on a strategy to prove Kepler's conjecture proposed by L. Fejes Tóth (1953), the first step of which is a quantitative version of the Newton-Gregory problem discussed in §1.1.

## Chapter 2

**Contents** 

# Dense lattice packings for $n \le 24$

2.1	Lattice terminology		16
	700 C 4 1 1 441	C71 11	40

2.2	The Cartesian lattice $\mathbb{Z}^n$		18
2.2	701 1 1	D 1 4- 1 1 D*	10

2.3	The checkerboard lattice $D_n$ and its dual $D_n^*$	 1
	2.2.1 The effect checkerhoord pecking D+	2

2.4	The zero-sum lattice $A_n$ and its dual $A_n^*$	21
	2.4.1 The glued zero-sum lattices $A^r$	24

	2.4.1	The glued zero-sum lattices $A'_n$	_
2.5	The $E_8$	(Gosset), E <sub>7</sub> , & E <sub>6</sub> lattices and their duals	4

There are many dense lattices more complex than the Cartesian lattice that offer superior uniformity
and nearest-neighbor configuration, as quantified by the standard metrics introduced in the Preface (namely,

packing density, covering thickness, mean-square quantization error, and kissing number). This section provides an overview of many of these lattices; the definitive comprehensive reference for this subject is Conway & Sloane (1998), to which the reader is referred for much more detailed discussion and further references on many of the topics discussed in this chapter. The subject of coding theory, reviewed in §5, is closely related to the subject of such dense lattice packings (see also §6). As mentioned in the Preface, the practical applications explored in Part II of this text leverage these constructions heavily.

#### 2.1 Lattice terminology

2.6

2.7

The notation  $L_n \cong M_n$  means that the lattices  $L_n$  and  $M_n$  are *equivalent* (when appropriately rotated and scaled) at the specified dimension n. Also note that the four most basic families of lattices introduced in this chapter, denoted  $\mathbb{Z}^n$ ,  $A_n$ ,  $D_n$ , and  $E_n$ , are often referred to as *root lattices* due to their relation to the root systems of Lie algebra.

There are three primary methods  $^{1}$  to define any given n-dimensional real lattice:

• As an *explicit description* of the points included in the lattice.

<sup>&</sup>lt;sup>1</sup>A convenient alternative method for building a cloud of lattice points near the origin is based on the stencil of nearest-neighbor points to the origin in the lattice, repeatedly shifting this stencil to each of the lattice points near the origin determined thus far in order to create additional lattice points in the cloud. Unfortunately, this simple alternative method does not work for all lattices, such as  $D_n^*$  and  $A_n^r$  (see §2.3 and 2.4).

#### 2.1. LATTICE TERMINOLOGY

- As an *integer linear combination* (that is, a linear combination with integer coefficients) of a set of n basis vectors  $\mathbf{b}^i$  defined in  $\mathbb{R}^{n+m}$  for  $m \ge 0$ ; for convenience, we arrange these basis vectors as the columns<sup>2</sup> of a *basis matrix*<sup>3</sup> B.
- As a *union of cosets*, or sets of nodal points, which themselves may or may not be lattices.

The standard forms of these definitions, as used throughout this chapter, make it straightforward to generalize application codes that can build easily upon any of the lattices so described.

Any real (or complex) lattice  $L_n$  has associated with it a dual lattice  $L_n^*$  defined such that

$$L_n^* = \left\{ \mathbf{x} \in \mathbb{R}^n \text{ (or } \mathbb{C}^n) : \mathbf{x} \cdot \bar{\mathbf{u}} \in \mathbb{Z} \text{ for all } \mathbf{u} \in L_n \right\}, \tag{2.1}$$

where  $\mathbb{Z}$  denotes the set of all integers, dot denotes the usual scalar product, and overbar denotes the usual complex conjugate. If B is a square basis matrix for  $L_n$ , then  $B^{-T}$  is a square basis matrix for  $L_n^*$ .

Unless specified otherwise, the word lattice in this paper implies a real lattice, defined in  $\mathbb{R}^n$ . However, note that it is straightforward to extend this work to complex lattices, defined in  $\mathbb{C}^n$ . To accomplish this extension, it is necessary to extend the concept of the integers, which are used to construct a lattice via the "integer" linear combination of the basis vectors in a basis matrix B, as described above. There are two primary such extensions:

• The *Gaussian integers*, defined as  $\mathscr{G} = \{a+bi : a,b \in \mathbb{Z}\}$  where  $i = \sqrt{-1}$ , which lie on a square array in the complex plane  $\mathbb{C}$ .

<sup>&</sup>lt;sup>2</sup>In the literature on this subject, it is more common to use a *generator matrix M* to describe the construction of lattices. The basis matrix convention B used here is related simply to the corresponding generator matrix such that  $B = M^T$ ; we find the basis matrix convention to be more natural in terms of its linear algebraic interpretation.

<sup>&</sup>lt;sup>3</sup>Note that integer linear combinations of the columns of most matrices do *not* produce lattices (as defined in the second paragraph of the "gentle introduction" of the Preface). The matrices listed in §2 as basis matrices are special in this regard. Note also that basis matrices are not at all unique, but the lattices constructed from alternative forms of them are equivalent; the forms of the basis matrices listed in §2 were selected based on their simplicity.

• The *Eisenstein integers*, defined as  $\mathscr{E} = \{a + b\omega : a, b \in \mathbb{Z}\}$  where  $\omega = (-1 + i\sqrt{3})/2$  [note that  $\omega^3 = 1$ ], which lie on a triangular array in the complex plane  $\mathbb{C}$ .

We may thus define three types of lattices from a basis matrix B:

- a real lattice, defined as a linear combination of the columns of B with integers as weights;
- a (complex)  $\mathcal{G}$  lattice, defined as a linear combination of the columns of B with Gaussian integers as weights; and
- a (complex)  $\mathscr E$  lattice, defined as a linear combination of the columns of B with Eisenstein integers as weights.

The special *n*-dimensional real,  $\mathscr{G}$ , and  $\mathscr{E}$  lattices formed by taking  $B = I_{n \times n}$  are denoted  $\mathbb{Z}^n$ ,  $\mathbb{Z}[i]^n$ , and  $\mathbb{Z}[\omega]^n$  respectively. Note also that, for any complex lattice with elements  $\tilde{\mathbf{z}} \in \mathbb{C}^n$ , there is a corresponding real lattice with elements  $\tilde{\mathbf{x}} \in \mathbb{R}^{2n}$  such that

$$\tilde{\mathbf{x}} = \begin{pmatrix} \Re{\{\tilde{z}_1\}} & \Im{\{\tilde{z}_1\}} & \dots & \Re{\{\tilde{z}_n\}} & \Im{\{\tilde{z}_n\}} \end{pmatrix}^T.$$
 (2.2)

The present sequence of papers focuses on the practical use of real lattice and nonlattice packings with n > 3. Thus, in the present Part I, we only make brief use of complex lattices to simplify certain constructions.

#### 2.2 The Cartesian lattice $\mathbb{Z}^n$

The Cartesian lattice,  $\mathbb{Z}^n$ , is defined  $\mathbb{Z}^n = \{(x_1, ..., x_n) : x_i \in \mathbb{Z}\}$ , and is constructed via integer linear combination of the columns of the basis matrix  $B = I_{n \times n}$ . The Cartesian lattice is self dual  $[(\mathbb{Z}^n)^* \cong \mathbb{Z}^n]$ .

## **2.3** The checkerboard lattice $D_n$ and its dual $D_n^*$

The checkerboard lattice,  $D_n$ , is an n-dimensional extension of the 3-dimensional face-centered cubic (FCC, a.k.a. cubic close packed) lattice. It is defined

$$D_n = \{ (x_1, \dots, x_n) \in \mathbb{Z}^n : x_1 + \dots + x_n = \text{even} \},$$
 (2.3a)

and may be constructed via integer linear combination of the columns of the  $n \times n$  basis matrix

$$B_{D_n} = \begin{pmatrix} -1 & 1 & & & 0 \\ -1 & -1 & 1 & & & \\ & & \ddots & \ddots & & \\ & & & -1 & 1 \\ 0 & & & & -1 \end{pmatrix}. \tag{2.3b}$$

(2.4a)

The dual of the checkerboard lattice, denoted  $D_n^*$  and reasonably identified as the *offset Cartesian lattice*, is an *n*-dimensional extension of the 3-dimensional *body-centered cubic (BCC)* lattice. It may be written as

 $D_n^* = D_n \cup ([1] + D_n) \cup ([2] + D_n) \cup ([3] + D_n) \cong \mathbb{Z}^n \cup ([1] + \mathbb{Z}^n),$ 

where the coset representatives [1], [2], and [3] are defined in this case such that

$$[1] = \begin{pmatrix} 1/2 \\ \vdots \\ 1/2 \\ 1/2 \end{pmatrix}, \quad [2] = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix}, \quad [3] = \begin{pmatrix} 1/2 \\ \vdots \\ 1/2 \\ -1/2 \end{pmatrix}.$$

The  $D_n^*$  lattice may also be constructed via integer linear combination of the columns of the  $n \times n$  basis matrix

$$B_{D_n^*} = \begin{pmatrix} 1 & 0 & 0.5 \\ 1 & 0.5 \\ & \ddots & \vdots \\ & 1 & 0.5 \\ 0 & & 0.5 \end{pmatrix}. \tag{2.4b}$$

It is important to recognize that, for  $n \ge 5$ , the contact graph of the  $D_n^*$  lattice is simply two disjoint nets given by the contact graphs of the  $\mathbb{Z}^n$  and shifted  $\mathbb{Z}^n$  sets of lattice points upon which  $D_n^*$  may be built [see (2.4a)]. Thus, as suggested by Conway & Sloane (1997), we introduce, for  $n \ge 4$ , a generalized net formed by connecting each node of the unshifted  $\mathbb{Z}^n$  set to the  $2^n$  nearest nodes on the shifted  $\mathbb{Z}^n$  set, and each node on the shifted  $\mathbb{Z}^n$  set to the  $2^n$  nearest nodes on the unshifted  $\mathbb{Z}^n$  set. The resulting net, of coordination number  $2^n$ , is uninodal, but is *not* a contact graph of the corresponding sphere packing.

#### **2.3.1** The offset checkerboard packing $D_n^+$

The packing  $D_n^+$ , reasonably identified as the *offset checkerboard packing*, is an *n*-dimensional extension of the 3-dimensional *diamond* packing, and is defined simply as

$$D_n^+ = D_n \cup ([1] + D_n); \tag{2.5}$$

note that  $D_n^+$  is a lattice packing only for even n, and that  $D_3^+$  is the *diamond packing* (for further discussion, see §4.4.1).

#### **2.4** The zero-sum lattice $A_n$ and its dual $A_n^*$

The zero-sum lattice,  $A_n$ , may be thought of as an *n*-dimensional extension of the 2-dimensional triangular lattice; in 3 dimensions,  $A_3 \cong D_3$ . It is defined

$$A_n = \{(x_0, \dots, x_n) \in \mathbb{Z}^{n+1} : x_0 + \dots + x_n = 0\},$$
 (2.6a)

and may be constructed via integer linear combination of the columns of the  $(n+1) \times n$  basis matrix

$$B_{A_n} = \begin{pmatrix} -1 & & & 0 \\ 1 & -1 & & \\ & \ddots & \ddots & \\ & & 1 & -1 \\ 0 & & & 1 \end{pmatrix}, \quad \text{with} \quad \mathbf{n}_{A_n} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{pmatrix}. \tag{2.6b}$$

Notice that  $A_n$  is constructed here via n basis vectors in n+1 dimensions. The resulting lattice lies in an n-dimensional subspace in  $\mathbb{R}^{n+1}$ ; this subspace is normal to the vector  $\mathbf{n}_{A_n}$ . An illustrative example is  $A_2$ , the triangular 2D lattice, which may conveniently be constructed on a plane in  $\mathbb{R}^3$  (see Figure 2.1).

Note that, starting from a (2D) triangular configuration of oranges or cannonballs (see Figure P.1a), one can stack additional layers of oranges in a trangular configuration on top, appropriately offset from the base

layer, to build up the (3D) FCC configuration mentioned previously (see Figure 1.2a). This idea is referred to as *lamination*, and will be extended further in §2.6 when considering the  $\Lambda_n$  and  $K_n$  families of lattices.

and  $R_2$  is the rectangular grid (not a lattice, nor even a nonlattice packing) obtained by stretching the  $\mathbb{Z}^2$ 

Also note that, in the special case of n = 2, the  $A_2$  lattice may also be written as

$$A_2 \cong R_2 \cup (\mathbf{a} + R_2), \quad \text{where} \quad \mathbf{a} = \begin{pmatrix} 1/2 \\ \sqrt{3}/2 \end{pmatrix}$$
 (2.6c)

lattice in the second element by a factor of  $\sqrt{3}$ .

matrix

The dual of the zero-sum lattice, denoted  $A_n^*$ , may be written as

$$A_n^* = \bigcup_{s=0}^n ([s] + A_n), \tag{2.7a}$$

where the n+1 coset representatives [s], for  $s=0,\ldots,n$ , are defined such that the k'th component of the vector [s] is

$$[s]_k = \begin{cases} \frac{s}{n+1} & k \le n+1-s, \\ \frac{s-n-1}{n+1} & \text{otherwise.} \end{cases}$$
 (2.7b)

The  $A_n^*$  lattice may be constructed via integer linear combination of the columns of the  $(n+1) \times n$  basis

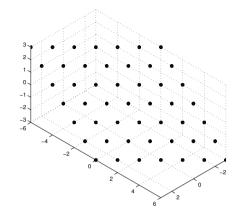


Figure 2.1: A cloud of points on the  $A_2$  lattice, defined on a plane in  $\mathbb{R}^3$ . Note that the normal vector  $\mathbf{n}_{A_2} = \begin{pmatrix} 1 & 1 & 1 \end{pmatrix}^T$  points directly out of the page in this view.

$$\begin{pmatrix} 1 & 1 & \cdots & 1 & \frac{-n}{n+1} \\ -1 & & 0 & \frac{-1}{n+1} \end{pmatrix}$$

$$B_{A_n^*} = \begin{pmatrix} 1 & 1 & \cdots & 1 & \frac{-n}{n+1} \\ -1 & & 0 & \frac{1}{n+1} \\ & -1 & & \frac{1}{n+1} \\ & & \ddots & \vdots \\ & & -1 & \frac{1}{n+1} \end{pmatrix}, \quad \text{with} \quad \mathbf{n}_{A_n^*} = \mathbf{n}_{A_n}. \tag{2.7c}$$

#### 2.4.1 The glued zero-sum lattices $A_n^r$

A related family of lattice packings, developed in §12 of Coxeter (1951) and reasonably identified as the glued zero-sum lattices  $A_n^r$ , is a family of lattices somewhere between  $A_n$  and  $A_n^*$  [as given in (2.7a)] defined via the union of r translates of  $A_n$  for  $n \ge 5$ :

$$A_n^r = A_n \cup ([s] + A_n) \cup ([2s] + A_n) \cup ... \cup ([(r-1)s] + A_n), \text{ where } r \cdot s = n+1,$$
 (2.8)

where the components of the "glue" vectors [s] are specified in (2.7b), and where r and s are integer divisors of (n+1) with 1 < s < n+1 and 1 < r < n+1, excluding the case  $\{r=2, s=3\}$  for n=5. The lattices  $A_9^5$ ,  $A_{11}^4$ ,  $A_{13}^7$ ,  $A_{14}^8$ ,  $A_{15}^8$ ,  $A_{17}^9$ ,  $A_{19}^{10}$ ,  $A_{20}^7$ , and  $A_{21}^{11}$  are found to have especially good covering thickness, with the last four currently the thinnest coverings available in their respective dimensions (see Baranovskii 1994, Anzin 2002, and Sikirić, Schürmann, & Vallentin 2008). Note also that  $A_7^2 \cong E_7$ ,  $A_7^4 \cong E_7^*$ , and  $A_8^3 \cong E_8$ , each of which is discussed further below.

Note finally that the contact graphs of some of the  $A_n^r$  lattices, such as  $A_9^5$  and  $A_{11}^4$ , are disjoint nets given by the contact graphs of the  $A_n$  and shifted  $A_n$  sets of lattice points upon which these glued zero-sum lattices are built [see (2.8)]. Thus, as in the case of  $D_n^*$  for n > 4 as discussed in §2.3, a *generalized net* may be formed by connecting each node of the unshifted  $A_n$  set to the nearest nodes on the shifted  $A_n$  set. Again, the resulting net is uninodal, but is not a contact graph of the corresponding sphere packing.

#### 2.5 The $E_8$ (Gosset), $E_7$ , & $E_6$ lattices and their duals

The Gosset lattice  $E_8 \cong E_8^*$ , which has a (remarkable) kissing number of  $\tau = 240$ , may be defined simply as

$$E_8 = D_8^+,$$
 (2.9a)

#### 2.5. THE $E_8$ (GOSSET), $E_7$ , & $E_6$ LATTICES AND THEIR DUALS

and may be constructed via integer linear combination of the columns of the  $8 \times 8$  basis matrix

$$B_{E_8} = \begin{pmatrix} 2 & -1 & & & 0 & 1/2 \\ 1 & -1 & & & & 1/2 \\ & 1 & -1 & & & 1/2 \\ & & 1 & -1 & & 1/2 \\ & & & 1 & -1 & & 1/2 \\ & & & & 1 & -1 & & -1/2 \\ & & & & 1 & -1 & -1/2 \\ & & & & & 1 & -1/2 \\ & & & & & & 1/2 \end{pmatrix}.$$

(2.9b)

(2.10a)

The lattice  $E_7$  is defined by restricting  $E_8$ , as constructed above, to a 7-dimensional subspace,

$$E_7 = \{(x_1, \dots, x_8) \in E_8 : x_1 + \dots + x_8 = 0\},$$

and may be constructed directly via integer linear combination of the columns of the  $8 \times 7$  basis matrix

$$B_{E_7} = \begin{pmatrix} -1 & & & 0 & 1/2 \\ 1 & -1 & & & 1/2 \\ & 1 & -1 & & & 1/2 \\ & & 1 & -1 & & & 1/2 \\ & & 1 & -1 & & & 1/2 \\ & & & 1 & -1 & & -1/2 \\ & & & 1 & -1 & -1/2 \\ & & & & 1 & -1/2 \\ 0 & & & & & -1/2 \end{pmatrix}, \quad \text{with} \quad \mathbf{n}_{E_7} = \begin{pmatrix} 1/2 \\ 1/2 \\ 1/2 \\ 1/2 \\ 1/2 \\ 1/2 \\ 1/2 \\ 1/2 \end{pmatrix}. \quad (2.10b)$$

The dual of the  $E_7$  lattice may be written as

$$E_7^* = E_7 \cup ([1] + E_7), \quad \text{where} \quad [1] = \begin{pmatrix} 1/4 \\ \vdots \\ 1/4 \\ -3/4 \\ -3/4 \end{pmatrix}, \quad (2.11a)$$

and may be constructed directly via integer linear combination of the columns of the  $8\times7$  basis matrix

$$B_{E_{7}^{*}} = \begin{pmatrix} -1 & & & 0 & -3/4 \\ 1 & -1 & & & -3/4 \\ & 1 & -1 & & 1/4 \\ & & 1 & -1 & & 1/4 \\ & & & 1 & -1 & & 1/4 \\ & & & & 1 & -1 & & 1/4 \\ & & & & & 1 & 1/4 \\ & & & & & 1/4 \end{pmatrix}, \quad \text{with} \quad \mathbf{n}_{E_{7}^{*}} = \mathbf{n}_{E_{7}}. \tag{2.11b}$$

The lattice  $E_6$  is defined by further restricting  $E_7$ , as defined in (2.10), to a 6-dimensional subspace,

$$E_6 = \{(x_1, \dots, x_8) \in E_7 : x_1 + x_8 = 0\},$$
 (2.12a)

and may be constructed directly via integer linear combination of the columns of the  $8 \times 6$  basis matrix

The dual of the  $E_6$  lattice may be written as

$$E_{6}^{*} = E_{6} \cup ([1] + E_{6}) \cup ([2] + E_{6}), \quad \text{where} \quad [1] = \begin{pmatrix} 0 \\ -2/3 \\ -2/3 \\ 1/4 \\ \vdots \\ 1/4 \\ 0 \end{pmatrix}, \quad [2] = -[1], \quad (2.13a)$$

and may be constructed directly via integer linear combination of the columns of the  $8 \times 6$  basis matrix

$$B_{E_6^*} = \begin{pmatrix} 0 & 0 & 1/2 \\ -1 & & 2/3 & 1/2 \\ 1 & -1 & 2/3 & 1/2 \\ 1 & -1 & -1/3 & 1/2 \\ & 1 & -1 & -1/3 & -1/2 \\ & & 1 & -1 & -1/3 & -1/2 \\ & & & -1/3 & -1/2 \\ 0 & & 0 & -1/2 \end{pmatrix}, \text{ with } N_{E^*} = N_E.$$
 (2.13b)

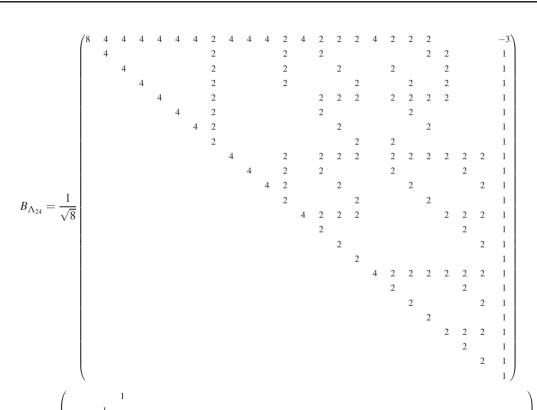
#### **2.6** The laminated lattices $\Lambda_n$ and the closely-related $K_n$ lattices

The lattices in the  $\Lambda_n$  and  $K_n$  families can be built up one dimension, or "laminate", at a time, starting from the integer lattice ( $\mathbb{Z} \cong \Lambda_1 \cong K_1$ ), to triangular ( $A_2 \cong \Lambda_2 \cong K_2$ ), to FCC ( $A_3 \cong D_3 \cong \Lambda_3 \cong K_3$ ), all the way up (one layer at a time) to the remarkable Leech lattice ( $\Lambda_{24} \cong K_{24}$ ). Both families of lattices may in fact be extended (but not uniquely) to at least n = 48.

The Leech lattice,  $\Lambda_{24}$ , is the unique lattice in n=24 dimensions with a (remarkable) kissing number of  $\tau=196,560$ . It may be constructed via integer linear combination of the columns of the  $24\times24$  basis matrix  $B_{\Lambda_{24}}$ , which is depicted below in the celebrated Miracle Octad Generator (MOG) coordinates (see Curtis 1976 and Conway & Sloane 1998). Further, as in the  $E_8\to E_7\to E_6$  progression described in §2.5, the  $\Lambda_n$  lattices for  $n=23,22,\ldots,1$  may all be constructed by restricting the  $\Lambda_{24}$  lattice to smaller and smaller subspaces via the normal vectors assembled in the matrix  $N_{\Lambda}$  depicted below<sup>4</sup>.

<sup>&</sup>lt;sup>4</sup>There are, of course, *many* equivalent constructions of  $\Lambda_1$  through  $\Lambda_{23}$  via restriction of  $\Lambda_{24}$ , and the available literature on the subject considers these symmetries at length. The convenient form of  $N_{\Lambda}$  depicted here was deduced, with some effort, from Figure 6.2 of Conway & Sloane (1998).

#### 2.6. THE LAMINATED LATTICES $\Lambda_N$ AND THE CLOSELY-RELATED $K_N$ LATTICES



Thus, the  $\Lambda_{23}$  lattice is obtained from the points of the  $\Lambda_{24}$  lattice in  $\mathbb{R}^{24}$  (which themselves are generated via integer linear combination of the columns of  $B_{\Lambda_{24}}$ ) which lie in the 23-dimensional subspace orthogonal to  $\mathbf{n}_{\Lambda_{23}}$ . Similarly, the  $\Lambda_{22}$  lattice is obtained from the points of the  $\Lambda_{24}$  lattice which lie in the 22-dimensional subspace orthogonal to both  $\mathbf{n}_{\Lambda_{23}}$  and  $\mathbf{n}_{\Lambda_{22}}$ , etc. Noting the block diagonal structure of  $N_{\Lambda}$ , it follows that  $\Lambda_n$  may be constructed using the basis matrix, denoted  $B_{\Lambda_n}$ , given by the  $n \times n$  submatrix in the upper-left corner of  $B_{\Lambda_{24}}$  for any  $n \in \mathbb{N}_1 = \{21, 20, 16, 9, 8, 5, 4\}$ . For the remaining dimensions,  $n \in \mathbb{N}_2 = \{19, 18, 17, 15, 14, 13, 12, 11, 10, 7, 6, 3, 2, 1\}$ ,  $\Lambda_n$  may be constructed via the appropriate restriction of the lattice generated by the next larger basis matrix in the set  $\mathbb{N}_1$ ; for example,  $\Lambda_{14}$  may be constructed in  $\mathbb{R}^{16}$  via restriction of the lattice generated by the basis matrix  $B_{\Lambda_{16}}$  to the subspace normal to the vectors (in  $\mathbb{R}^{16}$ ) given by the first 16 elements of  $\mathbf{n}_{\Lambda_{15}}$  and  $\mathbf{n}_{\Lambda_{14}}$ .

A similar sequence of lattices, denoted  $K_n$ , may be constructed via restriction of the Leech lattice (generated via  $B_{\Lambda_{24}}$ ) in a similar fashion (for details, see Figure 6.3 of Conway & Sloane 1998). Lattices from the  $\Lambda_n$  and/or  $K_n$  families have the maximal packing densities and kissing numbers amongst all lattices for the entire range considered here,  $1 \le n \le 24$ . Note that the  $\Lambda_n$  and  $K_n$  families are not equivalent in the range  $7 \le n \le 17$ , with  $\Lambda_n$  being superior to  $K_n$  by all four metrics introduced in the Preface at most values of n in this range, except for the narrow range  $11 \le n \le 13$ , where in fact  $K_n$  has a slight advantage. Note also that there is some flexibility in the definition of the lattices  $\Lambda_{11}$ ,  $\Lambda_{12}$ , and  $\Lambda_{13}$ ; the branch of the  $\Lambda_n$  family considered here is that which maximizes the kissing number  $\tau$  in this range of n, and thus the corresponding lattices are denoted  $\Lambda_{11}^{\max}$ ,  $\Lambda_{12}^{\max}$ , and  $\Lambda_{13}^{\max}$ . Note that  $K_{12}$  is referred to as the Coxeter-Todd lattice and  $\Lambda_{16}$ 

is referred to as the Barnes-Wall lattice.

#### 2.7 Numerically-generated lattices with thin coverings for n = 6 to 15

Recall from §2.1 that an *n*-dimensional real lattice may be defined as an integer linear combination of a set of *n* basis vectors  $\mathbf{b}^i$  defined in  $\mathbb{R}^{n+m}$  for  $m \ge 0$ ; that is, any lattice point may be written as

$$\mathbf{x} = y_1 \mathbf{b}^1 + y_2 \mathbf{b}^2 + \ldots + y_n \mathbf{b}^n = B \mathbf{y},$$

where the elements  $\{y_1, ..., y_n\}$  of the vector  $\mathbf{y}$  are taken as integers. The square of the distance of any lattice point from the origin is thus given by  $f(\mathbf{y}) = \mathbf{y}^T A \mathbf{y}$ , where  $A \triangleq B^T B$  is known as the *Gram matrix* associated with the lattice in question, and the function  $f(\mathbf{y})$  is referred to as the corresponding *quadratic form* [note that each term of  $f(\mathbf{y})$  is quadratic in the elements of  $\mathbf{y}$ ]. All of the lattices studied thus far, when scaled appropriately, are characterized by Gram matrices with *integer elements*, and thus their corresponding quadratic forms  $f(\mathbf{y})$  have integer coefficients (and are thus referred to as *integral quadratic forms*).

There is particular mathematical interest in discovering (or generating numerically) both lattice and non-lattice packings which minimize covering thickness and/or packing density. The numerical approach to this problem studied in Schürmann & Vallentin (2006) and Sikirić, Schürmann, & Vallentin (2008) has generated new lattices in dimensions n=6 to 15 with the thinnest covering thicknesses known amongst all lattices<sup>5</sup>. The lattice so generated in dimension 7 happens to correspond to an integral quadratic form, but the others, apparently, do not.

<sup>&</sup>lt;sup>5</sup>Gram matrices A corresponding to these 10 lattices (denoted  $L_7^{c1}, L_7^c, L_8^c, \dots, L_{15}^c$ ) are available at http://fma2.math.uni-magdeburg.de/~latgeo/covering\_table.html

<sup>(</sup>nonunique) basis matrices B corresponding to each of these lattices may be generated simply by taking the Cholesky decomposition of the corresponding Gram matrix, as  $A = B^T B$ .

#### **Chapter 3**

# Characteristics of exemplary lattice and nonlattice packings and nets

For all of the dense lattices described in §2, as well as for all of the rare packings and nets described in §4, Tables 3.1-3.2 list the known values of the packing density  $\Delta$ , the covering thickness  $\Theta$ , and the mean squared quantization error per dimension, G. Table 3.1 also lists the coordination sequence through k=10 of the corresponding net, as well as its local topological density  $td_{10}$ . If this net is a contact graph, the coordination number (that is, the first element of the coordination sequence) is equal to the kissing number of the corresponding packing; if this net is *not* a contact graph, it is marked with a G, and the kissing number  $\tau$  of the corresponding sphere packing is listed in parentheses.

#### CHAPTER 3. CHARACTERISTICS OF LATTICE AND NONLATTICE PACKINGS AND NETS

The other information appearing in Table 3.1 is described further in §4. Note that Table 3.1 alone has 8

columns and over 100 rows, with those results which we believe to be new denoted in italics. The original source of each of the several hundred existing results reported can not feasibly be spelled out here. Suffice it to say that the vast majority of those existing results related to lattices are discussed in Conway & Sloane (1998) and in the On-Line Encyclopedia of Integer Sequences<sup>1</sup>, where a large number of the original references are listed in detail. The vast majority of those existing results related to 3D nets (see §4), *including clear drawings of each* as well as detailed lists of original references, are given in the Reticular Chemistry Structure Resource<sup>2</sup>; for further discussion of this database and others, see O'Keeffe et al. (2008), Treacy et al. (2004), Blatov (2006), and Hyde et al. (2006). Note also that there are hundreds of new results reported in Tables 3.1 and 3.2, as denoted in italics; most of these are the result of painstaking numerical simulation, some of which

Note finally that there are a variety of (lattice-specific) ways to quantize to the nearest lattice point; for an introduction, see §6.

Available on the web at http://www.research.att.com/~njas/sequences/.

tooks weeks of CPU time (on a quad-core 3GHz Intel Xeon server) to complete.

<sup>&</sup>lt;sup>2</sup> Available on the web at, e.g., http://rcsr.anu.edu.au/nets/fcu, where "fcu" may be replaced by any of the lowercase boldface three-letter identifiers given in Table 3.1 and §4.

СН	APT	ER 3	. CI	IAF	AC	TER	ISTI	CS	OF	LA	[TI	CE.	AN.	D N	ON	LA	TI	CE I	PAC	KI	NG	S A	.ND	NE	TS	
	coordination sequence (through $k = 10$ )	<u>2</u> , 2, 2, 2, 2, 2, 2, 2, 2	<u>6</u> , 12, 18, 24, 30, 36, 42, 48, 54, 60	4, 8, 12, 16, 20, 24, 28, 32, 36, 40	3, 6, 9, 12, 15, 18, 21, 24, 27, 30	3, 4, 6, 8, 12, 14, 15, 18, 21, 22	<u>12</u> , 42, 92, 162, 252, 362, 492, 642, 812, 1002	<u>12</u> , 44, 96, 170, 264, 380, 516, 674, 852, 1052	8, 26, 56, 98, 152, 218, 296, 386, 488, 602	6, 18, 38, 66, 102, 146, 198, 258, 326, 402	4, 12, 30, 52, 80, 116, 156, 204, 258, 318	4, 12, 24, 42, 64, 92, 124, 162, 204, 252	4, 12, 25, 44, 67, 96, 130, 170, 214, 264	4, 10, 20, 34, 52, 74, 100, 130, 164, 202	4, 6, 12, 18, 36, 48, 60, 78, 108, 126	4, 6, 12, 17, 28, 38, 52, 64, 84, 104	<b>G</b> : 4, 12, 36, 72, 122, 188, 264, 354, 456, 570 ( $\tau = 8$ )	<b>G</b> : 4, 12, 30, 58, 94, 138, 190, 250, 318, 394 ( $\tau = 6$ )	4, 12, 28, 50, 76, 110, 148, 194, 244, 302	001 101 10 10 01 100 01	3, 6, 12, 24, 43, 64, 91, 124, 160, 202	3 6 12 24 38 56 77 102 120 160		3, 6, 12, 24, 35, 48, 69, 86, 108, 138	3, 4, 6, 8, 12, 16, 24, 32, 48, 54	Table 3.1a. (Continued on next page.)
	G	0.083333	0.080188	0.083333	0.09623	0.1652	0.078745	0.078745	0.078543	0.083333	0.08534	0.09114	0.09139	0.1092	0.1511	0.1943	0.08151	0.08333	0.08602	0.09129	0.09217	0.09706	0.09817	0.1072	0.1882	Table 3.
	Ð	1	1.2092	1.5708	2.4184	5.832	2.0944	2.0944	1.4635	2.7207	2.0405	2.7207	3.3068	8.781	9.1723	28.26	2.1549	2.7207	3.1416	3.0042	2.7251	4.3099	3.518	3.4281	9.739	
	$\nabla$	<del>-</del> 11	06906.0	0.78540	0.60460	29068.0	0.74048	0.74048	0.68017	0.52360	0.39270	0.34009	0.34009	0.2777	0.12354	0.1033	0.6046	0.52360	0.39270	0.2687	0.2551	0.2327	0.2207	0.1851	0.0555	

1/13	2077200	<u>24</u> , 144, 456, 1056, 2040, 3504, 5544, 8256, 11736, 16080	
; †	C000/0.0	<b>G</b> : 16, 80, 240, 544, 1040, 1776, 2800, 4160, 5904, 8080 ( $\tau$ = 24)	
780	0.078020	20, 110, 340, 780, 1500, 2570, 4060, 6040, 8580, 11750	CF
<u>829</u>	0.077559	10, 50, 150, 340, 650, 1110, 1750, 2600, 3690, 5050	IAP
348	0.08333	8, 32, 88, 192, 360, 608, 952, 1408, 1992, 2720	TE
829	0.08827	5, 20, 50, 110, 200, 340, 525, 780, 1095, 1500	R 3.
2.4	0.1221	5, 15, 35, 70, 125, 205, 315, 460, 645, 875	СН
.82	0.1398	5, 8, 20, 32, 80, 116, 170, 236, 380, 482	AR
335	0.08333	<b>G</b> : 4, 12, 36, 92, 200,384, 664, 1056, 1576, 2240 ( $\tau = 8$ )	AC'
855	0.08692	4, 12, 36, 84, 172, 292, 468, 692, 988, 1348	ГЕБ
814	0.09333	4, 12, 36, 74, 136, 228, 352, 518, 732, 994	RIST
458	0.09736	3, 6, 12, 24, 48, 90, 146, 230, 336, 478	TICS
226	0.075786	<u>40</u> , 370, 1640, 4930, 11752, 24050, 44200, 75010, 119720, 182002	OF
218	0.077647	30, 240, 1010, 2970, 7002, 14240, 26070, 44130, 70310, 106752	LA
286	0.075625	<b>G</b> : 32, 242, 992, 2882, 6752, 13682, 24992, 42242, 67232, 102002 ( $\tau = 1$	TJ.
889	0.07784	16, 120, 480, 1410, 3296, 6712, 12256, 20770, 33056, 50232	CE
243	0.076922	12, 72, 272, 762, 1752, 3512, 6372, 10722, 17012, 25752	AN
955	0.083333	10, 50, 170, 450, 1002, 1970, 3530, 5890, 9290, 14002	DΝ
223	0.08646	6, 30, 90, 240, 510, 1010, 1770, 2970, 4626, 7002	ION
6.4	0.1349	6, 21, 56, 126, 252, 461, 786, 1266, 1946, 2877	ΓLΑ
18:	0.1313	6, 10, 30, 50, 150, 230, 390, 570, 1050, 1420	ТТІ
955	0.08333	<b>G</b> : 4, 12, 36, 100, 258, 610, ? $(\tau = 10)$	CE
310	0.08666	4, 12, 36, 100, 244, 514, 980, 1682, 2724, 4162	PAC
884	0.08753	4, 12, 36, 100, 248, 522, 988, 1724, 2800, 4324	CKI
4.8	0.1350	3, 6, 12, 24, 48, 90, 168, 312, 556, 914	NG
825	0.09038	<b>G</b> : 4, 12, 36, 100, 268, ? $(\tau = 14)$	<b>S</b> A
910	28060.0	4, 12, 36, 100, 220, 428, 752, 1254, 1944, 2924	ND
61.	0.09605	3, 6, 12, 24, 48, 90, 168, 312, 532, 872	NE'

Table 3.1b. (Continued on next page.)

9	coordination sequence (through $k = 10$ )	Cl
0.074347	<u>72</u> , 1062, 6696, 26316, 77688, 189810, 405720, 785304, 1408104	ΙΑΡ
0.074244	54, 828, 5202, 20376, 60030, 146484, 312858, 605232, 1084806, 1830060	TEI
0.075591	60, 792, 4724, 18096, 52716, 127816, 271908, 524640, 938652, 1581432	R 3.
0.07459	32, 332, 1824, 6776, 19488, 46980, 99680, 192112, 343584, 578876	СН
0.077466	42, 462, 2562, 9492, 27174, 65226, 137886, 264936, 472626, 794598	AR
0.075120	<b>G</b> : 64, 728, 4032, 14896, 42560, 102024, 215488, 413792, 737856, 1240120 ( $\tau$ =	Æ
0.076490	14, 98, 462, 1596, 4410, 10374, 21658, 41272, 73206, 122570	ГER
i	32,?	ISI
0.08333	12, 72, 292, 912, 2364, 5336, 10836, 20256, 35436, 58728	ICS
0.08525	7, 42, 147, 462, 1127, 2562, 5047, 9492, 16317, 27174	S OI
0.14712	7, 28, 84, 210, 462, 924, 1715, 2996, 4977, 7924	F L
0.1259	7, 12, 42, 72, 252, 402, 777, 1182, 2457, 3492	ΛТТ
0.08333	<b>G</b> : 4, 12, 36, 100, ? $(\tau = 12)$	ICE
10980'0	4, 12, 36, 100, 276, 660, 1484, 2920, ?	AN
0.09322	4, 12, 36, 100, 276, 610, 1284, 2346, 4152, 6792	I DI
0.09479	3, 6, 12, 24, 48, 90, 168, 312, 580, 1046	101
0.073231	126. 2898, 25886, 133506, 490014, 1433810, 3573054, 7902594, 15942206, 2989	ΛΪ́ξα'
0.07273	64, 1092, 8064, 37842, 131328, 371940, 906816, 1976898, 3946048, 7344164	ГΤΙ
0.073116	56, 938, 7688, 39746, 150248, 455114, 1171928, 2668610, 5521880, 10585514	CE
9895200	84, 1498, 11620, 55650, 195972, 559258, 1371316, 2999682, 6003956, 11193882	PAC
9682200	56, 812, 5768, 26474, 91112, 256508, 623576, 1356194, 2703512, 5025692	CKI
0.07493	G: 128, 2186, 16256, 75938, 263552, 745418, 1817216, 3959426, 7902848, 1470 <sup>a</sup>	<b>A</b>
0.076187	16, 128, 688, 2746, 8752, 23536, 55568, 118498, 232976, 428752	S Al
i	i	ND
0.083333	14, 98, 462, 1666, 4942, 12642, 28814, 59906, 115598, 209762	NE
0.08442	8, 56, 224, 812, 2240, 5768, 12656, 26474, 49952, 91112	ΓS
i	8, 36, 120, 330, 792, 1716, 3432, 6434, 11432, 19412	
0.1214	8, 14, 56, 98, 392, 644, 1400, 2198, 5096, 7532	

12, 24, 48, 90, 168,		
0.09365 3, 6, 1		
0.010	3, 6, 12, 24, 48, 90, 168, 312, 580, 1046	Cl
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	<b>240</b> , 9120, 121680, 864960, 4113840, 14905440, 44480400, 114879360, 265422960, 561403680	НАРТ
0.075914 112, 2592, 67232416	112, 2592, 25424, 149568, 629808, 2100832, 5910288, 14610560, 32641008, 67232416	ER 3
0.077391 72, 13	72, 1332, 11832, 66222, 271224, 889716, 2476296, 6077196, 13507416, 2771	3.5 248 248
0.074735 <b>G</b> : 256, 656 171312160	<b>G</b> : 256, 6560, 65280, 384064, 1614080, 5374176, 15097600, 37281920, 83222 171312160 ( $\tau=16$ )	₽AR
0.075972 18, 16	18, 162, 978, 4482, 16722, 53154, 148626, 374274, 864146, 1854882	AC'
i		ГЕБ
0.083333 16,12	16, 128, 688, 2816, 9424, 27008, 68464, 157184, 332688, 658048	RIST
0.0838 9, 72,	72, 324, 1332, 4104, 11832, 28674, 66222, 136404, 271224	CIC
? 9,45,	9, 45, 165, 495, 1287, 3003, 6435, 12870, 24309, 43749	S Ol
0.1178 9, 16,	72, 128, 576, 968, 2340, 3768, 9648, 14716	F L
0.08333 <b>G</b> : 4,	<b>G</b> : 4, 12, 36, 100, 276, 724, ? $(\tau = 16)$	ΑТТ
0.08452 4, 12,	4, 12, 36, 100, 276, 724, ?	ICE
0.09206 4, 12,	4, 12, 36, 100, 276, 724, 1908, 4390, 9876, 19682	ΙΑì
0.09266 3, 6, 1	3, 6, 12, 24, 48, 90, 168, 312, 580, 1046	ND I
Table 3.1d. (Con emplary lattice a	Table 3.1d. (Continued from previous pages.) Characteristics of some exemplary lattice and uninodal nonlattice packings and nets through $n=8$ ,	NONL
ordered from der	ordered from dense to rare in each section. Values in italics are (as far	\TT
as we know) nev	as we know) new. At each n, bold double underlined values are proven	ΊC
to be optimal (n	to be optimal (maximum or minimum) amongst all packings, and bold	E P.
single underline	single underlined values are proven to be optimal amongst all <u>lattices.</u>	AC.
Bold values (wit,	Bold values (without underlines) are the best known values amongst all	KI
packings, and bo all lattices. The	backings, and bold undertuated values are the best known values amongst all lattices. The point symbol is provided for those nets with $\tau > 5$ : the	NGS
vertex symbol is	vertex symbol is provided for those nets with $\tau \le 4$ . Nets whose coordi-	AN
nation sequences	nation sequences are identified with a G are generalized nets, not contact	ID
graphs (see, e.g. the kissing numb	graphs (see, e.g., the second-to-last paragraph of $\S2.3$ ); in these cases, the kissing number $\tau$ is indicated in parentheses after the coordination se-	NETS

avence In all other cases the first element of the coordination seamence is

## $\Lambda_9$ 0.14577 9.0035 0.07206 272 $D_0^+$ 0.14577 4.3331 0.07110 144

CHAPTER 3. CHARACTERISTIGGE HATTICE AND ONLASTICE PACKINGS AND NETS

	1						
	$D_9^+$	0.14577	4.3331	0.07110	144		
	$D_9^*$	0.01288	8.6662	0.07469	18		
9	$A_9^*$	0.01268	4.3889	0.07582	20		
	$A_9^5$	0.08447	4.3402	0.07207	90		
	$L_9^c$	0.08149	4.2686	?	?		
	$\mathbb{Z}^9$	0.006442	126.81	0.08333	18		
	$\Lambda_{10}$	0.09202	12.409	0.07150	336		
	$D_{10}^{+}$	0.07969	7.7825	0.07081	180		
10	$A_{10}^{*}$	0.005128	5.2517	0.07570	22		
	$L_{10}^{c}$	0.02995	5.1545	?	?		
	$\mathbb{Z}^{10}$	0.002490	249.04	0.08333	20		
	K <sub>11</sub>	0.06043	?	?	432		
	$\Lambda_{11}^{max}$	0.05888	24.781	0.07116	438		
	$D_{11}^{+}$	0.04163	8.4072	?	220		
11	$A_{11}^{*}$	0.001974	6.2813	0.07562	24		
	$A_{11}^{4}$	0.04740	5.5983	0.07025	132		
	$L_{11}^{c}$	0.04124	5.5056	?	?		
	$\mathbb{Z}^{11}$	9.200e-4	491.40	0.08333	22		
	$K_{12}, K_{12}^*$	0.04945	17.783	0.07010	756		

CHAPTER 3. CHARACTERISTICS packing

AND NONLATTICE PACKINGS AND NETS

	$\Lambda_{16}, \Lambda_{16}^*$	0.01471	96.500	0.06830	4320		
16	$A_{16}^{*}$	9.116e-6	15.311	0.07549	34		
	$\mathbb{Z}^{16}$	3.591e-6	15,422	0.08333	32		
	$\Lambda_{17}$	0.008811	197.72	0.06822	5346		
1.7	$A_{17}^{*}$	2.807e-6	18.288	0.07549	36		
17	$A_{17}^{9}$	?	12.357	?	?		
	$\mathbb{Z}^{17}$	1.076e-6	30,936	0.08333	34		
	$\Lambda_{18}$	0.005928	301.19	0.06792	7398		
18	$A_{18}^{*}$	8.396e-7	21.841	0.07550	38		
	$\mathbb{Z}^{18}$	3.134e-7	62,158	0.08333	36		
	$\Lambda_{19}$	0.004121	607.62	0.06767	10668		
10	$A_{19}^{*}$	2.443e-7	26.082	0.07552	40		
19	$A_{19}^{10}$	?	21.229	?	?		
	$\mathbb{Z}^{19}$	8.892e-8	125,077	0.08333	38		
	$\Lambda_{20}$	0.003226	889.86	0.06731	17400		
20	$A_{20}^{*}$	6.924e-8	31.143	0.07553	42		
20	$A_{20}^{7}$	?	20.367	?	?		
	$\mathbb{Z}^{20}$	2.461e-8	252,020	0.08333	40		
	$\Lambda_{21}$	0.002466	1839.5	0.06701	27720		

#### **Chapter 4**

#### \_

## Rare nonlattice packings & nets for $n \le 8$

Contents	

4.1	Net terminology	43
4.2	2D nets	47

7.2	2D licts	
4.3	A List of Twelve "highly regular" uninodal 3D nets	48
1.1	Uninedal extension of some regular 2D note to higher dimensions	51

	8	0		
4.4.1	Extending <b>dia</b> : the $A_n^+$ and $D_n^+$ pack	ings	 	51
442	Augmenting $A^{+}$ : the $\widehat{A}^{+}$ necking			50

4.4.2	Augmenting $A_n^+$ : the $A_n^+$ packing						 									 	5
4.4.3	Extending <b>sod</b> : the ${}^{T}\!A_n^*$ packing .						 									 	4

#### 

We now turn our attention to the problem of infinite *rare* sphere packings, with packing density *lower* than that of the corresponding Cartesian packing, and the closely related problem of infinite nets. For n = 2,

4.4.4

4.4.5

CHAPTER 4. RARE NONLATTICE PACKINGS & NETS FOR N < 8

this problem is essentially trivial. For n=3, the richness of solutions to this problem is fascinating and, due to the intense interest in crystallographic structures with various desirable chemical properties, has been exhaustively studied and catalogued. For n>3, relatively few regular constructions are known, and it appears as if what academic interest there has been has not yet led to any applications of significance in science and engineering; Part III of this text intends to change this, thus motivating the present chapter.

Interest in n-dimensional space groups and symmetries dates back to the nineteenth century, with the work of Hessel, Bravais, Gadolin, Frankenheim, Barlow, Rodrigues, Möbius, Jordan, Sohncke, Fedorov,

zenberger (1980). Much of the related work in the field of geometry was developed by Coxeter (1970, 1973, 1974, 1987, 1989). Despite this intense interest, there are very few explicit constructions of regular rare sphere packings for n > 3 available today, outside of very short treatments of the subject by O'Keeffe (1991b) and Beukemann & Klee (1992), discussed below.

As mentioned in the abstract and explored in depth in Part III, certain emerging engineering applications

now motivate the further development and deployment of quasi-infinite n-dimensional nets, with a particular

#### 4.1. NET TERMINOLOGY

for the rapid spread of information in switchless computational interconnect systems with a reduced number of wires and, thus, reduced cost. In such systems, a logical network with n > 3 may easily be designed and built and, as we will see, there are significant potential benefits for so doing. We are thus motivated to revisit the problem of the design of structured nets with low coordination number. Note that none of the lattices discussed in §2 have a coordination number lower than that of the corresponding Cartesian lattice,  $\tau = 2n$ . However, for n = 3, there is a wide range of stable and unstable nonlattice packings that lead to such nets; as shown below, many of these packings and nets generalize naturally to higher dimensions.

focus on structured nets with low coordination number and high topological density. Such nets are well suited

#### 4.1 Net terminology

The terminology used to discuss 3D nets, most of which generalizes readily to the discussion of *n*-dimensional nets, has been clarified significantly over the last decade, and is now quite precise.

Recall first the measures defined in the Preface, including the *coordination number*, the *coordination sequence*, and a k-hop measure of *local topological density* given by the cumulative sum of all nodes reached within k hops from origin, denoted  $td_k$  (Tables 3.1 and 3.2 list this quantity for k = 10). O'Keeffe (1991a) defines another, sometimes preferred (see, e.g., Grosse-Kunstleve et al. 1996) measure of *global topological density*,  $td = \lim_{k \to \infty} td_k/k^n$ , which reveals the rate of growth of  $td_k$  with k in the limit of large k. [For a uninodal n-dimensional net, td may be found by representing the coordination sequence as an (n-1)'th-order polynomial in the number of hops k, then taking the leading coefficient of this polynomial and dividing

<sup>&</sup>lt;sup>1</sup>Recall, e.g., the "hypercube" computational interconnect system introduced several years ago; though designed with a logical network with n > 3, the hypercube, like most computational interconnect strategies deployed today, is significantly hampered by its inherent dependence on a Cartesian topology.

<sup>&</sup>lt;sup>2</sup>Or by approximating this coordination sequence as an (n-1)'th-order polynomial for large k, if such a polynomial does not fit exactly.

by n.] Despite some impressive efforts in representing coordination sequences with such polynomials (see, e.g., Conway & Sloane 1997, and the references contained therein), the measure td is currently unknown for most of the nets discussed here. As a matter of computational tractability, we thus resort in §3 to the tabulation of the local topological density measure,  $td_{10}$ , as this measure is much easier to compute.

Our attention in this text is focused almost exclusively on *equilibrium packings* (that is, on sphere packings which, if unperturbed, can bear compressive loads applied at the edges of a packing that is built out to fill a finite convex domain) and their corresponding *equilibrium nets* (which are constructed with tensile members connecting nearest-neighbor nodes, and can bear tensile loads applied at the edges of a finite convex domain)<sup>3,4</sup>. Equilibrium packings fall into two catagories: stable (that is, sphere packings which, if perturbed, oscillate about their equilibrium configurations, and return to these configurations if there is damping present) and unstable (that is, sphere packings which depart from equilibrium if perturbed); we consider both.

After years of conflicting terminology in the literature on nets, the concepts of *cycles*, *rings*, *strong rings*, *tilings*, *natural tilings*, *point symbols*, and *vertex symbols* have, in 3D, finally crystallized. The reader is referred to Blatov et al. (2009) and the references contained therein for description of this modern terminology, as well as numerous cautions concerning the conflicting nomenclatures adopted elsewhere in the published literature. In short:

• A cycle is a sequence of nodes in a net, connected by edges, such that the first and last nodes of the sequence

 $<sup>^3</sup>$ A family of structures with both tensile and compressive members, known as *tensegrity*, might be said to cover the gap between purely compressive sphere packings and purely tensile nets. The mathematical characterization of tensegrity systems in 3D is now precise, due largely to the work of Skelton & de Oliveira (2009). An interesting extension of the present study would be to generalize such tensegrity systems to n > 3 dimensions.

<sup>&</sup>lt;sup>4</sup>For the purpose of the applications studied in Parts II and III, we do not actually use the property of mechanical equilibrium of the corresponding structure; this property may rather be considered as a convenient means to an end when designing a regular packing or net. Several nets discussed in the literature (see, e.g., Wells 1977, page 80) are not equilibrium sphere packings, and might be interesting to consider further.

#### 4.1. NET TERMINOLOGY

issue further.

coincide, while none of the other nodes in the sequence appears more than once.

- A cycle sum, of cycles A and B, is the union of those edges in either A or B but not both.
- A ring is a cycle that is not the sum of two smaller cycles.
- A strong ring is a cycle that is not the sum of any number of smaller cycles.
- A *tiling* of  $\mathbb{R}^3$  by a 3D net is simply the dissection of 3D space into volumes whose faces, which in general may be curved (as *minimal surfaces*, like soap bubbles; see, e.g., Sadoc & Rivier 1999), are bounded by cycles of the net. A 3D net generally admits many tilings.

• The dual of a tiling is the unique new tiling obtained by placing a new vertex inside each original tile and

- connecting the vertices of adjacent tiles (that is, with shared faces) in the original tiling with edges. If a tiling and its dual are identical, the tiling is said to be *self-dual*. The dual of a dual is the original tiling.

   A *natural tiling* of  $\mathbb{R}^3$  by a 3D net is a tiling with the smallest possible tiles such that the tiles have the
- maximum combinatorial symmetry and all the faces of the tiles are strong rings. A 3D net often<sup>5</sup> admits a unique natural tiling. If a tiling and its dual are both natural, the pair is referred to as *natural duals*. If a natural tiling is self-dual, it is said to be *naturally self-dual*.
- The *point symbol* of a uninodal net, of the form  $A^a.B^b.C^c...$ , indicates that there are a pairs of edges touching the node at the origin with shortest cycles of length A, b pairs of edges touching the node at the origin with shortest cycles of length B (with B > A), etc. Note that the sum of the superscripts in a point symbol totals  $\tau(\tau 1)/2$ .
- The *vertex symbol* of a uninodal net, of the form  $A_a.B_b.C_c...$ , indicates that the first pair of edges touching the node at the origin has a shortest rings of length A, the second pair of edges touching the node at the origin has b shortest rings of length b, etc. If for any entry there is only 1 such shortest ring, the subscript

origin has b shortest rings of length B, etc. If for any entry there is only 1 such shortest ring, the subscript  $\overline{}^{5}$  Unfortunately, not all 3D nets have natural tilings, and some have multiple natural tilings; §3 of Blatov et al. (2007) discusses this

is suppressed; if for any entry there is no ring, a subscript \* is used. The entries are sorted such that smaller rings are listed first, and when two rings of the same size appear, the entry with the smaller subscript is listed first. In the special case of  $\tau = 4$ , the six entries of the vertex symbol are listed as three pairs of entries, with each pair of entries corresponding to opposite pairs of edges, and are otherwise again sorted from smallest to largest. Note that the number of entries in a vertex symbol is  $\tau(\tau - 1)/2$ .

The concepts of *cycles*, *rings*, *strong rings*, *point symbols*, and *vertex symbols* extend immediately to *n* dimensions; for practical considerations (specifically, because the number of entries in a vertex symbol gets unmanageable for large  $\tau$ ), we list the point symbol in Table 3.1 wherever  $\tau \geq 5$ , and the vertex symbol where  $\tau < 4$ . The extension of the tiling concept to *n* dimensions is more delicate, and is discussed further in §4.5.

Following Delgado-Friedrichs et al. (2003a,b), the *regularity* of a 3D net may now be characterized precisely. In short, consider a 3D net with p kinds of vertex and q kinds of edge and whose natural tiling is characterized by r kinds of face and s kinds of tile. Delgado-Friedrichs & Huson (2000) introduced a clear and self-consistent method for characterizing the regularity of such a net simply by forming the array pqrs: examining the 4-digit number so formed, referred to as the *transitivity* of the net, the most "regular" 3D nets are generally those with the smallest transitivity.

Finally, a *minimal net* is a net with the minimum possible number of vertices and edges in its primitive cell<sup>6</sup>; that is, upon deletion of any further edges in the primitive cell, the resulting net breaks into multiple disconnected structures. Beukemann & Klee (1992) establish that there are only 15 such minimal nets in 3D. Delgado-Friedrichs & O'Keeffe (2003) define a 3D net as *barycentric* if every vertex is placed in the center of gravity of its neighbors (to which it is connected by edges). Bonneau et al. (2004), in turn, establish that 7 of the 15 such minimal nets in 3D have *collisions*; that is, when arranged in barycentric fashion, the location of two or more vertices coincide (and, thus, the net is in a sense degenerate). Of the 8 remaining minimal nets

<sup>&</sup>lt;sup>6</sup>A *primitive cell* of a net is the smallest fundamental volume (e.g., hypercube) that, when repeated as an infinite array in all directions with zero spacing, generates the net.

without collision, five are uninodal.

#### **4.2 2D** nets

Consider first the development of uninodal 2D nets with low coordination number. Start from the triangular  $(A_2^* \cong A_2)$  lattice (see Figure P.1a), the corresponding net of which is an array of hexagons, and perform a red/black/blue coloring of the nodes such that no two nearest-neighbor nodes are the same color. If we retain only the red and black nodes, we are left with the *honeycomb packing* (see Figure P.1e), and the corresponding net is an array of hexagons. The coordination number of this stable sphere packing is  $\tau = 3$ , which is less than that of the 2D square packing  $(\tau = 4)$ ; this implies fewer wires in the corresponding computational interconnect application. Unfortunately, the topological density of this net is quite poor, with  $td_{10} = 166$  (that is, with information spreading from one node to only 165 others after a message progresses 10 hops in the network application). We are thus motivated to explore other ways of constructing structured (that is, easy-to-build and easy-to-navigate) nets with low coordination number (that is, with low cost) but high topological density (that is, with a fast spread of information).

Note that the honeycomb packing has a packing density which is less than that of the corresponding triangular and square lattices discussed previously (see Table P.1). If minimization of packing density is the goal<sup>7</sup>, then the honeycomb packing may be *augmented* by replacing every sphere with a set of three spheres in contact, each such set forming an equilateral triangle which touches the neighbors in exactly the same locations as the single sphere which it replaces in the original (unaugmented) packing (see, e.g., Heesch & Laves 1933, Figure 13). The packing density of the resulting stable *augmented honeycomb* packing is less than 2/3 that of the original honeycomb packing (see Table 3.1), and is the rarest uninodal sphere packing available in 2D.

 $<sup>^{7}</sup>$ Note that, for n > 3, the authors are actually unaware of any practical application, other than mathematical curiosity, for which

#### 4.3 A List of Twelve "highly regular" uninodal 3D nets

There are far too many 3D nets to review them all here. We thus identify a List of Twelve highly "regular" (as defined in §4.1, via their transitivity) uninodal 3D nets upon which we will focus our attention and which, following Delgado-Friedrichs et al. (2003a,b), we denote (listing from dense to rare):

- 1. **fcu**: face-centered cubic (FCC), 5. **nbo**: NbO, 9. **cds**: CdSO<sub>4</sub>,
- 2. **bcu**: body-centered cubic (BCC), 6. **dia**: diamond, 10. **bto**: B<sub>2</sub>O<sub>3</sub>,
- 3. **pcu**: cubic, 7. **sod**: sodalite, 11. **ths**: ThSi<sub>2</sub>,
- 4. **qtz**: quartz, 8. **qzd**: quartz dual, 12. **srs**: SrSi<sub>2</sub>. See Table 3.1 for the common names, associated packings, and key characteristics of each<sup>8</sup>. These twelve

nets have been studied thoroughly in the literature, including the landmark work of Wells (1977, 1979, 1983, 1984) and scores of important publications since, including Koch & Fischer (1995, 2006) and the numerous references contained therein; space does not allow a comprehensive review of this broad body of literature here, nor even a comprehensive analysis of these twelve well-studied nets. Suffice it to say here that included in our List of Twelve are the 5 *regular* nets (that is, of transitivity 1111), **bcu**, **pcu**, **nbo**, **dia**, and **srs**, and the 1 *quasiregular* net (of transitivity 1112), **fcu**, as well as 2 of the 14 *semiregular* nets (of transitivity 11*rs*), **qtz** and **sod** (both of which have transitivity 1121), as discussed in O'Keeffe et al. (2000) and Delgado-Friedrichs et al. (2003a,b). Also included in this list are the 5 uninodal minimal nets without collision, **pcu**, **dia**, **cds**, **srs**, and **ths**, the first 4 of which are naturally self-dual, as discussed in Bonneau et al. (2004, Table 1); note that **cds** is of transitivity 1221, and **ths** is of transitivity 12119. The remaining 2 nets on our List of Twelve,

**qzd** (transitivity 1211; see Delgado-Friedrichs et al. 2003c) and **bto** (transitivity 1221; see Blatov 2007), are included because of their close structural relationship to the others, as discussed further in §4.4. We also note

<sup>&</sup>lt;sup>8</sup>Again, clear drawings of each of these nets are available at http://rcsr.anu.edu.au/nets/fcu, where "fcu" may be replaced by any of the lowercase boldface three-letter identifiers given here.

#### 4.3. A LIST OF TWELVE "HIGHLY REGULAR" UNINODAL 3D NETS

considered here.

that four on our List of Twelve, qtz, qzd, bto, and srs, are chiral (that is, these nets twist in such a way that the nets and their reflections are not superposable).

The 12 remaining semiregular nets (of transitivity 11rs) of Delgado-Friedrichs et al. (2003b, Table 1) are

the next natural candidates in this taxonomy (hxg, crs, reo, and rhr might be of particular interest), perhaps followed by the 28 binodal edge-transitive nets (of transitivity 21rs) of Delgado-Friedrichs et al. (2006, Table 1) and the 3 binodal minimal nets without collision (of transitivity 2222, 2211, and 2321) of Bonneau et al. (2004, Table 1) [see also Delgado-Friedrichs & O'Keeffe (2007)]. Note that just half of the List of Twelve considered here (specifically, in order of frequency, dia, pcu, srs, ths, nbo, and cds) account for 66% of the 774 uninodal metal-organic frameworks (MOFs) tabulated in the Cambridge Structural Database (CSD), as reviewed by Ockwig et al. (2005), thus indicating the prevalence in nature of several of the structures

The idea of augmentation, introduced in §4.2, extends directly to many 3D nets in order to reduce packing density. For example, in the (stable) packings related to the dia and sod nets (discussed further in §4.4.1 and §4.4.3 respectively), both of which have coordination number 4, we may replace each sphere with a set of four spheres in contact, each such set of spheres forming a tetrahedron, creating what is referred to as the augmented diamond (dia-a) and augmented sodalite (sod-a) nets. In the case of the augmentation of the packing related to the dia net, each tetrahedral set touches the neighbors in exactly the same locations as the single sphere which it replaces in the original (unaugmented) packing (see Heesch & Laves 1933, Figure 12). In the case of the augmentation of the packing related to the sod net, as the angles between the 4 nearest neighbors of any node are not uniform in the sod net, each tetrahedral set is slightly larger than the single

sphere which they replace in the original (unaugmented) packing, and the contact points are slightly shifted (O'Keeffe 1991b); note that the packing associated with the sod-a net is the rarest uninodal stable 3D packing currently known. On the other hand, in the augmentation of the (unstable) packing related to the srs net, which has coordination number 3, we may replace each sphere with a set of three spheres in contact, each such set of spheres, as in the augmentation of the honeycomb packing, forming an equilateral triangle and touching the neighbors in exactly the same locations as the single sphere which it replaces in the original (unaugmented) packing (see Heesch & Laves 1933, Figure 10); note that the packing associated with the resulting **srs-a** net is the rarest uninodal unstable 3D packing known.

Comparing augmented honeycomb to honeycomb, **dia-a** (transitivity 1222) to **dia**, **sod-a** (transitivity 1332) to **sod**, and **srs-a** (transitivity 1221) to **srs**, it is seen that augmentation, while reducing the packing density  $\Delta$  (see Table 3.1), also significantly reduces both the topological density,  $td_{10}$ , and the regularity of the resulting net. Thus, the process of augmentation appears to be of little interest for the purpose of developing efficient computational interconnects. Note that Fischer (2005) and Dorozinski & Fischer (2006) show that the process of augmentation can be repeated indefinitely in order to obtain (non-uninodal) sphere packings of arbitrarily low packing density.

Finally, there are two other 3D nets which, though less regular than our List of Twelve, are worthy of "honorable mention": hexagonal close packing (hcp, transitivity 1232) and lonsdaleite (lon, transitivity 1222). As hinted by their identical packing densities (see Table 3.1), hcp is closely related to fcu, and lon is closely related to dia; curiously, both have slightly higher values of  $td_{10}$  than do their more regular cousins. The relations between these two pairs of packings is readily evident when they are considered as built up in layers, as introduced in the second paragraph of §2.4 and discussed further below.

The  $A_3$  lattice (a.k.a. FCC, corresponding to the **fcu** net) may be built up as an alternating series of three 2D triangular ( $A_2$ ) layers, offset from each other in the form abcabc..., with the nodes in one layer over the holes in the layer below; **hcp** is built up similarly, but with two alternating layers, offset from each other in the form abab...

Similarly, the sphere packings corresponding to the **dia** and **lon** nets may be built up as alternating series of approximately 2D honeycomb layers offset from each other. These honeycomb "layers" are in fact not quite 2D; if the nodes in a single layer are marked with an alternate red/black coloring, the red nodes are

#### 4.4. UNINODAL EXTENSION OF SOME REGULAR 3D NETS TO HIGHER DIMENSIONS

raised a bit and the black nodes lowered a bit. In both packings, the layers are offset from each other, with the lowered nodes in one layer directly over the raised nodes in the other. In the packing corresponding to the **dia** net, there are three such alternating layers stacked in the form abcabc...; in the packing corresponding to the **lon** net, there are two such alternating layers stacked in the form abab...

#### 4.4 Uninodal extension of some regular 3D nets to higher dimensions

The **fcu** net is based on the  $D_3 \cong A_3$  lattice, and thus may be extended to n dimensions in two obvious ways (that is, via  $A_n$  or  $D_n$ ). The **bcu** net is based on the  $D_3^* \cong A_3^*$  lattice, and thus may also be extended to n dimensions in two obvious ways (via  $A_n^*$  or  $D_n^*$ ). The **pcu** net is based on the  $\mathbb{Z}^3$  lattice, and thus extends to n dimension via  $\mathbb{Z}^n$ . This section explores how most of the other nets on the List of Twelve described above extend naturally to higher dimensions.

It is important to recall that the nets in the  $D_n^*$  case for n > 4 turn out to be a bit peculiar, as discussed further in §2.3; the  $T_n^{90}$  and  $T_n^{60}$  nets encountered in §4.4.7 are similar.

#### **4.4.1** Extending dia: the $A_n^+$ and $D_n^+$ packings

The **dia** net may be obtained from the well-known  $D_3^+$  packing defined in (2.5) (see also Sloane 1987), and thus extends naturally to n dimensions as  $D_n^+$ . However, there is an alternative construction of the **dia** net, described below and denoted  $A_n^+$ , which is equivalent to  $D_n^+$  for n=3 but extends to n dimensions differently. In fact, a third extension of the **dia** net to n dimensions, the  $V_n^{90}$  construction, is introduced in §4.4.6. These alternative extensions of the **dia** net to n dimensions, with low coordination number, are perhaps better suited than  $D_n^+$  for many practical applications. We thus stress that names such as "n-dimensional diamond" are parochial, as there are sometimes multiple "natural" n-dimensional extensions of a net related to a given three-dimensional crystalline structure (e.g.,  $D_n^+$ ,  $A_n^+$ , and  $V_n^{90}$ ). For n-dimensional nets in general, we thus strongly prefer names derived from a corresponding well-defined n-dimensional lattice or, when such a name

is not available, names evocative of their *n*-dimensional construction; this preference is in sharp contrast with the names suggested by O'Keeffe (1991b).

Recall the first paragraph of §4.2. Now start from a BCC  $(A_3^* \cong D_3^*)$  lattice and perform a red/black/blue/yellow coloring of the points such that no two nearest-neighbor points are the same color. If we retain only the red and black points, we are left with the diamond packing. The coordination number of this packing is  $\tau = 4$ , which is less than that of the 3D cubic packing  $(\tau = 6)$ , but also has a reduced topological density, as quantified by  $td_{10}$  (see Table 3.1). The diamond packing also has a packing density which is less than that of the corresponding FCC, BCC, and cubic lattices.

Note in general [see (2.7a)] that  $A_n^*$  may be defined as the union of n+1 shifted  $A_n$  lattices, which is analogous to the property [see (2.4a)] that  $D_n^*$  may be defined as the union of 4 shifted  $D_n$  lattices. Recall from (2.5) that  $D_n^+$ , which we referred to the *offset checkerboard packing*, was defined as the union of just 2 shifted  $D_n$  lattices, and generates the diamond packing in 3D (where  $D_3 \cong A_3$ ). Motivated by the previous paragraph and the first paragraph of §4.2, we are thus also keenly interested in the nonlattice packing considered in Table 1 of O'Keeffe (1991b), denoted here  $A_n^+$  and referred to as the *offset zero-sum packing*, and which is defined as the union of just 2 shifted  $A_n$  lattices [cf. (2.5), (2.7)]:

$$A_n^+ = A_n \cup ([1] + A_n), \text{ where } [1]_k = \begin{cases} \frac{1}{n+1} & k \le n, \\ \frac{-n}{n+1} & k = n+1. \end{cases}$$
 (4.1)

The coordination number of the regular uninodal nonlattice packing  $A_n^+$  is n+1, with these n+1 nearest neighbors forming a regular n-dimensional simplex [that is, a regular n-dimensional polytope with n+1 vertices—e.g., in n=3 dimensions, a tetrahedron]. The generalization of the honeycomb and diamond packings to higher dimensions given by  $A_n^+$  is significant, as it illustrates how a highly regular stable packing with coordination number lower than that of the corresponding Cartesian lattice may be extended to dimension n>3. Note also that the nonlattice packings  $A_n^+$  are distinct from the lattice packings  $A_n^+$  defined in (2.8),

#### 4.4. UNINODAL EXTENSION OF SOME REGULAR 3D NETS TO HIGHER DIMENSIONS

which are generated in a similar manner.

#### **4.4.2** Augmenting $A_n^+$ : the $\widehat{A}_n^+$ packing

The third paragraph of §4.3 discusses the augmentation of the  $A_3^+$  packing, replacing each sphere with a tetrahedral set of 4 smaller spheres. This idea extends immediately to the augmentation, in n dimensions, of the  $A_n^+$  packing discussed above, replacing each (n-dimensional) sphere with a regular n-dimensional simplex of n+1 smaller spheres.

#### **4.4.3** Extending sod: the ${}^{\mathsf{T}}\!A_n^*$ packing

packing, with the nodes of the net located at the *holes* of the packing rather than at the centers of the spheres of the packing. As noted by O'Keeffe (1991b), this construction extends immediately to the *n*-dimensional net formed by the Voronoï tesselation of space via the  $A_n^*$  packing. Constructing the  $A_n^*$  packing as defined in §2.4, the holes of this packing that are nearest to the origin (that is, in its Voronoï tesselation, the corners of the Voronoï cell which contains the origin) are given by the (n+1)! permutations of the vector (see Conway & Sloane, 1999, page 474):

The familiar sod net is formed by the edges of the Voronoï tesselation of space formed by the  $A_3^*$  (that is, BCC)

$$\frac{1}{2(n+1)} \begin{pmatrix} -n & -n+2 & -n+4 & \dots & n \end{pmatrix}^T.$$

These nodal points [which, like the lattice points of  $A_n^*$  itself, are defined in an (n+1)-dimensional space, but all lie in the *n*-dimensional subspace orthogonal to the vector  $\mathbf{n}_{A_n}$  defined in (2.6b)] are equidistant from their n+1 nearest neighbors, and form *permutohedra* (in 3D, *truncated octahedra*) which tile *n*-dimensional space.

Note that these nodal points themselves form a uninodal sphere packing with coordination number  $\tau = n + 1$ ; due to its relationship to the *tesselation* of space via the points of the  $A_n^*$  packing, we thus introduce the notation  ${}^{\mathsf{T}}A_n^*$  for this packing.

The **nbo** net, a subset of the **pcu** net, has an obvious uninodal extension to n dimensions with  $\tau = 4$ , which may be visualized as the contact graph formed by repeating a unit hypercube pattern as an infinite array with unit spacing (see Figure 5.3), where each hypercube itself has two paths which "snake" along the edges from the  $(0,0,\cdots,0,0)$  node to the  $(1,1,\cdots,1,1)$  node, one coordinate direction at a time; we thus suggest the symbol  $S_n$  to denote this construction. These two paths touch at the opposite corners of the unit hypercube:

```
path A: (0,0,\cdots,0,0), (0,0,\cdots,0,1), (0,0,\cdots,1,1), ..., (0,1,\cdots,1,1), (1,1,\cdots,1,1), and
path B: (0,0,\cdots,0,0), (1,0,\cdots,0,0), (1,1,\cdots,0,0), ..., (1,1,\cdots,1,0), (1,1,\cdots,1,1).
```

#### Extending ths and bto: the $Y_n^{90}$ and $Y_n^{60}$ constructions

illustrated in Figure 4.1a, starting with this Y stencil and adjoining translates of itself, tip to tip, builds up the honeycomb packing in 2D. Extending this idea to 3D, as illustrated in Figure 4.1b, we may "twist" the Y stencil by  $90^{\circ}$  at each level: starting with the basic Y stencil in, say, the  $e^{1}-e^{2}$  plane, we can shift to the right (in  $e^1$ ) and adjoin Y stencils twisted by 90° (that is, aligned in the  $e^1$ - $e^3$  plane), then shift to the right again and adjoin Y stencils twisted again (aligned in the  $e^1-e^2$  plane), etc. This construction forms the ths net in 3D, and extends immediately to dimension n > 3; we thus denote this construction  $Y_n^{90}$ .

The honeycomb packing  $A_2^+$ , of coordination number  $\tau = 3$ , contains a fundamental Y-shaped stencil. As

Instead of twisting the Y stencil by 90° at each step, we may instead twist it by 60°. This forms the bto net in 3D. As with the hcp versus fcu and lon versus dia nets in 3D, as described at the end of §4.3, there is a bit of flexibility in terms of the ordering of the the successive twists for n > 3. A highly regular net for odd n, which we denote  $Y_n^{(0)}$ , is formed by pairing off the dimensions after the first and alternating the twists as follows: starting with the basic Y stencil in, say, the  $e^1$ - $e^2$  plane, we continue by adjoining Y stencils in the  $e^1-e^4$  plane, then in the  $e^1-e^6$  plane, etc. We then adjoin Y stencils in the  $e^1-z_{23}^{60}$  plane, where  $z_{23}^{60}$  is the vector formed by rotating the  $e^2$  unit vector  $60^\circ$  in the direction towards  $e^3$ ; we continue by adjoining Y stencils in

#### 4.4. UNINODAL EXTENSION OF SOME REGULAR 3D NETS TO HIGHER DIMENSIONS

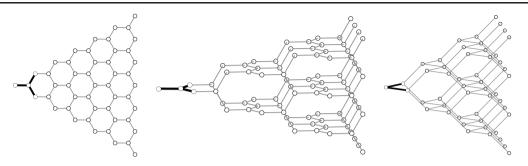


Figure 4.1: Construction of three rare packings: (left) the  $Y_2$  (honeycomb) net, (center) the  $Y_3^{90}$  (ths) net, and (right) the  $V_3^{90}$  (dia) net. All three constructions build from left to right in the above figures from a basic "Y" or "V" stencil, and have obvious extensions to higher dimensions.

the  $\mathbf{e}^1$ - $\mathbf{z}_{45}^{60}$  plane, then in the  $\mathbf{e}^1$ - $\mathbf{z}_{67}^{60}$  plane, etc. Next, we adjoin Y stencils in the  $\mathbf{e}^1$ - $\mathbf{z}_{23}^{120}$  plane, where  $\mathbf{z}_{23}^{120}$  is the vector formed by rotating the  $\mathbf{z}_{23}^{60}$  vector  $60^\circ$  further in the  $\mathbf{e}^2$ - $\mathbf{e}^3$  plane; we continue by adjoining Y stencils in the  $\mathbf{e}^1$ - $\mathbf{z}_{45}^{120}$  plane, then in the  $\mathbf{e}^1$ - $\mathbf{z}_{67}^{120}$  plane, etc., and repeat (that is, with stencils again aligned in the  $\mathbf{e}^1$ - $\mathbf{e}^2$  plane).

The  $Y_n^{90}$  and  $Y_n^{60}$  constructions have a parameter, denoted  $\alpha$  and defined as half of the angle between the two top branches of the Y stencil (thus,  $\alpha \to 0^\circ$  closes down the Y to an I, whereas  $\alpha \to 90^\circ$  opens up the Y to a T). The Voronoï volume of the  $Y_n^{90}$  and  $Y_n^{60}$  constructions may be written as simple functions of  $\alpha$  as

follows:

$$\begin{array}{l} \mathcal{V}_{\mathsf{Y}_n^{90}}(\alpha) = f_{\mathsf{Y}_n}(\alpha) \, \mathcal{V}_{\mathsf{Y}_n^{90}}(\bar{\alpha}) \\ \\ \mathcal{V}_{\mathsf{Y}_n^{60}}(\alpha) = f_{\mathsf{Y}_n}(\alpha) \, \mathcal{V}_{\mathsf{Y}_n^{60}}(\bar{\alpha}) \end{array} \right\} \quad \text{with} \quad \bar{\alpha} = 45^\circ, \quad f_{\mathsf{Y}_n}(\alpha) = (2 - \sqrt{2}) \, (1 + \cos \alpha) \, (\sqrt{2} \sin \alpha)^{n-1}.$$

This relation is plotted in Figure 4.2a. The characteristics of  $Y_n^{90}$  and  $Y_n^{60}$  reported in Table 3.1 are computed for  $\alpha = \cos^{-1}(1/n)$ , as marked with circles in Figure 4.2a, which maximizes the Voronoï volume and, thus, minimizes the packing density. An alternative natural choice is  $\alpha = 60$ , which results in barycentric constructions of  $Y_n^{90}$  and  $Y_n^{60}$ .

#### 4.4.6 Extending dia and qtz: the $V_n^{90}$ and $V_n^{60}$ constructions

The  $V_n^{90}$  and  $V_n^{60}$  constructions are defined in an identical manner as their  $Y_n^{90}$  and  $Y_n^{60}$  counterparts, with a V stencil replacing the Y stencil (see, e.g., Figure 4.1c), thus resulting in nets with coordination number  $\tau = 4$  instead of  $\tau = 3$ . These constructions lead to the dia and qtz nets in 3D.

As with the  $Y_n^{90}$  and  $Y_n^{60}$  construction, the  $V_n^{90}$  and  $V_n^{60}$  constructions have a parameter, denoted  $\alpha$  and defined as half of the angle between the two top branches of the V stencil. The Voronoï volume of the  $V_n^{90}$  and  $V_n^{60}$  constructions may be written as simple functions of  $\alpha$  as follows:

$$\frac{\mathcal{V}_{\mathsf{V}_n^{90}}(\alpha) = f_{\mathsf{V}_n}(\alpha) \, \mathcal{V}_{\mathsf{V}_n^{90}}(\bar{\alpha})}{\mathcal{V}_{\mathsf{V}_n^{60}}(\alpha) = f_{\mathsf{V}_n}(\alpha) \, \mathcal{V}_{\mathsf{V}^{60}}(\bar{\alpha})} \right\} \quad \text{with} \quad \bar{\alpha} = 45^\circ, \quad f_{\mathsf{V}_n}(\alpha) = 2^{n/2} \cos \alpha (\sin \alpha)^{n-1}.$$

This relation is plotted in Figure 4.2b. The characteristics of  $V_n^{90}$  and  $V_n^{60}$  reported in Table 3.1 are computed for  $\alpha = \cos^{-1}(1/\sqrt{n})$ , as marked with circles in Figure 4.2a, which maximize the Voronoï volumes and, thus, minimize the packing density. Note that the  $V_n^{90}$  and  $V_n^{60}$  constructions are barycentric for any  $\alpha$  in the range  $0 < \alpha < 90^{\circ}$ .

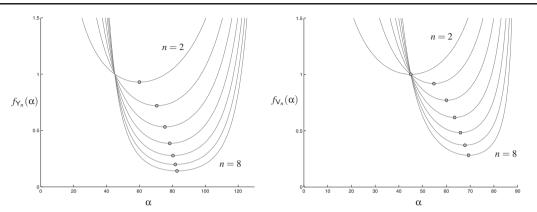


Figure 4.2: Variation of the Voronoï volume of the (left)  $Y_n^{90}$  &  $Y_n^{60}$  and (right)  $V_n^{90}$  &  $V_n^{60}$  packings as a function of  $\alpha$  for n = 2 to n = 8.

#### 4.4.7 Extending cds and qzd: the $T_n^{90}$ and $T_n^{60}$ constructions

The  $T_n^{90}$  and  $T_n^{60}$  constructions are defined in an analogous manner as their  $Y_n^{90}$ ,  $V_n^{90}$ ,  $Y_n^{60}$ , and  $V_n^{60}$  counterparts, and lead to the **cds** and **qzd** nets in 3D. The only difference now is that, instead of adjoining two new Y or V symbols on the tips of each Y or V symbol in the previous layer, we now adjoin a single new T symbol centered atop each T symbol in the previous layer, appropriately twisted; these constructions thus result in nets with coordination number  $\tau = 4$ . Note that the "horizontal" and "vertical" distances between nodes in these constructions are equal, and that these constructions are parameter free and barycentric.

Note that the  $x_1$  direction is special in the  $Y_n^{90}$ ,  $Y_n^{60}$ ,  $V_n^{90}$ ,  $V_n^{60}$ ,  $V_n^{90}$ , and  $V_n^{60}$  constructions. These constructions are configured in this way intentionally, in order to construct equilibrium packings; however, other variations are certainly possible. Note also that the  $Y_n^{60}$ ,  $V_n^{60}$ , and  $V_n^{60}$  constructions involve pairing off the dimensions after the first and rotating in each pair of dimensions  $V_n^{60}$ , at a time, in the manner described in  $V_n^{60}$ ,  $V_n^{60}$ , and  $V_n^{60}$  nets, respectively, as defined previously.

Note also that the  $Y_n^{90}$ ,  $V_n^{90}$ , and  $T_n^{90}$  constructions form square layers in the  $\mathbf{e}_2$ - $\mathbf{e}_3$  plane, the  $\mathbf{e}_4$ - $\mathbf{e}_5$  plane, etc., whereas the  $Y_n^{60}$ ,  $Y_n^{60}$ , and  $T_n^{60}$  constructions form triangular layers in these planes. In the resulting  $Y_n^{90}$ ,  $Y_n^{60}$ ,  $Y_n^{90}$ , and  $Y_n^{60}$  nets, there are, in fact, no edges of the net within these layers (that is, all of the edges connect nodes in different layers). On the other hand, in the resulting  $T_n^{90}$  and  $T_n^{60}$  nets, each node is connected via edges of the net to exactly two others (note: *not* four or six) within these layers. As with the peculiar  $D_n^*$  net discussed previously, the  $T_n^{90}$  and  $T_n^{60}$  constructions are, in fact, *not* contact graphs of the corresponding sphere packings  $T_n^{10}$  in order to form the  $T_n^{90}$  and  $T_n^{60}$  nets.

#### 4.4.8 Other extensions

Sections 4.4.1 through 4.4.7 summarize several uninodal families of n-dimensional extrapolations of some common 3D nets; most of these (unless indicated otherwise, via references to existing literature) are new. Note that O'Keeffe (1991b) mentions two other such extensions, one corresponding to the **lon** net and one corresponding to the **sod-a**, the latter of which is currently the rarest uninodal stable packing known for n > 3 (and which, consistent with the above developed naming conventions, we might suggest to identify as  ${}^{T}\widehat{A}_{n}^{*}$ ).

<sup>&</sup>lt;sup>10</sup>Note that there is a lower-symmetry form of **cds** in 3D with four nearest neighbors per node whose contact graph does generate the **cds** net; see Delgado-Friedrichs (2005, Figure 1). Lower symmetry forms of other  $T_n^{90}$  and  $T_n^{60}$  constructions, whose nets are contact graphs, might also exist.

Beukemann & Klee (1992, page 50) mentions two extensions of their own (at least, to n = 4), both related to the **dia** net. Judging from the vast assortment of distinct rare sphere packings and related nets available in 3D, there are certainly *many* more uninodal extensions to higher dimensions of regular rare 3D packings that are still awaiting discovery; we have focused our attention here on what appear to be several of the most regular. The regularity of n-dimensional nets for n > 3 is discussed further below.

#### **4.5** Regularity and transitivity of *n*-dimensional nets for n > 3

As reviewed in §4.1, the regularity of a 3D net is defined based on its transitivity, which in turn is based on the so-called natural tiling of the 3D net. The natural tiles of 3D nets have been thoroughly characterized in the literature for all of the most regular 3D nets available. In §4.4, we described uninodal extensions of several regular 3D nets to higher dimensions, and mentioned that many more such uninodal nets with n > 3 most certainly exist. The natural question to ask, then, is how the concepts of regularity and transitivity can be extended to higher dimensions, so that we may differentiate between these nets and identify those which are the most regular.

This question is difficult to visualize in dimensions higher than three, and requires a symbolic/numerical

description of the net to proceed. The net arising from the  $\mathbb{Z}^n$  lattice for  $n=4,5,\ldots$ , which is naturally tiled by n-dimensional hypercubes, is by far the easiest starting point. Denote first the symbols  $\{v,w,x,y,z\}$  as variables that range from 0 to 1. The 3D unit cube, denoted  $\{xyz\}$ , has six faces,  $\{xy0,xy1,x0z,x1z,0yz,1yz\}$ . Each face, in turn, has four edges; e.g.,  $\{0yz\}$  has edges  $\{0y0,0y1,00z,01z\}$ . Finally, each edge connects two nodes; e.g.,  $\{00z\}$  connects nodes  $\{000,001\}$ . The 4D unit hypercube,  $\{wxyz\}$ , has eight 3-faces, which we identify as  $\{wxy0,wxy1,wx0z,wx1z,w0yz,w1yz,0xyz,1xyz\}$ , each 3-face has six 2-faces, each 2-face has four edges, and each edge connects two nodes; etc.

In 3D, as reviewed in §4.1, the transitivity is based on the number of distinct nodes, edges, (2D) faces,

and (3D) tiles. By analogy, then, in 4D we may define the transitivity of a net based on the number of distinct nodes, edges, 2-faces, 3-faces, and (4D) tiles in the natural tiling. Similarly, in 5D, we may define the transitivity based on the number of distinct nodes, edges, 2-faces, 3-faces, 4-faces and (5D) tiles in the natural tiling; etc. Via this definition, the net derived from the  $\mathbb{Z}^4$  lattice has transitivity 11111, the net derived from the  $\mathbb{Z}^5$  lattice has transitivity 111111, etc.

For all of the other nets with n > 3 listed in Table 3.1, the computation of the transitivity remains an important unsolved problem. Note that, in a tiling corresponding to a 3D net, the (2D) faces of the (3D) tiles are, in general, minimal surfaces stretched over non-planar frames built from (1D) edges between several

nodal points defined in 3D. In a tiling corresponding to an n-dimensional net for n > 3, the 2-faces of the tiles are, in general, minimal surfaces stretched over nonplanar frames between several nodes defined in n dimensions. [Note that the computation of such minimal surfaces in n dimensions is straightforward using standard level set methods; see, e.g., Cecil (2005).] Several of these nonplanar 2-faces combine to form the boundaries of each 3-face, which itself is not confined to lie within a 3D subspace of the n-dimensional domain. Several of these 3-faces then combine to form the boundaries of each 4-face; etc.

Identification of such high-dimensional natural tilings is apparently a task that could be readily accomplished numerically, but is, in general, expected to be difficult to visualize.

#### **Chapter 5**

### **Coding theory**

Introduction

Contents

5.1	Introduction	UJ
5.2	Exemplary linear binary codes (LBCs)	72
	501 Dt	70

62

3.2	Excin	mary finear binary codes (LDCs)	• 1	-
	5.2.1	Binary single parity-check codes	. 7	2
	522	Ringry repetition codes	7	12

3.2.2	Binary repetition codes
5.2.3	Binary Hamming codes

5.2.4	Binary simplex codes
505	Enter ded binom Honorina and a

# 

77

	3.2.1	Biliary quadratic residue codes					
	5.2.8	Extending, puncturing, and shortening					
5.3	Exemplary linear ternary codes (LTCs)						
	5.3.1	Ternary single parity-check codes					
	5.3.2	Ternary repetition codes					
	5.3.3	Ternary Hamming codes					
	5.3.4	Ternary simplex codes					
	5.3.5	Ternary quadratic residue codes					
5.4	Exemp	plary linear quaternary codes (LQCs)					
	5.4.1	Quaternary single parity-check codes					
	5.4.2	Quaternary repetition codes					
	5.4.3	Quaternary Hamming codes					
	5.4.4	Quaternary simplex codes					

 Algebraic decoding
 88

 Cyclic form
 91

5.2.6

527

5.4.5

5.5.1

5.5.2.

5.5.3

5.5.4

Ripary quadratic residue codes

#### 5.1 Introduction

Though the lattices that arise from n-dimensional sphere packings have connections that permeate many foundational concepts in number theory and pure geometry, the list of successful direct applications in science and engineering of n-dimensional sphere packings with n > 3 is currently surprisingly short<sup>1</sup>; this list includes

- the numerical evaluation of integrals (Sloan & Kachoyan 1987),
- $\bullet \ \ the \ solution \ of \ the \ linear \ Diophantine \ inequalities \ that \ arise \ in \ integer \ linear \ programming \ (Schrijver \ 1986),$
- the characterization of crystals with curious five-fold symmetries (Janssen 1986),
  attempts at unifying the 4 fundamental forces (in 10, 11, or 26 dimensions) via superstring theory (Kaku 1990), and
- 1999), and
  the development of maximally effective numerical schemes to address an information-theoretic interference suppression problem known as the Witsenhausen counterexample (Grover, Sahai, & Park 2010).

Far and away the most elegant and practical application of n-dimensional sphere packings, however, is in the framing and understanding of error correcting codes (ECCs). The reader is referred to MacWilliams & Sloane (1977), Thompson (1983), Pless (1998), Conway & Sloane (1998), and Morelos-Zaragoza (2006) for some comprehensive reviews of this fascinating subject. A brief overview of this field is given here to emphasize the existing practical relevance of n-dimensional sphere packings with n > 3; we aim to augment this list of practical applications significantly in Parts II and III of this text, based heavily on the various aspects of n-dimensional sphere packing theory reviewed and extended in Part I.

To proceed, define  $\mathbf{F}_q$  [also denoted GF(q)] as the set of symbols in a *finite field* (a.k.a. *Galois field*) of order q, where  $q = p^a$  with p prime, and define  $\mathbf{F}_q^n$  as the set of all vectors of order n with elements

<sup>&</sup>lt;sup>1</sup>Notably, Conway & Sloane (1998, page 12) state: "A related application that has not yet received much attention is the use of these packings for solving *n*-dimensional *search* or *approximation* problems"; this is exactly the problem focused on in our Part II.

selected from  $\mathbf{F}_q$ . The cases of particular interest in this work are the *binary field*  $\mathbf{F}_2 = \{0,1\}$ , the *ternary field*  $\mathbf{F}_3 = \{0,1,2\}$ , and the *quaternary field*  $\mathbf{F}_4 = \{0,1,\omega,\bar{\omega}\}$ , where, as in §2.1,  $\omega = (-1+\mathrm{i}\sqrt{3})/2$  [note that  $\omega^2 = \bar{\omega}$ ,  $\bar{\omega}^2 = \omega$ , and  $\bar{\omega} \cdot \omega = 1$ ]. In a finite field  $\mathbf{F}_q$ , addition (+) and multiplication (·) are closed (that is, they map to elements within the field) and satisfy the usual rules: they are associative, commutative, and distributive, there is a 0 element such that a+0=a, there is a 1 element such that  $a\cdot 1=a$ , for each a there is an element (-a) such that a+(-a)=0, and for each  $a\neq 0$  there is an element  $a^{-1}$  such that  $a\cdot a^{-1}=1$ . If a is itself prime (e.g., if a0 e 2 or a1, then standard integer addition and multiplication mod a2 forms a finite field. If not (e.g., if a1, a bit more care is required in order to obtain closure within the finite field while respecting these necessary rules on addition and multiplication. For the cases considered in this section (specifically, a2, a3, and a3, addition and multiplication on a4, addition and multiplication as follows:

	+    0   1		. 11 (	) I 1			+	0	1	2		0	1	2
$\mathbf{F}_2$ :	0 0 1	=			=	$\mathbf{F}_3$ :	0	0	1	2	0	0	0	0
<b>r</b> 2.	1 1 0	0 0 0			13.	1	1	2	0	1	0	1	2	
	1    1   0	1	1 11 0	, 1			2	2	0	1	2	0	2	1
		+	0	1	ω	Ū		0	1	ω	Ū			
		0	0	1	ω	Ō	0	0	0	0	0			
	$\mathbf{F}_4$ :	1	1	0	Θ	ω	1	0	1	ω	Ū			
		ω	ω	Ō	0	1	ω	0	ω	Φ	1			
		Ø	Ō	ω	1	0	Ø	0	Ø	1	ω			

A vector in  $\mathbf{F}_q^n$  is a vector of length n with each element in  $\mathbf{F}_q$ . The *Hamming distance* between two such vectors is the number of elements that differ between them.

 $<sup>^2</sup>$ We limit our attention in the quaternary case to codes designed over the finite field  $\mathbf{F}_4$ ; though there is some attention in the literature to codes defined over  $\mathbb{Z}_4$  [that is, over the integers mod 4], codes defined over finite fields turn out to be, in a sense, more natural.

#### 5.1. INTRODUCTION

An  $[n,k]_q$  (if d is specified,  $[n,k,d]_q$ ) q-ary linear d code (LC) is defined via a set of d independent basis vectors  $\mathbf{v}^i \in \mathbf{F}_q^n$ . The d distinct codewords  $\mathbf{w}^i \in \mathbf{F}_q^n$  of the LC are given by all d-ary linear combinations of the basis vectors  $\mathbf{v}^i$  (that is, by all linear combinations with coefficients selected from  $\mathbf{F}_q$ , with addition and multiplication defined elementwise on  $\mathbf{F}_q$ ). The basis vectors  $\mathbf{v}^i$  are generally selected such the minimum distance d of the LC (that is, the minimum Hamming distance between any two resulting codewords) is maximized.

This work focuses on cases with q=2 [termed a linear binary code (LBC)], q=3 [termed a linear ternary code (LTC)], and q=4 [termed a linear quaternary code (LQC)]. In cases with q=2, which are common, we frequently write simply [n,k] or [n,k,d], dropping the q subscript. We denote by  $V_{[n,k]_q}$  (or  $V_{[n,k,d]_q}$ ) the  $n\times k$  basis matrix with the k basis vectors  $\mathbf{v}^i$  as columns, and by  $W_{[n,k]_q}$  (or  $W_{[n,k,d]_q}$ ) the  $n\times q^k$  codeword matrix with the  $q^k$  codewords  $\mathbf{w}^i$  as columns. Without loss of generality, we write  $V_{[n,k]_q}$  and a companion  $(n-k)\times n$  parity-check matrix  $H_{[n,k]_q}$  in the standard (a.k.a. systematic) form<sup>4</sup>

$$H_{[n,k]_q} = \begin{bmatrix} -P_{(n-k)\times k} & I_{(n-k)\times (n-k)} \end{bmatrix}, \quad V_{[n,k]_q} = \begin{bmatrix} I_{k\times k} \\ P_{(n-k)\times k} \end{bmatrix}, \quad \mathbf{w}^i = \begin{bmatrix} \mathbf{d}^i \\ \mathbf{b}^i \end{bmatrix}. \tag{5.1}$$

<sup>&</sup>lt;sup>3</sup>Nonlinear *q*-ary codes also appear in the literature, in which the valid codewords are *not* simply linear combinations of a set of basis vectors and must be enumerated differently. Such codes, which are related to nonlattice packings, are in general more difficult to decode than LCs, and are not considered further here.

<sup>&</sup>lt;sup>4</sup>In the literature on this subject, it is more common to use a "generator matrix" G to describe the construction of linear codes. The "basis matrix" convention V used here is related simply to the corresponding generator matrix such that  $V = G^T$ ; we find the basis matrix convention to be more natural in terms of its linear algebraic interpretation.

When written in systematic form, each of the data vectors  $\mathbf{w}^i$  block decomposes into its k data symbols  $\mathbf{d}^i$  and its r = n - k parity symbols  $\mathbf{b}^i$ ; note that r is sometimes called the redundancy of the code. Note also that  $H_{[n,k]_q}V_{[n,k]_q} = 0$  (on  $\mathbf{F}_q$ )<sup>6</sup>, which establishes that the basis vectors  $\mathbf{v}^i$  so constructed [and, thus, all of the resulting codewords  $\mathbf{w}^i$ ] each satisfy the parity-check equations,  $H_{[n,k]_q}\mathbf{w}^i = 0$  (on  $\mathbf{F}_q$ ), as implied by the rows of  $H_{[n,k]_q}$  and illustrated by the several examples given in systematic form in §5.2, §5.3, and §5.4. Note further that, for LBCs and LQCs, P = -P.

The key to designing a "good"  $[n,k]_q$  LC is to construct the *parity submatrix*  $P_{(n-k)\times k}$  in (5.1) in such a way that the minimum distance d of the resulting code is maximized for given values of n, k, and q. Indeed, the problem of designing a good binary error correcting code is essentially a finite sphere packing problem on  $\mathbf{F}_2$ ; thus the very close relationship between the design of error-correcting codes and the design of infinite dense sphere packings in  $\mathbb{R}^n$ , as discussed in §2.

For  $q=p^a$  with p prime, *conjugation* in  $\mathbf{F}_q$  (that is, for a scalar  $v \in \mathbf{F}_q$ ) is defined such that  $\bar{v}=v^p$ ; conjugation in  $\mathbf{F}_q^n$  (that is, for vectors  $\mathbf{v} \in \mathbf{F}_q^n$ ), as well as for matrices formed with a number of such vectors as columns, is performed elementwise. Any  $[n,k]_q$  linear code C has associated with it an  $[n,n-k]_q$  dual code  $C^{\perp}$  defined [cf. (2.1)] such that

$$C^{\perp} = \left\{ \mathbf{w} \in \mathbf{F}_q^n : \mathbf{w} \cdot \bar{\mathbf{u}} = 0 \text{ for all } \mathbf{u} \in C \right\}.$$
 (5.2)

The parity-check and codeword matrices of  $C^{\perp}$  may be written in systematic form as

$$H_{[n,n-k]_q}^{\perp} = \begin{bmatrix} \bar{P}^T & I_{(n-k)\times(n-k)} \end{bmatrix}, \quad V_{[n,n-k]_q}^{\perp} = \begin{bmatrix} I_{(n-k)\times(n-k)} \\ -\bar{P}^T \end{bmatrix}. \tag{5.3}$$

<sup>&</sup>lt;sup>5</sup>The word "bit", a portmanteau word for "binary digit", is generally reserved for the case with q=2; in the general case, we use the word "symbol" in its place.

<sup>&</sup>lt;sup>6</sup>The qualifiers "(on  $\mathbf{F}_q)$ " and "(mod q)" are used, as appropriate, to remind the reader that multiplication and addition in the equation indicated are performed elementwise on the finite field  $\mathbf{F}_q$ , as discussed above.

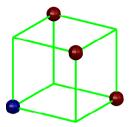
### 5.1. INTRODUCTION

where  $\bar{P}$  denotes conjugation in  $\mathbf{F}_q$  of each element of the parity submatrix P of the original  $[n,k]_q$  linear code C. Note that  $\bar{P}^T$  is of order  $k \times (n-k)$ , and, of course, that  $H^{\perp}_{[n,n-k]_q}V^{\perp}_{[n,n-k]_q}=0$  (on  $\mathbf{F}_q$ ). Note further that, for LBCs and LTCs,  $\mathbf{u}=\bar{\mathbf{u}}$  and  $P=\bar{P}$ .

A self-dual code C is a code for which the transpose of the codeword matrix V results in a new matrix H which is itself the parity-check matrix of a code which is equivalent to C, albeit not in systematic form.

Graphically, the codewords of an  $[n,k,d]_2$  LBC may be thought of as a carefully chosen subset of  $2^k$  of the  $2^n$  corners on a single n-dimensional unit hypercube, as illustrated for n=3 in Figure 5.1, whereas an  $[n,k,d]_3$  LTC may be thought of as a subset of  $3^k$  of the  $3^n$  gridpoints in a cluster of  $2^n$  unit hypercubes in n-dimensions, as illustrated for n=3 in Figure 5.2. For any q, d quantifies the minimum number of symbols which differ between any two codewords. It follows that:

- An LC with d=2 is single error detecting (SED) [see, e.g., Figures 5.1a and 5.2a]. In this case, the sum (on  $\mathbf{F}_q$ ) of the symbols in each transmitted codeword is zero, so if it is assumed that at most one symbol error occured and this sum is nonzero, then a symbol error in transmission occurred, whereas if it is zero, then a symbol error did not occur. However, if a symbol error in transmission occurred, the received (invalid) message is generally equidistant from multiple codewords, so it is not possible to correct the symbol error. Two or more symbol errors can cause the codeword to be misinterpreted.
- An LC with d=3 is single error correcting (SEC) [see, e.g., Figures 5.1b and 5.2b]. In this case, if it is again assumed that at most one symbol error in transmission occured, then if the received codeword is not a codeword, there is only one codeword that is unit Hamming distance away, so the single symbol error may in fact be corrected. Again, 2 or more symbol errors can cause the codeword to be misinterpreted.
- An LC with d = 4 is single error correcting and double error detecting (SECDED). In this case, if a single symbol error occurs, the received codeword will be unit Hamming distance away from a single codeword, and thus single symbol errors can be corrected. On the other hand, if two symbol errors occur, the received



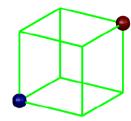
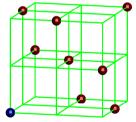


Figure 5.1: Valid codewords of (left) the (SED)  $[3,2,2]_2$  LBC, and (right) its dual, the (perfect, SEC)  $[3,1,3]_2$  LBC. The blue sphere denotes the origin, and d specifies the number of edges between any two codewords.



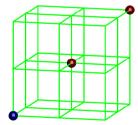


Figure 5.2: Valid codewords of (left) the (SED) [3,2,2]<sub>3</sub> LTC, and (right) its dual, the (SEC) [3,1,3]<sub>3</sub> LTC.

codeword is generally Hamming distance 2 away from multiple codewords, so double symbol errors can be detected but *not* corrected. Now, 3 or more symbol errors can cause the codewords to be misinterpreted.

# 5.1. INTRODUCTION

- An LC with d = 5 is double error correcting (DEC), with 3 or more symbol errors causing misinterpretation.
  An LC with d = 6 is double error correcting and triple error detecting (DECTED), with 4 or more symbol
- errors causing misinterpretation.
- An LC with d = 7 is triple error correcting (TEC), with 4 or more symbol errors causing misinterpretation.
  An LC with d = 8 is triple error correcting and quadruple error detecting (TECQED), with 5 or symbol
- errors causing misinterpretation.
  An LC with d = 9 is quadruple error correcting (QEC), with 5 or more symbol errors causing misinterpretation.
- The labels defined above are frequently used to quantify the error correction capability of an LC. Alternatively, if error correction is *not* attempted, then:
- An LC with d=2 is single error detecting, with 2 or more symbol errors causing misinterpretation.
- An LC with d=3 is double error detecting, with 3 or more symbol errors causing misinterpretation.
- An LC with d=4 is triple error detecting, with 4 or more symbol errors causing misinterpretation.
- An LC with d = 4 is unpic error detecting, with 4 or more symbol errors causing misinterpretation.

Error correcting algorithms are useful for a broad range of data transmission or data storage applications in which it is difficult or impossible to request that a corrupted codeword be retransmitted; algorithms which use such LCs for error detection only, on the other hand, may be used only when efficient handshaking is incorporated in a manner which makes it easy to request and resend any messages that might be corrupted during transmission.

An  $[n,k,d]_q$  LC is *perfect* if, for some integer t > 0, each possible *n*-dimensional *q*-ary codeword is of Hamming distance t or less from a single codeword (that is, there are no "wasted" codewords that are Hamming distance t + 1 or more from the codewords, and thus may not be corrected under the assumption that

at most t symbol errors have occured); note that a perfect code has odd d = 2t + 1 > 1. A remarkable proof by Tietäväinen (1973), which was based on related work by Van Lint, establishes that the *only* nontrivial perfect LCs are the  $[(q^m - 1)/(q - 1), (q^m - 1)/(q - 1) - m, 3]_q$  perfect q-ary Hamming codes and the  $[23, 12, 7]_2$  and  $[11, 6, 5]_3$  binary and ternary Golay codes, described further in §5.2 and §5.3.

An [n,k,d] LC is *quasi-perfect* if, for some integer t > 1, each possible n-dimensional q-ary codeword is either (a) of Hamming distance t-1 or less from a single codeword, and thus up to t-1 symbol errors may be corrected, or (b) of Hamming distance t from at least one codeword, and thus codewords with t symbol errors may be detected but not necessarily corrected (that is, there are no "wasted" codewords that are Hamming distance t+1 or more from a codeword, and thus may not be reconciled under the assumption that at most t symbol errors have occured); note that a quasi-perfect code has even d=2t>2.

Note finally, as illustrated for n = 3 in Figure 5.3, that a real lattice corresponding to an  $[n, k, d]_2$  LBC may often be constructed by forming a union of  $2^k$  cosets:

Construction A: 
$$\bigcup_{i=1}^{2^k} (\mathbf{w}_{[n,k,d]_2}^i + 2\mathbb{Z}^n), \tag{5.4a}$$

where the *coset representatives* in this construction,  $\mathbf{w}^i_{[n,k,d]_2}$  for  $i=1,\ldots,2^k$ , are the codewords of the  $[n,k,d]_2$  LBC under consideration and  $(\mathbf{w}+2\mathbb{Z}^n)$  denotes a  $\mathbb{Z}^n$  lattice scaled by a factor of 2 with all nodal points shifted by the vector  $\mathbf{w}$ ; thus, Construction A denotes the union of the nodal points in several such scaled and shifted  $\mathbb{Z}^n$  lattices. An alternative real lattice may sometimes be constructed via:

Construction B: 
$$\bigcup_{i=1}^{2^k} (\mathbf{w}_{[n,k,d]_2}^i + 2J) \text{ where } J = \left\{ \mathbf{x} \in \mathbb{Z}^n \middle| \left[ \sum_{i=1}^n x_i \right] \in 2\mathbb{Z} \right\},$$
 (5.4b)

where  $(2\mathbb{Z})$  denotes the even integers, and thus the last condition is sometimes written  $\sum_{i=1}^{n} x_i = 0 \pmod{2}$ .

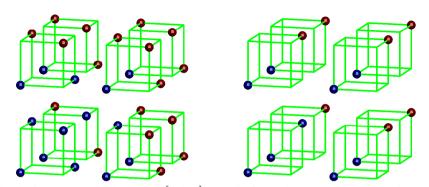


Figure 5.3: The lattice corresponding to an [n,k,d] LBC is formed by repeating the unit hypercube pattern given by the LBC (see, e.g., Figure 5.1) as an infinite array with unit spacing. In the above example, we illustrate this extension for (left) the face-centered cubic (FCC) lattice generated by the [3,2,2] LBC,  $D_3 = \bigcup_{i=1}^4 \left(\mathbf{w}^i_{[3,2,2]} + 2\mathbb{Z}^3\right)$ , and (right) the body-centered cubic (BCC) lattice generated by the [3,1,3] LBC,  $D_3^* = \bigcup_{i=1}^2 \left(\mathbf{w}^i_{[3,1,3]} + 2\mathbb{Z}^3\right)$ . The blue spheres, taken together, form a *primitive cell* that, repeated as an infinite array with *zero* spacing, tile (that is, fill) the domain.

In an analogous fashion, a complex lattice corresponding to an  $[n,k,d]_q$  LC may often be constructed by forming a union of  $q^k$  shifted and scaled n-dimensional  $\mathscr E$  lattices  $\mathbb Z[\omega]^n$  (see §2.1) such that

Construction 
$$A_{\mathscr{E}}^{\pi}: \bigcup_{i=1}^{q^k} (\mathbf{w}_{[n,k,d]_q}^i + \pi \mathbb{Z}[\omega]^n),$$
 (5.5a)

where, in the sequel, the multiplicative factor  $\pi$  takes two possible values (2 and  $\theta = \omega - \bar{\omega} = i\sqrt{3}$ ) and the coset representatives in this construction,  $\mathbf{w}^i_{[n,k,d]_q}$  for  $i=1,\ldots,q^k$ , are the codewords of the  $[n,k,d]_q$  LC under consideration. An alternative complex lattice may sometimes be constructed via:

Construction 
$$B_{\mathscr{E}}^{\pi}: \bigcup_{i=1}^{q^k} (\mathbf{w}_{[n,k,d]_q}^i + \pi J) \text{ where } J = \left\{ \mathbf{x} \in \mathbb{Z}[\omega]^n \middle| \left[ \sum_{i=1}^n x_i \right] \in \pi \mathscr{E} \right\},$$
 (5.5b)

where  $(\pi \mathcal{E})$  denotes the lattice of Eisenstein integers in the complex plane multiplied (that is, rotated and scaled) by the (possibly complex) factor  $\pi$ . Note the remarkable similarity in structure between the real constructions in (5.4a)-(5.4b) and the complex constructions in (5.5a)-(5.5b). Note also that real lattices corresponding to any of the complex lattices so constructed may easily be generated via (2.2).

# 5.2 Exemplary linear binary codes (LBCs)

We now summarize some of the families of LBCs available, presenting each in systematic form (5.1).

#### **5.2.1** Binary single parity-check codes

The simple [n, n-1, 2] binary single parity-check codes are SED, and include [2, 1, 2] (self-dual), [3, 2, 2], [4, 3, 2], [5, 4, 2], etc. Using such a code, for each (n-1) data bits to be transmitted, a parity bit is generated such that the sum (mod 2) of the data bits plus the parity bit is 0; when decoding, an error is flagged if this sum (mod 2) is 1. The [3, 2, 2] code illustrated in Figure 5.1a is given by

$$H_{[3,2,2]} = \begin{pmatrix} 1 & 1 & 1 \end{pmatrix}, \quad V_{[3,2,2]} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix}, \quad W_{[3,2,2]} = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix}. \tag{5.6}$$

<sup>&</sup>lt;sup>7</sup>As mentioned previously, when q = 2, we suppress the q subscript for notational clarity; we thus do this throughout §5.2.

# 5.2. EXEMPLARY LINEAR BINARY CODES (LBCS)

Other binary single parity-check codes have a partity submatrix P [see (5.3)] of similar form (a row of 1's). As seen for n = 3 in Figure 5.3a, via Construction A, the [n, n - 1, 2] binary single parity-check code generates the  $D_n$  lattice (see §2.3), which for n = 3 is FCC.

A single parity-check code (binary or otherwise), with d=2, can detect but not correct an error in an unknown position. However, it can correct an *erasure*; that is, the loss of data from a known position. A common application of this capability is in a RAID 5 system, a popular configuration for a relatively small *Redundant Array of Independent Disks*. In such a system, data is striped across n drives using a single parity check code; if any single drive fails, it can be recovered simply by achieving parity with the other disks.

#### **5.2.2** Binary repetition codes

The dual of the binary single parity-check codes are the simple [n, 1, n] binary repetition codes, which include [2, 1, 2] (SED, self-dual), [3, 1, 3] (SEC, perfect), [4, 1, 4] (SECDED), [5, 1, 5] (DEC), etc. This family of codes just repeats any given data bit n times; when decoding, one simply needs to determine which of the two codewords that the received code is nearest to. The [3, 1, 3] code illustrated in Figure 5.1b is given by

$$H_{[3,1,3]} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix}, \quad V_{[3,1,3]} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \quad W_{[3,1,3]} = \begin{pmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}.$$
 (5.7)

Other binary repetition codes have a partity submatrix of similar form (a column of 1's). As seen for n=3 in Figure 5.3b, via Construction A, the [n,1,n] binary repetition code generates the  $D_n^*$  lattice (see §2.3), which for n=3 is BCC. Via Construction B, on the other hand, the [8,1,8] binary repetition code generates the  $E_8$  lattice (see §2.5). Note also that the [3,2,2] binary single parity-check code with each bit in V repeated vertically m times leads to a [3m,2,2m] code, which may subsequently be rearranged into systematic form; taking m=4 and applying Construction B, the resulting [12,2,8] code, which is TECQED, generates the  $\Lambda_{12}^{\max}$  lattice (see §2.6).

#### **5.2.3** Binary Hamming codes

The  $[2^m - 1, 2^m - 1 - m, 3]$  binary Hamming codes are perfect and SEC, and include [3, 1, 3], [7, 4, 3], [15, 11, 3], [31, 26, 3], [63, 57, 3], [127, 120, 3], etc. For a given  $(2^m - 1 - m)$  data bits to be transmitted, each parity bit is generated such that the sum (mod 2) of a particular subset of the data bits plus that parity bit is 0. Note that, when decoding, the m parity bits may be used in a simple fashion to determine not only whether or not a single bit error occurred (which is true if one or more of these parity bits is nonzero), but if it did, which bit contains the error, as discussed further in §5.5. To illustrate, the venerable [7,4,3] code, with four data bits  $\{d_1,d_2,d_3,d_4\}$  and three parity bits  $\{b_1,b_2,b_3\}$ , is given by

$$H_{[7,4,3]} = \begin{pmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 \end{pmatrix}, \quad V_{[7,4,3]} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \end{pmatrix}, \quad \mathbf{w} = \begin{pmatrix} a_1 \\ d_2 \\ d_3 \\ d_4 \\ b_1 \\ b_2 \\ b_3 \end{pmatrix}, \tag{5.8a}$$

The parity-check matrix H of the [7,4,3] code has as columns all nonzero binary vectors of length (n-k)=3; when expressed in systematic form, the (n-k) columns of H corresponding to the identity matrix are shifted to the end, and the remaining k columns of H, in arbitrary order, make up the partity submatrix P. Other binary Hamming codes may be built up similarly. Via Construction A, the [7,4,3] binary Hamming code generates the  $E_7^*$  lattice (see §2.5).

# 5.2. EXEMPLARY LINEAR BINARY CODES (LBCS)

A Hamming code (binary or otherwise), with d=3, can only correct a single error in an unknown position. However, it can correct up to two *erasures* (cf. §5.2.1). A common application of this capability is in a RAID 6 system, a popular RAID configuration for large storage systems in data critical applications. In such a system, data may be striped across n drives using a Hamming code; if any single drive fails, it can be recovered using an appropriate parity check equation (that is, one of the parity check equations that takes that bit into account). If (while rebuilding the information on that disk, which might take a while if the disk is large) a *second* drive fails, then two useful equations may be derived from the (n-k) parity check equations: one that takes failed disk A into account but not failed disk B, and one that takes failed disk B into account but not failed disk A. By restoring parity in these two derived equations, the information on *both* drives may be rebuilt.

#### **5.2.4** Binary simplex codes

The dual of the binary Hamming codes are the  $[2^m - 1, m, 2^{m-1}]$  binary simplex codes [a.k.a. the binary maximum-length-sequence (MLS) codes], which include [3,2,2] (SED), [7,3,4] (SECDED), [15,4,8] (TEC-QED), etc. These codes are remarkable geometrically, as their codewords form a regular simplex. The [3,2,2] code is illustrated in Figure 5.1a; the [7,3,4] code is given by

$$H_{[7,3,4]} = \begin{pmatrix} 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & 1 \end{pmatrix}, \quad V_{[7,3,4]} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}.$$
 (5.9)

Other binary simplex codes have a partity submatrix given similarly by the transpose of the corresponding binary Hamming code. Via Construction A, the [7,3,4] binary simplex code generates the  $E_7$  lattice (see

§2.5). Via Construction B, the [15,4,8] binary simplex code generates the  $\Lambda_{15}$  lattice (see §2.6).

### 5.2.5 Extended binary Hamming codes

The  $[2^m, 2^m - 1 - m, 4]$  extended binary Hamming codes are quasi-perfect and SECDED, and include [4, 1, 4], [8, 4, 4] (self-dual), [16, 11, 4], etc. These codes are just binary Hamming codes (see §5.2.3) with an additional overall parity bit (see §5.2.1), and thus, assuming no more than 2 bit errors have occured, may be decoded similarly, as discussed further in §5.5. To illustrate, the venerable [8, 4, 4] code is given by

$$H_{[8,4,4]} = \begin{pmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}, \quad V_{[8,4,4]} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}.$$
 (5.10)

Other extended binary Hamming codes have a partity submatrix that may similarly be constructed by adding an overall parity bit to the corresponding binary Hamming code. Via Construction A, the [8,4,4] extended binary Hamming code again generates the  $E_8$  lattice.

# **5.2.6** Binary biorthogonal codes

The dual of the extended binary Hamming codes are the  $[2^m, m+1, 2^{m-1}]$  binary biorthogonal codes (a.k.a. Hadamard codes), and include [4,3,2] (SED), [8,4,4] (SECDED, self-dual), [16,5,8] (TECQED), [32,6,16], etc. The [32,6,16] code was used on the Mariner 9 spacecraft. These codes are distinguished by the characteristic that their codewords are mutually orthogonal [that is,  $\mathbf{w}^i \cdot \mathbf{w}^j = 0 \pmod{2}$  for  $i \neq j$ ]. Note that the [4,3,2] and [8,4,4] codes have already been discussed above. The binary biorthogonal codes each have a partity submatrix that is simply the transpose of the parity submatrix of the corresponding extended binary

# 5.2. EXEMPLARY LINEAR BINARY CODES (LBCS)

Hamming code, the construction of which is described in §5.2.5. Via Construction B, the [16,5,8] binary biorthogonal code generates the  $\Lambda_{16}$  lattice.

#### **5.2.7** Binary quadratic residue codes

The [n,(n+1)/2,d] binary quadratic residue codes are defined for all prime n for which there exists an integer 1 < x < n such that  $x^2 = 2 \pmod{n}$  [equivalently, for all prime n of the form  $n = 8m \pm 1$  where m is an integer], and include [7,4,3] (SEC, perfect, as introduced in §5.2.3), [17,9,5] (DEC), [23,12,7] (TEC, perfect, a.k.a. the binary Golay code), [31,16,7] (TEC), [41,21,9] (QEC), [47,24,11], etc. Adding an overall parity bit to these codes, the [n+1,(n+1)/2,d+1] extended binary quadratic residue codes include [8,4,4] (SECDED, quasi-perfect, self-dual, as introduced in §5.2.5), [18,9,6] (DECTED), [24,12,8] (TEC-

QED, quasi-perfect, self-dual, a.k.a. the extended binary Golay code), [32,16,8] (TECQED), [42,21,10], [48, 24, 12], etc. The venerable [24, 12, 8] extended binary Golay code, used by the Voyager 1 & 2 spacecraft, is given by

Note that P is symmetric. The [23,12,7] binary Golay code may be obtained by *puncturing* the [24,12,8] code listed above; that is, by eliminating any row of P (typically, the last).

Via Construction B, the [24, 12, 8] extended binary Golay code generates the *Leech half-lattice H*<sub>24</sub>, which may be joined with a translate of itself [that is,  $H_{24} + \mathbf{a}$  where  $a_1 = -3/2$  and  $a_k = 1/2$  for k = 2, ..., 24] to construct the  $\Lambda_{24}$  lattice.

Note that many of the binary codes introduced above fall within a larger family of codes collectively referred to as *Reed-Muller* codes, as illustrated in Figure 5.4.

The (perfect) binary Hamming and binary Golay codes may be *extended* to quasi-perfect codes by adding an overall parity bit, thereby increasing n by 1 and, in the case of these specific codes, increasing d by 1. A

# 5.2.8 Extending, puncturing, and shortening

code obtained by essentially the reverse of this process, removing a parity bit and thus reducing both n and d by 1, is sometimes said to be *punctured*. In contrast, a code obtained by removing  $\ell \geq 1$  data bits, thus reducing both n and k by  $\ell$ , is said to be *shortened*. A typical and common application is in error-correcting memory systems for computers, in which the data often comes naturally in blocks of 64 bits. Starting from the [127,120,3] binary Hamming code, one may eliminate 56 data bits to create a shortened [71,64,3] SEC code; alternatively, starting from the [128,120,4] extended binary Hamming code, one may eliminate 56 data bits to create a shortened [72,64,4] SECDED code. Many ECC Memory and RAID 6 storage systems are based on variants of such shortened binary Hamming codes, which are simple and fast to use. Note also that, via Construction B, the [21,9,8] code obtained by shortening the [24,12,8] extended binary Golay code by 3 data bits generates directly the  $\Lambda_{21}$  lattice.

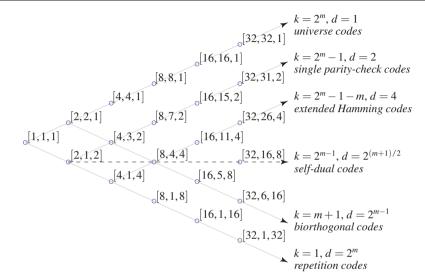


Figure 5.4: The family of  $[2^m, k, d]$  Reed-Muller binary codes for m = 0 to 5.

# 5.3 Exemplary linear ternary codes (LTCs)

We now summarize some of the families of LTCs available, presenting each in systematic form (5.1), noting that all have analogs in the binary setting.

#### 5.3.1 Ternary single parity-check codes

The  $[n, n-1, 2]_3$  ternary single parity-check codes are SED, and include  $[2, 1, 2]_3$  (self-dual),  $[3, 2, 2]_3$ ,  $[4, 3, 2]_3$ , etc. As illustrated for n = 3 in Figure 5.2a, the  $[3, 2, 2]_3$  code is given by

$$H_{[3,2,2]_3} = \begin{pmatrix} 1 & 1 & 1 \end{pmatrix}, \quad V_{[3,2,2]_3} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 2 & 2 \end{pmatrix}, \quad W_{[3,2,2]_3} = \begin{pmatrix} 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 \\ 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 0 & 2 & 1 & 2 & 1 & 0 & 1 & 0 & 2 \end{pmatrix}. \tag{5.12}$$

Other ternary single parity-check codes have a partity submatrix P [see (5.3)] of similar form (a row of 2's). Via Construction  $A_{\mathcal{E}}^{\theta}$ , the [3,2,2]<sub>3</sub> ternary single parity-check code generates the  $E_6^*$  lattice.

#### **5.3.2** Ternary repetition codes

The dual of the ternary single parity-check codes are the  $[n,1,n]_3$  ternary repetition codes, which include  $[2,1,2]_3$  (SED, self-dual),  $[3,1,3]_3$  (SEC),  $[4,1,4]_3$  (SECDED), etc. As illustrated for n=3 in Figure 5.2b, the  $[3,1,3]_3$  code is given by

$$H_{[3,1,3]_3} = \begin{pmatrix} 2 & 1 & 0 \\ 2 & 0 & 1 \end{pmatrix}, \quad V_{[3,1,3]_3} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \quad W_{[3,1,3]_3} = \begin{pmatrix} 0 & 1 & 2 \\ 0 & 1 & 2 \\ 0 & 1 & 2 \end{pmatrix}.$$
 (5.13)

Other ternary repetition codes have a partity submatrix of similar form (a column of 1's). Via Construction  $A_{\mathscr{E}}^{\theta}$ , the  $[3,1,3]_3$  ternary repetition code generates the  $E_6$  lattice. Via Construction  $B_{\mathscr{E}}^{\theta}$ , on the other hand, the  $[6,1,6]_3$  ternary repetition code generates the  $K_{12}$  lattice.

#### **5.3.3** Ternary Hamming codes

The  $[(3^m - 1)/2, (3^m - 1)/2 - m, 3]_3$  ternary Hamming codes are perfect and SEC, and include  $[4,2,3]_3$  (self-dual, a.k.a. the *tetracode*),  $[13,10,3]_3$ ,  $[40,36,3]_3$ , etc. To illustrate, the venerable  $[4,2,3]_3$  tetracode is given by

$$H_{[4,2,3]_3} = \begin{pmatrix} 1 & 1 & 1 & 0 \\ 1 & 2 & 0 & 1 \end{pmatrix}, \quad V_{[4,2,3]_3} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 2 & 2 \\ 2 & 1 \end{pmatrix}. \tag{5.14}$$

The parity-check matrix H of the  $[4,2,3]_3$  code has as columns those nonzero ternary vectors of length (n-k)=2 whose first nonzero entry is 1; when expressed in systematic form, the (n-k) columns of H corresponding to the identity matrix are shifted to the end, and the remaining k columns of H, in arbitrary order, make up the entries of -P. Other ternary Hamming codes may be built up similarly; for example, the  $[13,10,3]_3$  code is given by

Via Construction  $A_{\mathcal{E}}^{\theta}$ , the  $[4,2,3]_3$  tetracode again generates the  $E_8$  lattice.

### **5.3.4** Ternary simplex codes

The dual of the ternary Hamming codes are the  $[(3^m - 1)/2, m, 3^{m-1}]_3$  ternary simplex codes, which include  $[4,2,3]_3$  (SEC, perfect, self-dual),  $[13,3,9]_3$  (QEC),  $[40,4,27]_3$ , etc. These codes are remarkable geometrically, as their codewords are all equidistant from one another. Ternary simplex codes have a partity submatrix given by the negative transpose of the corresponding ternary Hamming code.

# 5.3.5 Ternary quadratic residue codes

The  $[n, (n+1)/2, d]_3$  ternary quadratic residue codes are defined for all prime n for which there exists an integer 1 < x < n such that  $x^2 = 3 \pmod{n}$  [equivalently, for all prime n of the form  $n = 12m \pm 1$ 

where m is an integer], and include  $[11,6,5]_3$  (DEC, perfect, a.k.a. the *ternary Golay code*),  $[13,7,5]_3$  (DEC),  $[23,12,8]_3$  (TECQED),  $[37,19,10]_3$ ,  $[47,24,14]_3$ , etc. Adding an overall parity bit to these codes, the  $[n+1,(n+1)/2,d+1]_3$  extended ternary quadratic residue codes include  $[12,6,6]_3$  (DECTED, quasi-perfect, self-dual, a.k.a. the extended ternary Golay code),  $[14,7,6]_3$  (DECTED),  $[24,12,9]_3$  (QEC),  $[38,19,11]_3$ ,  $[48,24,15]_3$ , etc. The venerable  $[12,6,6]_3$  extended ternary Golay code is given by

$$H_{[12,6,6]_3} = \begin{bmatrix} -P_{6\times6} & I_{6\times6} \end{bmatrix}, \quad V_{[12,6,6]_3} = \begin{bmatrix} I_{6\times6} \\ P_{6\times6} \end{bmatrix}, \quad P_{6\times6} = \begin{pmatrix} 0 & 1 & 1 & 1 & 1 & 1 \\ 1 & 0 & 1 & 2 & 2 & 1 \\ 1 & 1 & 0 & 1 & 2 & 2 \\ 1 & 2 & 1 & 0 & 1 & 2 \\ 1 & 2 & 2 & 1 & 0 & 1 \\ 1 & 1 & 2 & 2 & 1 & 0 \end{pmatrix}.$$
(5.16)

Note that  $P$  is symmetric. The  $\begin{bmatrix} 11,6,5 \end{bmatrix}_3$  ternary Golay code may be obtained by puncturing the  $\begin{bmatrix} 12,6,6 \end{bmatrix}_3$ 

code listed above. Via Construction  $B_{\mathcal{E}}^{\theta}$ , the [12,6,6]<sub>3</sub> extended ternary Golay code generates an intermediate lattice which may be joined with two translates of itself to generate the  $\Lambda_{24}$  lattice.

# 5.4 Exemplary linear quaternary codes (LQCs)

We now summarize some of the families of LQCs available, presenting each in systematic form (5.1).

#### 5.4.1 Quaternary single parity-check codes

The  $[n, n-1, 2]_A$  quaternary single parity-check codes are SED, and include  $[2, 1, 2]_A$  (self-dual)  $[3, 2, 2]_A$ 

$$\begin{split} H_{[3,2,2]_4} = \begin{pmatrix} 1 & 1 & 1 \end{pmatrix}, & V_{[3,2,2]_4} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix}, \\ W_{[3,2,2]_4} = \begin{pmatrix} 0 & 1 & \omega & \bar{\omega} & 0 & 1 & \omega & \bar{\omega} & 0 & 1 & \omega & \bar{\omega} \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & \omega & \omega & \omega & \bar{\omega} & \bar{\omega} & \bar{\omega} & \bar{\omega} \\ 0 & 1 & \omega & \bar{\omega} & 1 & 0 & \bar{\omega} & \bar{\omega} & \bar{\omega} & 1 & 0 \end{pmatrix}. \end{split}$$

Other quaternary single parity-check codes have a partity submatrix *P* of similar form.

#### **5.4.2** Quaternary repetition codes

The dual of the quaternary single parity-check codes are the  $[n,1,n]_4$  quaternary repetition codes, which include  $[2,1,2]_4$  (SED, self-dual),  $[3,1,3]_4$  (SEC),  $[4,1,4]_4$  (SECDED), etc. The  $[3,1,3]_4$  code is given by

$$H_{[3,1,3]_4} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix}, \quad V_{[3,1,3]_4} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \quad W_{[3,1,3]_4} = \begin{pmatrix} 0 & 1 & \omega & \bar{\omega} \\ 0 & 1 & \omega & \bar{\omega} \\ 0 & 1 & \omega & \bar{\omega} \end{pmatrix}.$$

Other quaternary repetition codes have a partity submatrix of similar form.

# **5.4.3** Ouaternary Hamming codes

The  $[(4^m - 1)/3, (4^m - 1)/3 - m, 3]_4$  quaternary Hamming codes are perfect and SEC, and include  $[5,3,3]_4$ ,  $[21,18,3]_4$ ,  $[85,81,3]_4$ , etc. To illustrate, the  $[5,3,3]_4$  code is given by

$$H_{[5,3,3]_4} = \begin{pmatrix} 1 & 1 & 1 & 1 & 0 \\ 1 & \omega & \bar{\omega} & 0 & 1 \end{pmatrix}, \quad V_{[5,3,3]_4} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 1 & 1 \end{pmatrix}. \tag{5.19}$$

(5.17)

(5.18)

The parity-check matrix H of the  $[5,3,3]_4$  code has as columns those nonzero quaternary vectors of length (n-k)=2 whose first nonzero entry is 1; when expressed in systematic form, the (n-k) columns of H

corresponding to the identity matrix are shifted to the end, and the remaining k columns of H, in arbitrary order, make up the entries of P. Other quaternary Hamming codes may be built up similarly.

### **5.4.4** Quaternary simplex codes

The dual of the quaternary Hamming codes are the  $[(4^m - 1)/3, m, 4^{m-1}]_4$  quaternary simplex codes, which include  $[5,2,4]_4$  (SECDED),  $[21,3,16]_4$ ,  $[85,4,64]_4$ , etc. These codes are remarkable geometrically, as their codewords are all equidistant from one another. Quaternary simplex codes have a partity submatrix given by the conjugate transpose of the corresponding quaternary Hamming code.

### 5.4.5 Quaternary quadratic residue codes

The  $[n, (n+1)/2, d]_4$  quaternary quadratic residue codes are defined for all prime n of the form  $n = 8m \pm 3$  where m is an integer, and include  $[5,3,3]_4$  (SEC, perfect, see §5.4.3),  $[11,6,5]_4$  (DEC),  $[13,7,5]_4$  (DEC),  $[19,10,7]_4$  (TEC),  $[29,15,11]_4$ ,  $[37,19,11]_4$ , etc. The related  $[n+1,(n+1)/2,d+1]_4$  extended quaternary quadratic residue codes include  $[6,3,4]_4$  (SECDED, quasi-perfect, self-dual, a.k.a. the hexacode),  $[12,6,6]_4$  (DECTED),  $[14,7,6]_4$  (DECTED, self-dual),  $[20,10,8]_4$  (TECQED),  $[30,15,12]_4$  (self-dual),  $[38,19,12]_4$ , etc. The venerable  $[6,3,4]_4$  hexacode is given by

$$H_{[6,3,4]_4} = \begin{pmatrix} 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & \omega & \bar{\omega} & 0 & 1 & 0 \\ 1 & \bar{\omega} & \omega & 0 & 0 & 1 \end{pmatrix}, \quad V_{[6,3,4]_4} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 1 & 1 \\ 1 & \omega & \bar{\omega} \\ 1 & \bar{\omega} & \bar{\omega} \end{pmatrix}. \tag{5.20}$$

Note that P is symmetric. The  $[5,3,3]_4$  quaternary quadratic residue code may be obtained by puncturing the  $[6,3,4]_4$  code listed above.

Via Construction  $A_{\epsilon}^2$ , the  $[6,3,4]_4$  hexacode generates the  $K_{12}$  lattice.

The  $\begin{bmatrix} 6 & 2 & 4 \end{bmatrix}$ , have code with  $A^3 = 64$  codewords is of particular importance due to the structured role.

### 5.4. EXEMPLARY LINEAR QUATERNARY CODES (LQCS)

Figure 5.5: A labelling of 16 points of the  $D_2$  lattice (due to Ungerboeck 1982). The  $A_{ijk}$  points have coordinates which are both even integers [e.g.,  $A_{000} = \begin{pmatrix} 0 & 0 \end{pmatrix}$ ], and the  $B_{ijk}$  points have coordinates which are both odd integers [e.g.,  $B_{000} = \begin{pmatrix} 1 & 1 \end{pmatrix}$ ]. The complete  $D_2$  lattice is formed by repeating this 2D pattern as an infinite array with unit spacing, as in Figure 5.3; note that each of the subsets of  $D_2$  corresponding to a particular label is itself an appropriate shift of a  $4D_2$  lattice (that is, the  $D_2$  lattice with the spacing quadrupled between the points).

 $2^{12} = 4096$  codewords **w**, and the corresponding  $\Lambda_{24}$  lattice. To construct the extended binary Golay code in this manner (see §11 of Conway & Sloane 1998), we may first arrange binary vectors of length 24 into  $4 \times 6$  arrays with binary entries. The sum of the bits (mod 2) in any row or column of this array gives its *parity*, which is said to be *even* if the bits sum to 0 and *odd* if the bits sum to 1. We then define the *projection* of

any binary vector  $\mathbf{d} \in \mathbf{F}_2^4$  onto a symbol  $x \in \mathbf{F}_4$  via the product  $x = \begin{pmatrix} 0 & 1 & \omega & \omega \end{pmatrix} \mathbf{d}$  (on  $\mathbf{F}_4$ ). The [24, 12, 8]

it plays in some convenient constructions of the [24, 12, 8] extended binary Golay code (see §5.2.7), with

extended binary Golay code is then given by the set of all w∈ F<sub>2</sub><sup>24</sup> such that, in the corresponding 4 × 6 array,
the parity of all of the columns matches the parity of the top row, and
the projection of the six columns of the array forms a codeword of the [6,3,4]<sub>4</sub> hexacode.

An alternative construction of the  $\Lambda_{24}$  lattice, due to Vardy & Be'ery (1993) and which also leverages claverly the [6, 3, 4], beyond is based on the Uncerbook (1982) portitioning of the  $D_2$  lattice (see §2.3)

into  $A_{ijk}$  and  $B_{ijk}$  subsets, as depicted in Figure 5.5. Binary vectors of length 24 are now constructed as  $2 \times 6$  arrays whose entries are points of  $D_2$ , labelled as shown. When considering a pair of such points [say,  $\mathbf{c} = (A_{i_1,j_1,k_1} \quad A_{i_2,j_2,k_2})^T$ ],

- the pair is said to be *even* or *odd* based on the sum (mod 2) of the indices  $\{i_1, j_1, i_2, j_2\}$ ,
- the index i<sub>1</sub> is known as the h-parity of the pair,
  the sum (mod 2) of k<sub>1</sub> and k<sub>2</sub> is known as the k-parity of the pair, and
- the *projection* of the pair is defined as above, based on the vector  $\mathbf{d} = \begin{pmatrix} i_1 & j_1 & i_2 & j_2 \end{pmatrix}^T$ .

The Leech lattice  $\Lambda_{24}$  is then given by the set of all  $\mathbf{u} \in \mathbb{Z}^{24}$  such that, in the corresponding  $2 \times 6$  array,

- all array entries are either points in the  $A_{ijk}$  subsets of  $D_2$  (referred to as a *type-A* array), or points in the  $B_{ijk}$  subsets of  $D_2$  (referred to as a *type-B* array),
- the overall *k* parity of the array [that is, the sum (mod 2) of the *k*-parity of the 6 pairs of points] is even if the array is type *A* and odd if the array is type *B*,
- the pairs of points in the 6 columns of the array are either all even (referred to as an *even* array) or all odd (referred to as an *odd* array),
- the overall h parity of the array [that is, the sum (mod 2) of the h-parity of the 6 pairs of points] is even if the array even and odd if the array is odd, and
- the projection of the six columns of the array forms a codeword of the  $[6,3,4]_4$  hexacode.

The union of all points corresponding to Type A arrays in this construction forms the *Leech half lattice*  $H_{24}$  mentioned in §5.2.7, whereas the union of all points corresponding to Type B arrays forms its translate,  $H_{24} + \mathbf{a}$ . The  $H_{24}$  lattice can be further decomposed into all points corresponding to even arrays, which forms the *Leech currents lattice*  $Q_{12}$  and all points corresponding to odd arrays, which forms the *Leech currents lattice*  $Q_{12}$  and all points corresponding to odd arrays, which forms

the Leech quarter lattice  $Q_{24}$ , and all points corresponding to odd arrays, which forms its translate,  $Q_{24} + \mathbf{b}$ . The  $\Lambda_{24}$  lattice is then given by the union of  $Q_{24}$ ,  $Q_{24} + \mathbf{b}$ ,  $Q_{24} + \mathbf{a}$ , and  $Q_{24} + \mathbf{a} + \mathbf{b}$ ; this construction is exploited in §6.1.5 when presenting a remarkably efficient algorithm for quantization from  $\mathbb{R}^{24}$  to  $\Lambda_{24}$ .

### .5 Decoding

The use of an  $[n,k,d]_q$  linear code (a.k.a. *linear block code*) in practice to communicate data over a noisy channel is straightforward:

- arrange the original data into *blocks* of length *k* over an *alphabet* of *q* symbols;
- code each resulting data vector  $\mathbf{d} \in \mathbf{F}_q^k$  into a longer codeword  $\mathbf{w} \in \mathbf{F}_q^n$  via  $\mathbf{w} = V_{[n,k,d]_q} \mathbf{d}$ ;
- transmit the corresponding codeword  $\mathbf{w} \in F_q^n$  over the noisy channel;
- receive the (possibly corrupted) message  $\hat{\mathbf{w}} \in F_q^n$  on the other end;
- *decode* the received message  $\hat{\mathbf{w}}$  leveraging  $H_{[n,k,d]_q}$ ; that is, find the most likely codeword  $\mathbf{w}$  corresponding to the received message  $\hat{\mathbf{w}}$ , and the data vector  $\mathbf{d}$  that generated it.

The decoding problem is quite rich; many creative schemes have been proposed over the years for decoding the various LCs that have been introduced thus far, as well as many others. This subject goes a bit beyond the scope of the present review, but we would be remiss if we didn't at least briefly introduce a few exemplary decoding strategies.

For the purpose of fast decoding of an LC, it is useful to consider convenient alternatives to the systematic form. If H and V are the parity-check and basis matrices of an  $[n,k,d]_q$  LC in systematic form, with HV=0 (on  $\mathbf{F}_q$ ), then an *equivalent* LC, possibly not in systematic form, is given by taking

$$\tilde{H} = HQ$$
 and  $\tilde{V} = Q^{-1}V$ . (5.21)

It follows immediately that, again,  $\tilde{H}\tilde{V}=0$  (on  $\mathbf{F}_q$ ). In the simplest such transformation, Q is a permutation matrix, and thus  $Q^{-1}=Q^T$ ; this transformation corresponds to reordering the rows of V and the corresponding columns of H (that is, reordering the data bits and parity bits in the corresponding LC). Other equivalent LCs may be constructed in this manner by relaxing the constraint that Q be a permutation matrix, effectively

taking linear combinations (on  $\mathbf{F}_q$ ) of the rows of V and the corresponding columns of H. Note further that reordering the columns of V and/or the rows of H leaves an LC unchanged.

#### 5.5.1 Algebraic decoding

Certain LBCs may be decoded quickly by arranging the columns of the parity-check matrix in a convenient order and examining the binary number given by the product of the parity-check matrix and the (possibly, corrupted) received message. To illustrate, consider the [7,4,3] binary Hamming code introduced in §5.2.3. Transforming as described above with

$$Q = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

results in a modified basis matrix  $\tilde{V}$ , and a modified parity-check matrix  $\tilde{H}$  arranged such that the columns of  $\tilde{H}$  appear in binary order:

$$\tilde{H}_{[7,4,3]} = \begin{pmatrix} 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{pmatrix}, \quad \tilde{V}_{[7,4,3]} = \begin{pmatrix} 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}, \quad \tilde{\mathbf{w}} = \begin{pmatrix} b_3 \\ b_2 \\ d_1 \\ b_1 \\ d_2 \\ d_3 \end{pmatrix}. \tag{5.22}$$

Taking the matrix  $\tilde{H}_{[2^m-1,2^m-1-m,3]}$  of a binary Hamming code arranged in such a fashion (in the above example, m=3) times (mod 2) any of the codewords  $\tilde{\mathbf{w}}$  (generated via  $\tilde{\mathbf{w}}=\tilde{V}_{[2^m-1,2^m-1-m,3]}\mathbf{d}$  where  $\mathbf{d}\in$ 

#### 5.5. DECODING

 $\mathbf{F}_2^{2^m-1-m}$ ) gives the zero vector. On the other hand, taking the matrix  $\tilde{H}_{[2^m-1,2^m-1-m,3]}$  times (mod 2) any invalid vector  $\hat{\mathbf{w}}$  gives the nonzero *syndrome vector*  $\mathbf{s}$ , of order m=n-k, which may be interpreted as a nonzero *m*-bit binary number called the *syndrome*, denoted s, of the received message. Conveniently, as a direct result of the structure of  $\tilde{H}$  used in this construction, the number s identifies precisely which bit of the received message vector  $\hat{\mathbf{w}}$ , arranged as shown above, must be flipped in order to determine the nearest codeword, thereby performing single error correction (SEC).

Consider now the class of  $[2^m, 2^m - 1 - m, 4]$  extended binary Hamming codes introduced in §5.2.5. Define the syndrome s as in the corresponding binary Hamming code of length  $(2^m - 1)$  as discussed above, neglecting the overall parity bit, and define p as the sum (mod 2) over all the bits, including the overall parity bit. There are zero bit errors if s = p = 0, there two bit errors (which may be detected but not uniquely corrected) if  $s \neq 0$  and p = 0, and there is a single bit error if p = 1 (in which case, if s = 0, this error is in the overall parity bit, and, if  $s \neq 0$ , this error is in one of the other bits and may be corrected based on s just as in the corresponding binary Hamming code). This strategy thus performs single error correction and double error detection (SECDED).

The extended binary Golay code introduced in §5.2.7 may be decoded via syndrome computation in a

similar fashion, though several more checks are involved, as the procedure performs triple error correction and quadruple error detection (TECQED) on the received message  $\hat{\mathbf{w}}$ . Recall the definitions of H, V, and  $P = P^T$  for the [24, 12, 8] extended binary Golay code in systematic form, as listed in (5.11). Note that  $V^TV = 0$ , and thus  $V^T$  serves as an alternative parity-check matrix for this code. Defining  $w_H(\mathbf{s})$  as the Hamming weight (that is, the number of nonzero elements) of the vector  $\mathbf{s}$ , and defining  $\mathbf{p}^i$  as the i'th column of P,  $\mathbf{e}^i$  as the i'th Cartesian unit vector, and 0 as the zero vector, we may decode  $\hat{\mathbf{w}}$  as follows:

set  $\mathbf{s} = V^T \hat{\mathbf{w}}$ , if  $w_H(\mathbf{s}) \le 3$  then set  $\mathbf{c} = [\mathbf{s}; 0]$ , flag= 0, return, end if (case A) set  $\mathbf{r} = P\mathbf{s}$ , if  $w_H(\mathbf{r}) \le 3$  then set  $\mathbf{c} = [0; \mathbf{r}]$ , flag= 0, return, end if (case B) for i = 1:12

if  $w_H(\mathbf{s} + \mathbf{p}^i) \le 2$  then set  $\mathbf{c} = [\mathbf{s} + \mathbf{p}^i; \mathbf{e}^i]$ , flag = 0, return, end if (case C) if  $w_H(\mathbf{r} + \mathbf{p}^i) \le 2$  then set  $\mathbf{c} = [\mathbf{e}^i; \mathbf{r} + \mathbf{p}^i]$ , flag = 0, return, end if (case D) end for flag=1; return (4 total errors, can not be uniquely corrected)

Upon return, assuming the received vector  $\hat{\mathbf{w}}$  has 4 or less bit errors, if flag = 0, then 3 or fewer errors are detected and the corrected vector is  $\mathbf{w} = \hat{\mathbf{w}} + \mathbf{c}$ , whereas if flag = 1, then 4 errors are detected and  $\hat{\mathbf{w}}$  can not be uniquely corrected. To verify this algorithm, noting that  $V^T\mathbf{w} = 0$  for any codeword  $\mathbf{w}$ , it is sufficient to analyze the algorithm for  $\mathbf{w} = 0$  only. Block partitioning  $\hat{\mathbf{w}} = \begin{bmatrix} \mathbf{x}; & \mathbf{y} \end{bmatrix}$ , consider the following 4 correctable cases:

Case A (0 parity bit errors, up to 3 data bit errors): Due to the structure of P, parity bit errors (that is,  $w_H(\mathbf{y}) \neq 0$ ) result in  $w_H(\mathbf{s}) \geq 6$ ; if  $w_H(\mathbf{s})$  is less than this, then  $\mathbf{y} = 0$  and  $\mathbf{s} = V^T \hat{\mathbf{w}} = I\mathbf{x} = \mathbf{x}$ .

Case B (0 data bit errors, up to 3 parity bit errors): Note that  $PV^T = H$ , and thus  $\mathbf{r} = H\hat{\mathbf{w}}$ . By an analogous argument as that used in Case A, due to the structure of P, data bit errors (that is,  $w_H(\mathbf{x}) \neq 0$ ) result in  $w_H(\mathbf{r}) \geq 6$ ; if  $w_H(\mathbf{s})$  is less than this, then  $\mathbf{x} = 0$  and  $\mathbf{r} = H\hat{\mathbf{w}} = I\mathbf{y} = \mathbf{y}$ .

Case C (1 parity bit error, up to 2 data bit errors): In this case, we individually check each of the (12) possible cases corresponding to a single parity bit error, essentially repeating the analysis of Case A, mutatis mutandis. That is, for each i, we consider the possibility that  $\mathbf{y} = \mathbf{e}^i$ , and thus  $\mathbf{s} = \mathbf{x} + \mathbf{p}^i$ , and check to see if  $w_H(\mathbf{x}) = w_H(\mathbf{s} + \mathbf{p}^i) \le 2$ .

Case D (1 data bit error, up to 2 parity bit errors): In this case, we individually check each of the (12) possible cases corresponding to a single data bit error, essentially repeating the analysis of Case B, mutatis mutandis (cf. Case C).

#### 5.5.2 Cyclic form

A *cyclic code* is an LC that may be transformed [via (5.21)] into a form in which all cyclic shifts of codewords are themselves also codewords. The basis matrix  $V = V_{n \times k}$  and parity-check matrix  $H = H_{(n-k)\times n}$  of any  $[n,k]_q$  cyclic code may be written in the standard form

$$H_{[n,k]_q} = \begin{pmatrix} h_k & h_{k-1} & \dots & h_0 & & & & 0 \\ & h_k & h_{k-1} & \dots & h_0 & & & \\ & & \ddots & \ddots & \ddots & \ddots & \\ 0 & & & h_k & h_{k-1} & \dots & h_0 \end{pmatrix}, \quad V_{[n,k]_q} = \begin{pmatrix} v_0 & & & & 0 \\ v_1 & v_0 & & & & \\ \vdots & v_1 & \ddots & & & \\ v_{n-k} & \vdots & \ddots & v_0 & & \\ & v_{n-k} & \ddots & v_1 & & \\ & & & \ddots & \vdots & \\ 0 & & & v_{n-k} \end{pmatrix}.$$
 (5.23)

A convenient construction which simplifies the analysis of an  $[n,k]_q$  cyclic code, as defined above, is the *cyclic shift* operator z. The use of this operator as discussed here is akin to its use in the Z-transform analysis of discrete-time linear systems, with the major difference being that it is used here in a cyclic context on  $\mathbf{F}_q$ : that is, arithmetic with polynomials in z and coefficients in  $\mathbf{F}_q$  is performed as usual, except that the coefficients of each power of z are combined via the arithmetic rules on  $\mathbf{F}_q$  (see the second paragraph of §5.1), and higher powers of  $z^k$  are simplified via the cyclic condition

$$z^n = 1. (5.24)$$

(5.25b)

In the deployment of an  $[n,k]_q$  cyclic code, the operator z appears in

the data polynomial

the basis polynomial 
$$v(z) = v_0 + v_1 z + ... + v_{n-k} z^{n-k}$$
, the codeword polynomial  $w(z) = w_0 + w_1 z + ... + w_{n-1} z^{n-1}$ , the received-message polynomial  $\hat{w}(z) = \hat{w}_0 + \hat{w}_1 z + ... + \hat{w}_{n-1} z^{n-1}$ , and the parity-check polynomial  $h(z) = h_0 + h_1 z + ... + h_k z^k$ .

The basis polynomial v(z) and parity-check polynomial h(z) are constructed in mutually-orthogonal manner that, taken together, enforces the cyclic condition (5.24):

$$v(z)h(z) = (z^n - 1),$$
 (5.25a)

 $d(z) = d_0 + d_1 z + ... + d_{k-1} z^{k-1}$ 

which may also be written

$$[v(z)h(z)] \bmod (z^n - 1) = 0;$$
(5.25b)

note that the mod command used in (5.25b) means that the polynomial  $[v(z)h(z)]$  is divided by the polynomial  $(z^n - 1)$  and the remainder is equal to 0. One such factorization of  $(z^n - 1)$  on  $\mathbb{F}$ , which exists for any  $n$  and

 $(z^n-1)$  and the remainder is equal to 0. One such factorization of  $(z^n-1)$  on  $\mathbf{F}_q$ , which exists for any n and q, is  $z^{n}-1=(z-1)(z^{n-1}+z^{n-2}+\ldots+z+1);$ (5.26)

this leads to the single parity check code  $[n, n-1, 2]_q$  if one takes v(z) = (z-1) and h(z) equal to the rest, and to the repetition code  $[n, 1, n]_q$  if one takes h(z) = (z - 1) and v(z) equal to the rest. Other cyclic codes over  $\mathbf{F}_q$  for prime q may be built from the unique irreducible factors of the polynomial  $(z^n - 1)$ , grouping these factors appropriately to form v(z) and h(z); a few such factorizations for various values of n are listed

#### 5.5. DECODING

in Table 5.1 for q = 2 (in which -1 = 1) and Table 5.2 for q = 3 (in which -1 = 2); others are easily found using Mathematica. Factoring  $(z^n - 1)$  over  $\mathbf{F}_4$  is more delicate, as the factorizations do not reduce to unique irreducible forms; one such factorization is listed in Table 5.3. Based on (5.25a) and such factorizations, a large number of cyclic codes may be constructed. However, only a few such codes have both favorable minimum distance d and an available simple error dectection/correction scheme; some such codes are listed in Table 5.4.

Given a data vector  $\mathbf{d} \in \mathbf{F}_q^k$ , the use of an LC in cyclic form is again straightforward:

- code d(z) into a codeword polynomial w(z) leveraging the basis polynomial v(z) [using nonsystematic coding, one simply takes w(z) = d(z) v(z);
- transmit the corresponding codeword  $\mathbf{w} \in F_a^n$  over the noisy channel;

• form a data polynomial d(z) with the k elements of **d** as coefficients;

- receive the (possibly corrupted) message  $\hat{\mathbf{w}} \in F_q^n$  on the other end;
- decode the corresponding  $\hat{w}(z)$  leveraging the parity-check polynomial h(z).

Cyclic coding. For the purpose of fast decoding, we now present two methods with which the basis polynomial v(z) may be leveraged to generate a codeword polynomial w(z) in systematic form [that is, rather than taking w(z) = d(z)v(z)]. By convention, the systematic form in the cyclic case usually shifts the k data symbols in d(z) to the end of the codeword, that is:

$$w(z) = b(z) + z^{n-k}d(z)$$

$$= b_0 + b_1 z + \dots + b_{n-k-1} z^{n-k-1} + d_0 z^{n-k} + d_1 z^{n-k+1} + \dots + d_{k-1} z^{n-1}.$$
(5.27)

If k/n < 0.5, a recursive approach may be used to determine the parity symbols in b(z). By (5.25b) and the fact that each valid codeword polynomial w(z) is itself a linear combination of the basis polynomials v(z), it is seen that

$$u(z) \mod (z^n - 1) = 0$$
 where  $u(z) \triangleq h(z) w(z) = u_0 + u_1 z + u_2 z^2 + \dots$ 

Initializing the last k symbols of w(z) as shown in (5.27), the remaining symbols of w(z) may thus be determined from the resulting convolution formulae for  $u_{n-1}$  through  $u_k$  as follows:

 $u_{n-1} = h_0 w_{n-1} + \ldots + h_k w_{n-k-1} = 0 \Rightarrow w_{n-k-1} = -[h_0 w_{n-1} + \ldots + h_{k-1} w_{n-k-2}]/h_k$ 

$$u_{n-2} = h_0 w_{n-2} + \dots + h_k w_{n-k-2} = 0 \quad \Rightarrow \quad w_{n-k-2} = -[h_0 w_{n-2} + \dots + h_{k-1} w_{n-k-3}] / h_k,$$

$$\vdots$$

$$u_k = h_0 w_k + \dots + h_k w_0 = 0 \quad \Rightarrow \quad w_0 = -[h_0 w_k + \dots + h_{k-1} w_1] / h_k.$$

If k/n > 0.5, a polynomial division approach to determine the parity symbols is more efficient. This is accomplished by writing the shift of the data symbols as some multiple of the basis polynomial v(z) plus a remainder r(z):

$$z^{n-k}d(z) = q(z)v(z) + r(z) \quad \Rightarrow \quad [z^{n-k}d(z)] \bmod v(z) = r(z),$$

where the mod command is interpreted as in (5.25b). Since the degree of v(z) is (n-k), the maximum degree of r(z) is (n-k-1). Calculating r(z) as shown above, taking b(z) = -r(z), and rearranging the above equations, it is seen that

$$w(z) = b(z) + z^{n-k}d(z) = q(z)v(z),$$

thus verifying that the polynomial w(z) so generated is in fact a valid codeword polynomial, as it is a multiple of the basis polynomial v(z).

### 5.5. DECODING

**Cyclic decoding.** In single parity-check codes, single symbol errors are flagged if  $h(z)\hat{w}(z) \neq 0$ . In repetition codes, the symbols of  $\hat{w}(z)$  may be corrected by simple majority vote.

Decoding of the binary Hamming and the extended binary Golay codes is introduced in §5.5.1. Such syndrome decoding methods extend easily to codes in cyclic form, in which the required syndrome computations are especially streamlined, as now shown. Note that any valid codeword polynomial w(z) is a multiple of the basis polynomial v(z); the *syndrome polynomial* s(z) of the received-message polynomial  $\hat{w}(z)$  is thus given by the remainder:

$$s(z) = \hat{w}(z) \bmod v(z).$$

Since the degree of v(z) is (n-k), the maximum degree of s(z) is (n-k-1), and thus the corresponding syndrome vector  $\mathbf{s}$  is of order m=(n-k), as expected [see discussion after (5.22)].

The polynomial multiplications and divisions involved in the cyclic coding and decoding algorithms described above are easy to code and efficient to calculate in either an *application-specific integrated circuit* (ASIC) or a *field-programmable gate array* (FPGA), in which repeated computations with shifted data may be performed quickly. The reduced storage associated with the vector representation of the basis matrix and the parity-check matrix in cyclic form help to facilitate such implementations.

#### 5.5.3 Shannon's theorem and turbo codes

The low-dimensional LBC, LTC, and LQC constructions given above are now supplanted by the more complex *turbo* codes for high performance coding applications such as 10GBase-T ethernet and deep space communication. Though these codes are generally much longer than the simple codes discussed above, they are built on the same fundamental principles, and achieve a coding efficiency over a noisy channel that is very close to the celebrated Shannon limit (Shannon 1949). For more information on such codes, the reader is referred to Gallager (1963), Berrou *et al.* (1993), and Moon (2005). Note also that the study of ternary and

quaternary codes is far more than a mathematical curiosity; new memory storage technology concepts leveraging, for example, DNA-based storage, with a four-character alphabet  $\{A, T, G, C\}$ , directly motivate the further development of non-binary error-correcting coding strategies.

#### 5.5.4 Soft-decision decoding

The type of decoding discussed in §5.5.1-5.5.3, in which the received vector  $\hat{\mathbf{w}}$  is assumed to be in  $\mathbf{F}_q^n$ , is known as *hard-decision decoding*.

Another formulation of the decoding problem assumes again that  $\mathbf{w} \in \mathbf{F}_q^n$ , but that  $\hat{\mathbf{w}} \in \mathbb{R}^n$ . The decoding problem in this case, called *soft-decision decoding*, is similar to that considered before (again, to find the most likely codeword  $\mathbf{w}$  corresponding to  $\hat{\mathbf{w}}$ , and the original data vector  $\mathbf{d}$  that generated it), but is now based on finding the codeword  $\mathbf{w}$  that minimizes the Euclidian distance to  $\hat{\mathbf{w}}$  rather than that which minimizes the Hamming distance.

For example, consider the soft-decision decoding of a binary parity check code. Assume that the transmitted codeword  $\mathbf{w} \in \mathbf{F}_2^n$  (that is, the symbols being transmitted are binary, and in this case rescaled to be  $\pm 1$ ) but that the received message  $\hat{\mathbf{w}} \in \mathbb{R}^n$  (that is, the symbols received are real). In this case, we may decode the received message by initially taking  $\mathbf{w} = \text{sign}(\hat{\mathbf{w}})$ . If the resulting decoded vector fails the parity check, we simply take the decision that we were least certain about (that is, the element of  $\hat{\mathbf{w}}$  that is closest to zero) and round it the other direction; this is known as *Wagner's decoding rule* (Silverman & Balser 1954).

Many soft-decision decoding algorithms are essentially generalizations of Wagner's decoding rule. Further, most soft-decision decoding algorithms may be framed as straightforward restrictions of a corresponding lattice quantization algorithm (see §6) to the appropriate subset of the lattice in question.

$$z^{5} - 1 = (z+1)(z^{4} + z^{3} + z^{2} + z + 1)$$

$$z^{7} - 1 = (z+1)(z^{3} + z + 1)(z^{3} + z^{2} + 1)$$

$$z^{15} - 1 = (z+1)(z^{2} + z + 1)(z^{4} + z + 1)(z^{4} + z^{3} + 1)(z^{4} + z^{3} + z^{2} + z + 1)$$

$$z^{23} - 1 = (z+1)(z^{11} + z^{9} + z^{7} + z^{6} + z^{5} + z + 1)(z^{11} + z^{10} + z^{6} + z^{5} + z^{4} + z^{2} + 1)$$

Table 5.1. Unique irreducible factors of  $(z^n - 1)$  over  $\mathbf{F}_2$  for various values of n.

$$z^{4} - 1 = (z+2)(z+1)(z^{2}+1)$$

$$z^{11} - 1 = (z+2)(z^{5} + 2z^{3} + z^{2} + 2z + 2)(z^{5} + z^{4} + 2z^{3} + z^{2} + 2)$$

$$z^{13} - 1 = (z+2)(z^{3} + 2z + 2)(z^{3} + z^{2} + 2)(z^{3} + z^{2} + z + 2)(z^{3} + 2z^{2} + 2z + 2)$$

Table 5.2. Unique irreducible factors of  $(z^n - 1)$  over  $\mathbb{F}_3$  for various values of n.

$$z^5 - 1 = (z^2 + \omega z + 1)(z^3 + \omega z^2 + \omega z + 1)$$

Table 5.3. A useful (though nonunique) factorization of  $(z^5 - 1)$  over  $\mathbf{F}_4$ ; note that Table 5.1 provides an alternative factorization of  $(z^5 - 1)$  over  $\mathbf{F}_2$  which is also valid over  $\mathbf{F}_4$ .

. / \

code	description	v(z)	h(z)					
$[n, n-1, 2]_2$	§5.2.1	z+1	$z^{n-1} + z^{n-2} + \dots + z + 1$					
$[n,1,n]_2$	§5.2.2	$z^{n-1} + z^{n-2} + \ldots + z + 1$	z+1					
$[7,4,3]_2$	§5.2.3	$z^3 + z + 1$	$z^4 + z^2 + z + 1$					
$[15, 11, 3]_2$	§5.2.3	$z^4 + z + 1$	$z^{11} + z^8 + z^7 + z^5 + z^3 + z^2 + z + 1$					
$[31, 26, 3]_2$	§5.2.3	$z^5 + z^2 + 1$	$(z^{31}-1)/(z^5+z^2+1)$ over $\mathbf{F}_2$					
$[63,57,3]_2$	§5.2.3	$z^6 + z + 1$	$(z^{63}-1)/(z^6+z+1)$ over $\mathbf{F}_2$					
[127 120 3]	85 2 3	$a^{7} + a^{3} + 1$	$(z^{127} - 1)/(z^7 + z^3 + 1)$ over <b>F</b> <sub>2</sub>					

# Chapter 6

# Further connections between lattice theory and coding theory

# **Contents**

6.1	Quantization onto lattices	. 100
	6.1.1 Quantization to $\mathbb{Z}^n$	100

# 6.2.1 Cases with $n \le 8$ 106 6.2.2 Cases with n > 8 106

CHAPTER 6 FURTHER CONNECTIONS BETWEEN LATTICE AND CODING THEORY

# **6.1** Quantization onto lattices

We now introduce some methods for quantization from an arbitrary point  $\mathbf{x}$  in  $\mathbb{R}^n$  onto a point  $\tilde{\mathbf{x}}$  on a discrete lattice, which may be defined via integer linear combination of the columns of the corresponding basis matrix B. The solution to this problem is lattice specific, and is thus treated lattice by lattice in the subsections below. Note that §6.1.1 through §6.1.4 are adapted from Conway & Sloane (1998), and §6.1.5 is adapted from Vardy

& Be'ery (1993). Note also that we neglect the problem of scaling of the lattices in this discussion, which is trivial to implement in code. **6.1.1 Ouantization to**  $\mathbb{Z}^n$ 

### Quantize to $\mathbb{Z}^n$ simply by rounding each element of **x** to the nearest integer.

6.1.4

6.1.5

# **6.1.2** Quantization to $D_n$

Quantize to  $D_n$  by rounding **x** two different ways:

- Round each element of  $\mathbf{x}$  to the nearest integer, and call the result  $\hat{\mathbf{x}}$ .
- Round each element of x to the nearest integer, and call the result x.
  Round each element of x to the nearest integer except that element of x which is furthest from an integer, and round that element the wrong way (that is, round it down instead of up, or up instead of down); call

the result  $\hat{\hat{\mathbf{x}}}$ .

Compute the sum *s* of the individual elements of  $\hat{\mathbf{x}}$ ; the desired quantiziation is  $\tilde{\mathbf{x}} = \hat{\mathbf{x}}$  if is *s* is even, and  $\tilde{\mathbf{x}} = \hat{\hat{\mathbf{x}}}$  if *s* is odd.

### **6.1.3** Quantization to $A_n$

The  $A_n$  lattice is defined in an n-dimensional subspace C of Y =  $\mathbb{R}^{n+1}$ . The subspace C is spanned by the n columns of the corresponding basis matrix  $B_{A_n}$ , and the orthogonal complement of C is spanned by the vector  $\mathbf{n}_{A_n}$ . Thus, the nearest point in the subspace,  $\mathbf{y}_{\mathsf{C}} \in \mathsf{C}$ , to any given point  $\mathbf{y} \in \mathsf{Y}$  is given by

$$\mathbf{y}_{\mathsf{C}} = \mathbf{y} - (\mathbf{y}, \mathbf{n}_{A_n}) \cdot \mathbf{n}_{A_n}.$$

An orthogonal basis  $\hat{B}_{A_n}$  of C may easily be determined from  $B_{A_n}$  via Gram Schmidt orthogonalization. With this orthogonal basis, the vectors  $\mathbf{x} \in \mathbb{R}^n$  comprising the  $A_n$  lattice may be related to the corresponding vectors  $\mathbf{y}_C \in C \subset Y$  (that is, on an n-dimensional subspace of  $\mathbb{R}^{n+1}$ ) via the equation

$$\mathbf{y}_{\mathsf{C}} = \hat{B}_{A_n} \mathbf{x}. \tag{6.1a}$$

Thus, starting from some point  $\mathbf{x} \in \mathbb{R}^n$  but not yet quantized onto the lattice, we can easily determine the corresponding (n+1)-dimensional vector  $\mathbf{y}_{\mathsf{C}}$  which lies within the *n*-dimensional subspace  $\mathsf{C}$  of  $\mathbb{R}^{n+1}$  via (6.1a). Given this value of  $\mathbf{y}_{\mathsf{C}} \in \mathsf{C}$ , we now need to quantize onto the lattice. We may accomplish this with the following simple steps:

- Round each component of  $\mathbf{y}_{\mathsf{C}}$  to the nearest integer, and call the result  $\hat{\mathbf{y}}$ . Define the deficiency  $\Delta = \sum_{i} \hat{y}_{i}$ , which quantifies the orthogonal distance of the point  $\hat{\mathbf{y}}$  from the subspace  $\mathsf{C}$ .
- If  $\Delta = 0$ , then  $\tilde{\mathbf{y}} = \hat{\mathbf{y}}$ . If not, define  $\mathbf{d} = \mathbf{y}_{\mathsf{C}} \hat{\mathbf{y}}$ , and distribute the integers  $0, \dots, n$  among the indices  $i_0, \dots, i_n$  such that

$$-1/2 \le d(\hat{y}_{i_0}) \le d(\hat{y}_{i_1}) \le \ldots \le d(\hat{y}_{i_n}) \le 1/2.$$

# CHAPTER 6 FURTHER CONNECTIONS BETWEEN LATTICE AND CODING THEORY

If  $\Delta > 0$ , then nudge  $\hat{\mathbf{y}}$  back onto the C subspace by defining  $\tilde{y}_{i_k} = \begin{cases} \hat{y}_{i_k} - 1 & k < \Delta, \\ \hat{y}_{i_k} & \text{otherwise.} \end{cases}$ If  $\Delta < 0$ , then nudge  $\hat{\mathbf{y}}$  back onto the C subspace by defining  $\tilde{y}_{i_k} = \begin{cases} \hat{y}_{i_k} + 1 & k > n + \Delta, \\ \hat{y}_{i_k} & \text{otherwise.} \end{cases}$ 

Back in *n*-dimensional parameter space, the quantized value  $\tilde{y} \in C$  corresponds to

$$\tilde{\mathbf{x}} = \hat{B}_{A_n}^T \tilde{\mathbf{y}}.\tag{6.1b}$$

### **Quantization to the union of cosets**

The dual lattices  $D_n^*$  and  $A_n^*$ , the triangular lattice  $A_2$ , and the packing  $D_n^+$  (including the lattice  $E_8 \cong E_8^* \cong D_8^+$ ) are described via the union of simple, real cosets in (2.4a), (2.7a), (2.6c), and (2.5), respectively. The lattices  $E_7$  and  $E_7^*$  may be built via the union of simple, real cosets via Construction A [see (5.4a)], with coset representatives  $\mathbf{w}_{[n,k,d]}^{i}$  defined in (5.8) and (5.9) respectively. To quantize a lattice described in such a manner (as a union of simple cosets), one may quantize to each coset independently, then select from these individual quantizations that lattice point which is nearest to the original point  $\mathbf{x}$ .

The lattices  $E_6$  and  $E_6^*$  may be built via the union of complex cosets [which are scaled and shifted complex  $\mathscr{E}$  lattices  $\mathbb{Z}[\omega]^3$ ] via Construction  $A_{\mathscr{E}}^{\pi}$  [see (5.5a)], with coset representatives  $\mathbf{w}_{[n,k,d]}^i$  given in (5.13) and (5.12) respectively. Following Conway & Sloane (1984), to discretize a point **x** to coset *i* in these cases:

- Determine the complex vector  $\mathbf{z} \in \mathbb{C}^3$  corresponding to  $\mathbf{x} \in \mathbb{R}^6$ . Shift and scale such that  $\hat{\mathbf{z}} = (\mathbf{z} \mathbf{a}_i)/\theta$ .
- Determine the real vector  $\hat{\mathbf{x}} \in \mathbb{R}^6$  corresponding to  $\hat{\mathbf{z}} \in \mathbb{C}^3$ . Quantize the first, second, and third pairs of elements of  $\hat{\mathbf{x}}$  to the real triangular  $A_2$  lattice to create the quantized vector  $\hat{\tilde{\mathbf{x}}}$ .
- Determine the complex vector  $\hat{\mathbf{z}} \in \mathbb{C}^3$  corresponding to  $\hat{\mathbf{x}} \in \mathbb{R}^6$ . Unscale and unshift such that  $\tilde{\mathbf{z}} = \theta \hat{\mathbf{z}} + \mathbf{a}_i$ .
- Determine the real vector  $\tilde{\mathbf{x}} \in \mathbb{R}^6$  corresponding to  $\tilde{\mathbf{z}} \in \mathbb{C}^3$ .

#### **6.1.5** Quantization to $\Lambda_{24}$

We now jump to the Leech lattice in dimension n = 24. Recall from §2.6 that the best lattices in dimensions n = 9 to n = 23 may all be determined as lower-dimensional cross-sections of  $\Lambda_{24}$ ; once the (difficult) n = 24 case is mastered, quantization to these intermediate dimensions is relatively straightforward.

Efficient quantization to  $\Lambda_{24}$  is a problem that received intense scrutiny in the 1980s and early 1990s. The best algorithm described in the literature, due to Vardy & Be'ery (1993), is based on the construction of  $\Lambda_{24}$  described in the last paragraph of §5.4.5, and essentially represents a culmination of the previous efforts that led to it. This remarkable algorithm requires only about 3000 to 3600 floating-point operations and comparisons, and a comparable number of integer operations and comparisons, to compute the point of the  $\Lambda_{24}$  lattice that is closest to any given point  $\mathbf{r} \in \mathbb{R}^{24}$ . The algorithm leverages effectively many of the fundamental symmetries inherent in  $\Lambda_{24}$ , including its close relationships with both carefully-chosen subsets of the  $D_2$  lattice (Figure 5.5) as well as the  $[6,3,4]_4$  hexacode (§5.4.5).

Though it was proposed in 1993, the logic inherent to this algorithm is so intricate that, as of the writing of this review, an executable version of it did not appear to be readily available in the literature. We have thus written an efficient Fortran 90 implementation of this algorithm, which is available online at:

http://renaissance.ucsd.edu/software/DecodeLeech.tgz

This implementation is thoroughly commented, and is written in a notation consistent with that of Vardy & Be'ery (1993). Thus, in addition to being a useful code for new practical applications of the Leech lattice in science and engineering, it is hoped that this executable code can itself be a helpful guide in the understanding of this complex algorithm.

<sup>&</sup>lt;sup>1</sup>Our implementation of this algorithm executes in about 0.3 milliseconds on a 2008 vintage laptop (2.53GHz Intel Core 2 Duo), which is sufficiently fast for many applications. It is also trivial to parallelize this code efficiently over four separate computational threads, as quantization to each Leech quarter lattice is handled independently.

# CHAPTER 6. FURTHER CONNECTIONS BETWEEN LATTICE AND CODING THEORY using the notation introduced at the end of 85.4.5, this algorithm first splits the problem of

In short, using the notation introduced at the end of §5.4.5, this algorithm first splits the problem of quantizating a point  $\mathbf{r} \in \mathbb{R}^{24}$  to the nearest  $\Lambda_{24}$  point into two subproblems:

- quantizing to  $H_{24}$ ; that is, when forming the original vector  $\mathbf{r} \in \mathbb{R}^{24}$  into a  $2 \times 6$  array of points  $\mathbf{r}_{hn} \in \mathbb{R}^2$  for h = 0, 1 and  $n = 0, \dots, 5$ , quantizing each  $\mathbf{r}_{hn}$  to the best  $A_{ijk}$  points in the Ungerboeck partitioning of  $D_2$  such that the overall k parity of the array is even, while the projection of the  $2 \times 6$  array of points forms a codeword of the  $[6,3,4]_4$  hexacode; and
- quantizing to  $H_{24} + \mathbf{a}$ ; that is, quantizing to the best  $B_{ijk}$  points in the Ungerboeck partitioning of  $D_2$  such that the overall k parity of the array is odd, while, again, the projection of the  $2 \times 6$  array of points forms a codeword of the  $[6,3,4]_4$  hexacode.

The best of the two lattice points selected by these subproblems is then returned.

During the execution of each of these two subproblems, the closest point to  $\mathbf{r}_{hn}$  in each  $A_{ijk}$  family (in the even overall k parity case) or in each  $B_{ijk}$  family (in the odd overall k parity case) is first identified, and the *squared Euclidian distance* (SED) to each of these points is calculated. For each i and j, the "preferred" value of k (that is, the one that leads to the least SED for that point) is determined, and the SED penalty  $\delta$  for chosing the other value of k is computed. The algorithm then further splits the quantization to  $H_{24}$  (and, similarly to  $H_{24} + \mathbf{a}$ ) into two smaller sub-subproblems:

- quantizing to Q<sub>24</sub>; that is, to arrays with the specified overall k parity such that, additionally, the overall h parity is even; and
   quantizing to Q<sub>24</sub> + b; that is, to arrays with the specified overall k parity such that, additionally, the overall
- h parity is odd.

The best of the two lattice points selected by these sub-subproblems is then returned.

The quantization to  $Q_{24}$  and its 3 translates is, in turn, decomposed into 5 distinct steps:

1. Only two sets of indices  $\{i_0, j_0, i_1, j_1\}$  project to each symbol  $p \in \mathbb{F}_4$ ; in this step, for each symbol p and for

each column n of the  $2 \times 6$  array, we identify the "preferred representation" as that set which, when taken together with their corresponding preferred values of  $k_0$  and  $k_1$ , minimize the SED of the column, and the other set, referred to as the "non-preferred representation"; we also calculate the SED penalty associated with chosing the non-preferred representation. Conveniently, it turns out that the preferred representation and the non-preferred representation necessarily have opposite h parity.

- 2. The three lists of penalties associated with changing the column-wise *k* parities (case 0), the column-wise *h* parities (case 1), or both (case 2) are then sorted (our implementation uses mergesorts, due to their cache efficiency; heapsorts or quicksorts are viable alternatives).
- The SED for each preferred "block" (that is, each pair of columns) is then computed.
   For each of the 64 codewords of the hexacode [see (5.20)], we then find the smallest possible correction(s)
- to the set of preferred representations such that the total k parity and the total h parity match the specified values required for the particular translate of Q<sub>24</sub> being considered (of 4 possible cases). This step leverages the sorted lists computed in step 2.
  5. For each of 16 sets of symbols [given by w<sub>0</sub> ∈ F<sub>4</sub> and w<sub>1</sub> ∈ F<sub>4</sub>], calculate the total SED of corrected representations, determined in step 4, corresponding to the 4 valid codewords of the hexacode [given by
- representations, determined in step 4, corresponding to the 4 valid codewords of the hexacode [given by  $w_2 \in \mathbf{F}_4$  and  $\{w_3, w_4, w_5\}$  selected according to  $V_{[6,3,4]_4}$  defined in (5.20)]. We then find the minimum total SED amongst these 16 corrected representations, and return the corresponding lattice point.

# 6.2 Enumerating nearest-neighbor lattice points

In the practical use of lattices in engineering applications, one occasionally needs to generate a list of all lattice points that are nearest neighbors to a given lattice point. It is sufficient to generate a list of all lattice points that are nearest neighbors of the origin, then to shift these points as necessary to the vicinity of any other lattice point. The present section describes two methods to generate such lists of nearest neighbors on a

CHAPTER 6 FURTHER CONNECTIONS BETWEEN LATTICE AND CODING THEORY

lattice. **6.2.1** 

### Cases with $n \leq 8$

Noting first (see §2.1) that a basis matrix B of an n-dimensional lattice might itself have more than n rows, the following algorithm is found to be effective for all lattices up to about n = 8:

- 0. Initialize p = 1.
- integers  $\{-p, ..., 0, ..., p\}$ , and that *all possible vectors* that can be created in such a fashion, except the origin, are present (without duplication) in this distribution.

  2. Compute the distance of each transformed point  $\tilde{\mathbf{y}}^i = B\tilde{\mathbf{z}}^i$  in this distribution from the origin, and eliminate

1. Define a distribution of points  $\tilde{z}^i$  such that each element of each of these vectors is selected from the set of

- 2. Compute the distance of each transformed point y = BZ in this distribution from the origin, and eminiate those points in the distribution that are farther from the origin than the minimum distance computed in the set.
  3. Count the number of points remaining in the distribution. If this number equals the (known) kissing number
- of the lattice under consideration, as listed in Tables 3.1-3.2, then determine an orthogonal  $\hat{B}$  from B via Gram Schmidt orthogonalization, set  $\tilde{\mathbf{x}}^i = \hat{B}^T \tilde{\mathbf{y}}^i$  for all i, and exit; otherwise, increment p and repeat from step 1. Though this simple algorithm is not at all efficient, for  $n \le 8$  it really need not be, as the nearest neighbor

Though this simple algorithm is not at all efficient, for  $n \le 8$  it really need not be, as the nearest neighbor distribution is identical around every lattice point, and thus this algorithm need only be run once for any given lattice.

## **6.2.2** Cases with n > 8

For n > 8, the algorithm described above is prohibitively expensive. We thus focus here on an efficient manner of obtaining the 196,560 nearest neighbors to the origin of the Leech lattice  $\Lambda_{24}$ , then on the restriction of this set of neighbors, one dimension at a time, down to n = 9.

#### 6.2. ENUMERATING NEAREST-NEIGHBOR LATTICE POINTS

To proceed, it is first necessary to enumerate the codewords of the binary Golay code following the approach described in §5.2.7. Recall that the basis matrix of the binary Golay code has dimension  $24 \times 12$ ; thus, the  $2^{12} = 4096$  codewords of the binary Golay code follow immediately as a binary linear combination (that is, as a linear combination, mod 2, with binary coefficients) of the columns of this matrix.

Then, in order to identify all of the nearest neighbors of the Leech lattice, we may proceed (following Conway & Sloane 1998) by constructing three distinct sets of points:

- The first set, consisting of 98,304 points, is obtained using the binary Golay codewords discussed above. Construct first a  $24 \times 24$  matrix A with -3 everywhere along the main diagonal and 1 everywhere else. Then, take each codeword of the binary Golay code, one at a time, replace each 0 with -1, and perform elementwise multiplication of this modified codeword to each column of A, thereby generating 24 points for each of the  $2^{12}$  binary Golay codewords, or  $2^{12} \cdot 24 = 98,304$  points.
- The next set, consisting of 1,104 points, is composed of vectors with 22 zero elements and two elements that are either 4 or -4. As there are 276 ways to select the locations of the nonzero elements, and  $2^2 = 4$  valid ways to populate them, we obtain  $2^2 \cdot 276 = 1,104$  points.
- The third set, consisting of 97,152 points, is obtained using the 759 vectors of the Witt design, which are just the 759 binary Golay codewords (discussed above) of weight 8. Note that each of these vectors has 8 ones and 16 zeros. Construct an  $8 \times 128$  matrix C such that each element of each column is either a 2 or -2, with an even number of minus signs in each column (note that there are  $2^7 = 128$  such columns possible). We then distribute the elements in each of the 128 columns of C into each of 8 positions where the ones sit in each of the 759 vectors of the Witt design, thereby obtaining the remaining  $128 \cdot 759 = 97,152$  points.

The 98,304 + 1,104 + 97,152 = 196,560 points so generated are the nearest neighbors to the origin of  $\Lambda_{24}$ . Then, throwing out those points  $\mathbf{z}$  for which  $\mathbf{z} \cdot \mathbf{n}_{\Lambda_{23}} \neq 0$  (see §2.6) leaves the 93,150 neighbors of  $\Lambda_{23}$ ; additionally throwing out those points  $\mathbf{z}$  for which  $\mathbf{z} \cdot \mathbf{n}_{\Lambda_{22}} \neq 0$  leaves the 49,896 neighbors of  $\Lambda_{22}$ ; etc.

# Part II **Derivative-free optimization**

7	Extending lattice theory for coordinated derivative-free optimization	

7	Extending lattice theory for coordinated derivative-free optimization	109
8	Kriging interpolation	129

8	Kriging interpolation	129
9	Global ontimization leveraging Kriging-based interpolation	130

9	Global optimization leveraging Kriging-based interpolation	139
10	Lattice-based derivative-free ontimization via global surrogates (LARDOGS)	147

10	Lattice-based derivative-free optimization via global surrogates (LABDOGS)	147
11	Optimization via incomplete function evaluations (αDOGS)	161

11	Optimization via incomplete function evaluations (\alpha DOGS)	161
12	Ontimization in the presence of complex constraint boundaries (latticeMADS)	163

# **Chapter 7**

# **Extending lattice theory for coordinated derivative-free optimization**

### **Contents**

7.1	Introduction to	o derivative-free op	otimization	 
	7.1.1 (7)	. 1 6 16	the second second second	 110

7.1.4	The surrogate management framework (SMF)	114
7.1.5	Framing the search for a uniform simplex as a discrete Thomson problem	116

CHAPTER 7 EXTENDING LATTICE THEORY FOR DERIVATIVE-FREE OPTIMIZATION

The minimization of computationally expensive, high-dimensional functions is often most efficiently per-

# 7.1 Introduction to derivative-free optimization

7.2

formed via gradient-based optimization algorithms such as nonlinear conjugate gradients and L-BFGS-B. In complex systems for which an accurate computer model is available, the gradient required by such algorithms may often be found via adjoint analysis. However, when the function in question is not sufficiently smooth to leverage gradient information effectively during its optimization (see, e.g., Figure 7.1), a derivative-free approach is necessary. Such a scenario is evident, for example, when optimizing a finite-time-average approximation of an infinite-time-average statistic of a chaotic system such as a turbulent flow. Such an approximation may be determined via simulation or experiment. The truncation of the averaging window used to determine this approximation renders derivative-based optimization strategies ill suited, as the truncation

proximation may be determined via simulation or experiment. The truncation of the averaging window used to determine this approximation renders derivative-based optimization strategies ill suited, as the truncation error, though small, is effectively decorrelated from one flow simulation/experiment to the next. This effective decorrelation of the truncation error is reflected by the exponential growth, over the entire finite time horizon considered, of the adjoint field related to the optimization problem of interest in the simulation-based setting. As a result, derivative-free algorithms are often required for the optimization of nonsmooth scalar func-

tions in n dimensions. The core idea of all efficient algorithms for problems of this type is to keep function

#### 7.1. INTRODUCTION TO DERIVATIVE-FREE OPTIMIZATION

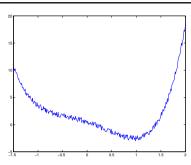


Figure 7.1: Prototypical nonsmooth optimization problem for which local gradient information is ill suited to accelerate the optimization algorithm.

evaluations far apart until convergence is approached. *Generalized pattern search (GPS)* algorithms, a modern class of methods particularly well suited to such problems, accomplish this by coordinating the search with an underlying grid which is refined, and coarsened, as appropriate.

One of the most efficient subclasses of GPS algorithms, known as the *surrogate management framework* (*SMF*; see Booker *et al.* 1999), alternates between an exploratory *search* over an interpolating function which summarizes the trends exhibited by existing function evaluations, and an exhaustive *poll* which checks the

function on neighboring points to confirm or confute the local optimality of any given *candidate minimum point* (*CMP*) on the underlying grid. The original SMF algorithm implemented a GPS step on an underlying Cartesian grid, augmented with a Kriging-based surrogate search. Rather than using the *n*-dimensional Cartesian grid (the typical choice), Part II of this text suggests the use of lattices derived from *n*-dimensional

sphere packings. As reviewed in Part I, such lattices are significantly more uniform and have many more nearest neighbors than their Cartesian counterparts. Both of these facts make them far better suited for coordinating GPS algorithms<sup>1</sup>, as demonstrated in a variety of numerical tests presented later in Part II.

### 7.1.1 The inherent role of uniform simplexes in derivative-free optimization

One of the earliest derivative-free optimization approaches to appear in the literature is the downhill simplex method (see Spendley, Hext, & Himsworth 1962 and Nelder & Mead 1965). The downhill simplex method is inherently based on an iterative, amoeba-like evolution (moving one point at a time) of a set of n+1points in n dimensions towards the minimum of a (possibly, nonsmooth) function. A large body of literature appeared after the original introduction of this method, much of which was aimed at heuristic strategies designed to keep the evolving simplex as regular as possible as the iteration proceeds, while expanding or contracting as appropriate. The grid-based methods considered in the present work are fundamentally different, so we will not dwell on such grid-free methods in this introduction. However, it is worth noting the inherent dependence on the regularity an evolving *simplex* (that is, on an *n*-dimensional polytope with n+1vertices) in this classical method, and an analogous focus in the present work on the identification (see §7.3) and characterization (see §7.2 and 7.5) of a maximally-uniform simplex (referred to in the present work as a minimum positive basis) around the best point encountered thus far as the iteration proceeds, referred to in the present work as a candidate minimum point. The role of the simplex in both cases is essentially identical: to identify the best direction to move next using a minimum number of new function evaluations.

<sup>&</sup>lt;sup>1</sup>In fact, as mentioned previously, Conway & Sloane (1998, p. 12) state: "A related application that has not yet received much attention is the use of these packings for solving n-dimensional search or approximation problems"; this is exactly the focus of Part II.

candidate minimum point (CMP).

### 7.1.2 Global convergence via a dumb method: exhaustive sampling (ES)

discussed above, turbulent flow simulations or experiments), a derivative-free optimization algorithm which converges to within an accurate tolerance of the global minimum of a nonconvex function of interest with a minimum number of function evaluations is desired. It is noted that, in the general case, proof of convergence of an optimization algorithm to a global minimum is possible only when, in the limit of a large number of function evaluations N, the function evaluations become dense in the feasible region of parameter space (Torn & Zilinskas, 1987). Though the algorithms developed in the present work, when implemented properly, satisfy this condition, so do far inferior approaches, such as a rather unintelligent algorithm which we call exhaustive sampling (ES), which simply covers the feasible parameter space with a grid, evaluates the function at every gridpoint, refines the grid by a factor of two, and repeats until terminated. Thus, a guarantee of global convergence is not sufficient to establish the efficiency of an optimization algorithm. If function evaluations are relatively expensive, and thus only a relatively small number of function evaluations can ultimately be afforded, effective heuristics for rapid convergence are perhaps even more important than rigorous proofs of the behavior of the optimization algorithm in the limit of large N, a limit that might actually be argued to be of limited relevance when function evaluations are expensive. Given that such algorithms are often used in applications in which only a few hundred function evaluations can be afforded, careful attention to such

Due to the often significant expense associated with performing repeated function evaluations (for example, as

# heuristics forms an important foundation for the present study. 7.1.3 Successive polling (SP) and generalized pattern search (GPS) algorithms

If, for the moment, we give up on the goal of global convergence, the perhaps simplest grid-based derivative-free optimization algorithm, which we call *successive polling (SP)*, proceeds as follows:

free optimization algorithm, which we call *successive polling (SP)*, proceeds as follows:

• Start with a coarse grid and evaluate the function at some starting point on this grid, identified as the first

### CHAPTER 7. EXTENDING LATTICE THEORY FOR DERIVATIVE-FREE OPTIMIZATION

- Then, poll (that is, evaluate) the function values on gridpoints which neighbor the CMP in parameter space, at a sufficient number of gridpoints to *positively span*<sup>2</sup> the feasible neighborhood of the CMP [this step ensures convergence, as discussed further in Torczon 1997, Booker *et al.* 1999, and Coope & Price 2001]. When polling:
- (a) If any poll point is found to have a function value less than that of the CMP, immediately consider this new point the new CMP and terminate the present poll step.
- (b) If no poll points are found to have function values less than that of the CMP, refine the grid by a factor of two.
- Initiate a new poll step, either (a) around the new CMP or (b) around the old CMP on the refined grid, and repeat until terminated.

Though the basic SP algorithm described above, on its own, is not very efficient, there are a variety of effective techniques for accelerating it. All grid-based schemes which effectively build on this basic SP idea are classified as *generalized pattern search* (*GPS*) algorithms.

### 7.1.4 The surrogate management framework (SMF)

The most efficient subclass of GPS algorithms, known as the Surrogate Management Framework (SMF; see Booker *et al.*, 1999), leverages inexpensive interpolating "surrogate" functions (often, Kriging interpolations are used) to summarize the trends of the existing function evaluations, and to provide suggested new regions of parameter space in which to perform one or more additional function evaluation(s) between each poll step. SMF algorithms thus alternate beween two steps:

<sup>&</sup>lt;sup>2</sup>That is, such that any feasible point in the neighborhood of the CMP can be reached via a *linear combination with non-negative coefficients* of the vectors from the CMP to the poll points.

# 7.1. INTRODUCTION TO DERIVATIVE-FREE OPTIMIZATION

chapters that follow, then discuss how they may be combined.

has a value greater than that of the CMP.

above.

There is substantial flexibility during the search step described above. An effective search is essential for an efficient SMF algorithm. In the case that the search behaves poorly and fails to return improved function values, the SMF algorithm essentially reduces to the SP algorithm. If, however, the surrogate-based search is effective, the SMF algorithm will converge to a minimum far faster than a simple SP-based minimization. As

the search and poll steps are essentially independent of each other, we will discuss them each in turn in the

(i) Search over the inexpensive interpolating function to identify, based on the existing function evaluations, the most promising gridpoint at which to perform a new function evaluation. Perform a function evaluation at this point, update the interpolating function, and repeat. The search step may be terminated either when it returns a gridpoint at which the function has already been evaluated, or when the function, once evaluated,

(ii) Poll the neighborhood of the new CMP identified by the search algorithm, following rules (a) and (b)

coarsened appropriately in order to explore the vicinity of this new CMP efficiently (that is, with a coarse grid first, then refined as necessary). Note also that the interpolating surrogate function of the SMF may be used to *order* the function evaluations of the poll step, such that those poll points which are most likely to have a function value lower than that of the CMP are evaluated first. By so doing, the poll steps will, on average, terminate sooner, and the computational cost of the overall algorithm may be reduced further.

Note that if the search produces a new CMP which is several gridpoints away from the previous function evaluations, which occasionally happens when exploring functions with multiple minima, the grid may be

To the best of our knowledge, all previous GPS and SMF implementations have been coordinated using Cartesian grids. However, like in the game of checkers (contrast "American" checkers with "Chinese" checkers).

Cartesian grids are not the only choice for discretizing parameter space. Other structured choices arising from

### CHAPTER 7. EXTENDING LATTICE THEORY FOR DERIVATIVE-FREE OPTIMIZATION

*n*-dimensional sphere packing theory (see Tables 7.1 and 7.2, and further characterizations in §3) are significantly more uniform and have many more nearest neighbors, especially as the dimension of the problem in question is increased; both of these properties suit these alternative lattices well for coordinating grid-based optimization algorithms.

Part I of this study consisely summarizes n-dimensional sphere packing theory, describing almost everything one needs to know about lattices up to dimension n=24 in order to use them effectively in practical engineering applications. To extend the lattice theory described in Part I of this text in order to coordinate a derivative-free optimization, a few additional component algorithms are needed, which are described in the remainder of §7. For simplicity, Part II focuses on the use of just two such lattices, the zero-sum lattice  $A_n$ , which is an n-dimensional analog of the 2-dimensional hexagonal lattice and the 3-dimensional face-centered-cubic lattice, and the Gosset lattice  $E_8$ , which is an 8-dimensional analog of the 3-dimensional diamond packing, and is especially uniform; both of these lattices are described completely in §2. The utility of other lattices in this setting will be explored in future work.

### 7.1.5 Framing the search for a uniform simplex as a discrete Thomson problem

Thomson (1904), in his study of the structure of the atom, is credited with being the first to address the problem<sup>3</sup>: "Where should k inimical dictators settle on a planet in order to be as far away from each other as possible?" This question extends naturally to n-dimensional planets, and has received significant attention in the years since Thomson's original paper. The question is readily answered numerically by assigning an identical "charge" to each of n identical "particles", restricting particle motion to the surface of the sphere, and iteratively moving each particle (with some damping applied) in the direction of the force caused by the other particles (projected onto the sphere) until all particles come to equilibrium. The precise solution

<sup>&</sup>lt;sup>3</sup>This curious problem, articulated by Meschkowski (1960) in terms of inimical dictators (see also L. Fejes Toth 1971), assumes that all locations on the planet's surface are equally desirable, and that the inimical dictators all cooperate.

# 7.1. INTRODUCTION TO DERIVATIVE-FREE OPTIMIZATION

which is proportional to the inverse square of this distance. The setting based on other distance measures (e.g., measured along the surface of the sphere instead of along a straight line) and other power laws are referred to as generalized Thomson problems; in particular, the case based on the p'th power in the limit that  $p \to \infty$  (that is, the max value) was studied in Tammes (1930), in his study of the boundaries of pollen grains.

reached is a function of the distance metric and power law used when computing the force between any two particles; in the electrostatic setting, Thomson used the Euclidian distance between the particles, and a force

In this chapter, we generalize this classical question in two ways, and introduce a new metric to characterize the solution found:

- First, the locations where the particles are allowed to settle are restricted to a discrete set of points on a sphere, which are specified as the nearest-neighbor lattice points to the CMP.
- Next, we allow some the particles' locations on the sphere to be specified (that is, fixed) in advance, and only move the remaining (free) particles to arrive at the best solution possible.
- Finally, the new metric we introduce is a check of whether or not the distribution produced by numerical solution of the resulting "discrete Thomson problem" forms a *positive basis* of the feasible neighborhood of the CMP; that is, in the case with no active constraints (cf. §7.4), whether or not all points on the unit sphere around the CMP can be reached via a linear combination *with non-negative coefficients* of the vectors from the CMP to the optimized particle locations.

After developing a method to test for a positive basis, the remainder of this section develops three efficient algorithms to iterate on this "discrete Thomson problem" until a positive basis is found. To accomplish this, these algorithms first solve the discrete Thomson problem numerically for n + m particles where m = 1. If the optimization algorithm succeeds in producing a positive basis, the algorithm exits; otherwise, m is increased by one and the process repeated until a positive basis is determined. The resulting algorithm is leveraged heavily during the poll step of the lattice-based SMF algorithms developed later in Part II.

### 7.2 Testing for a positive basis

Given a subset of the nearest-neighbor lattice points, we will at times need an efficient test to determine whether or not the vectors to these points from the CMP form a positive basis of the feasible domain around the CMP. Without loss of generality, we will shift this problem so that the CMP corresponds to the origin in the discussion that follows.

A set of vectors  $\{\tilde{\mathbf{x}}^1, \dots, \tilde{\mathbf{x}}^k\}$  for k > n+1 is said to positively span  $\mathbb{R}^n$  if any point in  $\mathbb{R}^n$  may be

reached via a linear combination of these vectors with non-negative coefficients. Since the 2n basis vectors  $\{\mathbf{e}^1,\dots,\mathbf{e}^n,-\mathbf{e}^1,\dots,-\mathbf{e}^n\}$  positively span  $\mathbb{R}^n$ , a convenient test for whether or not the vectors  $\{\tilde{\mathbf{x}}^1,\dots,\tilde{\mathbf{x}}^k\}$  positively span  $\mathbb{R}^n$  is to determine whether or not each vector in the set  $\mathsf{E}=\{\mathbf{e}^1,\dots,\mathbf{e}^n,-\mathbf{e}^1,\dots,-\mathbf{e}^n\}$  can be reached by a positive linear combination of the vectors  $\{\tilde{\mathbf{x}}^1,\dots,\tilde{\mathbf{x}}^k\}$ . That is, for each vector  $\mathbf{e}\in\mathsf{E}$ , a solution  $\mathbf{z}$ , with  $z_i\geq 0$  for  $i=1,\dots,k$ , to the equation  $\tilde{X}\mathbf{z}=\mathbf{e}$  is sought, where  $\tilde{X}=(\tilde{\mathbf{x}}^1,\dots,\tilde{\mathbf{x}}^k)$ . If such a  $\mathbf{z}$  exists for each vector  $\mathbf{e}\in\mathsf{E}$ , then the vectors  $\{\tilde{\mathbf{x}}^1,\dots,\tilde{\mathbf{x}}^k\}$  positively span  $\mathbb{R}^n$ ; if such a  $\mathbf{z}$  does not exist, then the vectors  $\{\tilde{\mathbf{x}}^1,\dots,\tilde{\mathbf{x}}^k\}$  do not positively span  $\mathbb{R}^n$ .

Thus, testing a set of vectors to determine whether or not it positively spans  $\mathbb{R}^n$  reduces simply to testing for the existence of a solution to 2n well-defined *linear programs* in standard form. Techniques to perform such tests, such as Matlab's linprog algorithm, are well developed and readily available. Further, if a set of k vectors positively spans  $\mathbb{R}^n$ , it is a simple matter to check whether or not this set of vectors is also a positive basis of  $\mathbb{R}^n$ , if such a check is necessary, simply by checking whether or not any subset of k-1 vectors chosen from this set also positively span  $\mathbb{R}^n$ . Note that a positive basis with k vectors will necessarily have k in the range  $n+1 \le k \le 2n$ ; the case with k=n+1 is referred to as a *minimal* positive basis, and the case with k=2n is referred to as a *maximal* positive basis.

# 7.3 Selecting a positive basis from nearest-neighbor lattice points

In §6 of Part I, we described how to enumerate all points which are nearest neighbors of the origin of a lattice (and thus, with the appropriate shift, all points which are nearest neighbors of any CMP on the lattice). In §7.2 above, we described how to test a subset of such points to see if the vectors from the origin to these points form a positive basis around the CMP. We now present a general algorithm to solve the problem of selecting a positive basis from the nearest-neighbors of the CMP using a minimal number of new poll points, while creating the maximum achievable angular uniformity between the vectors from the CMP to each of these points (that is, while minimizing the skewness of the resulting poll set). Note in Figure 7.2 that, as the number of nearest neighbors increases, the flexibility in solving this (apparently, NP-hard) problem also increases, though a perfectly distributed minimal positive basis (using n+1 points) is not always available. Ideally, for m=1, the solution to the discrete Thomson problem will produce a positive basis with good angular uniformity; if it does not, we may successively increment m by one and try again until we succeed in producing a positive basis. We have studied three algorithms for solving this problem:

angular uniformity; if it does not, we may successively increment m by one and try again until we succeed in producing a positive basis. We have studied three algorithms for solving this problem: Algorithm A. If the kissing number  $\tau$  of the lattice under consideration is relatively large (that is, if  $\tau \gg n$ ; for example, for the Leech lattice  $\Lambda_{24}$ ), then a straightforward algorithm can first be used to solve Thomson's problem on a continuous sphere in n dimensions. This can be done simply and quickly by fixing  $q \geq 0$  repulsive particles at the prespecified lattice points, and initializing n+m-q free repulsive particles on the sphere randomly. Then, at each iteration, a straightforward force-based algorithm may be used to move each free particle along the surface of the sphere a small amount in the direction that the other particles are tending to push it, and iterating until the set of particles approaches an equilibrium. The free particle that is nearest to a nearest-neighbor lattice point around the CMP is then moved to said lattice point and fixed there, and the remaining free particles adjusted until they reach a new equilibrium. This adjust/fix/adjust/fix sequence is repeated until all particles are fixed at lattice points.

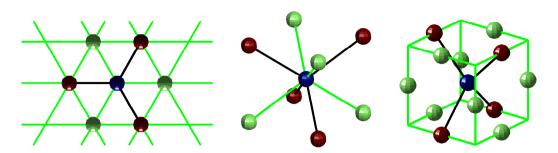


Figure 7.2: Various minimal positive bases (shown in red) around the origin (shown in blue) in the (left) triangular, (center) BCC, and (right) FCC lattices. Note that the triangular and BCC lattices each have two perfectly distributed minimal positive bases. In contrast, there are several choices for selecting a minimal positive basis in the FCC lattice, but none is perfectly distributed.

Algorithm B. If the kissing number  $\tau$  of the lattice under consideration is relatively small (that is, if  $\tau$  is not well over an order of magnitude larger than n), then it turns out to be more expedient to solve the discrete Thomson problem directly. To accomplish this, again taking the q presepecified repulsive particles as fixed, we initialize n + m - q free repulsive particles randomly on n + m - q nearest-neighbor lattice points around the CMP and then, at each iteration, move the two or three<sup>4</sup> free particles that are furthest from equilibrium in

<sup>&</sup>lt;sup>4</sup>Moving more than two or three particles at a time in this algorithm makes each iteration computationally intensive, and has little impact on overall convergence of the algorithm, whereas moving only one at a time is found to significantly impede convergence to the optimal solution.

# 7.3. SELECTING A POSITIVE BASIS FROM NEAREST-NEIGHBOR LATTICE POINTS

the force-based model described above (that is, those free particles which have the highest force component projected onto the surface of the sphere) into new positions selected from the available locations in such a way as to minimize the maximum force (projected onto the sphere) over the entire set of (fixed and free) particles. Though each iteration of this algorithm involves an exhaustive search for placing the two or three free particles in question, it converges quickly when  $\tau$  is O(100) or less.

Algorithm C. For intermediate kissing numbers  $\tau$ , a hybrid approach may be used: a "good" initial distribution may be found using Algorithm A, then this distribution may be refined using Algorithm B.

In each of these algorithms, to minimize the number of new function evaluations required at each poll step, a check is first made to determine whether any previous function evaluations have already been performed on the nearest-neighbor lattice points around the CMP. If so, then particles are fixed at these locations, while the remaining particles are adjusted via one of the three algorithms described above. By so doing, previously-calculated function values may be used with maximum effectiveness during the polling procedure. When performing the poll step of a surrogate-based search, in order to orient the new poll set favorably (and, on average, exit the poll step quickly), a particle may also be fixed at the nearest neighbor point with the lowest value of the surrogate function; when polling, this poll point is thus evaluated first.

The iterative algorithms described above, though in practice quite effective, are not guaranteed to converge from arbitrary initial conditions to a positive basis for a given value of m, even if such a positive basis exists. To address this issue, if the algorithm used fails to produce a positive basis, the algorithm may be repeated using a new random starting distribution. Our numerical tests indicate that this repeated random initialization scheme usually generates a positive basis within a few initializations when such a positive basis indeed exists.

Since at times, for a given m, there exists no configuration of the free particles on the nearest-neighbor lattice points that produces a positive basis, particularly when the previous function evaluations being leveraged are poorly configured, the number of new random initializations is limited to a prespecified value. Once this value

# CHAPTER 7. EXTENDING LATTICE THEORY FOR DERIVATIVE-FREE OPTIMIZATION

is reached, m is increased by one and the process repeated. As the cost of each function evaluation increases, the user can increase the number of random initializations attempted using one of the above algorithms for each value of m in order to avoid the computation of extraneous poll points that might in fact be unnecessary if sufficient exploration by the discrete Thomson algorithm described above is performed.

Numerical tests have demonstrated the efficacy of this rather simple strategy, which reliably generates a positive basis while keeping computational costs to a minimum even when leveraging a relatively poor configuration of previous function evaluations and when working in relatively high dimension n. Additionally, the algorithm itself is independent of the lattice being used; the only inputs to the algorithm are the dimension of the problem, the locations of the nearest-neighbor lattice points, and the identification of those nearest-neighbor lattice points for which previous function evaluations are available.

When implementing a global search in n dimensions, or even when implementing a local search on a function

### 7.4 Implementation of feasible domain boundaries

which is ill-defined for certain nonphysical values of the parameters (such as negative concentrations of chemicals), it is important to restrict the optimization algorithm to look only over a prespecified "feasible" region of parameter space. For simplicity, the present work assumes rectangular constraints on this feasible domain (that is, simple upper and lower bounds on each parameter value). An efficient n-dimensional lattice with packing radius  $\rho_n$  is used to quantize the interior of the feasible domain, efficient (n-1)-dimensional lattices with packing radius  $\rho_{n-1} = \rho_n/2$  are used to quantize the portions of the boundary of the feasible domain with one active constraint (that is, the "faces"), efficient (n-2)-dimensional lattices with packing radius  $\rho_{n-2} = \rho_n/4$  are used to quantize the portions of the boundary of the feasible domain with two active constraints (that is, the "edges"), etc. The present section describes how to search over the boundaries of the feasible domain, and how to move on and off of these boundaries as appropriate, while carefully restricting all function evaluations to the interior and boundary lattices in order to coordinate an efficient search.

#### 7.4. IMPLEMENTATION OF FEASIBLE DOMAIN BOUNDARIES

We distinguish between two scenarios in which the polling algorithm as described thus far must be adjusted to avoid violating the (n-1)-dimensional boundaries<sup>5</sup> of the feasible domain. In the first scenario, the CMP is relatively far (that is, greater than  $\rho_n$  but less than  $2\rho_n$ ) from the boundary of the feasible domain, and thus one or more of the poll points as determined by one of the algorithms proposed in §7.3 might land slightly outside this boundary. In this scenario, an effective remedy is simply to *eliminate* all lattice points which land outside of the feasible domain from the list of potential poll points, and then to *augment* this restricted list of potential poll points with all lattice points on the nearby (n-1)-dimensional constraint surface which are less than  $2\rho_n$  from the CMP. From this modified list of potential poll points, the poll set may be selected in the usual fashion using one of the algorithms described in §7.3.

In the second scenario, the CMP is relatively close (that is, less than  $\rho_n$ ) to the boundary of the feasible domain. In this scenario, it is most effective simply to shift the CMP onto the nearest lattice point on the (n-1)-dimensional constraint surface. With the CMP on the feasible domain boundary, each poll step explores a minimum positive basis selected on the lattice quantizing the (n-1)-dimensional boundary and, in addition, polls an additional lattice point on the interior of the feasible domain to allow the algorithm to move back off this constraint boundary. Ideally, this additional point would be located on a inward-facing vector normal to the (n-1)-dimensional feasible domain boundary a distance  $\rho_n$  from the CMP; we thus choose the interior lattice point closest to this location.

Multiple active constraints are handled in an analogous manner (see Figure 7.3). In an n-dimensional optimization problem with  $p \ge 2$  active constraints, the CMP is located on an active constraint "surface" of dimension n-p. An efficient (n-p)-dimensional lattice with packing radius  $\rho_{n-p} = \rho_n/2^p$  is used to quantize this active constraint surface, and a poll set is constructed by creating a positive basis selected from the points neighboring the CMP within the (n-p)-dimensional active constraint surface, together with p additional points located on the (n-p+1)-dimensional constraint surfaces neighboring the CMP. Ideally,

<sup>&</sup>lt;sup>5</sup>That is, the portions of the boundary with a single active constraint.

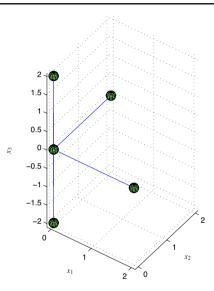


Figure 7.3: A scenario in which a CMP at  $\mathbf{x} = \begin{pmatrix} 0 & 0 & 0 \end{pmatrix}^T$  sits on an (n-2) = 1-dimensional edge of an n = 3-dimensional feasible region with bounds  $x_1 \ge 0$  and  $x_2 \ge 0$ . Note that the feasible neighborhood of this edge is positively spanned by the nearest neighbors on the integer lattice, and that two additional vectors are added to the poll set to facilitate moving off of each of these active constraint boundaries.

### 7.5. QUANTIFYING THE SKEWNESS OF MINIMAL POSITIVE BASES

these p additional points would be located on vectors normal to the (n-p)-dimensional active constraint surface a distance  $\rho_{n-p+1} = \rho_n/2^{p-1}$  from the CMP; we thus choose the lattice points on the (n-p+1)-dimensional feasible domain boundaries closest to these locations.

In practice, it is found that, once an optimization routine moves onto  $p \ge 1$  feasible domain boundaries, it only somewhat infrequently moves back off. To account for this, the p additional poll points mentioned in the previous paragraph are polled *after* the other poll points forming the positive basis within the (n-p)-dimensional active constraint surface.

### 7.5 Quantifying the skewness of minimal positive bases

A final relevant metric of a lattice that relates to the performance of the corresponding lattice-based optimization is the deviation from perfect uniformity of the best minimal positive basis available on nearest-neighbor lattice points. The best nearest-neighbor minimal positive basis skewness of a lattice, s, is thus now defined as the ratio between the largest and the smallest angles between any two vectors in the best minimal positive basis available on nearest-neighbor lattice points, minus one. Therefore, s = 0 indicates a perfectly uniform minimal positive basis on nearest-neighbor lattice points, as exhibited by  $A_2$  (see Figure 7.2a) and  $A_3^*$  (Figure 7.2b). In constrast,  $A_3$  through  $A_8$  all have s = 0.3333 (see, e.g.,  $A_3$  in Figure 7.2c).

Surprisingly, the best nearest-neighbor minimal positive basis skewness of  $E_8$  is s=1; one might initially expect it to be much smaller than this (indeed, one might hope that it would be fairly close to s=0) due to the relatively large kissing number ( $\tau=240$ ) of this n=8 lattice. Interestingly, the best nearest-neighbor positive basis of  $E_8$  when using n+2 points (that is, instead of a minimal positive basis with n+1 points) is perfectly uniform. The tests reported later in Part II thus use n+2 points instead of n+1 points when polling on the  $E_8$  lattice.

A minimal positive basis on nearest-neighbor lattice points doesn't even exist on the  $\mathbb{Z}^n$  lattice (indeed, a positive basis on nearest neighbors of the  $\mathbb{Z}^n$  lattice requires a full 2n points). This was, in fact, a matter

### CHAPTER 7. EXTENDING LATTICE THEORY FOR DERIVATIVE-FREE OPTIMIZATION

of significant inconvenience in previous work when using the Cartesian lattice as the default choice for such problems, as using a maximal positive basis rather than a minimal positive basis essentially doubles the cost of each complete poll step for large n. When developing a minimal positive basis for the  $\mathbb{Z}^n$  lattice, it is thus common (see, e.g., Booker et al. 1999) to select the Cartesian unit vectors  $\mathbf{e}^1$  through  $\mathbf{e}^n$  and one additional "oddball" vector in the  $(-1,-1,\ldots,-1)$  direction which is  $\sqrt{n}$  longer. Note the "clustering" of the Cartesian unit vectors in directions generally opposite to the oddball vector. To quantify, the skewness of this minimal positive basis is  $\cos^{-1}(-1/\sqrt{n})/(\pi/2) - 1$ , which in dimensions n = 2 through 8 is given by 0.5, 0.3918, 0.3333, 0.2952, 0.2677, 0.2468, and 0.2301. Note that, while the skewness of the angular distribution of this minimal positive basis actually decreases gradually as the dimension of the problem increases, the ratio in lengths of the vectors to the nearest-neighbor lattice points and the oddball vector in this basis increases like  $\sqrt{n}$  (that is, from 1.4142 in n=2 to 2.8284 in n=8). This is quite unfortunate, as it leads to a peculiar nonisotropic behavior of the optimization algorithm over parameter space (for further discussion on this point, see the sixth paragraph of §10.1). The tests reported later in Part II use this peculiar minimum positive basis, with a long oddball vector, when polling on the  $\mathbb{Z}^n$  lattice.

We now have all of the ingredients necessary to coordinate SMF algorithms, as introduced in §7.1, with any of the lattices listed in Tables 3.1-3.2 of Part I, while both reusing previous function evaluations as effeciently as possible as well as respecting sharp bounds on the feasible region of parameter space.

7.5. QUANTIFYING THE SKEWNESS OF MINIMAL POSITIVE BASES

n	lattice	name	Δ	Θ	G	τ
	$A_2$	hexagonal	0.90690	1.2092	0.080188	6
2	$\mathbb{Z}^2$	square	0.78540	1.5708	0.083333	4
	$A_3$	face-centered cubic (FCC)	0.74048	2.0944	0.078745	12
3	$A_3^*$	body-centered cubic (BCC)	0.68017	1.4635	0.078543	8
	$\mathbb{Z}^3$	cubic	0.52360	2.7207	0.083333	6
	$E_8$	Gosset	0.25367	4.0587	0.071682	240
	$D_8$		0.12683	32.470	0.075914	112
	$A_8$	zero-sum	0.08456	32.993	0.077391	72
8	$D_8^*$		0.03171	8.1174	0.074735	16
	$A_8^*$		0.02969	3.6658	0.075972	18
	$\mathbb{Z}^8$	Cartesian	0.01585	64.939	0.083333	16

Table 7.1. Characteristics of select distinct lattices in dimensions 2, 3, and 8, ordered from dense to rare (for a more complete characterization, see Tables 3.1 and 3.2 of Part I). Listed (see Part I) are the packing density,  $\Delta$ , covering thickness,  $\Theta$ , mean squared quantization error per dimension, G, and kissing number,  $\tau$ . Note that  $\mathbb{Z}^n$  is significantly outperformed in every standard metric in every dimension n > 1 by the available alternatives.

CHAPTER 7. EXTENDING LATTICE THEORY FOR DERIVATIVE-FREE OPTIMIZATION

	$A_2$	$A_3$	$D_4$	$D_5$	$E_6$	<i>E</i> <sub>7</sub>	$E_8$	K <sub>12</sub>	$\Lambda_{16}$	$\Lambda_{24}$
$f_{\Delta}$	1.155	1.414	2	2.83	4.62	8	16	152	4096	$1.68e^{7}$
$f_{\tau}$	1.5	2	3	4	6	9	15	31.5	135	4095

Table 7.2. The densest, most uniform lattices available in several dimensions, and two factors quantifying the degree to which these lattices are better than the corresponding Cartesian grid in the same dimension;  $f_{\Delta}$  denotes the factor of improvement in the packing density, an indication of the uniformity of the lattice, and  $f_{\tau}$  denotes the factor of improvement in the kissing number, an indication of the flexibility available in selecting a positive basis from the nearest neighbors on the lattice. Note that the improvements becoming especially pronounced as the dimension n is increased.

# **Chapter 8**

# **Kriging interpolation**

Co	nte	nts

8.2	Notation of statistical description	131
8.3	Statistical modeling assumptions of the ordinary Kriging model	132
0.4		

### 8.1 Interpolation - basic concepts

The purpose of the search step of an SMF algorithm (see §7.1) is to interpolate, and extrapolate, the trends exhibited by the existing function evaluations in order to suggest new regions of parameter space, perhaps far from the CMP, where the function value is anticipated, with some reasonable degree of probability, to be lower than that of the CMP. There are a variety of possibile ways of accomplishing this; we leverage here the Kriging interpolation strategy (Krige 1951; Matheron 1963; Jones 2001; Rasmussen & Williams 2006).

The problem of interpolation is the problem of drawing a smooth curve through data points in order to

estimate the function values in regions where the function itself has not yet been computed. The problem of interpolation, thus, necessarily builds on some hypothesis that models the function behavior in order to "connect the dots". The most common such model is a mechanical one, based on a thin piece of wood, or "spline", that is "bent" in order to touch all the data points; this mechanical model leads directly to the mathematical algorithm known as cubic spline interpolation. A perhaps equally valid hypothesis, which forms the foundation for the Kriging interpolation strategy, is to *model the underlying function as a realization, with maximum likelihood, of some stochastic process*. The stochastic model used in this approach is selected to be general enough to model a broad range of functions reasonably well, yet simple enough to be fairly inexpensive to tune appropriately based on the measured data. There are many such stochastic models which one can select; the simple stochastic model considered here leads to the easy-to-use interpolation strategy commonly referred to as ordinary Kriging.

### 8.2 Notation of statistical description

To begin, consider N points  $\{\mathbf{x}^1, \dots, \mathbf{x}^N\}$ , at which the function will ultimately be evaluated, and model the function's value at these N points with the random vector

$$\mathbf{f} = \begin{pmatrix} f(\mathbf{x}^1) \\ \vdots \\ f(\mathbf{x}^N) \end{pmatrix} = \begin{pmatrix} f_1 \\ \vdots \\ f_N \end{pmatrix}.$$

To proceed further, we need a clear statistical framework to describe this random vector.

The cumulative distribution function (CDF) of the random vector  $\mathbf{f}$  denoted  $d_{\mathbf{r}}(\mathbf{f})$  is a mapping from

The cumulative distribution function (CDF) of the random vector  $\mathbf{f}$ , denoted  $d_{\mathbf{f}}(\underline{\mathbf{f}})$ , is a mapping from  $\underline{\mathbf{f}} \in \mathbb{R}^n$  to the real interval [0,1] that monotonically increases in each of the components of  $\underline{\mathbf{f}}$ , and is defined

$$d_{\mathbf{f}}(\underline{\mathbf{f}}) = P(f_1 \leq f_1, f_2 \leq f_2, \dots, f_n \leq f_n),$$

where  $\underline{\mathbf{f}}$  is some particular value of the random vector  $\mathbf{f}$  and P(S) denotes a probability measure that the conditions stated in S are true. In the scalar case, for example,  $d_f(1) = 0.6$  means that it is 60% likely that the random variable f satisfies the condition  $f \leq 1$ . For a random vector  $\mathbf{f}$  whose CDF is modelled as being differentiable everywhere, the probability density function (PDF)  $p_{\mathbf{f}}(\mathbf{f}') \geq 0$  is a scalar function of  $\mathbf{f}'$  defined such that

$$d_{\mathbf{f}}(\underline{\mathbf{f}}) = \int_{-\infty}^{\underline{f}_1} \int_{-\infty}^{\underline{f}_2} \cdots \int_{-\infty}^{\underline{f}_n} p_{\mathbf{f}}(\mathbf{f}') df_1' df_2' \cdots df_n' \quad \Leftrightarrow \quad p_{\mathbf{f}}(\mathbf{f}') = \frac{\partial^n d_{\mathbf{f}}(\underline{\mathbf{f}})}{\partial f_1 \partial f_2 \cdots \partial f_n} \Big|_{\underline{\mathbf{f}} = \underline{\mathbf{f}}'}.$$

For small  $|\Delta \mathbf{f}'|$ , the quantity  $p_{\mathbf{f}}(\mathbf{f}')\Delta f_1'\Delta f_2'\cdots\Delta f_n'$  represents the probability that the random vector  $\mathbf{f}$  takes some value within a small rectangular region centered at the value  $\mathbf{f}'$  and of width  $\Delta f_i'$  in each coordinate

direction  $\mathbf{e}_i$ . Note that the integral of  $p_{\mathbf{f}}(\mathbf{f}')$  over all possible values of  $\mathbf{f}'$  is unity, that is

$$\int_{\mathbb{R}^n} p_{\mathbf{f}}(\mathbf{f}') d\mathbf{f}' = 1.$$

The expected value of a function  $\mathbf{g}(\mathbf{f})$  of a random vector  $\mathbf{f}$  is given by

$$\mathcal{E}\left\{\mathbf{g}(\mathbf{f})\right\} = \int_{\mathbb{D}^n} \mathbf{g}(\mathbf{f}') \, p_{\mathbf{f}}(\mathbf{f}') \, d\mathbf{x}'.$$

The expected value may be interpreted as the average of the quantity in question over many realizations. In particular, the mean  $\bar{\mathbf{f}}$  and covariance  $P_{\mathbf{f}}$  of the random vector  $\mathbf{f}$  are defined as

$$\bar{\mathbf{f}} \triangleq \mathcal{E}\left\{\mathbf{f}\right\} = \int_{\mathbb{R}^n} \mathbf{f}' \, p_{\mathbf{f}}(\mathbf{f}') \, d\mathbf{f}', \qquad P_{\mathbf{f}} \triangleq \mathcal{E}\left\{\left(\mathbf{f} - \bar{\mathbf{f}}\right) \left(\mathbf{f} - \bar{\mathbf{f}}\right)^T\right\} = \int_{\mathbb{R}^n} \left(\mathbf{f}' - \bar{\mathbf{f}}\right) \left(\mathbf{f}' - \bar{\mathbf{f}}\right)^T \, p_{\mathbf{f}}(\mathbf{f}') \, d\mathbf{f}'.$$

# 8.3 Statistical modeling assumptions of the ordinary Kriging model

The PDF of the random vector  $\mathbf{f} = \mathbf{f}_{n \times 1}$  in this analysis is modelled as Gaussian, and is thus restricted to the generic form

$$p_{\mathbf{f}}(\mathbf{f}') = \frac{1}{(2\pi)^{n/2}|P_{\mathbf{f}}|^{1/2}} \exp \frac{-(\mathbf{f}' - \overline{\mathbf{f}})^T P_{\mathbf{f}}^{-1}(\mathbf{f}' - \overline{\mathbf{f}})}{2},$$
(8.1a)

where the covariance  $P_{\mathbf{f}}$  is modelled as a constant  $\sigma^2$ , referred to as the variance, times a correlation matrix R whose  $\{i, j\}$ 'th component  $r_{ij}$  is given by a model of the correlation of the random function f between points  $\mathbf{x}^i$  and  $\mathbf{x}^j$ , where this correlation model  $r(\cdot, \cdot)$  itself decays exponentially with the distance between points  $\mathbf{x}^i$ 

and  $\mathbf{x}^{j}$ ; that is,

for some yet-to-be-determined constant 
$$\mu$$
. There is extensive debate in the recent literature (see, e.g., Isaaks & Srivastava 1989; Rasmussen & Williams 2006) on the statistical modeling assumptions one should use in a Kriging model of this sort. It is straightforward to extend the present investigation to incorporate less

Gaussian model (8.1a) is itself modelled as uniform over all of its components:  $\mathbf{\bar{f}} \triangleq \mu \mathbf{1} \tag{8.1c}$ 

for some yet-to-be-determined constants  $\sigma^2$ ,  $\theta_{\ell} > 0$ , and  $0 < p_{\ell} \le 2$  for  $\ell = 1, ..., n$ . The mean  $\bar{\mathbf{f}}$  in the

 $P_{\mathbf{f}} \triangleq \sigma^2 R$ , where  $r_{ij} \triangleq r(\mathbf{x}^i, \mathbf{x}^j)$  and  $r(\mathbf{x}, \mathbf{y}) \triangleq \prod_{\ell=1}^n \exp\left(-\theta_\ell |x_\ell - y_\ell|^{p_\ell}\right)$ 

(8.1b)

(8.2)

Adjusting the coefficients of the model based on the data

If the vector of observed function values is

$$\mathbf{f}^o = egin{pmatrix} f_1^o \ dots \ g_0 \end{pmatrix},$$

then the PDF corresponding to this observation in the statistical model proposed in (8.1) can be written as

$$p_{\mathbf{f}}(\mathbf{f}^{o}) = \frac{1}{(2\pi)^{n/2} (\sigma^{2})^{n/2} |R|^{1/2}} \exp \frac{-(\mathbf{f}^{o} - \mu \mathbf{1})^{T} R^{-1} (\mathbf{f}^{o} - \mu \mathbf{1})}{2\sigma^{2}}.$$
 (8)

The process of Kriging modeling boils down to selecting the parameters  $\sigma^2$ ,  $\theta_\ell$ ,  $p_\ell$ , and  $\mu$  in the statistical model proposed in (8.1) to maximize the PDF evaluated for the function values actually observed,  $\mathbf{f} = \mathbf{f}^o$ , as given in (8.2).

Maximizing  $p_{\mathbf{f}}(\mathbf{f}^o)$  is equivalent to minimizing the negative of its log. Thus, for simplicity, consider

$$J = -\log[p_{\mathbf{f}}(\mathbf{f}^o)] = \frac{n}{2}\log(2\pi) + \frac{n}{2}\log(\sigma^2) + \frac{1}{2}\log(|R|) + \frac{(\mathbf{f}^o - \mu \mathbf{1})^T R^{-1}(\mathbf{f}^o - \mu \mathbf{1})}{2\sigma^2}.$$
 (8.3)

Setting the derivatives of J with respect to  $\mu$  and  $\sigma^2$  equal to zero and solving, the optimal values of  $\mu$  and  $\sigma^2$  are determined immediately:

$$\mu = \frac{\mathbf{1}^T R^{-1} \mathbf{f}^o}{\mathbf{1}^T R^{-1} \mathbf{1}}, \qquad \sigma^2 = \frac{(\mathbf{f}^o - \mu \mathbf{1})^T R^{-1} (\mathbf{f}^o - \mu \mathbf{1})}{n}.$$
 (8.4)

With these optimal values of  $\mu$  and  $\sigma^2$  applied, noting that the last term in (8.3) is now constant, what remains to be done is to minimize

$$J_1 = \frac{n}{2}\log(\sigma^2) + \frac{1}{2}\log(|R|) \tag{8.5}$$

with respect to the remaining free parameters  $\theta_{\ell}$  and  $\theta_{\ell}$ , where  $\sigma^2$  is given as a function of R in (8.4) and R, in turn, is given as a function of the free parameters  $\theta_{\ell}$  and  $\theta_{\ell}$  in (8.1b). This minimization must, in general, be performed numerically. However, the function  $\theta_{\ell}$  is smooth in the parameters  $\theta_{\ell}$  and  $\theta_{\ell}$ , so this optimization

<sup>&</sup>lt;sup>1</sup>To simplify this optimization,  $p_{\ell}$  may be specified by the user instead of being determined via optimization; this is especially appropriate to do when the number of function evaluations N is relatively small, and thus there is not yet enough data to determine both the  $\theta_{\ell}$  and  $p_{\ell}$  uniquely. If this approach is followed,  $p_{\ell} = 1$  or 2 are natural choices; the case with  $p_{\ell} = 1$  is referred to as an Ornstein-Uhlenbeck process, whereas the case with  $p_{\ell} = 2$  is infinitely differentiable everywhere.

#### 8.5. USING THE TUNED STATISTICAL MODEL TO PREDICT NEW FUNCTION VALUES

may be performed efficiently with a standard gradient-based algorithm, such as the nonquadratic conjugate gradient algorithm, where the gradient itself, for simplicity, may easily be determined via a simple finite difference or complex-step derivative approach.

Note that, after each new function evaluation, the Kriging parameters adjust only slightly, and thus the previously-converged values of these parameters form an excellent initial guess for this gradient-based optimization algorithm. Note also that, while performing this optimization, the determinant of the correlation matrix occasionally reaches machine zero. To avoid the numerical difficulty that taking the log of zero would otherwise induce, a small  $[O(10^{-6})]$  term may be added to the diagonal elements of R. By so doing, the Kriging predictor does not quite have the value of the sampled data at each sampled point; however, it remains quite close, and the algorithm is made numerically robust [Booker *et al*, 1999].

#### 8.5 Using the tuned statistical model to predict new function values

Once the parameters of the stochastic model have been tuned as described above, the tuned Kriging model facilitates the computationally inexpensive prediction of the function value at any new location  $\bar{\mathbf{x}}$ . To perform this prediction, consider now the N+1 points  $\{\mathbf{x}^1,\ldots,\mathbf{x}^N,\bar{\mathbf{x}}\}$ , and model the function's value at these N+1 points with the vector

$$\bar{\mathbf{f}} = \begin{pmatrix} \mathbf{f} \\ f(\bar{\mathbf{x}}) \end{pmatrix} = \begin{pmatrix} \mathbf{f} \\ \bar{f} \end{pmatrix},$$

where  $\mathbf{f}$  is the  $N \times 1$  random vector considered previously and  $\bar{f}$  is the random scalar modeling the function at the new point. Analogous statistical assumptions as laid out in (8.1) are again applied, with the correlation matrix now written as

matrix now written as 
$$\bar{R} = \begin{bmatrix} R & \bar{\mathbf{r}} \\ \bar{\mathbf{r}}^T & 1 \end{bmatrix}, \qquad P_{\bar{\mathbf{f}}} \triangleq \sigma^2 \bar{R}, \tag{8.6}$$

(8.7)

where R is the  $N \times N$  correlation matrix considered previously and, consistent with this definition, the vector  $\bar{\mathbf{r}}$  is constructed with components

$$ar{r}_i = r(\mathbf{x}^i, \bar{\mathbf{x}}), \quad \text{where} \quad r(\mathbf{x}, \mathbf{y}) \triangleq \prod_{\ell=1}^n \exp\left(-\theta_\ell |x_\ell - y_\ell|^{p_\ell}\right).$$

Following Jones (2001), note by the matrix inversion lemma that  $\bar{R}^{-1}$  may be written

 $= \ldots = C_2 \cdot \exp \frac{-[\overline{f'} - \widehat{f}]^T [\overline{f'} - \widehat{f}]}{2s^2},$ 

$$\bar{R}^{-1} = \begin{bmatrix} R & \bar{\mathbf{r}} \\ \bar{\mathbf{r}}^T & 1 \end{bmatrix}^{-1} = \begin{bmatrix} R^{-1} + R^{-1}\bar{\mathbf{r}}(1 - \bar{\mathbf{r}}^T R^{-1}\bar{\mathbf{r}})^{-1}\bar{\mathbf{r}}^T R^{-1} & -R^{-1}\bar{\mathbf{r}}(1 - \bar{\mathbf{r}}^T R^{-1}\bar{\mathbf{r}})^{-1} \\ -(1 - \bar{\mathbf{r}}^T R^{-1}\bar{\mathbf{r}})^{-1}\bar{\mathbf{r}}^T R^{-1} & (1 - \bar{\mathbf{r}}^T R^{-1}\bar{\mathbf{r}})^{-1} \end{bmatrix}.$$
 (8.7)

Keeping the paramter values  $\sigma^2$ ,  $\theta_\ell$ ,  $p_\ell$ , and  $\mu$  as tuned previously, we now examine the variation of the PDF in the remaining unknown random variable,  $\bar{f}$ . Substituting (8.6) and (8.7) into a PDF of the form (8.1a), we may write

$$p_{\bar{\mathbf{f}}}(\bar{\mathbf{f}}') = C_1 \cdot \exp \frac{-(\bar{\mathbf{f}}' - \mu \mathbf{1})^T \bar{R}^{-1} (\bar{\mathbf{f}}' - \mu \mathbf{1})}{2\sigma^2} = C_1 \cdot \exp \frac{-\left[\frac{\mathbf{f}' - \mu \mathbf{1}}{\bar{f}' - \mu}\right]^T \bar{R}^{-1} \left[\frac{\mathbf{f}' - \mu \mathbf{1}}{\bar{f}' - \mu}\right]}{2\sigma^2}$$
(8.8)

#### 8.5. USING THE TUNED STATISTICAL MODEL TO PREDICT NEW FUNCTION VALUES

where, with a minor amount of algebraic rearrangement, the mean and variance of this scalar Gaussian distribution modeling the random scalar  $\bar{f}$  work out to be<sup>2</sup>

$$\hat{f}(\bar{\mathbf{x}}) = \mathcal{E}\{f(\bar{\mathbf{x}})\} = \mathcal{E}\{\bar{f}\} = \mu + \mathbf{r}^T R^{-1} (\mathbf{f}^o - \mu \mathbf{1}),$$

$$s^2(\bar{\mathbf{x}}) = \mathcal{E}\{[f(\bar{\mathbf{x}}) - \hat{f}]^2\} = \mathcal{E}\{[\bar{f} - \hat{f}]^2\} = \sigma^2 (1 - \mathbf{r}^T R^{-1} \mathbf{r}).$$
(8.9b)

Equations (8.9a)-(8.9b) give the final formulae for the Kriging predictor,  $\hat{f}(\bar{x})$ , and its associated uncertainty,  $s^2(\bar{\mathbf{x}})$ . When applied numerically to a representative test problem, as expected, the Kriging predictor function, which we denote  $\hat{f}(\bar{\mathbf{x}})$ , interpolates [that is, it goes through every observed function value at points  $\bar{\mathbf{x}} = \mathbf{x}^1$ to  $\bar{\mathbf{x}} = \mathbf{x}^N$ ], whereas the uncertainty function, denoted  $s^2(\bar{\mathbf{x}})$ , is zero at each sampled point, and resembles a Gaussian "bump" between these sampled points, as seen in Figure 8.1. Note that, once the parameters of the statistical model have been determined, as described in §8.4, the formula (8.9a)-(8.9b) for the Kriging pre-

dictor  $\hat{f}(\bar{\mathbf{x}})$  and its corresponding uncertainty  $s^2(\bar{\mathbf{x}})$  at any test point  $\bar{\mathbf{x}}$  is computationally quite inexpensive<sup>3</sup>.

$$s^{2}(\bar{\mathbf{x}}) = \sigma^{2} \left( 1 - \mathbf{r}^{T} R^{-1} \mathbf{r} + \frac{\left( 1 - \mathbf{r}^{T} R^{-1} \mathbf{r} \right)^{2}}{\mathbf{1}^{T} D^{-1} \mathbf{1}} \right). \tag{8.9b'}$$

Which formula [(8.9b) or (8.9b')] is used in the present model is ultimately a matter of little consequence as far as the overall derivativefree optimization algorithm is concerned; we thus prefer the form given in (8.9b) due to its computational simplicity.

<sup>&</sup>lt;sup>2</sup>An alternative interpretation of this process models the constant  $\mu$  itself as a stochastic variable rather than as a constant. Following this line of reasoning ultimately gives the same formula for the predictor  $\hat{f}(\bar{x})$  as given in (8.9a), and a slightly modified formula for its associated uncertainty,

<sup>&</sup>lt;sup>3</sup>Note that, for maximum efficiency,  $R^{-1}$  should be saved between function evaluations and reused for every new computation of  $\hat{f}$ and  $s^2$  required.

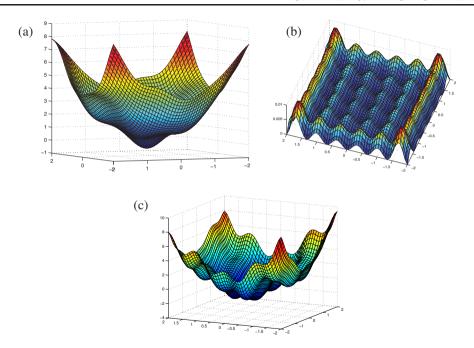


Figure 8.1: (a) The Kriging predictor,  $\hat{f}(\mathbf{x})$ , and (b) its associated uncertainty,  $s^2(\mathbf{x})$ , for a perturbed quadratic bowl sampled on a square grid of  $7 \times 7$  points. (c) The corresponding  $J(\mathbf{x}) = \hat{f}(\mathbf{x}) - c \cdot s^2(\mathbf{x})$  search function used for a global search in two dimensions (see §9).

## Global optimization leveraging Kriging-based interpolation

The previous chapter reviewed the Kriging interpolation strategy which, based on a sparse set of observed function values  $f^o(\mathbf{x}^i)$  for  $i=1,\ldots,N$ , develops a function predictor  $\hat{f}(\mathbf{x})$  and a model of the uncertainty  $s^2(\mathbf{x})$  associated with this prediction for any given set of parameter values  $\mathbf{x}$ . Leveraging this Kriging model, an efficient search algorithm can now be developed for the derivative-free optimization algorithm summarized in §7.1.

The effectiveness of the various Kriging-based search strategies which one might propose may be tested by applying them repeatedly to simple test problems via the following procedure:

#### CHAPTER 9. GLOBAL OPTIMIZATION LEVERAGING KRIGING-BASED INTERPOLATION

- a search function  $J(\mathbf{x})$  is first developed based on a Kriging model fit to the existing function evaluations,
- a gradient-based search is used to minimize this (computationally inexpensive, smoothly-varying) search function,
- the function  $f(\mathbf{x})$  is sampled at the point  $\tilde{\mathbf{x}}$  which minimizes the search function  $\mathbf{x}$ ,
- the Kriging model is updated, and the search is repeated.

search in order to keep function evaluations far apart.

In the present work, we consider a scalar test problem with multiple minima,  $f(x) = \sin(x) + x^2$ , on the interval  $x \in [-10, 10]$ , and use four starting points to initialize the search x = -10, x = -5.2, x = 6, and x = 10. Ineffective search strategies will not converge to the global minimum of f(x) in this test, and may not even converge to a local minimum. More effective search strategies converge to the global minimum following this approach, and the number of function evaluations required for convergence indicates the effectiveness of the search strategy used.

Perhaps the most "obvious" strategy to use in such problems is simply fitting a Kriging model to the known data, then searching the Kriging predictor itself,  $J(\mathbf{x}) = \hat{f}(\mathbf{x})$ , for its minimum value. This simple approach has been implemented in a variety of examples with reasonably good results (see Booker *et al*, 1999). However, as shown clearly in Figure 9.1, this approach can easily break down. The Kriging predictor does not necessarily model the function accurately, and its minimization fails to guarantee convergence to even a local minimum of the function  $f(\mathbf{x})$ . This observed fact can be motivated informally by identifying the Kriging predictor as an *interpolating* function which only under extraördinary conditions predicts a function value significantly lower than all of the previously-computed function values; under ordinary conditions, a strategy of minimizing the predictor will thus often stall in the vicinity of the previously-evaluated points.

To avoid the shortcomings of a search defined solely by the minimization of the predictor, another strategy

<sup>&</sup>lt;sup>1</sup>For the moment, to focus our attention on the behavior of the search algorithm itself, no underlying grid is used to coordinate the

explored by Booker *et al* (1999) is to evaluate the function at *two* points in parameter space during the search: one point chosen to minimize the predictor, and the other point chosen to maximize the predictor uncertainty. Such a heuristic provides a guarantee of global convergence, as the seach becomes dense in the parameter space as the total number of function evaluations, *N*, approaches infinity (see §7.1.2). However, this approach generally does not converge quickly as compared with the improved methods described below, as the extra search point has no component associated with the predictor, and is thus often evaluated in relatively "poor" regions of parameter space.

We are thus motivated to develop a more flexible strategy to explore slightly away from the minima of

the predictor. To achieve this, consider the minimization of  $J(\mathbf{x}) = \hat{f}(\mathbf{x}) - c \cdot s^2(\mathbf{x})$ , where c is some constant (see Cox & John 1997 and Jones 2001). A search coordinated by this function will tend to explore regions of parameter space where both the predictor of the function value is relatively low and the uncertainty of this prediction in the Kriging model is relatively high. With this strategy, the search is driven to regions of higher uncertainty, with the  $-c \cdot s^2(\mathbf{x})$  term in  $J(\mathbf{x})$  tending to cause the algorithm to explore away from previously evaluated points. Additionally, minimizing  $\hat{f}(\mathbf{x}) - c \cdot s^2(\mathbf{x})$  allows the algorithm to explore the vicinity of multiple local minima in successive iterations in order to determine, with an increasing degree of certainty, which local "bowl" in fact has the deepest minimum. The parameter c provides a natural means to "tune" the degree to which the search is driven to regions of higher uncertainty, with smaller values of c focusing the search more on refining the vicinity of the lowest function value(s) already found, and larger values of c focusing the search more on exploring regions of parameter space which are still relatively poorly sampled. This parameter may tuned based on knowledge of the function being minimized: if the function is suspected to have multiple minima, c can be made relatively large to ensure a more exploratory search, whereas if the function is suspected of having a single minimum, c can be made relatively small to ensure a more focused search in the vicinity of the CMP. For an appropriate intermediate value of c, the resulting algorithm is often quite effective at both global exploration and local refinement of the minimum, as illustrated in Figure 9.2.

#### CHAPTER 9. GLOBAL OPTIMIZATION LEVERAGING KRIGING-BASED INTERPOLATION

for a two-dimensional problem in Figure 8.1c. Note also that, in the spirit of Booker *et al* (1999) [who effectively suggested, in the present notation, exploring based on both c=0 and  $c\to\infty$  at each search step], one can perform a search using multiple but finite values of c at each search step, returning a set of points designed to focus, to varying degrees, on the competing objectives of global exploration and local refinement. If at each search step k at least one point is included which minimizes  $\hat{f}(\mathbf{x}) - c_k \cdot s^2(\mathbf{x})$  for a value of  $c_k$  which itself approaches  $\infty$  as  $k\to\infty$ , then the search drives at least some new function evaluations sufficiently far from the existing points that the function evaluations eventually become dense over the feasible domain, thus guaranteeing global convergence. Thus, an  $\hat{f}(\mathbf{x}) - c \cdot s^2(\mathbf{x})$  search, when used properly, can indeed be used in a clobally convergence transport.

The strategy of searching  $J(\mathbf{x}) = \hat{f}(\mathbf{x}) - c \cdot s^2(\mathbf{x})$  also extends naturally to multiple dimensions, as illustrated

itself approaches  $\infty$  as  $k \to \infty$ , then the search drives at least some new function evaluations sufficiently far from the existing points that the function evaluations eventually become dense over the feasible domain, thus guaranteeing global convergence. Thus, an  $\hat{f}(\mathbf{x}) - c \cdot s^2(\mathbf{x})$  search, when used properly, can indeed be used in a globally convergent manner.

Minimizing  $J(\mathbf{x}) = \hat{f}(\mathbf{x}) - c \cdot s^2(\mathbf{x})$  is not the only strategy to take advantage of the estimate of the uncertainty of the predictor provided by the Kriging model. Another effective search strategy involves maximizing the probability of achieving a target level of improvement below the current CMP; this is called the *maximum likelihood of improvement (MLI)* approach [see Kushner 1964, Stuckman 1988, Perttunen 1991, Elder 1992,

the probability of achieving a target level of improvement below the current CMP; this is called the *maximum likelihood of improvement (MLI)* approach [see Kushner 1964, Stuckman 1988, Perttunen 1991, Elder 1992, and Mockus 1994]. If the current CMP has a function value  $f_{\min}$ , then this search strategy seeks that  $\mathbf{x}$  for which the probability of finding a function value  $f(\mathbf{x})$  less than some prespecified target value  $f_{\text{target}}$  [that is, for which  $f(\mathbf{x}) \leq f_{\text{target}} < f_{\min}$ ] is maximized in the Kriging model. If  $f(\mathbf{x})$  is known to be a positive function, a typical target value in this approach is  $f_{\text{target}} = (1 - \delta) f_{\min}$ , where  $\delta$  may be selected somewhere in the range of 0.01 to 0.2. As for the parameter c discussed in the previous paragraph, the parameter  $\delta$  in this strategy tunes the degree to which the search is driven to regions of higher uncertainty, with smaller values of  $\delta$  focusing the search more on refining the vicinity of the lowest function value(s) already found, and larger values of  $\delta$  focusing the search more on exploring regions of parameter space which are still relatively poorly sampled. As seen in Figure 9.3, the MLI search offers performance similar to the  $\hat{f}(\mathbf{x}) - c \cdot s^2(\mathbf{x})$  method dis-

cussed previously. In contrast with the  $\hat{f}(\mathbf{x}) - c \cdot s^2(\mathbf{x})$  approach, even for a fixed (finite) value of  $\delta$ , the MLI

approach eventually drives the function evaluations far enough away from existing points that the function evaluations eventually become dense over the feasible domain, thus guaranteeing global convergence. Thus, the MLI approach is inherently globally convergent.

Even more sophisticated search strategies can also be proposed, as reviewed elegantly by Jones (2001). However, the simplicity, flexibility, and performance given by the strategy of minimizing  $J(\mathbf{x}) = \hat{f}(\mathbf{x}) - c \cdot \mathbf{x}$ 

 $s^2(\mathbf{x})$  renders this approach as adequate for our testing purposes here. Since both the  $J(\mathbf{x}) = \hat{f}(\mathbf{x}) - c \cdot s^2(\mathbf{x})$  search function and the MLI search function are inexpensive to compute, continuous, and smooth, but in general have multiple minima, an efficient gradient-based search,

initialized from several well-selected points in parameter space, may be used to to minimize them. As the uncertainty  $s^2(\mathbf{x})$  goes to zero at each sample point,  $J(\mathbf{x})$  will tend to dip between each sample point. Thus, a search is initialized on  $2n \cdot N$  total points forming a positive basis near (say, at a distance of  $\rho_n/2$ ) to each of the N sample points, and each of these starting points is marched to a local minima of the search function using an efficient gradient-based search (which is constrained to remain within the feasible domain of x). The lowest point of the paths so generated will very likely be the global minima of the search function. For simplicity, the necessary gradients for this search may be computed via a simple second-order central finite difference scheme applied to the Kriging model, though more sophisticated and efficient approaches are also

possible.

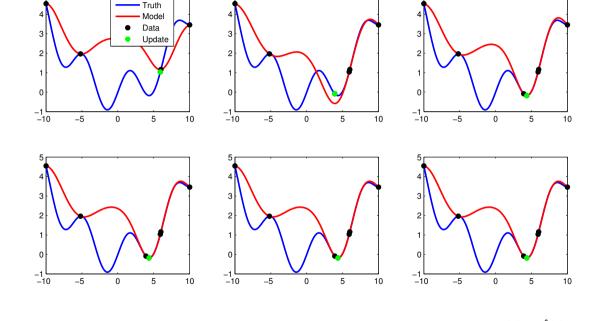


Figure 9.1: Convergence of a search algorithm based on minimizing the Kriging predictor,  $J(\mathbf{x}) = \hat{f}(\mathbf{x})$ , at each iteration. This algorithm does not necessarily converge to even a local minimum, and in this example has stalled, far from the global minimum, after six iterations.

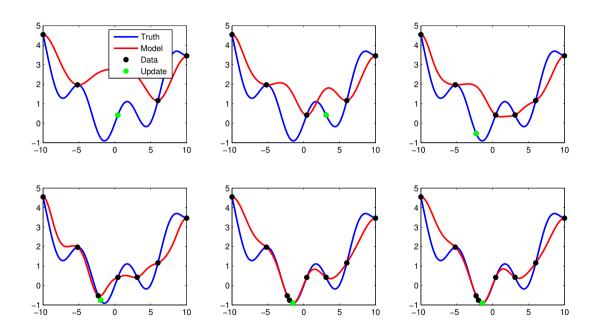


Figure 9.2: Convergence of a search algorithm based on minimizing the search function  $J(\mathbf{x}) = \hat{f}(\mathbf{x}) - c \cdot s^2(\mathbf{x})$  at each iteration, taking c = 1. Note that the global minimum is found after just a few iterations. However, global convergence is not guaranteed.

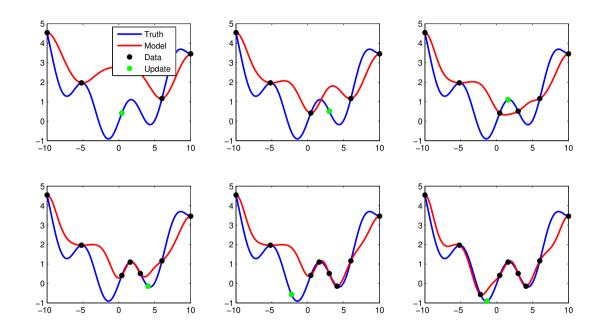


Figure 9.3: MLI search with a target T = 10%. Note convergence to global minimum, as well as exploratory nature of the search which guarantees global convergence.

## Lattice-based derivative-free optimization via global surrogates

#### **Contents**

10.1	SP applied to randomly-shifted quadratic bowls	149
10.2	SP applied to randomly-shifted Rosenbrock functions	152

§7) while leveraging a Kriging interpolant (see §8) to perform an efficient global search based on the search function  $J(\mathbf{x}) = \hat{f}(\mathbf{x}) - c \cdot s^2(\mathbf{x})$  (see §9). The full algorithm has been implemented in an efficient numerical code, dubbed Checkers, and is tested in this section in n = 2 to n = 8 dimensions using the  $\mathbb{Z}^n$ ,  $A_n$ , and  $E_8$  lattices to coordinate the search, and is applied here to:

Putting everything together, we now develop and test what we identify as the *Lattice Based Derivative-free Optimization via Global Surrogates (LABDOGS)* algorithm. This algorithm consists of an SMF-based optimization (see §7.1) coordinated by uniform *n*-dimensional lattices (see Part I and further extensions in

 $f_{\mathcal{Q}}(\mathbf{x}) = (\mathbf{x} - \mathbf{x}^o)^T A(\mathbf{x} - \mathbf{x}^o)$ 

• randomly shifted quadratic bowls

• randomly shifted Rosenbrock functions:  

$$f_R(\mathbf{x}) = \sum_{i=0}^{n-1} \{ [1 - (x_i - x_i^o)]^2 + (-1)^n 500 [(x_{i+1} - x_{i+1}^o) - (x_i - x_i^o)^2]^2 \},$$

• the Branin function:

$$f_B(\mathbf{x}) = [1 - 2x_2 + 0.05\sin(4\pi x_2) - x_1]^2 + [x_2 - 0.5\sin(2\pi x_1)]^2$$

• and the "T<sub>1</sub>" function:

$$f_M(\mathbf{x}) = \sin(5x_1) + \sin(5x_2) + 0.02[(5x_1 + 1.5)^2 + (5x_2 + 1.5)^2].$$

Note that the first two test functions are *n*-dimensional and have unique minima, whereas the last two test functions are 2-dimensional and have multiple minima. Much further testing remains to be done, and will be reported in future work.

#### 10.1 SP applied to randomly-shifted quadratic bowls

To test the hypothesis that the efficiency of a pattern search is significantly affected by the packing efficiency and/or the nearest-neighbor distribution of the lattices which coordinate it, a large number of SP optimizations were first performed on randomly-shifted quadratic bowls to gather and compare statistical data on the performance of  $\mathbb{Z}^n$ -based,  $A_n$ -based, and  $E_8$ -based SP optimizations. The positive-definite matrices A>0 and offsets  $\mathbf{x}^o$  defining the quadratic bowls to be minimized, as well as the starting points used in the searches, were selected at random for every set of tests, and the initial  $\mathbb{Z}^n$ ,  $A_n$ , and  $E_8$  lattices were scaled such that the initial number of points per unit volume of parameter space was identical.

The  $\mathbb{Z}^n$ -based,  $A_n$ -based, and  $E_8$ -based SP algorithms were run from the same starting points on the same quadratic test functions to the same level of convergence. Note that several of the significant built-in acceleration features of the full Checkers code were in fact turned off for this baseline comparison. Most notably, complete polls were performed (that is, the poll steps were not terminated immediately upon finding a lower CMP), and no attempt was made to reuse previously-computed points when forming each successive poll set, or to orient optimally any given poll set. In fact, the angular distribution of the poll set around the CMP was fixed from one step to the next in these initial tests.

Two quantitative measures of the relative efficiency of the optimization algorithms to be tested are now defined. The metric p is defined as the *percentage of runs* in which the lattice-based algorithm requires fewer function evaluations than does the  $\mathbb{Z}^n$ -based algorithm to converge 99.99% of the way from the initial value of  $J(\mathbf{x})$  to the optimal value of  $J(\mathbf{x})$  [which, in these test problems, is easy to compute analytically]. The metric r is defined as the *ratio of the average number of function evaluations* required for the lattice-based algorithm to converge 99.99% of the way from the initial value of  $J(\mathbf{x})$  to the optimal value of  $J(\mathbf{x})$  divided by the average number of function evaluations needed for the  $\mathbb{Z}^n$ -based algorithm to converge the same amount.

The p and r measures described above (averaged over 5000 runs for each value of n) were calculated in

## CHAPTER 10. LATTICE-BASED DERIVATIVE-FREE OPT. VIA GLOBAL SURROGATES

n	2	3	4	5	6	7	8
p	74.77	81.32	84.03	84.53	84.43	84.56	85.28
r	0.4290	0.4161	0.3273	0.3585	0.3150	0.3345	0.3060

Table 10.1. Performance comparison between the  $A_n$ -based SP algorithm and the  $\mathbb{Z}^n$ -based SP algorithm applied to randomly shifted quadratic bowls for n=2 to 8. It is seen that the  $A_8$ -based SP algorithm outperformed the  $\mathbb{Z}^8$ -based SP algorithm 85% of the time, and on average required 30% as many function evaluations to reach the same level of convergence.

n	8
p	90.65
r	0.1554

Table 10.2. Performance comparison between the  $E_8$ -based SP algorithm and the  $\mathbb{Z}^8$ -based SP algorithm applied to randomly shifted quadratic bowls. It is seen that the  $E_8$ -based SP algorithm outperformed the  $\mathbb{Z}^8$ -based SP algorithm 91% of the time, and on average required 17% as many function evaluations to reach the same level of convergence, thus offering nearly twice the performance of  $A_n$ .

the case of the  $A_n$  lattice (for n=2 to n=8) and the  $E_8$  lattice, and are reported in Tables 10.1 and 10.2. Note that values of p over 50% and values of r less than 1 indicate that, on average, the lattice-based SP algorithm outperforms the  $\mathbb{Z}^n$ -based SP algorithm, with p quantifying how often and r quantifying how much.

Note in Table 3.1 that the "best" lattice in n = 2 and n = 3, according to several standard metrics, is  $A_n$ ; however, as the dimension of the problem increases, several other lattices become available, and that by n = 8 the  $E_8$  lattice appears to be the best choice. This observation is consistent with the numerical results

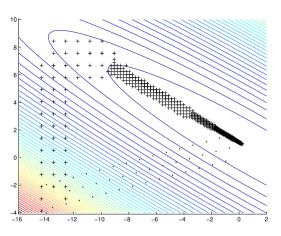


Figure 10.1: Typical paths taken by the  $A_2$ -based SP algorithm (dots) and the  $\mathbb{Z}^2$ -based SP algorithm (+) on a randomly-shifted quadratic bowl.

reported in Tables 10.1 and 10.2, which indicates that the  $A_n$ -based optimizations provided a consistent and substantial improvement over the  $\mathbb{Z}^n$ -based optimizations over the entire range n = 2 to 8, and that, in n = 8, the  $E_8$ -based optimization significantly outperformed the  $A_8$ -based optimization.

The mechanism by which the lattice-based SP algorithms outperform the  $\mathbb{Z}^n$ -based SP algorithm on quadratic test problems is now examined in detail. As described previously, the  $\mathbb{Z}^n$  minimal positive basis vectors are distributed with poor angular uniformity and can not be selected on nearest-neighbor lattice points.

descent direction is configured somewhere approximately midway between the oddball vector and one of the Cartesian unit vectors), the search path must "zig-zag" to move towards the actual minimum. If the local curvature of the function is small compared to the current lattice spacing, then the search algorithm must take several steps in a rather poor direction before it must eventually turn back down the "valley floor", as illustrated by the path of the  $\mathbb{Z}^n$ -based SP algorithm in Figure 10.1. Once in this valley, the lattice spacing must be diminished such that each step of the "zig-zag" path required to proceed down the valley floor in fact decreases the function; this leads to otherwise unnecessary lattice refinement and thus very slow progress by the SP algorithm. This effect is exacerbated when the vectors of the poll set are of substantially different length, as the entire set of vectors must be scaled down until movement along the direction of the longest poll vector during this zig-zagging motion still decreases the function. This leads to the poor convergence behavior demonstrated by the  $\mathbb{Z}^n$ -based SP algorithm along the narrow valley floor of the quadratic bowl indicated in Figure 10.1. Of course, the present arguments are statistical in nature, and in specific cases either the  $A_n$ -based SP algorithm or the  $\mathbb{Z}^n$ -based SP algorithm will sometimes get "lucky" and converge remarkably quickly. However, it is clear that the optimal descent direction at any given iteration is more likely to be "far" from

When the optimal descent direction is poorly approximated by these n+1 vectors (such as when the optimal

#### 10.2 SP applied to randomly-shifted Rosenbrock functions

the poll vectors when the poll set is distributed with poor angular uniformity.

The  $A_n$ -based and  $\mathbb{Z}^n$ -based SP algorithms were also applied to a randomly-shifted Rosenbrock function in a similar fashion. Figure 10.2 demonstrates a typical case, indicating the respective rates of convergence of the two SP algorithms. The  $A_n$ -based SP algorithm demonstrates a substantially improved convergence rate compared to the  $\mathbb{Z}^n$ -based SP algorithm.

These results demonstrate that the efficiency of the SP portion of a pattern search can be substantially improved simply by implementing a more efficient lattice to discretize parameter space.

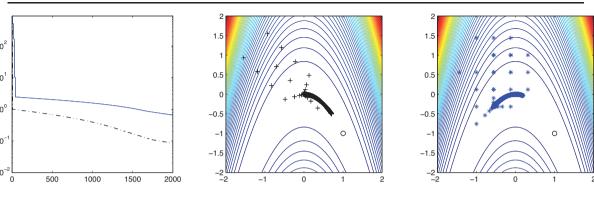


Figure 10.2: A sample SP minimization comparing the  $A_n$ -based case (dash-dot line at left and black + at center) with the  $\mathbb{Z}^n$ -based case (solid line at left and blue \* at right) on a randomly shifted Rosenbrock function. Note the superior convergence rate of the  $A_n$ -based approach (as illustrated in the convergence plot at left), resulting in further progress toward the minimum at [1,-1] (as illustrated in the subfigures at center and right).

#### 10.3 LABDOGS applied to randomly shifted Rosenbrock functions

To test the hypothesis that the efficiency of the full LABDOGS algorithm is significantly affected by the choice of the lattices which coordinate it, a more demanding test than a quadratic bowl is required. We thus consider here the application of the full LABDOGS algorithm to randomly shifted Rosenbrock functions. The "valley" in which the minimum of the Rosenbrock function lies is narrow, curved, and relatively flat

### CHAPTER 10. LATTICE-BASED DERIVATIVE-FREE OPT. VIA GLOBAL SURROGATES

(that is, with a vanishing second derivative) along the bottom. This makes it a difficult test case for any SMF-like algorithm to approximate with a surrogate function of sufficient accuracy to be particularly useful along the valley floor, other than simply to indicate where the function evaluations are currently relatively sparse. In other words, both the search and poll components of the LABDOGS algorithm are put to the test when searching along the valley floor of the Rosenbrock function.

Two comparisons of the efficiencies of the  $A_n$ -based and  $\mathbb{Z}^n$ -based LABDOGS algorithms (using c=5)

applied to randomly shifted Rosenbrock functions are reported here. As in the SP tests described previously, the initial  $A_n$  and  $\mathbb{Z}^n$  lattices were scaled appropriately so as to be of the same initial density.

Recall in the SP tests the metric p, which quantified *how often* the lattice-based method outperformed the Cartesian-based method, and the metric r, which quantifying *how much* the lattice-based method outperformed the Cartesian-based method. In this section, we use two similar metrics,  $\bar{p}$  and  $\bar{r}$ , but now terminate

each optimization after a particular number of iterations rather than after convergence to a given percentage of the (known) optimal solution. Specifically, the metric  $\bar{p}$  is defined as the percentage of runs in which the  $A_n$ -based LABDOGS algorithm converged further than did the  $\mathbb{Z}^n$ -based LABDOGS algorithm after 300 function evaluations, whereas the metric  $\bar{r}$  is defined as the ratio of the average function value to which the  $A_n$ -based LABDOGS algorithm converged after 300 function evaluations divided by the average function value to which the  $\mathbb{Z}^n$ -based LABDOGS algorithm converged after 300 function evaluations. The results for n=2 to 5 (averaged over 200 runs for n=2, 3, and 4, and 100 runs for n=5) are reported in Table 10.3. Note that values of  $\bar{p}$  over 50% and values of  $\bar{r}$  less than 1 indicate that, on average, the lattice-based LABDOGS algorithm outperforms the  $\mathbb{Z}^n$ -based LABDOGS algorithm consistently and significantly outperforms the  $\mathbb{Z}^n$ -based LABDOGS algorithm consistently and significantly outperforms the  $\mathbb{Z}^n$ -based LABDOGS algorithm.

Figure 10.3 compares the convergence of the  $A_n$ -based and  $\mathbb{Z}^n$ -based LABDOGS algorithms on a representative realization of the Rosenbrock function in n = 6. The convergence of the two algorithms are similar

#### 10.3. LABDOGS APPLIED TO RANDOMLY SHIFTED ROSENBROCK FUNCTIONS

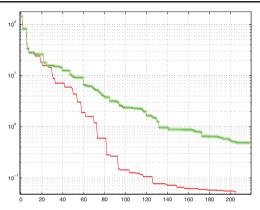


Figure 10.3: Convergence of the Checkers code using  $A_n$  (red) vs  $\mathbb{Z}^n$  (green), on an n = 6 Rosenbrock function.

in behavior during the first 20 iterations, during which they share a nearly identical search, with the differences between the two becoming more and more apparent as convergence is approached. Initially, the poll steps return much smaller improvements than the search steps. Once the surrogate model adequately represents the walls of the Rosenbrock function, thereby identifying the "valley floor", the search becomes less effective, and both algorithms rely more heavily on the polling algorithm to identify the minimum.

CHAPTER 10. LATTICE-BASED DERIVATIVE-FREE OPT. VIA GLOBAL SURROGATES

n	2	3	4	5
$\bar{p}$	64.0	56.0	63.0	68.0
ī	0.651	0.699	0.773	0.758

Table 10.3. Performance comparison between the  $A_n$ -based LABDOGS algorithm and the  $\mathbb{Z}^n$ -based LABDOGS algorithm applied to randomly shifted Rosenbrock functions. For n=2, it is seen that the  $A_n$ -based SP algorithm outperformed the  $\mathbb{Z}^n$ -based SP algorithm about 64% of the time, and on average converged to a function value 65% better using the same number of function evaluations.

#### 10.4 LABDOGS applied to Branin and $T_1$

Thus far, only functions with unique minima have been explored. As the LABDOGS algorithm has the capability to locate and explore multiple local minima in an attempt to identify and refine an estimate of the global minimum, some searches were performed on two test functions with multiple minima, Branin and  $T_1$ , to demonstrate this capability.

On the interval -2 < x < 2, -2 < y < 2, the Branin function has five local minima. As seen in Figure 10.4, with the search parameter c=2, the LABDOGS algorithm does an excellent job of locating and exploring all of these local minima, eventually converging to an accurate estimate of the global minimum. With c=10000, the search tends to be more "space-filling", acting at each step to reduce the maximum uncertainty of the Kriging surrogate. It is clearly evident that, as the number of function evaluations gets large in the c=10000 case, this search will tend to explore nearly uniformly over the entire feasible domain. [In the limit that c is infinite, the function evaluations become dense as  $N \to \infty$ , thereby assuring global convergence.] However, for a small number of total function evaluations N [which should be the primary problem of interest if function evaluations are expensive!], the strategy with smaller c in fact identifies and refines the estimate of the global

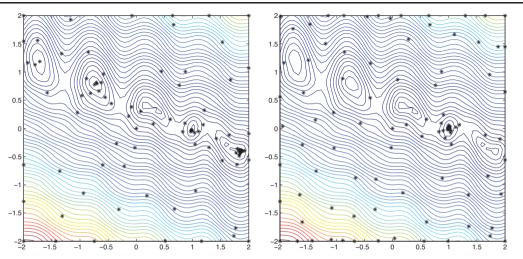


Figure 10.4: Points evaluated by the LABDOGS algorithm when exploring the Branin function (with multiple minima), with (left) c=2 and (right) c=10000. Note the more "focused" sampling when c is small and the more "exploratory" sampling when c is large.

minimum point much sooner, as the case with large c wastes a lot of computational effort reducing the uncertainty of the surrogate in areas predicted to have poor function values.

Similar behavior can be seen for the  $T_1$  test function in Figure 10.5. Initially, the algorithm happens upon

#### CHAPTER 10. LATTICE-BASED DERIVATIVE-FREE OPT. VIA GLOBAL SURROGATES

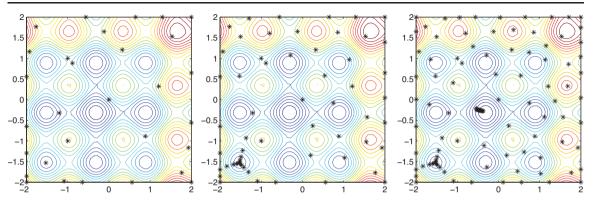


Figure 10.5: Points evaluated by the LABDOGS algorithm when exploring the  $T_1$  function (with multiple minima) with c=1000 after (left) 30 function evaluations, (center) 60 function evaluations, and (right) 100 function evaluations. Note (after 30 function evaluations) that the LABDOGS algorithm initially identifies and converges to a local minimum near the lower-left corner. Ultimately (after 100 function evaluations), the LABDOGS algorithm successfully identifies a refined estimate of the global minimum.

the local minimum in the lower-left corner of the feasible domain. With its exploratory function evaluations, however, the algorithm ultimately identifies and refines its estimate of the global minimum.

While these results indicate encouraging global exploration, further testing of the LABDOGS algorithm

#### 10.5. LABDOGS PERFORMANCE SUMMARY

on nonconvex functions is certainly warranted, particularly in high-dimensional problems. In particular, further refinement of the algorithm to provide the most robust combination of "focused" and "exploratory" sampling remains to be performed; however, the present results clearly demonstrate the capability and flexibility of the LABDOGS algorithm to strike this balance while maintaining maximum computational efficiency.

#### 10.5 LABDOGS Performance Summary

This chapter proposes a new algorithm, dubbed LABDOGS, for derivative-free optimization formed via the tight integration of

• the efficient SMF algorithm (see §7.1) for a surrogate-based search coordinated by an underlying grid,

- in order to keep function evaluations far apart until convergence is approached,
- a uniform "grid" selected from those available in lattice theory (see Part I and further extensions in §7) to coordinate such an optimization algorithm, in order to reduce the average quantization error of a grid of a given density and to better distribute the poll points during the poll step, and
- a highly effective search algorithm, leveraging a Kriging interpolant (see §8) to construct the search function  $J(\mathbf{x}) = \hat{f}(\mathbf{x}) c \cdot s^2(\mathbf{x})$  combining both the function predictor and a model of its associated uncertainty, in order to provide a flexible combination of global exploration and local refinement during the search (see §9).

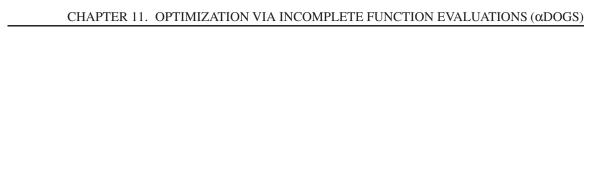
The numerical results achieved via this algorithm, as reported in this chapter, indicate effective convergence of the resulting algorithm on a range of benchmark optimization problems, and reveal a clear advantage for using an efficient lattice derived from an n-dimensional sphere packing to coordinate such a search, rather than the heretofore default choice,  $\mathbb{Z}^n$ , which is simply untenable in light of the clear advantages of using alternative lattices which are, quantifiably, both more uniform and have a more favorable distribution of nearest neighbors, especially as the dimension of the optimization problem is increased.

#### CHAPTER 10. LATTICE-BASED DERIVATIVE-FREE OPT. VIA GLOBAL SURROGATES

The flexible numerical code we have developed which implements this algorithm, dubbed Checkers, has been written from scratch, and each subroutine of the code has been scrutinized to maximize its overall efficiency for systems with expensive function evaluations.

## evaluations (\alpha DOGS)

**Optimization via incomplete function** 



# **Optimization in the presence of complex constraint boundaries (latticeMADS)**

## Part III

## Structured computational interconnects

<b>A</b>	

13 Joe chap 114 Joe chap 2

15 Joe chap 3

16 Joe chap 4

165

167

169

171

### C ------

Joe chap 1

### \_

Joe chap 2

## Joe chap 3

**Chapter 16** 

Joe chap 4

#### **References (Part I)**

Anzin (2002) On the density of a lattice covering for n = 11 and n = 14 (in Russian), *Uspekhi Mat. Nauk* 57, no. 2, 187188; English translation in *Math. Surveys* 57, 407409.

Aste, T, & Weaire, D (2008) The pursuit of the perfect packing. Taylor & Francis.

Baranovskii, EP (1994) The perfect lattices  $\Gamma(\mathfrak{A}^n)$ , and the covering density of  $\Gamma(\mathfrak{A}^n)$ . European Journal of Combinatorics **15**, 317-323.

Belitz, P, & Bewley, T (2011) New Horizons in Sphere Packing Theory, Part II: Lattice Based Derivative Free Optimization via Global Surrogates. *SIAM Journal on Optimization*, submitted.

Berrou, C, Glavieux, A, & Thitimajshima, P (1993) Near Shannon limit error-correcting coding and decoding: Turbo-Codes, In *ICC'93*, Geneva, Switzerland. 1064-1070.

Beukemann, A, & Klee, WE (1992) Minimal nets. Z. Krist. 201, 37-51.

Blahut, RE (2003) Algebraic codes for data transmission. Cambridge.

Blatov, VA (2006) Multipurpose crystallochemical analysis with the program package TOPOS. *IUCr CompComm Newsletter* 7, 438.

Blatov, VA (2007) Topological relations between three-dimensional periodic nets. I. Uninodal nets. *Acta Cryst.* **A63**, 329-343.

Blatov, VA, Delgado-Friedrichs, O, O'Keeffe, M, & Proserpio, DM (2007) Three-periodic nets and tilings: natural tilings for nets. *Acta Cryst.* **A63**, 418-425.

- Blatov, VA, O'Keeffe, M, & Proserpio, DM (2009) Vertex-, face-, point-, Schläfli-, and Delaney-symbols in nets, polyhedra and tilings: recommended terminology. *CrystEngComm.*, 2010, DOI: 10.1039/b910671e.
- Blichfeldt, HF (1935) The minimum values of positive quadratic forms in six, seven and eight variables. *Math. Z.* **39**, 115.
- Bonneau, C, Delgado-Friedrichs, O, O'Keeffe, M, & Yaghi, OM (2004) Three-periodic nets and tilings: minimal nets. *Acta Cryst.* **A60**, 517-520.
- Brown, H, Bülow, R, Neubüser, J, Wondratschek, H, & Zassenhaus, H (1978) Crystallographic groups of 4-dimensional space. Wiley.
- Cecil, T (2005) A numerical method for computing minimal surfaces in arbitrary dimension. J. Comp. Phys. 206, 650-660.
- Cessna, J, & Bewley, T (2011) New Horizons in Sphere Packing Theory, Part III: Structured Computational Interconnects. *SIAM Journal on Computing*, submitted.
- Clark, WE (1930) The Aryabhatiya of Aryabhata: An Ancient Indian Work on Mathematics and Astronomy. University of Chicago Press.

Cohn, H, & Kumar, A (2009) Optimality and uniqueness of the Leech lattice among lattices, Annals of Mathematics 170,

- 1003-1050.
- Conway, JH, & Sloane, NJA (1984) On the Voronoï regions of certain lattices. SIAM J. Alg. Disc. Meth. 5, 294-302.
- Conway, JH, & Sloane, NJA (1997) Low-dimensional lattices. VII Coordination sequences. *Proc. R. Soc. Lond. A* **453**, 2369-2389.

Conway, JH, & Sloane, NJA (1998) Sphere Packings, Lattices, and Groups, Springer.

Coxeter, HSM (1951) Extreme forms. Canadian Journal of Mathematics, 3, 391-441.

Coxeter, HSM (1970) *Twisted Honeycombs*. Regional Conference Series in Mathematics, No. 4, American Mathematical Society, Providence.

Coxeter, HSM (1973) Regular Polytopes. Dover.

Coxeter, HSM (1974) Regular Complex Polytopes. Cambridge.

Coxeter, HSM (1987) *Projective Geometry*. Springer.

Coxeter, HSM (1989) *Introduction to Geometry*. Wiley.

Croft, Falconer, & Guy (1991) Unsolved Problems in Geometry. Springer.

Curtis, RT (1976) A new combinatorial approach to M24. Math. Proc. Camb. Phil. Soc. 79, 25-42.

Delgado-Friedrichs, O, Foster, MD, O'Keeffe, M, Proserpio, DM, Treacy, MMJ, & Yaghi, OM (2005) What do we know about three-periodic nets? *Journal of Solid State Chemistry* **178**, 2533-2554.

Delgado-Friedrichs, O, & Huson, DH (2000) 4-Regular Vertex-Transitive Tilings of E<sup>3</sup>. *Discrete & Computational Geometry* **24**, 279-292.

Delgado-Friedrichs, O, & O'Keeffe, M (2003) Identification of and symmetry computation for crystal nets. *Acta Cryst.* **A59**, 351-360.

Delgado-Friedrichs, O, & O'Keeffe, M (2007) Three-periodic nets and tilings: face-transitive tilings and edge-transitive

- nets revisited. Acta Cryst. A63, 344-347.
- Delgado-Friedrichs, O, O'Keeffe, M, & Yaghi, OM (2003a) Three-periodic nets and tilings: regular and quasiregular nets. *Acta Cryst.* **A59**, 22-27.
- Delgado-Friedrichs, O, O'Keeffe, M, & Yaghi, OM (2003b) Three-periodic nets and tilings: semiregular nets. *Acta Cryst.* **A59**, 515-525.
- Delgado-Friedrichs, O, O'Keeffe, M, & Yaghi, OM (2003c) The CdSO<sub>4</sub>, rutile, cooperite, and quartz dual nets: interpenetration and catenation. *Solid State Sciences.* **5**, 73-78.
- Delgado-Friedrichs, O, O'Keeffe, M, & Yaghi, OM (2006) Three-periodic nets and tilings: edge-transitive binodal structures. *Acta Cryst.* **A62**, 350-355.
- Dorozinski, TE, & Fischer, W (2006) A novel series of sphere packings with arbitrarily low density. *Z. Kristallogr.* **221**, 563-566.
- Fejes Tóth, L (1972) Lagerungen in der Ebene auf der Kugel und im Raum. Springer.
- Fejes Tóth, L (1975) Research problem 13. *Period. Math. Hungar.* **6**, 197-199.
- Fischer, W (2005) On sphere packings of arbitrarily low density. Z. Kristallogr. 220, 657662.

Fejes Tóth, L (1959) Verdeckung einer Kugel durch Kugeln. Publ. Math. Debrecen 6, 234-240.

- Friedman, E (2009) Erich's Packing Center: http://www2.stetson.edu/~efriedma/packing.html.
- Gallager, RG (1963) Low Density Parity Check Codes. Monograph, M.I.T. Press.

- Gandini, PM, & Willis, JM (1992) On finite sphere packings. *Math Pannon* 3, 19-29.
- Gandini, PM, & Zucco, A (1992) On the sausage catastrophe in 4-space. Mathematicka 39, 274-278.
- Gauss, CF (1831) Untersuchungen über die Eigenscahften der positiven ternären quadratischen Formen von Ludwig August Seber, *Göttingische gelehrte Anzeigen*, 1831 Juli 9, also published in *J. Reine Angew. Math.* **20** (1840), 312320, and *Werke*, vol. 2, Königliche Gesellschaft der Wissenschaften, Göttingen, 1876, pp. 188196.
- Gensane, T (2004) Dense packings of equal spheres in a cube. Electronic J. Combinatorics 11 (1), #R33, 1-17.
- Grosse-Kunstleve, RW, Brunner, GO, & Sloane, NJA (1996) Algebraic description of coordination sequences and exact topological densities for zeolites. *Acta Cryst.* A52, 879-889.
- Grover, P, Sahai, A, & Park, SY (2010) The finite-dimensional Witsenhausen counterexample. *IEEE Transactions on Automatic Control*, submitted. Draft available as: arXiv:1003.0514v1.
- Hales, S (1727) *Vegetable Staticks*, London: Printed for W. and J. Innys; and T. Woodward. Hales, TC (2005) A proof of the Kepler conjecture. *Annals of Mathematics*, **162**, 10651185.
- Hales, TC (2006), Historical overview of the Kepler conjecture, *Discrete Comput Geom* **36**, 520.
- Harriot, T (1614) *De Numeris Triangularibus Et inde De Progressionibus Artithmeticis: Magisteria magna*, in Berry, J, and Stedall, J, editors (2009) *Thomas Harriot's Doctrine of Triangular Numbers: the 'Magisteria Magna'*, published by the European Mathematical Society.
- Heath, TL (1931) A Manual of Greek Mathematics. Oxford University Press (republished in 2003 by Dover).
- Heesch, H, & Laves, RT (1933) Über dünne Kugelpackeungen, Z. Krist. 85, 443-453.

Hyde, ST, Delgado-Friedrichs, O, Ramsden, SJ, Robins, V (2006) Towards enumeration of crystalline frameworks: The 2D hyperbolic approach. *Solid State Sci.* **8**, 740752.

Janssen, T (1986) Crystallography of quasi-crystals. Acta Cryst. A42, 261-271.

Kaku, M (1999) Introduction to Superstring and M-Theory. Springer-Verlag.

Kelvin, WT (1887) On the Division of Space with Minimum Partitional Area, *Philosophical Magazine* **24**, No. 151, p. 503.

Kepler, J (1611) Strena seu de nive sexangula (The six-cornered snowflake).

Koch, E, & Fischer, W (1995) Sphere packings with three contacts per sphere and the problem of the least dense sphere packing. *Z. Krist.* **210**, 407414.

Koch, E, & Fischer, W (2006) Sphere packings and packings of ellipsoids. In *International Tables for Crystallography*,

Vol C, Sec 9.1.1, pp 746-751.

Korkine, A, & Zolotareff, G (1873) Sur les formes quadratiques. *Math. Ann.* **6**, 366389.

Korkine, A, & Zolotareff, G (1875) Sur les formes quadratiques positives, Math. Ann. 11, 242292.

Lagrange, JL (1773) Recherches darithmétique. Nov. Mem. Acad. Roy. Sc. Bell Lettres Berlin In Œuvres, 3, Gauthier-Villars, Paris, 18671892, pp. 693758.

Leech, J (1956) The problem of the thirteen spheres. *Math. Gazette* **40**, 22-23.

MacWilliams, FJ, & Sloane, NJA (1977) The Theory of Error Correcting Codes. North Holland.

- Meschkowski, H (1960) Ungelöste und unlösbare Probleme der Geometrie, Braunsehweig.
- Milnor, J (1976) Hilberts problem 18: on crystallographic groups, fundamental domains, and on sphere packings, in *Mathematical Developments Arising from Hilbert Problems*, Proceedings of Symposia in Pure Mathematics, vol. 28, AMS, Providence, RI, pp. 491506.
- Moon, T (2005) Error Correction Coding, Mathematical Methods and Algorithms. Wiley.
- Morelos-Zaragoza, RH (2006) The Art of Error Correcting Coding. Wiley.
- Ockwig, NW, Delgado-Friedrichs, O, O'Keeffe, M, & Yaghi, OM (2005) Recticular Chemistry: Occurance and Taxonomy of Nets and Grammar for the Design of Frameworks. *Acc. Chem. Res.* **38**, 176-182.
- O'Keeffe, M (1991a) Dense and rare four-connected nets. Z. Krist. 196, 21-37
- O'Keeffe, M (1991b) N-dimensional diamond, sodalite, and rare sphere packings. Acta Cryst. A47, 748-753.
- O'Keeffe, M (2008) Three-periodic nets and tilings: regular and related infinite polyhedra. Acta Cryst. A64, 425-429.
- O'Keeffe, M, Eddaoudi, M, Li, H, Reineke, T, & Yaghi, OM (2000) Frameworks for Extended Solids: Geometrical Design Principles. *Journal of Solid State Chemistry* **152**, 3-20.
- O'Keeffe, M, Peskov, MA, Ramsden, SJ, & Yaghi, OM (2008) The Recticular Chemistry Structure Resource (RCSR) Database of, and Symbols for, Crystal Nets. *Acc. Chem. Res.* **41**, 1782-1789.
- Pless, V (1998) Introduction to the Theory of Error-Correcting Codes. Wiley.
- Rukeyser, M (1972) The Traces of Thomas Hariot. Random House.

- Sadoc, JF, & Rivier, N. editors (1999) Foams and Emulsions, Kluwer,
- Schrijver, A (1986) Thoery of linear and integer programming. Wiley.
- Schürmann, A (2006) On packing spheres into containers; about Kepler's finite sphere packing problem. Documenta Mathematica 11, 393406.
- Schürmann, A, & Vallentin, F (2006) Computational approaches to lattice packing and covering problems, Discrete Comput. Geom. 35, 73-116.
- Schütte, K. & van der Waerden, BL (1953) Das Problem der dreizehn Kugeln, Math Annalen 125, 325-334.
- Schwarzenberger, RLE (1980) N-dimensional crystallography. Pitman.
- Shannon, CE (1949) Communication in the presence of noise. Proc. Institute of Radio Engineers, 37, 1021.
- Sikirić, MD, Schürmann, A, & Vallentin, F (2008) A generalization of Voronoi's reduction theory and its application. Duke Mathematical Journal 142, 127-164.
- Silverman, RA, & Balser, M (1954) Coding for a constant data rate source. IRE Trans. Inform. Theory IT-4, 50-63.
- Skelton, RE, & de Oliveira, MC (2009) Tensegrity Systems. Springer.
- Sloan, IH, & Kachoyan PJ (1987) Lattice methods for multiple integration: theory, error analysis and examples. SIAM J. Num. Anal. 24, 116-128.
- Sloane, NJA (1987) Theta-Series and Magic Numbers for Diamond and Certain Ionic Crystal Structures, J. Math. Phys.
- **28**, 1653-1657.

Szpiro, G (2003) Kepler's conjecture: how some of the greatest minds in history helped solve one of the oldest math problems in the world. Wiley.

Thompson, TM (1983) From error-correcting codes through sphere packings to simple groups. The Mathematical Association of America.

Thue, A (1892) Om nogle geometrisk taltheoretiske Theoremer, Forand. Skand. Natur. 14, 352353.

Tietäväinen, A (1973) On the nonexistence of perfect codes over finite fields, SIAM J. Appl. Math. 24, 88-96.

Treacy, MMJ, Rivin, I, Balkovsky, E, Randall, KH, Foster, MD (2004) Enumeration of periodic tetrahedral frameworks. II. Polynodal graphs. *Microporous Mesoporous Mater.* **74**, 121132.

Ungerboeck, G (1982) Channel coding with multilevel/phase signals. *IEEE Trans. Inform. Theory*, **IT-28**, 55-67.

Vardy A & Be'ery Y (1993) Maximum likelihood decoding of the Leech lattice, *IEEE Trans. Inform. Theory* **39**, 1435-1444.

Weaire D & Phelan R (1994) A counterexample to Kelvin's conjecture on minimal surfaces, *Phil. Mag. Lett.* **69**, 107110.

Wells, AF (1977) Three-Dimensional Nets and Polyhedra. Wiley.

Wells, AF (1979) Further Studies of Three-dimensional Nets. ACA Monograph No. 8.

Wells, AF (1983) Six New Three-Dimensional 3-Connected Nets 4.n<sup>2</sup>. Acta Cryst. **B39**, 652-654.

Wells, AF (1984) Structural inorganic chemistry. Oxford University Press.

Willis, JM (1983) Research problem 35. Period. Math. Hungar. 14, 312-314.

Zong, C (1999) Sphere Packings. Springer.

#### **Additional References (Part II)**

Booker, A, Dennis, JR, Frank, P, Serafini, D, Torczon, V, & Trosset, M (1999) A rigorous framework for optimization of expensive functions by surrogates. *Structural and Multidisciplinary Optimization* 17, 113.

- Conway, JH, & Sloane, NJA (1998) Sphere Packings, Lattices, and Groups, Springer.
- Coope, ID, & Price, CJ (2001) On the convergence of grid-based methods for unconstrained optimization. SIAM J. Optim., 11, 859869.
- Cox, DD & John, S (1997) SDO: A statistical method for global optimization. In *Multidisciplinary Design Optimization: State of the Art* (edited by Alexandrov, N, & Hussaini, MY), 315329. SIAM.
- Elder, JF IV (1992) Global Rd optimization when probes are expensive: the GROPE algorithm. *Proceedings of the 1992 IEEE International Conference on Systems, Man, and Cybernetics,* **1**, 577582, Chicago.

Fejes Tóth, L (1971) Perfect distribution of points on a sphere. Periodica Mathematica Hungarica 1, 25-33.

- Isaaks, EH, & Srivastav, RM (1989) An Introduction to Applied Geostatistics. Oxford.
- Jones, DR (2001) A Taxonomy of Global Optimization Methods Based on Response Surfaces. Journal of Global Opti-
- mization 21, 345-383.

Krige, DG (1951) A statistical approach to some mine valuations and allied problems at the Witwatersrand. Master's thesis of the University of Witwatersrand, South Africa.

Kushner, HJ (1964) A new method of locating the maximum point of an arbitrary multipeak curve in the presence of noise. *Journal of Basic Engineering*, **86**, 97106.

Matheron, G (1963) Principles of geostatistics. *Economic Geology* **58**, 1246-1266.

Meschkowski, H (1960) Ungelöste und unlösbare Probleme der Geometrie, Braunsehweig.

Mockus, J (1994) Application of Bayesian approach to numerical methods of global and stochastic optimization. *Journal of Global Optimization*, **4**, 347365.

Nelder, JA, & Mead, R (1965) A simplex method for function minimization Computer Journal 7, 308313.

Perttunen, C (1991) A computational geometric approach to feasible region division in constrained global optimization. *Proceedings of the 1991 IEEE Conference on Systems, Man, and Cybernetics.* 

Rasmussen, CE, & Williams, CKI (2006) *Gaussian processes for machine learning*. MIT Press.

lutionary Operation *Technometrics* **4**, 441-461.

Stuckman, BE (1988) A global search method for optimizing nonlinear systems. *IEEE Transactions on Systems, Man*,

Spendley, W, Hext, GR, & Himsworth, FR (1962) Sequential Application of Simplex Designs in Optimisation and Evo-

and Cybernetics, 18, 965977.

Tammes, PML (1930) On the origin of number and arrangement of the places of exit on the surface of pollen-grains. *Recueil des travaux botaniques néerlandais*, **27**, 1-84.

Thomson, JJ (1904) On the Structure of the Atom: an Investigation of the Stability and Periods of Oscillation of a number of Corpuscles arranged at equal intervals around the Circumference of a Circle; with Application of the Results to the Theory of Atomic Structure. *Philosophical Magazine Series* 6, 7 (39), 237–265.

Torczon, V (1997) On the convergence of pattern search algorithms. SIAM J. Optim., 7, 1-25.

Torn, A, & Zilinskas, A (1987) Global Optimization, Springer.

# **Additional References (Part III)**