



**HAL**  
open science

## **Incremental construction of nested designs based on two-level fractional factorial designs**

Rodrigo Cabral Farias, Luc Pronzato, Maria-João Rendas

► **To cite this version:**

Rodrigo Cabral Farias, Luc Pronzato, Maria-João Rendas. Incremental construction of nested designs based on two-level fractional factorial designs. J. Pilz, V.B. Melas, A. Bathke. Statistical Modeling and Simulation for Experimental Design and Machine Learning Applications, Springer, pp.77-110, 2023, 978-3-031-40054-4. <10.1007/978-3-031-40055-1\_5>. <hal-02483004>

**HAL Id: hal-02483004**

**<https://hal.science/hal-02483004v1>**

Submitted on 18 Feb 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

# Chapter 1

## Incremental construction of nested designs based on two-level fractional factorial designs

Rodrigo Cabral-Farias, Luc Pronzato and Maria-João Rendas

**Abstract** The incremental construction of nested designs having good spreading properties over the  $d$ -dimensional hypercube is considered, for values of  $d$  such that the  $2^d$  vertices of the hypercube are too numerous to be all inspected. A greedy algorithm is used, with guaranteed efficiency bounds in terms of packing and covering radii, using a  $2^{d-m}$  fractional-factorial design as candidate set for the sequential selection of design points. The packing and covering properties of fractional-factorial designs are investigated and a review of the related literature is provided. An algorithm for the construction of fractional-factorial designs with maximum packing radius is proposed. The spreading properties of the obtained incremental designs, and of their lower dimensional projections, are investigated. An example with  $d = 50$  is used to illustrate that their projection in a space of dimension close to  $d$  has a much higher packing radius than projections of more classical designs based on Latin hypercubes or low discrepancy sequences.

### 1.1 Introduction

We consider the incremental construction of designs with large packing radius in the  $d$ -dimensional hypercube, using the coffee-house rule of [20] and [21, Chap. 4]: each new point introduced maximises the distance to its nearest neighbour in the current design. This simple algorithm is known to guarantee an efficiency of 50% in terms of packing and covering radii, for each design size along the construction. Intuitively, when  $d$  is large, the first points selected are vertices of the hypercube, and we shall provide arguments that validate

---

R. Cabral-Farias, L. Pronzato, M.-J. Rendas  
Laboratoire I3S, Université Côte d’Azur - CNRS, 2000 route des Lucioles, 06903 Sophia Antipolis, France, e-mail: [cabral@i3s.unice.fr](mailto:cabral@i3s.unice.fr), [Luc.Pronzato@cnrs.fr](mailto:Luc.Pronzato@cnrs.fr), [rendas@i3s.unice.fr](mailto:rendas@i3s.unice.fr)

this intuition. However, when  $d$  is very large, it is impossible to inspect all vertices and select one at every iteration. We show that restriction of the search to fractional factorial designs having a large enough covering radius does not entail any loss of performance up to some design size: an example shows that designs of size up to  $2^{15} + 1 = 32\,769$ , with 50% packing and covering efficiencies, can be constructed in this way when  $d = 50$ . The packing and covering properties of these designs when projected on smaller dimension subspaces are investigated. Transformation rules based on rescaling are proposed to generate designs that populate the interior of the hypercube. Numerical computations indicate that the designs obtained have slightly larger covering radii than more classical space-filling designs based on (non-incremental) Latin hypercubes or (incremental) Sobol' low discrepancy sequence, but have significantly larger packing radii.

The paper is organised as follows. Section 1.2 sets notation and recalls the definitions of packing and covering radii and the incremental construction of designs based on the coffee-house rule. The main properties of two-level fractional factorial designs are recalled in Sections 1.3 and 1.4 to make the paper self-contained. Their spreading properties are investigated in Sections 1.5 (packing radius) and 1.6 (covering radius). An algorithm is given in Sect. 1.5 for the construction of fractional factorial designs with large covering radii. Section 1.7 studies the restriction of the coffee-house rule to two-level fractional factorial designs, and shows that the 50% packing and covering efficiencies are preserved when the fractional factorial design has minimum Hamming distance at least  $d/4$ . A rescaling rule is proposed to generate incremental designs not concentrated on the vertices of the hypercube, and properties of projections on smaller dimensional subspaces are investigated. An example in dimension  $d = 50$  illustrates the presentation. Section 1.8 briefly concludes.

## 1.2 Greedy coffee-house design

Let  $\mathcal{X}$  denote a compact subset of  $\mathbb{R}^d$  with nonempty interior; throughout the paper we consider the case where  $\mathcal{X}$  is the  $d$ -dimensional hypercube  $\mathcal{C}_d = [-1, 1]^d$ . Denote by  $\mathbf{X}_k = \{\mathbf{x}_1, \dots, \mathbf{x}_k\}$  a  $k$ -point design when the ordering of the  $\mathbf{x}_i$  is not important, and by  $\mathbf{X}_k = [\mathbf{x}_1, \dots, \mathbf{x}_k]$  the ordered sequence; for  $1 \leq k_1 \leq k_2$ ;  $\mathbf{X}_{k_1:k_2}$  denotes the design formed by  $[\mathbf{x}_{k_1}, \mathbf{x}_{k_1+1}, \dots, \mathbf{x}_{k_2}]$ , with  $\mathbf{X}_{1:k} = \mathbf{X}_k$ . The  $j$ th coordinate of a design point  $\mathbf{x}_i$  is denoted by  $\{\mathbf{x}_i\}_j$ ,  $j = 1 \dots, d$ ;  $\|\mathbf{x}\| = \left(\sum_{i=1}^d \{\mathbf{x}\}_i^2\right)^{1/2}$  denotes the  $\ell_2$  norm of the vector  $\mathbf{x} \in \mathbb{R}^d$ ,  $\|\mathbf{x}\|_1 = \sum_{i=1}^d |\{\mathbf{x}\}_i|$  (respectively,  $\|\mathbf{x}\|_\infty = \max_{i=1, \dots, d} |\{\mathbf{x}\}_i|$ ) is its  $\ell_1$  (respectively,  $\ell_\infty$ ) norm. For any  $\mathbf{x} \in \mathbb{R}^d$  and any  $k$ -point design  $\mathbf{X}_k$  in  $\mathcal{X}$  we denote  $d(\mathbf{x}, \mathbf{X}_k) = \min_{i=1, \dots, k} \|\mathbf{x} - \mathbf{x}_i\|$ . For  $\mathbf{x}$  and  $\mathbf{x}'$  two vectors of same size,  $\mathbf{z} = \mathbf{x} \circ \mathbf{x}'$  denotes their Hadamard product, with components

$\{\mathbf{z}\}_i = \{\mathbf{x}\}_i \{\mathbf{x}'\}_i$ .  $\mathcal{B}(\mathbf{x}, r)$  denotes the closed ball with centre  $\mathbf{x}$  and radius  $r$ . For  $\mathcal{A}$  a finite set,  $|\mathcal{A}|$  is the number of elements in  $\mathcal{A}$ .

Space-filling design aims at constructing a set  $\mathbf{X}_k$  of points in  $\mathcal{X}$ , with given cardinality  $k$ , that “fill”  $\mathcal{X}$  in a suitable way; see, e.g., [25, 26]. Two measures of performance are standard. The *covering radius* of  $\mathbf{X}_k$  is defined by

$$\text{CR}(\mathbf{X}_k) = \max_{\mathbf{x} \in \mathcal{X}} d(\mathbf{x}, \mathbf{X}_k). \quad (1.1)$$

It corresponds to the smallest  $r$  such that the  $k$  closed balls of radius  $r$  centred at the  $\mathbf{x}_i$  cover  $\mathcal{X}$ .  $\text{CR}(\mathbf{X}_k)$  is also called the dispersion of  $\mathbf{X}_k$  [22, Chap. 6] and corresponds to the minimax-distance criterion [12] used in space-filling design; small values are preferred. Another widely used geometrical criterion of spreadness is the *packing radius*

$$\text{PR}(\mathbf{X}_k) = \frac{1}{2} \min_{\mathbf{x}_i, \mathbf{x}_j \in \mathbf{X}_k, \mathbf{x}_i \neq \mathbf{x}_j} \|\mathbf{x}_i - \mathbf{x}_j\|. \quad (1.2)$$

$\text{PR}(\mathbf{X}_k)$  is also called *separating radius*, it corresponds to the largest  $r$  such that the  $k$  open balls of radius  $r$  centred at the  $\mathbf{x}_i$  do not intersect;  $2 \text{PR}(\cdot)$  corresponds to the maximin-distance criterion [12] often used in computer experiments; large values are preferred. We may also consider the combined measure given by the mesh ratio

$$\tau(\mathbf{X}_k) = \frac{\text{CR}(\mathbf{X}_k)}{\text{PR}(\mathbf{X}_k)},$$

with  $\tau(\mathbf{X}_k) \geq 1$  for any design  $\mathbf{X}_k$  when  $\mathcal{X}$  is convex, since the  $k$  balls  $\mathcal{B}(\mathbf{x}_i, \text{PR}(\mathbf{X}_k))$  cannot cover  $\mathcal{X}$ .

When the objective is to construct a sequence  $\mathbf{X}_k = [\mathbf{x}_1, \dots, \mathbf{x}_k]$  such that  $\text{PR}(\mathbf{X}_k)$  is reasonably large, and/or  $\text{CR}(\mathbf{X}_k)$  is reasonably small, for all  $k \in \{2, \dots, n\}$ , the following greedy algorithm is called coffee-house design ([20], [21, Chap. 4]). See also [14] for an early suggestion.

#### Algorithm 1 (Coffee-house)

- 0) Select  $\mathbf{x}_1 \in \mathcal{X}$ , set  $\mathbf{S}_1 = \{\mathbf{x}_1\}$  and  $k = 1$ .
- 1)  $\text{for } k = 1, 2, \dots$  **do**
  - find  $\mathbf{x}^* \in \text{Arg max}_{\mathbf{x} \in \mathcal{X}} d(\mathbf{x}, \mathbf{S}_k)$ , set  $\mathbf{S}_{k+1} = \mathbf{S}_k \cup \{\mathbf{x}^*\}$ .

The point  $\mathbf{x}^*$  can be obtained by a Voronoi tessellation of  $\mathcal{X}$  (when  $d$  is small enough) or a MCMC method, see [25]. Note that the choice of  $\mathbf{x}^*$  is not necessarily unique. The construction is much easier when a finite candidate set  $\mathcal{X}_n$  with  $n$  points is substituted for  $\mathcal{X}$  at Step 1. In the paper we show that when  $\mathcal{X} = \mathcal{C}_d$ , a well-chosen  $\mathcal{X}_n$  yields a drastic simplification of calculations

for very large  $d$  but does not entail any loss of performance for the greedy algorithm. For a given order of magnitude of the anticipated number of design points to be used, we informally define the notions of small, large and very large dimension  $d$  as follows: small  $d$  are such that the construction of designs with  $2^d$  points may be considered; large  $d$  correspond to situations where the greedy construction above with a candidate set  $\mathcal{X}_n$  containing all  $2^d$  vertices of  $\mathcal{C}_d$  is conceivable; very large  $d$  cover cases where exploration of all  $2^d$  vertices of  $\mathcal{C}_d$  is unfeasible. For instance,  $\mathcal{C}_{50}$  has more than  $10^{15}$  vertices, a situation considered in Sect. 1.7.

Let  $\text{CR}_n^* = \min_{\mathbf{X}_n} \text{CR}(\mathbf{X}_n)$  denote the minimum covering radius for an  $n$ -point design in  $\mathcal{X}$ ,  $n \geq 1$ , and  $\text{PR}_n^* = \max_{\mathbf{X}_n} \text{PR}(\mathbf{X}_n)$  denote the maximum packing radius,  $n \geq 2$ . The following property is a consequence of [8]. Note that the *efficiencies*  $\text{CR}_k^*/\text{CR}(\mathbf{S}_k)$  and  $\text{PR}(\mathbf{S}_k)/\text{PR}_k^*$  belong to  $[0, 1]$  by definition; large values are preferred for both.

**Theorem 1** *The sequence of designs  $\mathbf{S}_k$  constructed with Algorithm 1 satisfies*

$$\frac{\text{CR}_k^*}{\text{CR}(\mathbf{S}_k)} \geq \frac{1}{2} \quad (k \geq 1) \quad \text{and} \quad \frac{\text{PR}(\mathbf{S}_k)}{\text{PR}_k^*} \geq \frac{1}{2} \quad (k \geq 2). \quad (1.3)$$

Moreover,  $\tau(\mathbf{S}_k) \leq 2$  for all  $k \geq 2$ .

*Proof.* By construction, for all  $k \geq 1$ ,  $\text{PR}(\mathbf{S}_{k+1}) = d(\mathbf{x}_{k+1}, \mathbf{S}_k)/2 = \text{CR}(\mathbf{S}_k)/2$ . Therefore, for all  $k \geq 2$ ,  $\tau(\mathbf{S}_k) = \text{CR}(\mathbf{S}_k)/\text{PR}(\mathbf{S}_k) = 2 \text{PR}(\mathbf{S}_{k+1})/\text{PR}(\mathbf{S}_k) \leq 2$ . Also, from the pigeonhole principle, for *any pair* of  $k$  and  $(k+1)$ -point designs  $\mathbf{X}_k$  and  $\mathbf{X}'_{k+1}$ , one of the ball  $\mathcal{B}(\mathbf{x}_i, \text{CR}(\mathbf{X}_k))$  with  $\mathbf{x}_i$  in  $\mathbf{X}_k$  contains two points  $\mathbf{x}'_i$  and  $\mathbf{x}'_j$  of  $\mathbf{X}'_{k+1}$ , which implies  $\text{CR}(\mathbf{X}_k) \geq \text{PR}(\mathbf{X}'_{k+1})$ . Therefore, for the greedy construction we have in particular  $\text{CR}_k^* \geq \text{PR}(\mathbf{S}_{k+1}) = \text{CR}(\mathbf{S}_k)/2$  and  $\text{PR}_{k+1}^* \leq \text{CR}(\mathbf{S}_k) = 2 \text{PR}(\mathbf{S}_{k+1})$ .  $\square$

In the rest of the paper we take  $\mathcal{X} = \mathcal{C}_d$ . Take  $\mathbf{x}_1 = \mathbf{0}_d$ , the null vector of dimension  $d$ , which corresponds to the centre of  $\mathcal{C}_d$ . The design  $\mathbf{S}_1 = [\mathbf{x}_1]$  has thus minimum covering radius, with  $\text{CR}(\mathbf{S}_1) = \sqrt{d}$ . When applying Algorithm 1, for all  $k$  such that  $\text{CR}(\mathbf{S}_k) = \sqrt{d}$ , the point  $\mathbf{x}^*$  chosen at Step 1 is then necessarily a vertex of  $\mathcal{C}_d$ ; that is,  $\mathbf{x}_k \in \{-1, 1\}^d$  for  $k = 2, \dots, k_*(d)$ , where  $k_*(d)$  is the first  $k$  such that  $\text{CR}(\mathbf{S}_k) < \sqrt{d}$ . (Note that it implies that  $\text{CR}_n^* \geq \sqrt{d}/2$  for all  $n \leq k_*(d) - 1$ .) This simple property has the important consequence that the greedy construction of a design  $\mathbf{S}_n$  via Algorithm 1 initialised at  $\mathbf{x}_1 = \mathbf{0}_d$  can restrict its attention to the set of vertices of  $\mathcal{C}_d$ , provided that  $n \leq k_*(d)$ . For  $d \leq 4$ , since any pair of distinct vertices of the hypercube are at distance at least  $2 \geq \sqrt{d}$ , Algorithm 1 sequentially (and indifferently) selects the vertices of  $\mathcal{C}_d$  until they are exhausted, and  $k_*(d) = 2^d + 1$ . For larger  $d$ , the behaviour depends on the order in which vertices are selected in the first iterations; that is, on the particular choices of  $\mathbf{x}^*$  made at Step 1. The largest values of  $k_*(d)$  obtained for  $d$  up to 8 are indicated in Table 1.1. We shall see in Sect. 1.7.1 that  $k_*(d)$  is large enough

for practical applications when  $d$  gets very large; see (1.12). Occasionally, Algorithm 1 may still take  $\mathbf{x}^*$  among vertices of  $\mathcal{C}_d$  when  $k \geq k_*(d)$ , that is, when  $\text{CR}(\mathbf{S}_k) < \sqrt{d}$ ; this again depends on the first choices made for  $\mathbf{x}^*$ . Denote by  $k_{NV}(d)$  the first  $k > 1$  such that  $\mathbf{x}^*$  chosen at Step 1 is not a vertex of  $\mathcal{C}_d$  (with necessarily  $k_{NV}(d) \geq k_*(d)$ ); the largest values of  $k_{NV}(d)$  that we have obtained are also indicated in Table 1.1.

The difficulty is that the inspection of all  $2^d$  vertices of  $\mathcal{C}_d$  is unpractical for very large  $d$ . The main objective of the paper is therefore to propose a method for selecting a subset  $\mathcal{X}_n$  of  $2^d$  vertices of  $\mathcal{C}_d$  on which Algorithm 1 can be applied, ensuring that  $\max_{\mathbf{x} \in \mathcal{X}_n} d(\mathbf{x}, \mathbf{S}_k) = \text{CR}(\mathbf{S}_k) = \sqrt{d}$  for all  $k \leq n$ , with  $n$  large enough to allow the construction of designs of practical size. The method relies on the notion of fractional factorial design, the basic properties of which are recalled in the next two sections. Their spreading properties in terms of packing and covering radii are then investigated in Sections 1.5 and 1.6. We prefer not to call those designs “space-filling” since they are supported on the vertices of  $\mathcal{C}_d$ ; they nevertheless satisfy the bounds of Th. 1.

**Table 1.1** First  $k$  such that  $\text{CR}(\mathbf{S}_k) < \sqrt{d}$  and first  $k > 1$  such that  $\mathbf{x}_k$  is not a vertex of  $\mathcal{C}_d$ .

$d$	2	3	4	5	6	7	8
$k_*(d)$	5	9	17	17	33	65	129
$k_{NV}(d)$	5	9	17	33	65	129	133

### 1.3 Two-level fractional factorial designs

This section only gives a brief summary of the topic; one may refer to [2, 3] for a thorough and illuminating exposition.

#### 1.3.1 Half fractions: $m = 1$

A  $2^d$  factorial (or *full factorial*) design is formed by the  $2^d$  vertices of  $\mathcal{C}_d$ ; each design point  $\mathbf{x}_i$  is such that  $\{\mathbf{x}_i\}_j \in \{-1, 1\}$ ,  $i = 1, \dots, 2^d$ ,  $j = 1, \dots, d$ . The notation used for a  $2^d$  factorial design is illustrated in Table 1.2(a) for the case  $d = 3$ . The coordinates of design points correspond to factors and are denoted by lowercase letters\*.

A (regular)  $2^{d-m}$  fractional factorial design is obtained by setting  $d - m$  coordinates (sometimes called *basic factors*) of the  $2^{d-m}$  design points at values given by a  $2^{d-m}$  factorial design, the other  $m$  coordinates being defined

\* The design in Table 1.2(a) is listed in what is called *standard order*.

by *generating equations*, or *generators*, that explain how they are obtained (calculated) from the basic factors. Without any loss of generality, we can suppose that the basic factors correspond to the first  $d - m$  coordinates. Table 1.2(b) shows the  $2^{4-1}$  fractional factorial design  $\mathbf{X}_{2^{4-1}}^{(a)}$  obtained from the generating equation  $\{\mathbf{x}\}_4 = d = abc$ . By the product of two factors we mean the entrywise (Hadamard) product of the corresponding columns in the design viewed as a  $n \times d$  matrix. Since all  $\{\mathbf{x}\}_i$  belong to  $\{-1, 1\}$  for  $\mathbf{x}$  in  $\mathbf{X}_{2^{d-m}}$ , this implies in particular that the product of a factor by itself gives a vector with all components equal to 1, which we denote by  $\mathbf{1}$ . The equation  $d = abc$  is thus equivalent to  $\mathbf{1} = abcd$ , called *defining relation*. Changing the generating equation to  $d = ab$  gives another  $2^{4-1}$  fractional factorial design  $\mathbf{X}_{2^{4-1}}^{(b)}$ , presented in Table 1.2(c). Both designs are called a *half fraction* of the full factorial design with  $d = 4$ . Since  $d = ab$  in Table 1.2(c),  $\{\mathbf{x}\}_4 = \{\mathbf{x}\}_1\{\mathbf{x}\}_2$  for all  $\mathbf{x}$  in  $\mathbf{X}_{2^{4-1}}^{(b)}$ , and this design does not allow us to estimate separately the main effect of  $\{\mathbf{x}\}_4$  and the interaction  $\{\mathbf{x}\}_1\{\mathbf{x}\}_2$ ; these effects are said *confounded*, or *aliased*. The equation  $d = ab$  also implies  $a = bd$  and  $b = ad$ , showing that the effects of  $\{\mathbf{x}\}_1$  and  $\{\mathbf{x}\}_2\{\mathbf{x}\}_4$  are confounded, as well as those of  $\{\mathbf{x}\}_2$  and  $\{\mathbf{x}\}_1\{\mathbf{x}\}_4$ . We say that this design has resolution  $R = III$  (notation with a Roman numeral is traditional): no  $p$  factor effect is confounded with any other effect containing less than  $R - p$  factors,  $p = 0, \dots, R$ . For the design in Table 1.2(b), we get  $a = bcd$ ,  $b = acd$ ,  $c = abd$  and of course  $d = abc$  which is the generating equation. Here none of the main and 2 factor interaction effects are confounded, and the design has resolution  $R = IV$ . In general, designs of high resolution are preferable. When  $m = 1$ , the highest possible resolution  $R = d$  is obtained for the half fraction with defining relation  $\{\mathbf{x}\}_d = \prod_{i=1}^{d-1} \{\mathbf{x}\}_i$  (unique up to a sign change and a permutation of variables).

(a) a $2^3$ factorial design	(b) a $2^{4-1}$ fractional factorial design	(c) another $2^{4-1}$ fractional factorial design
$\mathbf{X}_8 \{\mathbf{x}\}_1 = a \{\mathbf{x}\}_2 = b \{\mathbf{x}\}_3 = c$	$\mathbf{X}_8 a b c d = abc$	$\mathbf{X}_8 a b c d = ab$
$\mathbf{x}_1$ -1 -1 -1	$\mathbf{x}_1$ -1 -1 -1 -1	$\mathbf{x}_1$ -1 -1 -1 1
$\mathbf{x}_2$ 1 -1 -1	$\mathbf{x}_2$ 1 -1 -1 1	$\mathbf{x}_2$ 1 -1 -1 -1
$\mathbf{x}_3$ -1 1 -1	$\mathbf{x}_3$ -1 1 -1 1	$\mathbf{x}_3$ -1 1 -1 -1
$\mathbf{x}_4$ -1 -1 1	$\mathbf{x}_4$ -1 -1 1 1	$\mathbf{x}_4$ -1 -1 1 1
$\mathbf{x}_5$ 1 1 -1	$\mathbf{x}_5$ 1 1 -1 -1	$\mathbf{x}_5$ 1 1 -1 1
$\mathbf{x}_6$ 1 -1 1	$\mathbf{x}_6$ 1 -1 1 -1	$\mathbf{x}_6$ 1 -1 1 -1
$\mathbf{x}_7$ -1 1 1	$\mathbf{x}_7$ -1 1 1 -1	$\mathbf{x}_7$ -1 1 1 -1
$\mathbf{x}_8$ 1 1 1	$\mathbf{x}_8$ 1 1 1 1	$\mathbf{x}_8$ 1 1 1 1

**Table 1.2** (a): a  $2^3$  factorial design; (b) and (c): two  $2^{4-1}$  fractional factorial designs, with  $\{\mathbf{x}\}_1 = a$ ,  $\{\mathbf{x}\}_2 = b$ ,  $\{\mathbf{x}\}_3 = c$  and  $\{\mathbf{x}\}_4 = d$ .

### 1.3.2 Several generators

#### Defining relations

A  $2^{d-m}$  fractional factorial design with  $m > 1$  requires more than one generating equation, and the construction of suitable designs with high resolution has motivated intensive research since the pioneering papers [2, 3]. To ensure that the resolution is larger than  $II$ , the generating equations are chosen independent, which means that a generator cannot be obtained by multiplying together two other generators. It implies that there are no repetitions within the columns of the design table. If this were not the case, two main effects would be confounded since we would have  $\{\mathbf{x}\}_i = \{\mathbf{x}\}_j$  for some  $i, j \in \{1, \dots, d\}$  and all  $\mathbf{x}$  in the design. The generators are called *principal* when there are only positive signs in the defining relations. When multiplying the  $m$  generating equations by the  $d$  independent factors and by themselves in all possible ways, we obtain the *complete set of defining relations*. *Principal defining relations* are obtained from principal generators. The complete set of principal defining relations defines a unique *fraction*, that is a unique design up to  $2^m$  sign changes in the variables defined by the generators. The set contains  $2^m$  defining relations (including the trivial one  $\mathbb{1} = \mathbb{1}$ ), each one having the form  $\mathbb{1}$  equals the product of a subset of factors, called *word*. For example, the  $2^{6-2}$  design with independent factors  $a, b, c$  and  $d$  and generating equations  $e = abcd$  and  $f = acd$  has defining relations  $\mathbb{1} = abcde$ ,  $\mathbb{1} = acdf$  and  $\mathbb{1} = bef$ , the latter being obtained by multiplying the first two since  $(abcde) \times (acdf) = a^2bc^2d^2ef = bef$ .

#### Resolution

The resolution of the design is given by the shortest word length within the complete set of defining relations, here  $R = III$  (since  $\mathbb{1} = bef$ ). Another choice of generating equations may yield a different set of defining relations and a design with different resolution. For instance, choosing  $e = abc$  and  $f = acd$  in a  $2^{6-2}$  design yields the complete set of defining relations  $\mathbb{1} = abce = acdf = bdef$ , and the resolution is now  $IV$ . To identify the resolution of a design, the notation  $2_R^{d-m}$  is used. For example, the design with six variables and generating equations  $e = abcd$  and  $f = acd$  is denoted  $2_{III}^{6-2}$ . Designs  $2_{III}^{d-m}$  (with  $n = 2^{d-m}$  points in  $d$  variables) can be constructed for  $d$  up to  $n-1$  and  $n$  a power of 2 (designs with  $d = n-1$  are called *saturated*, see Sect. 1.3.3), designs  $2_{IV}^{d-m}$  can be constructed for  $d$  up to  $n/2$  and  $n$  a power of 2; generators for  $2_{III}^{7-4}$ ,  $2_{III}^{15-11}$ ,  $2_{III}^{31-26}$ ,  $2_{IV}^{7-3}$ ,  $2_{IV}^{8-4}$ ,  $2_{IV}^{16-11}$  can be found in [2]. The nonregular designs of [23] have resolution  $III$  and allow the exploration of  $n-1$  variables for  $n$  a multiple of 4. Generators for designs  $2_V^{8-2}$  and  $2_V^{11-4}$  are given in [3]. The *Matlab* function `fracfactgen.m` implements

the algorithm of [6] for the construction of a  $2^{d-m}$  fractional factorial design of prescribed resolution (when it exists).

A design with resolution  $R$  contains a full factorial design in any subset of  $R - 1$  variables. Omitting  $p$  variables from a  $2_R^{d-m}$  design with resolution  $R$  produces a design of resolution  $R$  in  $d - p$  variables but with  $n = 2^{d-m}$  points. All words containing characters associated with the dropped variables must be removed from the set of defining relations. The resulting design may duplicate some design points, and a more economical design with similar word pattern in the defining relations may exist in general. Bounds on the maximum resolution attainable for a  $2^{d-m}$  design are given in [7].

### Word length pattern

The *word length pattern*  $\mathcal{A}(\mathbf{X}_n)$  of a  $2^{d-m}$  design with resolution  $R$  is defined by the distribution of word lengths in the complete set of defining relations,

$$\mathcal{A}(\mathbf{X}_n) = [1, 0, \dots, 0, A_R(\mathbf{X}_n), A_{R+1}(\mathbf{X}_n), \dots, A_d(\mathbf{X}_n)],$$

with  $A_k(\mathbf{X}_n)$  denoting the number of words of length  $k$  ( $A_0(\mathbf{X}_n) = 1$  since the word  $\mathbf{1}$  is always present and  $\sum_{k=0}^d A_k(\mathbf{X}_n) = 2^m$ ). Among two designs  $\mathbf{X}_n^{(a)}$  and  $\mathbf{X}_n^{(b)}$  having the same (maximum) resolution  $R$ , the paper [7] recommends to select the one with *minimum aberration*: let  $i_*$  be the smallest  $i \geq 1$  such that  $A_i(\mathbf{X}_n^{(a)}) \neq A_i(\mathbf{X}_n^{(b)})$ , then  $\mathbf{X}_n^{(a)}$  is preferred to  $\mathbf{X}_n^{(b)}$  if  $A_{i_*}(\mathbf{X}_n^{(a)}) < A_{i_*}(\mathbf{X}_n^{(b)})$ , and  $\mathbf{X}_n^{(b)}$  is preferred to  $\mathbf{X}_n^{(a)}$  otherwise\*\*. The construction of a minimum aberration design can thus be viewed as the sequential minimisation of the  $A_i(\mathbf{X}_n)$  for  $i \geq 1$ . A minimum aberration  $2^{d-m}$  design has necessarily generators that contain all  $d$  variables [7]; lists of generators are tabulated in [33].

### 1.3.3 Minimum size

**Proposition 1** *The size of a  $2^{d-m}$  design necessarily satisfies  $n = 2^{d-m} \geq d + 1$ . The designs for which equality holds have resolution III.*

*Proof.* Since the  $m$  generators must be independent and each one must involve at least 2 of the  $d - m$  basic factors, we get

$$m \leq \sum_{k=2}^{d-m} \binom{d-m}{k} = 2^{d-m} - (d+1-m),$$

---

\*\* It may happen, though rarely, that two designs with different defining relations have exactly the same word length pattern; the minimum aberration criterion then does not provide any preference.

that is,  $n = 2^{d-m} \geq d + 1$ . Designs for which equality holds are those that use all possible independent generators (without any loss of generality, we only consider principal generators). They cannot have resolution  $R$  larger than  $III$  since there are generators defined as product of two basic factors, and thus defining relations involving words of length 3. We prove by contradiction that they cannot have resolution  $II$ .

If the design has resolution  $II$ , it means that one of the defining relations has been obtained by multiplying two relations  $\mathbb{1} = \underline{w}$  and  $\mathbb{1} = \underline{z}$ , with words  $\underline{w}$  and  $\underline{z}$  that only differ by two letters, say  $a, b$ . There are two possibilities: either  $\underline{w} = \underline{ta}$  and  $\underline{z} = \underline{tb}$ , or  $\underline{w} = \underline{tab}$  and  $\underline{z} = \underline{t}$ . In both case, the multiplication  $\underline{w} \times \underline{z}$  gives the defining relation  $\mathbb{1} = ab$ , which cannot exist since the generators are independent.  $\square$

Proposition 1 gives a lower bound on the number of points for a given dimension  $d$ ; it also gives an upper bound on the number of generators that can be used for a given  $d$ ,

$$m \leq m^*(d) = \lfloor d - \log_2(d + 1) \rfloor. \tag{1.4}$$

That is, for  $2^k \leq d < 2^{k+1}$  we can construct  $2^{d-m}$  fractional factorial designs with  $m \in \{1, 2, \dots, m^*(d) = d - k - 1\}$ . Values of  $d$ ,  $m$  and  $n$  for minimum-size  $2^{d-m}$  designs with  $d + 1 = 2^{d-m}$ , called *saturated designs*, for  $d$  up to  $d = 255$  are given in Table 1.3.

**Table 1.3** Saturated designs.

$d$	3	7	15	31	63	127	255
$m$	1	4	11	26	57	120	247
$n$	4	8	16	32	64	128	256

## 1.4 Two-level factorial designs and error correcting codes

### 1.4.1 Definitions and properties

The construction of a two-level factorial design  $\mathbf{X}_n$  possesses strong similarities with the construction of an error correcting code  $\mathbf{C}_n$  with binary alphabet  $\{0, 1\}$ : design points correspond to codewords in  $\mathbf{C}_n$  and  $d$  is the length of the code, with  $n = |\mathbf{C}_n|$  (and  $n = 2^{d-m}$  for a fractional factorial design). Associating levels 1 and  $-1$  to symbols 0 and 1, respectively, we obtain that the product rule used in Sect. 1.3 corresponds now to addition modulo 2, and the codes corresponding to fractional factorial designs, which are obtained through generating equations, are linear.

Since  $\{\mathbf{x}_i\}_j \in \{-1, 1\}$  for each design point and any  $j \in \{1, \dots, d\}$ , the the Hamming distance  $d_H(\mathbf{x}_i, \mathbf{x}_j)$ , which counts the number of components that differ between two design points  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , satisfies

$$d_H(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{4} \|\mathbf{x}_i - \mathbf{x}_j\|^2 = \frac{1}{2} \|\mathbf{x}_i - \mathbf{x}_j\|_1.$$

The minimum distance of  $\mathbf{C}_n$ ,  $\rho_H(\mathbf{C}_n)$ , is defined as the minimum Hamming distance between two codewords in  $\mathbf{C}_n$  and we shall write  $\rho_H(\mathbf{X}_n) = \rho_H(\mathbf{C}_n)$  with  $\mathbf{C}_n$  the code associated with  $\mathbf{X}_n$ . More generally,  $\rho_H(\mathbf{X}_n) = \min_{\mathbf{x}_i, \mathbf{x}_j \in \mathbf{X}_n, \mathbf{x}_i \neq \mathbf{x}_j} d_H(\mathbf{x}_i, \mathbf{x}_j)$  for any design  $\mathbf{X}_n$  supported on the vertices of  $\mathcal{C}_d$ . Therefore,

$$\rho_H(\mathbf{X}_n) = \text{PR}^2(\mathbf{X}_n).$$

Similarly, the (Hamming) covering radius  $\text{CR}_H(\mathbf{X}_n)$  of a two-level fractional factorial design  $\mathbf{X}_n$  corresponds to the covering radius  $\text{CR}_H(\mathbf{C}_n)$  of the associated code, and we define more generally

$$\text{CR}_H(\mathbf{X}_n) = \max_{\mathbf{x} \in \{-1, 1\}^d} \min_{\mathbf{x}_i \in \mathbf{X}_n} d_H(\mathbf{x}, \mathbf{x}_i). \quad (1.5)$$

Several results from coding have their counterpart in design theory. Suppose that  $\rho_H(\mathbf{X}_n) \geq 2k + 1$  for some  $k \in \mathbb{N}$ . For each of the  $n$  design points  $\mathbf{x}_i$ , there are  $\binom{d}{\ell}$  points in  $\{-1, 1\}^d$  that are at distance  $\ell$  from  $\mathbf{x}_i$ . Since the  $n$  Hamming balls centred at the  $\mathbf{x}_i$  with radii  $k$  do not intersect, we obtain the *sphere-packing bound*, see, e.g., [32, Th. 20.1]:  $n \sum_{\ell=0}^k \binom{d}{\ell} \leq 2^d$ . For a fractional factorial design  $\mathbf{X}_n$  with  $n = 2^{d-m}$ , it gives

$$2^m \geq \sum_{\ell=0}^k \binom{d}{\ell}. \quad (1.6)$$

Note that  $\text{CR}_H(\mathbf{X}_n) \geq \lfloor \rho_H(\mathbf{X}_n)/2 \rfloor$ . When  $\rho_H(\mathbf{X}_n) = 2k + 1$  and equality is reached in (1.6), all points in  $\{-1, 1\}^d$  are at Hamming distance at most  $k$  to exactly one design point in  $\mathbf{X}_{2^{d-m}}$ , which corresponds to the notion of *perfect code*.

Delete now the  $p - 1$  last coordinates of each  $\mathbf{x}_i \in \mathbf{X}_n$ , with  $p = \rho_H(\mathbf{X}_n)$ . The  $n$  points that are obtained belong to  $\{-1, 1\}^{d-(p-1)}$  and are all distinct. Therefore, their number  $n$  is less than  $2^{d-p+1}$ , which gives the *Singleton bound* ([29], [32, Th. 20.2]):  $n \leq 2^{d-p+1}$ . For a fractional factorial design  $\mathbf{X}_{2^{d-m}}$ , we obtain

$$\rho_H(\mathbf{X}_{2^{d-m}}) \leq m + 1. \quad (1.7)$$

Another result from coding theory gives an upper bound on the size  $n$  of a design  $\mathbf{X}_n$  supported on  $\{-1, 1\}^d$  when  $\rho_H(\mathbf{X}_n)$  is large: Plotkin bound ([16, Th. 5.5.2], [24]) states that

$$n \leq \left\lfloor \frac{\rho_H(\mathbf{X}_n)}{\rho_H(\mathbf{X}_n) - d/2} \right\rfloor \quad (1.8)$$

when  $\rho_H(\mathbf{X}_n) > d/2$ .

Besides the value of the packing radius  $\text{PR}(\mathbf{X}_n)$ , the distribution of the distances  $\|\mathbf{x}_i - \mathbf{x}_j\|$ , or  $d_H(\mathbf{x}_i, \mathbf{x}_j)$ , between pairs of design points is also of interest. This is particularly true in the present context where there exist many pairs of points at the same distance since all design points are vertices of the hypercube. In [12] a design  $\mathbf{X}_n^*$  is called maximin-distance optimal when it maximises  $\text{PR}(\mathbf{X}_n)$  and minimises the number of pairs of points at distance  $2 \text{PR}(\mathbf{X}_n^*)$ . That definition is extended as follows in [19]. For a given design  $\mathbf{X}_n$ , consider the list  $[d_1, d_2, \dots, d_q]$  of intersite distances sorted in decreasing order, with  $d_1 = 2 \text{PR}(\mathbf{X}_n)$  and  $1 \leq q \leq n(n-1)/2$ . Denote by  $\mathcal{J}(\mathbf{X}_n) = [J_1, \dots, J_q]$  the associated counting list defined by  $J_k = |(i, j) : \|\mathbf{x}_i - \mathbf{x}_j\| = d_k, \mathbf{x}_i, \mathbf{x}_j \in \mathbf{X}_n|$ ,  $k = 1, \dots, q$ . In [19], a design is called maximin-distance optimal if it maximises  $d_1$ , and among all such designs minimises  $J_1$ , maximises  $d_2$ , and among all such designs minimises  $J_2 \dots$  and so on. Following [35], we call (Hamming) *distance distribution* of a design  $\mathbf{X}_n$  supported on  $\{-1, 1\}^d$  the list  $\mathcal{B}(\mathbf{X}_n) = [B_0(\mathbf{X}_n), B_1(\mathbf{X}_n), \dots, B_d(\mathbf{X}_n)]$  where

$$B_k(\mathbf{X}_n) = \frac{1}{n} |(i, j) : d_H(\mathbf{x}_i, \mathbf{x}_j) = k, \mathbf{x}_i, \mathbf{x}_j \in \mathbf{X}_n|, \quad k = 0, \dots, d \quad (1.9)$$

(so that  $\sum_{k=0}^d B_k(\mathbf{X}_n) = n$  and  $B_0(\mathbf{X}_n) = 1$  when all points are distinct).

Let  $\mathbf{X}_n$  be a  $2^{d-m}$  fractional factorial design;  $\mathbf{X}_n$  is balanced, i.e., each value  $+1$  and  $-1$  appears equally often for each factor, and for any  $\mathbf{x}_i \in \mathbf{X}_n$ ,  $\sum_{\mathbf{x}_j \in \mathbf{X}_n, j \neq i} d_H(\mathbf{x}_i, \mathbf{x}_j) = \sum_{\mathbf{x}_j \in \mathbf{X}_n, j \neq i} \sum_{k=1}^d d_H(\{\mathbf{x}_i\}_k, \{\mathbf{x}_j\}_k) = nd/2$ . Therefore,  $\sum_{k=1}^d k B_k(\mathbf{X}_n) = nd/2$ , and interpreting  $B_k(\mathbf{X}_n)/(n-1)$  as a weight on  $k$ , we get

$$\rho_H(\mathbf{X}_n) = \min\{k \in \{1, \dots, d\} : B_k(\mathbf{X}_n) > 0\} \leq \frac{nd}{2(n-1)}. \quad (1.10)$$

Let  $p$  denote the number of generators written as the product of an odd number of basic factors ( $p \geq 0$ ). For any  $\mathbf{x}_i \in \mathbf{X}_n$ , the design point  $\mathbf{x}_j$  obtained by changing the signs of the  $d-m$  basic factors is at Hamming distance  $d_H(\mathbf{x}_i, \mathbf{x}_j) = d-m+p$  from  $\mathbf{x}_i$ ; that is,

$$B_{d-m+p}(\mathbf{X}_n) \geq 1.$$

In particular, it implies that  $\rho_H(\mathbf{X}_n) \leq d-m+p$ . Also, since each point has at most one point at Hamming distance  $d$ , when  $p = m$  we have  $B_d(\mathbf{X}_n) \leq 1$  and thus  $B_d(\mathbf{X}_n) = 1$ ; see for example the designs  $\mathbf{X}_{16}^{(a1)}$ ,  $\mathbf{X}_{16}^{(a2)}$ ,  $\mathbf{X}_{32}^{(a)}$  and  $\mathbf{X}_{32}^{(b)}$  of Table 1.4.

Due to the equivalence between Hamming and Euclidean distances for a  $2^{d-m}$  design, design selection based on maximin-distance optimality in the sense of [19] sequentially minimises the  $B_k(\mathbf{X}_n)$  for  $k \geq 1$ ; it is similar to selection by the minimum aberration criterion of [7] applied to the distance distribution instead of the word length pattern. In [15] minimum aberration designs are called *maximin word length*.

As noticed in [35], MacWilliams' theorem, see, e.g., [32, Th. 20.3], implies that the distance distribution  $B_k(\mathbf{X}_n)$  and the word length pattern  $\mathcal{A}(\mathbf{X}_n)$  of a given  $2^{d-m}$  design  $\mathbf{X}_n$  are related by

$$A_j(\mathbf{X}_n) = \frac{1}{n} \sum_{k=0}^d B_k(\mathbf{X}_n) P_j(k; d, 2), \quad j = 0, \dots, d, \quad (1.11)$$

$$B_j(\mathbf{X}_n) = n 2^{-d} \sum_{k=0}^d A_k(\mathbf{X}_n) P_j(k; d, 2), \quad j = 0, \dots, d,$$

where the  $P_j(x; d, s)$  are the Krawtchouk polynomials defined by

$$P_j(x; d, s) = \sum_{i=0}^j (-1)^i (s-1)^{j-i} \frac{\Gamma(x+1)}{\Gamma(x+1-i)\Gamma(i+1)} \\ \times \frac{\Gamma(d+1-x)}{\Gamma(d+i+1-x-j)\Gamma(j+1-i)},$$

so that  $P_j(k; d, 2) = \sum_{i=0}^j (-1)^i \binom{k}{i} \binom{d-k}{j-i}$ .

Several extensions of the results above, in various directions, are present in the literature. Let us mention a few. Fractional factorial designs with  $s$  levels, with  $s$  any prime number, are considered in [35], together with designs where different factors may have different numbers of levels, and the notion of generalised minimum aberration is introduced; see also [4]. Space-filling properties of fractional factorial designs with more than two levels are studied in [36], where it is shown that the generalised minimum aberration designs of [35] have good performance in terms of maximin distance for the  $\ell_1$  norm when allowing permutations of factor levels. Starting from an initial  $s$ -level balanced design  $\mathbf{X}_n$ , where each level appears exactly  $n/s$  times for each one of the  $d$  factor, [34] shows how to construct a design  $\mathbf{X}'_n$  with  $d$  factors at  $qs$  levels, for  $n$  divisible by  $qs$  ( $\mathbf{X}'_n$  is a Latin hypercube design when  $q = n/s$ ). When  $s > 2$ , the space-filling properties of  $\mathbf{X}'_n$  (measured by the maximin distance for the  $\ell_1$  norm) can be improved by level permutation, using the approach in [36]. Following the approach of [18], properties of  $2^{d-m}$  designs for prediction with a Gaussian process model defined on the vertices  $\{-1, 1\}^d$  of the hypercube  $[-1, 1]^d$  are investigated in [15]; a practical conclusion is that maximin word length (minimum aberration) designs often coincide with maximin distance designs, but not always. The paper [1] shows how to decompose

a minimum aberration  $2^{d-m}$  design into layers containing two points each, in such a way that the resulting design has suitable space-filling properties. The construction of two-level factorial designs having small covering radius (1.5) is considered in [11] (note, however, that  $\text{CR}_H(\mathbf{X}_n)$  is not necessarily an adequate measure of the space-filling properties of  $\mathbf{X}_n$  over the full hypercube  $[-1, 1]^d$ ); a few general properties are given, and the construction of minimum-size covering designs having  $\text{CR}_H(\mathbf{X}_n) = 1$  and minimum-size designs with  $\text{CR}_H(\mathbf{X}_n) = 2$  is detailed for  $d \leq 7$  (with rather intensive computer search for  $d = 7$ ). The centred  $L_2$ -discrepancy  $\text{CL}_2(\mathbf{X}_n)$  of [9] is a popular measure of uniformity of a design  $\mathbf{X}_n$ . For a  $2^{d-m}$  fractional factorial design,  $\text{CL}_2(\mathbf{X}_n)$  is a function of the  $A_i(\mathbf{X}_n)$  in the word-length pattern  $\mathcal{A}(\mathbf{X}_n)$  [5]; see also [31] for related results. A relation between  $\text{CL}_2(\mathbf{X}_n)$  and the distance distribution  $\mathcal{B}(\mathbf{X}_n)$  is established in [30] for more general balanced designs (with  $n$  runs and  $d$  factors, each one taking  $s$  levels, and, for each factor, each level appearing equally often).

### 1.4.2 Examples

The  $2^{4-1}_{IV}$  design  $\mathbf{X}_8$  of Table 1.2(b), with generator  $d = abc$ , has word length pattern  $\mathcal{A} = [1\ 0\ 0\ 0\ 1]$ ; its distance distribution is  $\mathcal{B} = [1\ 0\ 6\ 0\ 1]$ ; it reaches the bound (1.7) since  $\rho_H(\mathbf{X}_8) = 2 = m + 1$ . The design in Table 1.2(c) with  $d = ab$  has resolution *III*,  $\mathcal{A} = [1\ 0\ 0\ 1\ 0]$  and  $\mathcal{B} = [1\ 1\ 3\ 3\ 0]$ ; it is thus worse than previous one both in terms of aberration and maximin distance.

Other examples with more factors are presented in Table 1.4.  $\mathbf{X}_{16}^{(a1)}$  (respectively,  $\mathbf{X}_{16}^{(a2)}$ ) is better than  $\mathbf{X}_{16}^{(b1)}$  (respectively,  $\mathbf{X}_{16}^{(b2)}$ ) both in terms of resolution and maximin distance.  $\mathbf{X}_{16}^{(a2)}$  reaches the bound (1.6), it corresponds to a perfect code of length 7, distance 3 and covering radius 1; see, e.g., [32, p. 215]. The three  $2^{7-2}_{IV}$  designs  $\mathbf{X}_{32}^{(a)}$ ,  $\mathbf{X}_{32}^{(b)}$  and  $\mathbf{X}_{32}^{(c)}$  are those in Table 1 of [7]; they all have resolution *IV* and the hierarchy  $(a) \prec (b) \prec (c)$  is respected both in terms of aberration and maximin distance, where  $(a) \prec (b)$  means that  $(b)$  is preferable to  $(a)$  for the criterion considered. The word length pattern of  $\mathbf{X}_{64}^{(a)}$  is (slightly) better than that of  $\mathbf{X}_{64}^{(b)}$ , but  $\mathbf{X}_{64}^{(b)}$  does better than  $\mathbf{X}_{64}^{(a)}$  in terms of maximin distance. The two  $2^{11-5}_{IV}$  designs  $\mathbf{X}_{64}^{(c)}$  and  $\mathbf{X}_{64}^{(d)}$  are given in Table 3 of [15];  $\mathbf{X}_{64}^{(c)}$  has minimum aberration but is worse than  $\mathbf{X}_{64}^{(d)}$  in terms of maximin distance.

**Table 1.4**  $2^{d-m}$  designs, generators, word length patterns, distance distributions and covering radii.

design	generators	$\mathcal{A}$	$\mathcal{B}$	$\text{CR}_H$
$2_{IV}^{6-2} \mathbf{X}_{16}^{(a1)}$	$abc, acd$	[1 0 0 0 3 0 0]	[1 0 3 8 3 0 1]	1
$2_{III}^{6-2} \mathbf{X}_{16}^{(b1)}$	$abcd, acd$	[1 0 0 1 1 1 0]	[1 0 4 6 3 2 0]	2
$2_{IV}^{7-3} \mathbf{X}_{16}^{(a2)}$	$bcd, abd, acd$	[1 0 0 0 7 0 0 0]	[1 0 0 7 7 0 0 1]	1
$2_{III}^{7-3} \mathbf{X}_{16}^{(b2)}$	$bcd, abd, abcd$	[1 0 0 2 3 2 0 0]	[1 0 1 6 5 2 1 0]	2
$2_{IV}^{7-2} \mathbf{X}_{32}^{(a)}$	$abc, bcd$	[1 0 0 0 3 0 0 0]	[1 1 3 11 11 3 1 1]	1
$\mathbf{X}_{32}^{(b)}$	$abc, ade$	[1 0 0 0 2 0 1 0]	[1 0 6 9 9 6 0 1]	1
$\mathbf{X}_{32}^{(c)}$	$abcd, abce$	[1 0 0 0 1 2 0 0]	[1 0 5 12 7 4 3 0]	1
$2_{IV}^{11-5} \mathbf{X}_{64}^{(a)}$	$abcde, abcdf, abcef, abdef, cdef$	[1 0 0 0 6 12 8 0 1 4 0 0]	[1 0 1 0 14 24 6 8 9 0 1 0]	2
$\mathbf{X}_{64}^{(b)}$	$abcd, abce, acdf, cdef, abcdef$	[1 0 0 0 7 9 6 6 2 1 0 0]	[1 0 0 4 11 18 15 8 4 2 1 0]	2
$\mathbf{X}_{64}^{(c)}$	$cde, bde, abcdf, abce, adef$	[1 0 0 0 4 14 8 0 3 2 0 0]	[1 0 0 2 14 22 8 6 9 2 0 0]	2
$\mathbf{X}_{64}^{(d)}$	$cdef, adef, abef, abcf, bcdf$	[1 0 0 0 5 10 10 5 0 0 0 1]	[1 0 0 0 25 0 27 0 10 0 1 0]	3

## 1.5 Maximin-distance properties of two-level factorial designs

### 1.5.1 Neighbouring pattern and distant-sites pattern

For any design  $\mathbf{X}_n$  supported on the  $2^d$  vertices of  $\mathcal{C}_d$  and any  $\mathbf{x}_i \in \mathbf{X}_n$ , we call neighbouring pattern of  $\mathbf{x}_i$  the counting list  $\mathcal{L}(\mathbf{x}_i; \mathbf{X}_n) = [1, I_1(\mathbf{x}_i; \mathbf{X}_n), \dots, I_d(\mathbf{x}_i; \mathbf{X}_n)]$  with  $I_k(\mathbf{x}_i; \mathbf{X}_n) = |\{j : d_H(\mathbf{x}_i, \mathbf{x}_j) = k, \mathbf{x}_j \in \mathbf{X}_n\}|$ . Similarly, we call distant-sites pattern the list  $\bar{\mathcal{L}}(\mathbf{x}_i; \mathbf{X}_n) = [0, \bar{I}_1(\mathbf{x}_i; \mathbf{X}_n), \dots, \bar{I}_d(\mathbf{x}_i; \mathbf{X}_n)]$  with  $\bar{I}_k(\mathbf{x}_i; \mathbf{X}_n) = |\{j : d_H(\mathbf{x}_i, \mathbf{x}_j) = k, \mathbf{x}_j \notin \mathbf{X}_n\}|$ .  $2^{d-m}$  fractional factorial designs satisfy the following property.

**Proposition 2** *All design points  $\mathbf{x}_i$  of a  $2^{d-m}$  fractional factorial design have the same neighbouring pattern and the same distant-sites pattern.*

*Proof.* Take any  $\mathbf{x} \in \mathbf{X}_n$ ; without any loss of generality we suppose that basic factors correspond to the first  $d - m$  coordinates, and we denote by  $\underline{\mathbf{x}}$  the corresponding part of  $\mathbf{x}$ . The remaining  $m$  components are constructed from the generators that define the design; we can write  $\{\mathbf{x}\}_{d-m+k} = g_k(\underline{\mathbf{x}})$ , with  $g_k(\underline{\mathbf{x}})$  equal to the product of some components of  $\underline{\mathbf{x}}$ ,  $k = 1, \dots, m$ . We collect those  $m$  components in a vector  $\mathbf{g}(\underline{\mathbf{x}})$  and write  $\mathbf{x} = (\underline{\mathbf{x}}, \mathbf{g}(\underline{\mathbf{x}}))$ .

Suppose that there exist  $\mathbf{x}_j \in \mathbf{X}_n$  such that  $d_H(\mathbf{x}, \mathbf{x}_j) = k$ . We first show that for any  $\mathbf{x}' \in \mathbf{X}_n$  there also exists a  $\mathbf{x}'_j \in \mathbf{X}_n$  such that  $d_H(\mathbf{x}', \mathbf{x}'_j) = k$ . Using the same notation as above, we can write  $\mathbf{x}' = (\underline{\mathbf{x}}', \mathbf{g}(\underline{\mathbf{x}}'))$ , and, since

$\mathbf{x}' \in \mathbf{X}_n$ ,  $\underline{\mathbf{x}}' = \mathbf{z} \circ \underline{\mathbf{x}}$  with  $\mathbf{z}$  a  $(d - m)$ -dimensional vector with components in  $\{-1, 1\}$ . Therefore,

$$\mathbf{x}' = (\mathbf{z} \circ \underline{\mathbf{x}}, \mathbf{g}(\mathbf{z} \circ \underline{\mathbf{x}})) = (\mathbf{z} \circ \underline{\mathbf{x}}, \mathbf{g}(\mathbf{z}) \circ \mathbf{g}(\underline{\mathbf{x}})) = (\mathbf{z}, \mathbf{g}(\mathbf{z})) \circ \mathbf{x}.$$

The vector  $\mathbf{x}'_j = (\mathbf{z} \circ \underline{\mathbf{x}}_j, \mathbf{g}(\mathbf{z} \circ \underline{\mathbf{x}}_j)) = (\mathbf{z}, \mathbf{g}(\mathbf{z})) \circ \mathbf{x}_j$  also belongs to  $\mathbf{X}_n$  (since the first  $d - m$  coordinates of design points in  $\mathbf{X}_n$  form a  $2^{d-m}$  factorial design), and satisfies  $d_H(\mathbf{x}', \mathbf{x}'_j) = d_H(\mathbf{x}, \mathbf{x}_j)$ .

To conclude the proof that all design points have the same neighbouring pattern, we only need to show that if  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are two distinct points in  $\mathbf{X}_n$ , say with  $d_H(\mathbf{x}, \mathbf{x}_i) = d_H(\mathbf{x}, \mathbf{x}_j) = k$ , then  $\mathbf{x}'_i = (\mathbf{z} \circ \underline{\mathbf{x}}_i, \mathbf{g}(\mathbf{z} \circ \underline{\mathbf{x}}_i))$  and  $\mathbf{x}'_j = (\mathbf{z} \circ \underline{\mathbf{x}}_j, \mathbf{g}(\mathbf{z} \circ \underline{\mathbf{x}}_j))$  are distinct points in  $\mathbf{X}_n$  satisfying  $d_H(\mathbf{x}', \mathbf{x}'_i) = d_H(\mathbf{x}', \mathbf{x}'_j) = k$ . The equality between distances has already been proved; the points are distinct since  $\mathbf{x}'_i = (\mathbf{z}, \mathbf{g}(\mathbf{z})) \circ \mathbf{x}_i \neq (\mathbf{z}, \mathbf{g}(\mathbf{z})) \circ \mathbf{x}_j = \mathbf{x}'_j$ .

Denote  $I'_k(\mathbf{x}_i) = |\{j : d_H(\mathbf{x}_i, \mathbf{x}_j) = k, \mathbf{x}_j \in \{-1, 1\}^d\}|$ . Since  $\bar{I}_k(\mathbf{x}_i; \mathbf{X}_n) = I'_k(\mathbf{x}_i) - I_k(\mathbf{x}_i; \mathbf{X}_n)$  and  $I'_k(\mathbf{x}_i) = I'_k(\mathbf{x}_j)$  for any  $\mathbf{x}_i$  and  $\mathbf{x}_j$  in  $\mathbf{X}_n$ , all design points have also the same distant-sites pattern.  $\square$

This property explains why division by  $n$  in the definition (1.9) of distance distribution yields integer values for the  $B_k(\mathbf{X}_n)$ : we have  $\mathcal{L}_k(\mathbf{x}_i; \mathbf{X}_n) = \mathcal{B}(\mathbf{X}_n)$  for any  $2^{d-m}$  fractional factorial design  $\mathbf{X}_n$  and any  $\mathbf{x}_i \in \mathbf{X}_n$ . A straightforward consequence is we do not need to consider all pairs of points in  $\mathbf{X}_n$  to construct the distance distribution, but only the distances between one point and the  $n - 1$  others. In particular, this point can be taken as  $\mathbf{1}_d$ , the  $d$ -dimensional vector with all components equal to 1 (provided that the design is constructed with principal generators with non negative signs, which we assume throughout the paper). As an illustration, below we consider the distance distribution of fractional factorial designs with  $n = d + 1$ , see Sect. 1.3.3, which is very peculiar.

**Proposition 3** *Saturated  $2^{d-m}$  fractional factorial designs ( $n = d + 1$ ) are maximin-distance optimal; their distance distribution satisfies  $B_0(\mathbf{X}_n) = 1$ ,  $B_{(d+1)/2}(\mathbf{X}_n) = n - 1$  and  $B_i(\mathbf{X}_n) = 0$  for  $i > 0, i \neq (d + 1)/2$ .*

*Proof.* From Prop. 2, we only need to consider the distance between one particular point, which we denote  $\mathbf{x} = (\underline{\mathbf{x}}, \mathbf{g}(\underline{\mathbf{x}}))$ , and other points  $\mathbf{x}' \in \mathbf{X}_n$ ,  $\mathbf{x}' = (\underline{\mathbf{x}}', \mathbf{g}(\underline{\mathbf{x}}'))$ . We show that  $d_H(\mathbf{x}, \mathbf{x}') = 2^{d-m-1} = n/2 = (d + 1)/2$  when  $d_H(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = 1, 2, \dots, m$ .

Suppose that  $d_H(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = 1$ , let  $a$  be the basic factor that changes between  $\underline{\mathbf{x}}$  and  $\underline{\mathbf{x}}'$ . The number of generators that contain  $a$  is

$$n_a = \sum_{k=1}^{d-m-1} \binom{d-m-1}{k} = 2^{d-m-1} - 1,$$

since there remains  $d - m - 1$  factors available and each defining relation contains at least two factors. It gives  $d_H(\mathbf{g}(\underline{\mathbf{x}}), \mathbf{g}(\underline{\mathbf{x}}')) = 2^{d-m-1} - 1$ , and thus  $d_H(\mathbf{x}, \mathbf{x}') = 2^{d-m-1}$ .

Suppose now that  $d_H(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = 2$ , with  $a$  and  $b$  the modified factors. The number of generators containing  $a$  and not containing  $b$  is

$$n_{a\bar{b}} = \sum_{k=1}^{d-m-2} \binom{d-m-2}{k} = 2^{d-m-2} - 1,$$

since now there only remains  $d-m-2$  factors available. We also need to count generators that contain  $b$  and not  $a$ , which gives  $d_H(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = 2 + 2(2^{d-m-2} - 1) = 2^{d-m-1}$ .

The same calculation can be repeated when  $d_H(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = p$ , with factors  $a_1, \dots, a_p$  being modified, for any  $p \leq d-m$ . Suppose first that  $p$  is odd. There are  $2^{d-m-p} - 1$  generators with  $a_1$  alone (without  $a_2, \dots, a_p$ ),  $2^{d-m-p}$  with  $a_1 a_2 a_3$  alone (without  $a_4, \dots, a_p$ ), etc., and  $2^{d-m-p}$  with all the  $a_i$ ,  $i = 1, \dots, p$ . It gives

$$\begin{aligned} d_H(\underline{\mathbf{x}}, \underline{\mathbf{x}}') &= p + p(2^{d-m-p} - 1) + \binom{p}{3} 2^{d-m-p} + \binom{p}{5} 2^{d-m-p} + \dots + 2^{d-m-p} \\ &= \left[ \binom{p}{1} + \binom{p}{3} + \binom{p}{5} + \dots + \binom{p}{p} \right] 2^{d-m-p} = 2^{p-1} 2^{d-m-p} = 2^{d-m-1}. \end{aligned}$$

Suppose now that  $p$  is even. Similar calculation gives

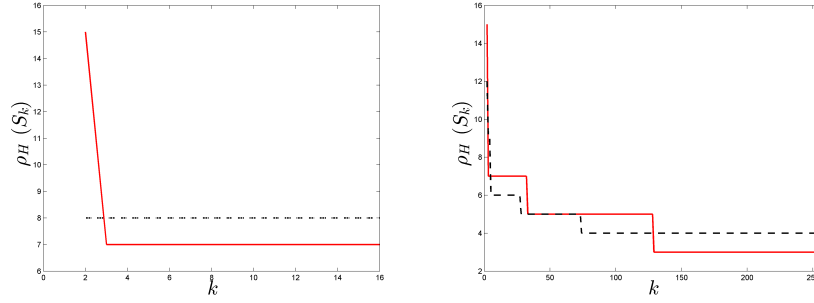
$$\begin{aligned} d_H(\underline{\mathbf{x}}, \underline{\mathbf{x}}') &= \left[ \binom{p}{1} + \binom{p}{3} + \binom{p}{5} + \dots + \binom{p}{p-1} \right] 2^{d-m-p} \\ &= 2^{p-1} 2^{d-m-p} = 2^{d-m-1}. \end{aligned}$$

Therefore,  $d_H(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = 2^{d-m-1} = n/2 = (d+1)/2$  for any  $\underline{\mathbf{x}}' \in \mathbf{X}_n$ ,  $\underline{\mathbf{x}}' \neq \underline{\mathbf{x}}$  (note that it gives equality in the upper bound (1.10)). Plotkin bound (1.8) indicates that the size  $n$  of a design  $\mathbf{X}_n$  supported on  $\{-1, 1\}^d$  and such that  $\rho_H(\mathbf{X}_n) > d/2$ , with  $d$  odd, is at least  $d+1$ , showing that saturated  $2^{d-m}$  designs are maximin-distance optimal among all designs supported on  $\{-1, 1\}^d$ . The  $n$  design points of a saturated design are vertices of a regular simplex in  $\mathcal{C}_d$  with (Euclidean) edge length  $\sqrt{2(d+1)}$  and form a maximin-optimal design in  $\mathcal{C}_d$ .  $\square$

The application of Algorithm 1 to the candidate set  $\mathcal{X}_n$  defined by a  $2^{d-m}$  design  $\mathbf{X}_n$  with  $n = d+1$ , initialised at any  $\underline{\mathbf{x}}_i \in \mathbf{X}_n$ , ensures that  $\rho_H(\mathbf{S}_k) = (d+1)/2 = 2^{d-m-1}$  for all  $k = 2, \dots, n$  (in fact, the property is true for any sequential selection of points within  $\mathbf{X}_n$ ). As the example below illustrates, the performance achieved in terms of  $\rho_H$  may be superior to those obtained when the candidate set is  $\{-1, 1\}^d$ , the set of vertices of  $\mathcal{C}_d$  (i.e., the full factorial  $2^d$  design). Note that we have  $\rho_H(\mathbf{S}_k) = \mathbf{CR}_H(\mathbf{S}_{k-1})$ ,  $2 \leq k \leq |\mathbf{X}_n|$ , when applying Algorithm 1 to a candidate set  $\mathcal{X}_n \subseteq \{-1, 1\}^d$ ; see the proof of Th. 1.

*Example 1.* The left panel of Fig. 1.1 presents the evolution of  $\rho_H(\mathbf{S}_k)$  for  $d = 15$  and  $n = 16 = 2^4$  when the candidate set in Algorithm 1 is the full  $2^d$  factorial design (red solid line) and a  $2^{d-m}$  fractional factorial design with  $m = 11$  (black dashed line). In the second case,  $\rho_H(\mathbf{S}_k) = (d + 1)/2$  for all  $k = 2, \dots, 16$ , which for  $k \geq 3$  is larger than the value obtained in the first case where all the  $2^{15} = 32\,768$  vertices are used as candidates.

The fact that restricting the set of candidate points to a subset of  $\{-1, 1\}^d$  may be beneficial is further illustrated on the right panel of Fig. 1.1. There, the red solid line corresponds again to the candidate set given by the full  $2^d$  factorial design, whereas a  $2^{d-m}$  design with  $m = 7$  and  $\rho_H(\mathbf{X}_n) = 4$  is used for the black dashed-line curve ( $n = 2^8 = 256$ ).



**Fig. 1.1** Evolution of  $\rho_H(\mathbf{S}_k)$  when Algorithm 1 is applied to the candidate set  $\mathcal{X}_n$  given by the  $2^d$  full factorial design (red solid line) and when  $\mathcal{X}_n = \mathbf{X}_n$  is a  $2^{d-m}$  fractional factorial design (black dashed line);  $d = 15$ , the algorithm is initialised at a  $\mathbf{x}_i \in \mathcal{X}_n$ . Left:  $m = 11, n = d + 1 = 16, k = 2, \dots, 16$ . Right:  $m = 7, n = 256, k = 2, \dots, 250$ .

In the Appendix, we give conditions on the choice of generators that provide guarantees on the minimum Hamming distance  $k$  of a fractional factorial design  $\mathbf{X}_n$ , i.e.,  $\rho_H(\mathbf{X}_n) \geq k$ , for  $k = 2, 3$  and  $4$ . However, the derivation of such conditions gets cumbersome when  $k \geq 5$ , and in the next section we present an algorithm for the optimal selection of  $m$  generators among all  $2^{d-m} - (d - m) - 1$  possible generators having length at least 2.

### 1.5.2 Optimal selection of generators by simulated annealing

**SA algorithm for the maximisation of  $\rho_H$**

- 0) Construct the set  $\mathcal{G}$  all  $2^{d-m} - (d-m) - 1$  generators (of length  $\geq 2$ ), choose an initial set  $G$  of  $m$  distinct generators in  $\mathcal{G}$ , construct the corresponding design  $\mathbf{X}_n$  and its distance distribution  $\mathcal{B}(\mathbf{X}_n)$ ; set  $\mathbf{X}_n^* = \mathbf{X}_n$  and  $k = 1$ .
- 1) Select a random generator  $g$  within  $G$  and a random generator  $g'$  within  $\mathcal{G} \setminus G$ . Construct  $G' = G \setminus \{g\} \cup \{g'\}$  and its associated design  $\mathbf{X}'_n$  and distance distribution  $\mathcal{B}(\mathbf{X}'_n)$ .
- 2) Compute  $i^* = \min\{i : B_i(\mathbf{X}'_n) \neq B_i(\mathbf{X}_n^*)\}$  and  $\delta^* = B_{i^*}(\mathbf{X}'_n) - B_{i^*}(\mathbf{X}_n^*)$ ; if  $\delta^* \leq 0$ , set  $\mathbf{X}_n^* = \mathbf{X}'_n$ .
- 3) Compute  $i^+ = \min\{i : B_i(\mathbf{X}'_n) \neq B_i(\mathbf{X}_n)\}$ , and  $P = \min\left\{\exp\left[-\frac{B_{i^+}(\mathbf{X}'_n) - B_{i^+}(\mathbf{X}_n)}{T_k}\right], 1\right\}$ ; accept the move  $\mathbf{X}_n = \mathbf{X}'_n$  and  $G = G'$  with probability  $P$ .
- 4) if  $k = K$ , stop; otherwise  $k \leftarrow k + 1$ , return to 1.

A logarithmic decrease of the “temperature”  $T_k$ , such as  $T_k = 2^{d-m} / \log(k)$  ensures global asymptotic convergence to a maximin-distance optimal fractional factorial design. In practice, we wish to stop the algorithm after a number of iterations  $K$  which is not excessively large. Numerical experimentation indicates that a faster decrease of  $T_k$ , yielding a behaviour close to that of a simple descent method, is often suitable. For instance, we take  $T_k = 2^{d-m} / k^{4/5}$  in Example 2 (Sect. 1.7).

**Remark 1** *In the applications we have in mind,  $n = 2^{d-m}$  is relatively small even when  $d$  is large (it must nevertheless satisfy the bound  $n \geq d + 1$  of Prop. 1), and the set  $\mathcal{G}$  remains of reasonable size. A similar simulated annealing algorithm can be used for the construction of minimum aberration designs. However, the construction of the word length pattern for a  $2^{d-m}$  design requires the calculation of  $2^m - 1$  defining relations, which becomes prohibitively computationally demanding when  $m$  is large (a consequence of  $n$  being reasonably small). For instance, for  $d = 50$  and  $m = 35$  (the situation in Example 2 below), we have  $2^{d-m} - (d - m) - 1 = 37\,752$  whereas  $2^m - 1 > 3.43 \times 10^{10}$ . Computation of the word length pattern from the distance distribution using (1.11) may then be advantageous.*

## 1.6 Covering properties of two-level factorial designs

In this section, we investigate the covering properties of fractional factorial designs measured by the Hamming covering radius  $\text{CR}_H$  defined by (1.5). One should notice that the best design in terms of maximin distance is not

always the best one in terms of covering radius, compare for instance  $\mathbf{X}_{64}^{(a)}$  and  $\mathbf{X}_{64}^{(d)}$  in Table 1.4.

### 1.6.1 Bounds on $\text{CR}_H(\mathbf{X}_n)$

We already know that  $\text{CR}_H(\mathbf{X}_n) \geq \lfloor \rho_H(\mathbf{X}_n)/2 \rfloor$  for any design on  $\{-1, 1\}^d$ ; see Sect. 1.4.1.  $\text{CR}_H(\mathbf{X}_n)$  also satisfies the following property.

**Proposition 4** *The Hamming covering radius  $\text{CR}_H(\mathbf{X}_n)$  of a  $2^{d-m}$  fractional factorial design  $\mathbf{X}_n$  satisfies*

$$\text{CR}_H(\mathbf{X}_n[d-m+1:d]) \leq \text{CR}_H(\mathbf{X}_n) \leq m,$$

with  $\mathbf{X}_n[d-m+1:d]$  denoting the projection of  $\mathbf{X}_n$  on the  $m$  dimensional space defined by the non-basic factors (i.e., those constructed through generators).

*Proof.* We use the same notation as in Sect. 1.5.1 and, for  $\mathbf{x}_i \in \mathbf{X}_n$ , we denote  $\mathbf{x}_i = (\underline{\mathbf{x}}_i, \mathbf{g}(\underline{\mathbf{x}}_i))$ ; also, we split any  $\mathbf{x} \in \mathbb{R}^d$  into  $\mathbf{x} = (\underline{\mathbf{x}}, \bar{\mathbf{x}})$ , with  $\underline{\mathbf{x}} \in \mathbb{R}^{d-m}$  and  $\bar{\mathbf{x}} \in \mathbb{R}^m$ . For any  $\mathbf{x} \in \{-1, 1\}^d$ , there exists  $j = j(\mathbf{x})$  such that  $\underline{\mathbf{x}}_j = \underline{\mathbf{x}}$ , so that

$$\begin{aligned} \min_{\mathbf{x}_i \in \mathbf{X}_n} d_H(\mathbf{x}, \mathbf{x}_i) &\leq d_H(\mathbf{x}, \mathbf{x}_j) = d_H(\underline{\mathbf{x}}, \underline{\mathbf{x}}_j) + d_H(\bar{\mathbf{x}}, \mathbf{g}(\underline{\mathbf{x}}_j)) \\ &= d_H(\bar{\mathbf{x}}, \mathbf{g}(\underline{\mathbf{x}}_j)) = d_H(\bar{\mathbf{x}}, \mathbf{g}(\underline{\mathbf{x}})). \end{aligned}$$

We also have  $d_H(\mathbf{x}, \mathbf{x}_i) = d_H(\underline{\mathbf{x}}, \underline{\mathbf{x}}_i) + d_H(\bar{\mathbf{x}}, \mathbf{g}(\underline{\mathbf{x}}_i)) \geq d_H(\bar{\mathbf{x}}, \mathbf{g}(\underline{\mathbf{x}}_i))$ . Therefore,

$$\begin{aligned} \max_{\mathbf{x} \in \{-1, 1\}^d} \min_{\mathbf{x}_i \in \mathbf{X}_n} d_H(\bar{\mathbf{x}}, \mathbf{g}(\underline{\mathbf{x}}_i)) &\leq \text{CR}_H(\mathbf{X}_n) = \max_{\mathbf{x} \in \{-1, 1\}^d} \min_{\mathbf{x}_i \in \mathbf{X}_n} d_H(\mathbf{x}, \mathbf{x}_i) \\ &\leq \max_{\mathbf{x} \in \{-1, 1\}^d} d_H(\bar{\mathbf{x}}, \mathbf{g}(\underline{\mathbf{x}})). \end{aligned}$$

Finally,  $\max_{\mathbf{x} \in \{-1, 1\}^d} \min_{\mathbf{x}_i \in \mathbf{X}_n} d_H(\bar{\mathbf{x}}, \mathbf{g}(\underline{\mathbf{x}}_i)) = \text{CR}_H(\mathbf{X}_n[d-m+1:d])$  and  $\max_{\mathbf{x} \in \{-1, 1\}^d} d_H(\bar{\mathbf{x}}, \mathbf{g}(\underline{\mathbf{x}})) = m$  conclude the proof.  $\square$

### 1.6.2 Calculation of $\text{CR}_H(\mathbf{X}_n)$

The direct calculation of  $\text{CR}_H(\mathbf{X}_n)$  through  $\max_{\mathbf{x} \in \{-1, 1\}^d} \min_{\mathbf{x}_i \in \mathbf{X}_n} d_H(\mathbf{x}, \mathbf{x}_i)$  is unfeasible for large  $d$ . Exploiting Prop. 2, a possible alternative consists in restricting the set of candidates to the  $\binom{d}{\delta}$  points at a given Hamming distance  $\delta$  from an arbitrary point  $\mathbf{x}_1$  in  $\mathbf{X}_n$ , for an increasing sequence  $\delta_i$ , initialised at  $\delta_1 = \lfloor \rho_H(\mathbf{X}_n)/2 \rfloor$ . We thus compute  $C(\mathbf{X}_n, \mathbf{x}_1, \delta) = \max_{\mathbf{x}: d_H(\mathbf{x}, \mathbf{x}_1) = \delta} \min_{\mathbf{x}_i \in \mathbf{X}_n} d_H(\mathbf{x}, \mathbf{x}_i)$  for  $\delta = \delta_1, \delta_1 + 1 \dots$  and stop at the first

$\delta_i$  when  $C(\mathbf{X}_n, \mathbf{x}_1, \delta)$  starts decreasing; the value  $\delta_{i-1}$  equals  $\text{CR}_H(\mathbf{X}_n)$ . Although it requires less computations than direct calculation, this approach is still too costly for large  $d$  unless  $\text{CR}_H(\mathbf{X}_n)$  is very small (meaning that  $n$  is very large) or very large (meaning that  $n$  is very small). In Example 2 considered below, with  $d = 50$  and  $n = 32\,768$ , we have  $\text{CR}_H(\mathbf{X}_n) = 13$  and the construction is unapplicable. Therefore, hereafter we present a simple local ascent algorithm for searching a distant point from  $\mathbf{X}_n$ , which we initialise at a design point.

The construction relies on the search of a point in  $\{-1, 1\}^d$  at maximum Hamming distance from  $\mathbf{X}_n$ . From Prop. 2, we only need to consider moves from an arbitrary point of  $\mathbf{X}_n$ . The order of inspection of the  $d$  factors in the `for` loop of Step 1 may be randomised.

### Algorithmic construction of a lower bound on $\text{CR}_H(\mathbf{X}_n)$

- 0) Set  $\mathbf{x} = \mathbf{x}_1 \in \mathbf{X}_n$ ,  $\Delta = 0$ , `continue` = 1
- 1) while `continue` = 1
  - Try successively all points at Hamming distance 1 from  $\mathbf{x}$ :
  - `for`  $i = 1, \dots, d$ 
    - set  $\{\mathbf{x}'\}_j = \{\mathbf{x}\}_j$  for  $j \neq i$  and  $\{\mathbf{x}'\}_i = -\{\mathbf{x}\}_i$ , compute  $\Delta' = \min_{\mathbf{x}_i \in \mathbf{X}_n} d_H(\mathbf{x}', \mathbf{x}_i)$ .
    - if  $\Delta' > \Delta$ , set  $\mathbf{x} = \mathbf{x}'$ ,  $\Delta = \Delta'$  and break the `for` loop.
    - otherwise, if  $i = d$ , set `continue` = 0 (all possible moves have been unsuccessfully exhausted).
- 2) Return  $\Delta$ , which forms a lower bound on  $\text{CR}_H(\mathbf{X}_n)$ .

**Remark 2** *Convergence to a point at maximum distance from  $\mathbf{X}_n$  is not guaranteed. The algorithm can be modified to incorporate a simulated annealing scheme that accepts moves such that  $\Delta' < \Delta$  with some probability: at Step 1, in the `for` loop we then set  $\mathbf{x} = \mathbf{x}'$ ,  $\Delta = \Delta'$  with probability  $\min\{\exp[(\Delta' - \Delta)/T_k], 1\}$  for some decreasing temperature profile  $T_k$ , do not break the loop and never set `continue` = 0; the algorithm is stopped when the number of iterations reaches a predefined bound.*

## 1.7 Greedy constructions based on fractional factorial designs

### 1.7.1 Base designs

We first consider a specialisation of the  $n$  first iterations of Algorithm 1 to the case where  $\mathcal{X} = \mathcal{C}_d$  and the candidate set at Step 1 is a  $2^{d-m}$  fractional factorial design  $\mathbf{X}_n$ .

#### Algorithm 2

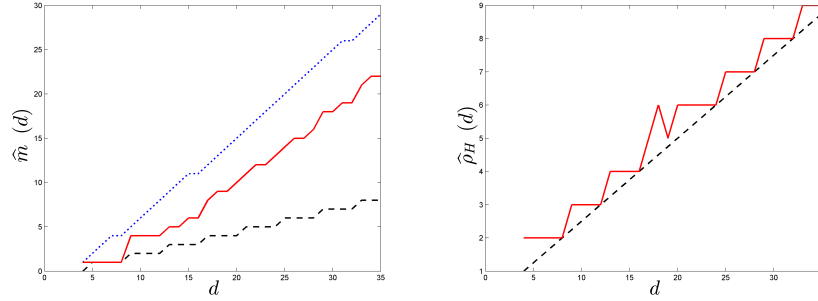
- 0) Construct a  $2^{d-m}$  fractional factorial design  $\mathbf{X}_n$ ; set  $\mathbf{S}_1 = \{\mathbf{0}\}$  and  $k = 1$ .  
 1) **for**  $k = 1, \dots, n$  **do**  
     find  $\mathbf{x}^* = \arg \max_{\mathbf{x} \in \mathbf{X}_n} d(\mathbf{x}, \mathbf{S}_k)$ , set  $\mathbf{S}_{k+1} = \mathbf{S}_k \cup \{\mathbf{x}^*\}$ .

Note that the distances  $d(\mathbf{x}, \mathbf{S}_k)$ ,  $\mathbf{x} \in \mathbf{X}_n$ , can be computed recursively as  $d(\mathbf{x}, \mathbf{S}_k) = \min\{d(\mathbf{x}, \mathbf{S}_{k-1}), \|\mathbf{x} - \mathbf{x}_k\|\}$  and that the generation of  $\mathbf{S}_k$  for  $k \leq n + 1$  has complexity  $\mathcal{O}(knd)$ . Algorithm 2 also satisfies the following property.

**Proposition 5** *If  $\rho_H(\mathbf{X}_n) \geq d/4$ , then, for  $k \leq n + 1$ , the design  $\mathbf{S}_k$  constructed by Algorithm 2 could also have been generated by Algorithm 1 initialised at  $\mathbf{S}_1 = \{\mathbf{0}\}$  and it satisfies the bounds of Th. 1.*

*Proof.* Since  $\rho_H(\mathbf{X}_n) \geq d/4$ , for any  $k \leq n$  and any  $\mathbf{x}_i \in \mathbf{X}_n \setminus \mathbf{S}_{2:k}$  we have  $d_H(\mathbf{x}_i, \mathbf{S}_{2:k}) \geq d/4$ ; that is,  $d(\mathbf{x}_i, \mathbf{S}_{2:k}) \geq \sqrt{d}$ . Therefore,  $d(\mathbf{x}_i, \mathbf{S}_k) = \sqrt{d} = \max_{\mathbf{x} \in \mathcal{C}_d} d(\mathbf{x}, \mathbf{S}_k)$ . The restriction to the set  $\mathbf{X}_n$  at Step 1 thus entails no loss of performance and Th. 1 applies.  $\square$

Proposition 3 shows that minimum-size  $2^{d-m}$  designs  $\mathbf{X}_n$  with  $n = d + 1$  satisfy  $\rho_H(\mathbf{X}_n) = (d + 1)/2$ . We shall not provide an explicit construction ensuring the existence of fractional factorial designs satisfying  $\rho_H(\mathbf{X}_n) \geq d/4$  for all values of  $d$  (see, however, Remark 3). Instead, for each  $d = 4, \dots, 35$ , using the algorithm of Sect. 1.5.2 we have searched the smallest  $m = \hat{m}(d)$  (that is, the largest possible design size  $2^{d-\hat{m}(d)}$ ) for which we can find a design  $\mathbf{X}_n$  with minimum Hamming distance at least  $d/4$ . We denote by  $\hat{\rho}_H(d) \geq d/4$  the distance we have obtained. The left panel of Fig. 1.2 shows  $\hat{m}(d)$  (red solid line), together with the upper bound  $m^*(d)$  given by (1.4) (blue dotted line) and the lower bound  $m_*(d) = \lceil d/4 - 1 \rceil$  implied by the Singleton bound (1.7) (black dashed line); the right panel presents  $\hat{\rho}_H(d)$ . For instance, for  $d = 35$ , we can construct a design with  $n = 2^{35-22} = 8192$  points and minimum Hamming distance 9. A construction with  $d = 50$  and  $m = 35$  will be considered in Example 2. Note that the value  $k_*(d)$  of Sect. 1.2 satisfies



**Fig. 1.2** Left:  $\widehat{m}(d)$  (red solid line),  $m^*(d) = \lfloor d - \log_2(d+1) \rfloor$  (blue dotted line),  $m_*(d) = \lfloor d/4 - 1 \rfloor$  (black dashed line). Right:  $\widehat{\rho}_H(d)$  (red solid line), the black dashed line corresponds to  $d/4$ .

$$k_*(d) \geq 2^{d-\widehat{m}(d)}. \quad (1.12)$$

**Remark 3** Let  $d_*$  be a dimension satisfying  $d_* + 1 = 2^{d_*-m}$ ; see Table 1.3. The corresponding minimum-size  $2^{d_*-m}$  fractional factorial design  $\mathbf{X}_n$  satisfies  $\rho_H(\mathbf{X}_n) = (d_* + 1)/2$ . By removing any  $d_* - d$  factors from  $\mathbf{X}_n$ , with  $d \geq \lceil 2(d_* - 1)/3 \rceil$ , we obtain a design  $\mathbf{X}'_n$  in  $[-1, 1]^d$  such that  $\rho_H(\mathbf{X}'_n) \geq \rho_H(\mathbf{X}_n) - (d_* - d) \geq d/4$ . However, these designs have too few points to be of practical interest for computer experiments.

For all  $k \leq n - 1$  the choice of  $\mathbf{x}^*$  at Step 1 of Algorithm 2 is arbitrary; in particular, if this choice is randomised, an unlucky selection may thus yield  $\rho_H(\mathbf{S}_{2:k}) = \rho_H(\mathbf{X}_n)$  for all  $k = 3, \dots, n + 1$ . This weakness can be overcome through a slight modification of Step 1, yielding the following algorithm.

### Algorithm 3

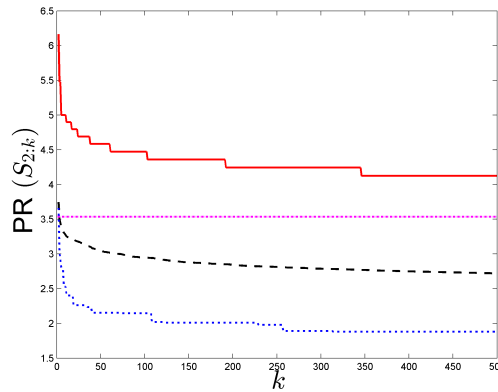
- 0) Construct a  $2^{d-m}$  fractional factorial design  $\mathbf{X}_n$  with  $\rho_H(\mathbf{X}_n) \geq d/4$ ; set  $\mathbf{S}_2 = \{\mathbf{0}, \mathbf{x}_2\}$  and  $k = 2$ , with  $\mathbf{x}_2$  an arbitrary point in  $\mathbf{X}_n$ .
- 1)  $\text{for } k = 2, \dots, n$  **do**  
 find  $\mathbf{x}^* = \arg \max_{\mathbf{x} \in \mathbf{X}_n} d(\mathbf{x}, \mathbf{S}_{2:k})$ , set  $\mathbf{S}_{k+1} = \mathbf{S}_k \cup \{\mathbf{x}^*\}$ .

Note that all  $\mathbf{x}_i$  in  $\mathbf{X}_n$  satisfy  $d(\mathbf{x}_i, \mathbf{S}_1) = \sqrt{d} = \arg \max_{\mathbf{x} \in \mathbf{X}_n} d(\mathbf{x}, \mathbf{S}_1)$  and have the same neighbouring pattern, see Sect. 1.2 and Prop. 2.

*Example 2.* For  $d = 50$  and  $m = 35$ , the algorithm of Sect. 1.5.2 yields a design  $\mathbf{X}_n$  of  $n = 2^{15} = 32\,768$  points, with resolution IV ( $A_4(\mathbf{X}_n) = 2$ ),  $\rho_H(\mathbf{X}_n) = 13 > d/4$  ( $B_{13}(\mathbf{X}_n) = 2$ ) and the algorithm of Sect. 1.6.2 gives  $\text{CR}_H(\mathbf{X}_n) \geq 13$ . Algorithm 3 generates a sequence of nested designs  $\mathbf{S}_k$  that satisfy the efficiency bounds (1.3) for all  $k \leq n + 1 = 32\,769$ . The construction

is very fast since there are only  $n = 32\,768$  points in  $\mathcal{X}_n = \mathbf{X}_n$  to be considered at Step 1 of Algorithm 3 (to be compared with the  $2^d > 1.1258 \times 10^{15}$  vertices of  $\mathcal{C}_d$ ). Figure 1.3 presents the evolution of the packing radius  $\text{PR}(\mathbf{S}_{2:k})$  as a function of  $k$  for Algorithm 3 (red solid line), for  $k = 3, \dots, 500$  (the value 500 is rather arbitrarily, chosen in agreement with the “10  $d$ ” rule of [17]). When including  $\mathbf{x}_1 = \mathbf{0}$ , it satisfies  $\text{PR}(\mathbf{S}_k) = \sqrt{d}/2 \simeq 3.5355$  for  $k \leq 32\,769$ . The curve in black dashed line (middle) is obtained when Algorithm 1 is applied to the candidate set  $\mathcal{X}_n$  given by the first  $2^{19}$  points of Sobol’ sequence; the blue dotted line (bottom) corresponds to designs given by the first  $k$  points of this Sobol’ sequence.

**Fig. 1.3** Evolution of  $\text{PR}(\mathbf{S}_{2:k})$  in Algorithm 3 with  $\mathcal{X}_n$  given by a  $2^{50-35}$  design (red solid line, top), of  $\text{PR}(\mathbf{S}_k)$  in Algorithm 1 with  $\mathcal{X}_n$  given by the first  $n = 2^{19}$  points of Sobol’ sequence  $(\mathbf{X}_i^S)_i$  in  $\mathcal{C}_d$  (black dashed line, middle), and of  $\text{PR}(\mathbf{X}_i^S)$ ,  $i = 2, \dots, 500$  (blue dotted line, bottom). The horizontal line indicates the value  $\sqrt{d}/2$  ( $d = 50$ ).



The complexity of Algorithm 3 is only linear in  $k$  and grows like  $\mathcal{O}(knd)$ . If necessary, it can be further reduced for large  $n$  by first constructing nested half-designs from  $\mathbf{X}_n$ , following ideas similar to those in [1]. Theorem 2 in their paper shows that only  $2^{d-m} - 1$  different half-designs need to be considered when starting from an arbitrary  $2^{d-m}$  fractional factorial design (note that here those half-designs must be compared in terms of their  $\rho_H$  values, whereas aberration is used in [1]).

### 1.7.2 Rescaled designs

A fractional factorial design  $\mathbf{X}_n$  has all its points on the vertices of  $\mathcal{C}_d$ , which is advantageous in terms of packing radius in the full dimensional space. However, performance in terms of prediction/interpolation of an unknown function by a non-parametric model (in particular with kriging) is more related to  $\text{CR}(\mathbf{X}_n)$  [12], and it may then be beneficial to have design points inside  $\mathcal{C}_d$ . In [1], the iterative decomposition of  $\mathbf{X}_n$  into half-designs is used

to construct multi-layer designs having two points on each layer. Here, since the design points in  $\mathbf{S}_k$  constructed with Algorithm 3 are ordered, for a fixed  $K \leq n$  we can directly apply a scaling procedure to  $\mathbf{S}_K$  to obtain a design  $\tilde{\mathbf{S}}_K$  with points inside  $\mathcal{C}_d$ .

First, following [1], and to avoid having pairs of points too close together, we impose that  $\tilde{\mathbf{S}}_K$  has no point in an hypercube  $[-a, a]^d$ ,  $0 < a < 1$ , except the centre  $\mathbf{x}_1 = \mathbf{0}$ . We choose  $a$  by setting the ratio  $r$  of the volume of the neglected empty hypercube to the volume of  $\mathcal{C}_d$ ; that is,  $a = a_r = r^{1/d}$ .

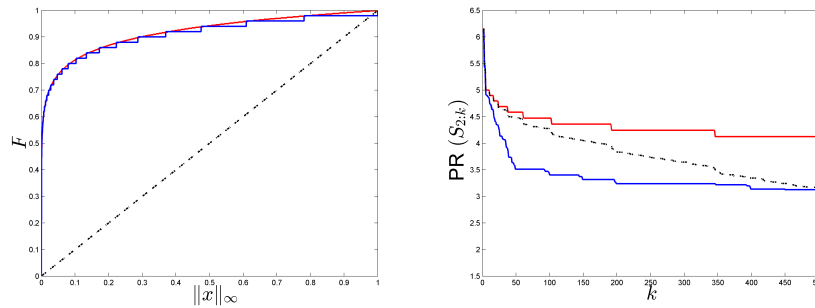
Next, we need to choose how we rescale the points  $\mathbf{x}_i$  of  $\mathbf{S}_K$ , for  $i = 2, \dots, K + 1$ , in order to obtain a suitable distribution of distances to the centre for the  $\ell_\infty$  norm. We shall denote by  $\tilde{\mathbf{x}}_i = \beta_{i,r,K,\gamma} \mathbf{x}_i$  the rescaled design points,  $i = 1, \dots, K + 1$ , with  $\beta_{1,r,K,\gamma} = 0$  and

$$\beta_{i,r,K,\gamma} = \left[ 1 - \frac{(i-2)(1-a_r^\gamma)}{K-1} \right]^{1/\gamma}, \quad i = 2, \dots, K + 1. \quad (1.13)$$

Here  $\gamma$  is a scalar in  $[1, d]$ : linear scaling with  $\gamma = 1$  yields a design with points more densely distributed close to the centre  $\mathbf{0}$  than near the boundary of  $\mathcal{C}_d$ ; when  $\gamma = d$ , the empirical distribution of the  $\|\tilde{\mathbf{x}}_i\|_\infty$  converges to the uniform distribution on  $[0, 1]$ , obtained for points  $\mathbf{x}$  uniformly distributed in  $\mathcal{C}_d$ , as  $K$  tend to infinity. Values of  $\gamma$  between 1 and  $d$  provide behaviours between these two extreme cases. When  $N$  points are needed, with  $K + 1 < N \leq n + 1$ , the rescaling procedure can be applied periodically, using

$$\beta_{i,r,K,\gamma} = \left\{ 1 - \frac{[(i-2) \bmod K](1-a_r^\gamma)}{K-1} \right\}^{1/\gamma}. \quad (1.14)$$

*Example 2 (continued).* The left panel of Fig. 1.4 shows the empirical cumulative distribution function (cdf) of the  $\|\tilde{\mathbf{x}}_i\|_\infty$  for the design obtained with Algorithm 3 for  $\gamma = 1$ ,  $r = 10^{-6}$  and  $K = 500$  (red solid line). When  $\gamma = d$ , the empirical cdf is visually confounded with the dashed-line diagonal, which corresponds to points uniformly distributed in  $\mathcal{C}_d$ . The black dashed line (middle) on the right panel of Fig. 1.4 presents the evolution of  $\text{PR}(\tilde{\mathbf{S}}_{2:k})$  after linear rescaling of the designs  $\mathbf{S}_k$  obtained by Algorithm 3; the top curve (red solid line) is identical to that on Fig. 1.3 and correspond to  $\mathbf{S}_{2:k}$ . Periodic rescaling of  $\mathbf{S}_k$  with  $\gamma = 1$ ,  $r = 10^{-6}$  and  $K = 50$  in (1.14) yields the two curves in blue solid line, for the cdf (left) and for the evolution of  $\text{PR}(\tilde{\mathbf{S}}_{2:k})$  (right). With an horizon  $N = 500$  and  $K = 50$ , there are 10-uples of points with the same  $\ell_\infty$  norm, which explains the stair-case shape of the cdf observed on the left panel. The faster decrease of the scaling factor yields a faster decrease of the packing radius on the right panel.



**Fig. 1.4** Linear rescaling with  $\gamma = 1$  and  $r = 10^{-6}$  in (1.13) when the  $\mathbf{x}_i$  are generated with Algorithm 3 in Example 2. Left: empirical cdf of the  $\|\tilde{\mathbf{x}}_i\|_\infty$  for  $K = 500$  (top solid line in red) and  $K = 50$  (blue stair-case solid line); when  $\gamma = d$  and  $K = 500$  the cdf is confounded with the dashed-line diagonal. Right: same as on Fig. 1.3 for the red solid line (top), evolution of  $\text{PR}(\tilde{\mathbf{S}}_{2;k})$  after linear rescaling of designs  $\mathbf{S}_k$  given by Algorithm 3 for  $K = 500$  (black dashed line) and  $K = 50$  (blue solid line, bottom).

### 1.7.3 Projection properties

Space-filling designs in the full  $d$ -dimensional space do not necessarily have good properties when projected on an axis-aligned sub-space with dimension  $d' < d$ . In this section, we compare the projection properties of designs  $\mathbf{S}_k$  generated with Algorithm 3 with those of Sobol' sequences and Latin hypercube designs, for a fixed  $k$ .

Sobol' sequence is a particular low discrepancy  $(t, s)$  sequence, see, e.g., [22, Chap. 4], which permits the fast generation of designs  $\mathbf{X}_k^S$  having good space-filling properties when  $k$  is a power of 2, also for large  $d$ . A Latin hypercube (Lh) design  $\mathbf{X}_k^{Lh}$  with  $k$  points in  $[-1, 1]^d$  has the  $k$  levels  $2i/(k-1) - 1$ ,  $i = 0, \dots, k-1$ , for each of the  $d$  factors, but this does not ensure good space-filling properties in the full  $d$ -dimensional space. In [19], maximin-distance optimal Lh designs are constructed by simulated annealing. A different space-filling criterion is considered in [13], whose optimisation yields so-called maximum projection designs. In the continuation of Example 2 (with  $d = 50$ ) presented below, we use the ESE algorithm of [10] to construct a maximin-distance optimal Lh design.

For any  $d' \in \{1, \dots, d\}$  and any  $r \in \{1, \dots, \binom{d}{d'}\}$ , let  $P_{d',r}$  denote one of the  $\binom{d}{d'}$  distinct projections on an axis-aligned  $d'$  dimensional sub-space. For any  $k$ -point design  $\mathbf{X}_k = \{\mathbf{x}_1, \dots, \mathbf{x}_k\}$  we then denote by  $P_{d',r}(\mathbf{X}_k)$  the corresponding design for the  $d'$  factors associated with  $P_{d',r}$ ; that is,  $P_{d',r}(\mathbf{X}_k) = \{P_{d',r}(\mathbf{x}_1), \dots, P_{d',r}(\mathbf{x}_k)\}$ , and consider the following criteria:

$$\text{CR}_{d'}(\mathbf{X}_k) = \max_{r=1,\dots,(d')} \max_{\mathbf{x} \in [-1,1]^{d'}} d(\mathbf{x}, P_{d',r}(\mathbf{X}_k)), \quad (1.15)$$

$$\text{PR}_{d'}(\mathbf{X}_k) = \frac{1}{2} \min_{r=1,\dots,(d')} \min_{\mathbf{x}_i, \mathbf{x}_j \in \mathbf{X}_k, \mathbf{x}_i \neq \mathbf{x}_j} \|P_{d',r}(\mathbf{x}_i) - P_{d',r}(\mathbf{x}_j)\|. \quad (1.16)$$

When applied to designs generated by Algorithm 3, we obtain the following properties for  $\text{CR}_{d'}(\mathbf{S}_k)$  and  $\text{PR}_{d'}(\mathbf{S}_k)$ .

**Proposition 6** *Let  $\mathbf{S}_k$  be a design generated by Algorithm 3,  $2 \leq k \leq n$ . We have*

$$\text{PR}_{d'}(\mathbf{S}_k) = [\min\{\max\{d' - d + \rho_H(\mathbf{S}_{2:k}), 0\}, d'/4\}]^{1/2} \quad (1.17)$$

for any  $d' \in \{1, \dots, d\}$ ,

$$\sqrt{\frac{d'}{2} - \frac{[d' \bmod 2]}{4}} \leq \text{CR}_{d'}(\mathbf{S}_k) \leq \sqrt{d'} \quad \text{for any } d' \in \{1, \dots, d\} \quad (1.18)$$

and  $\text{CR}_{d'}(\mathbf{S}_k) = \sqrt{d'}$  for  $d' \geq \frac{4}{3}[d - \rho_H(\mathbf{S}_{2:k})]$ .

*Proof.* We first prove (1.17). For any  $d' \leq d_0 = d - \rho_H(\mathbf{S}_{2:k})$  there exists at least one subset of  $d'$  factors,  $i_1, \dots, i_{d'}$  say, and two points  $\mathbf{x}_j \neq \mathbf{x}_\ell$  in  $\mathbf{S}_{2:k}$  such that  $d_H(\{\mathbf{x}_j\}_{i_1, \dots, i_{d'}}, \{\mathbf{x}_\ell\}_{i_1, \dots, i_{d'}}) = 0$ . Therefore,  $\|P_{d',r}(\mathbf{x}_j) - P_{d',r}(\mathbf{x}_\ell)\| = 0$  for some projection  $P_{d',r}$ , and  $\text{PR}_{d'}(\mathbf{S}_{2:k}) = 0 = \text{PR}_{d'}(\mathbf{S}_k)$ . Take now  $d' = d_0 + 1$ . On the one hand, there exist two points  $\mathbf{x}_j \neq \mathbf{x}_\ell$  in  $\mathbf{S}_{2:k}$  that satisfy  $d_H(P_{d',r}(\mathbf{x}_j), P_{d',r}(\mathbf{x}_\ell)) \leq 1$ , which implies that  $\rho_H(P_{d',r}(\mathbf{S}_{2:k})) \leq 1$ . On the other hand, the existence of a projection  $P_{d',r}$  such that  $\rho_H(P_{d',r}(\mathbf{S}_{2:k})) = 0$  would imply  $\rho_H(\mathbf{S}_{2:k}) \leq d - d' = \rho_H(\mathbf{S}_{2:k}) - 1$ . Therefore,  $\min_{r=1,\dots,(d')} \rho_H(P_{d',r}(\mathbf{S}_{2:k})) = 1$ . Proceeding in the same way, by induction we get  $\min_{r=1,\dots,(d')} \rho_H(P_{d',r}(\mathbf{S}_{2:k})) = d' - d_0$  for any  $d' \in \{d_0, d_0 + 1, \dots, d\}$ . We have thus obtained  $\text{PR}_{d'}(\mathbf{S}_{2:k}) = [\max\{d' - d + \rho_H(\mathbf{S}_{2:k}), 0\}]^{1/2}$ . Since  $\mathbf{S}_k = \mathbf{S}_{2:k} \cup \{\mathbf{0}\}$ , we obtain  $\text{PR}_{d'}(\mathbf{S}_k) = \min\{\text{PR}_{d'}(\mathbf{S}_{2:k}), \sqrt{d'}/2\}$ , which gives (1.17).

Now we prove (1.18).  $\mathbf{S}_k$  contains the origin, and furthest points from the origin are vertices of the projected hypercube, therefore  $\text{CR}_{d'}(\mathbf{S}_k) \leq \sqrt{d'}$ . Since  $k \leq n$ , there exists  $\mathbf{x}_i \in \mathbf{X}_n \setminus \mathbf{S}_{2:k}$  such that  $d_H(\mathbf{x}_i, \mathbf{S}_{2:k}) \geq \rho_H(\mathbf{S}_{2:k}) \geq d/4$ . Take any  $d'$  such that  $(4/3)[d - \rho_H(\mathbf{S}_{2:k})] \leq d' \leq d$ . For any projection  $P_{d',r}$ , we have

$$\begin{aligned} d_H[P_{d',r}(\mathbf{x}_i), P_{d',r}(\mathbf{S}_{2:k})] - \frac{d'}{4} &\geq d_H(\mathbf{x}_i, \mathbf{S}_{2:k}) + d' - d - \frac{d'}{4} \\ &\geq \frac{3}{4}d' + \rho_H(\mathbf{S}_{2:k}) - d \geq 0. \end{aligned}$$

Therefore,  $\text{CR}[P_{d',r}(\mathbf{S}_{2:k})] \geq \sqrt{d'}$  and  $\text{CR}[P_{d',r}(\mathbf{S}_k)] = \text{CR}_{d'}(\mathbf{S}_k) = \sqrt{d'}$ .

Consider finally the favourable case where, for every projection  $P_{d',r}$ ,  $P_{d',r}(\mathbf{S}_k)$  contains the  $2^{d'}$  full factorial design. For  $d'$  even, with  $d' = 2p$ , the point  $\mathbf{x}$  with  $p$  coordinates at 0 and the other  $p$  at 1 satisfies  $d(P_{d',r}(\mathbf{x}))$ ,

$P_{d',r}(\mathbf{S}_k) = \sqrt{d'/2} \leq \text{CR}_{d'}(\mathbf{S}_k)$ . When  $d'$  is odd, with  $d' = 2p + 1$ , consider the point  $\mathbf{x}$  with  $p$  coordinates at 0,  $p$  at 1, and the last one equal to  $1/2$ ; we have  $d(P_{d',r}(\mathbf{x}), P_{d',r}(\mathbf{S}_k)) = \sqrt{p+1/4} = \sqrt{d'/2 - 1/4} \leq \text{CR}_{d'}(\mathbf{S}_k)$ .  $\square$

**Remark 4** *The lower bound on  $\text{CR}_{d'}(\mathbf{S}_k)$  in (1.18) is very optimistic in general. However, when  $\mathbf{S}_{2:k}$  contains a  $2^{d-m}$  fractional factorial design with resolution  $R \geq d' + 1$ , then each projected designs  $P_{d',r}(\mathbf{S}_{2:k})$  contains a full factorial design, see Sect. 1.3.2, and the bound becomes accurate.*

*For any projection  $P_{d',r}$ , there are  $2^{d'}$  distinct points  $P_{d',r}(\mathbf{x}_i)$  at most when  $\mathbf{x}_i$  varies in  $\mathbf{S}_{2:k}$  (which has  $k - 1$  elements), so that  $\text{PR}[P_{d',r}(\mathbf{S}_k)] = \text{PR}[P_{d',r}(\mathbf{S}_{2:k})] = 0$  when  $2^{d'} < k - 1$ . One may note that this case is already covered by (1.17). Indeed, Singleton bound, see Sect. 1.4.1, implies that the  $k - 1$  points of any  $P_{d',r}(\mathbf{S}_{2:k})$  with  $d'' = d - [\rho_H(\mathbf{S}_{2:k}) - 1]$  are all distinct; therefore,  $k - 1 \leq 2^{d-\rho_H(\mathbf{S}_{2:k})+1}$ , and  $2^{d'} < k - 1$  implies  $d' \leq d - \rho_H(\mathbf{S}_{2:k})$ .*

When  $d$  is large, we cannot compute the values of  $\text{CR}_{d'}(\mathbf{X}_k)$  and  $\text{PR}_{d'}(\mathbf{X}_k)$  in (1.15) and (1.16) exactly, and we shall consider the following approximations that use  $q$  projections at most, instead of  $\binom{d}{d'}$ :

$$\begin{aligned} \widehat{\text{CR}}_{d'}(\mathbf{X}_k) &= \max_{r=1, \dots, \min\{q, \binom{d}{d'}\}} \max_{\mathbf{x} \in \mathcal{X}_{d',Q}} d(\mathbf{x}, P_{d',r}(\mathbf{X}_k)), \\ \widehat{\text{PR}}_{d'}(\mathbf{X}_k) &= \frac{1}{2} \min_{r=1, \dots, \min\{q, \binom{d}{d'}\}} \min_{\mathbf{x}_i, \mathbf{x}_j \in \mathbf{X}_k, \mathbf{x}_i \neq \mathbf{x}_j} \|P_{d',r}(\mathbf{x}_i) - P_{d',r}(\mathbf{x}_j)\|, \end{aligned}$$

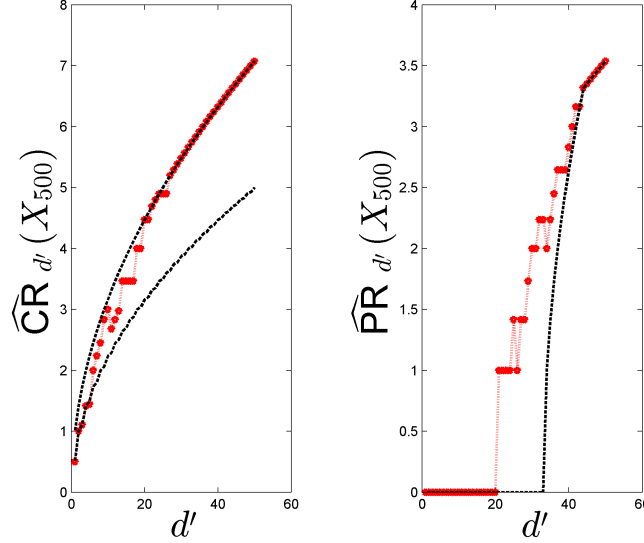
where  $\mathcal{X}_{d',Q}$  is a finite set of  $Q$  points in  $[-1, 1]^{d'}$ . The  $q$  projections are chosen randomly without repetition.  $\widehat{\text{CR}}_{d'}(\mathbf{X}_k)$  gives an optimistic (under) estimation of  $\text{CR}_{d'}(\mathbf{X}_k)$  due to the substitution of a finite set  $\mathcal{X}_{d',Q}$  for  $\mathcal{C}_{d'}$  and to the use of  $q$  random projections only. When  $d' \geq (4/3)[d - \rho_H(\mathbf{S}_{2:k})]$ ,  $\max_{\mathbf{x} \in \mathcal{C}_d} d(P_{d',r}(\mathbf{x}), P_{d',r}(\mathbf{S}_k)) = \sqrt{d'}$  for all projections  $P_{d',r}$ ; see the proof of Prop. 6. Therefore, for such  $d'$ , values of  $\widehat{\text{CR}}_{d'}(\mathbf{S}_k)$  smaller than  $\sqrt{d'}$  are only due to the substitution of  $\mathcal{X}_{d',Q}$  for  $\mathcal{C}_d$ .  $\widehat{\text{PR}}_{d'}(\mathbf{X}_k)$  over-estimates  $\text{PR}_{d'}(\mathbf{X}_k)$  due to the restriction to  $q$  projections, but  $\widehat{\text{PR}}_{d'}(\mathbf{S}_k) = 0$  when  $2^{d'} \leq k - 2$ ; see Remark 4.

Equation (1.17) indicates that  $\text{PR}_{d'}(\mathbf{S}_k) = 0$  for the designs obtained with Algorithm 3 when  $d' \leq d_0 = d - \rho_H(\mathbf{S}_{2:k})$ . Although the rescaling procedure of Sect. 1.7.2 prevents the exact coincidence of projected design points,  $\text{PR}_{d'}(\tilde{\mathbf{S}}_k)$  remains very close to zero when  $d' \leq d_0$  for rescaled designs. As the example below will illustrate, the performances in terms of  $\text{PR}_{d'}$  are thus much worse than those of more classical designs based on Lh and Sobol' sequences for small  $d'$ . They are much better, however, for  $d'$  close to  $d$ . The example also illustrates that rescaling decreases  $\text{PR}_{d'}$  for those large  $d'$ , but has the benefit of slightly improving (decreasing)  $\text{CR}_{d'}$ .

*Example 2 (continued).* We take  $q = 100$ , so that  $\min \left\{ q, \binom{d}{d'} \right\} = q$  for  $d' = 2$  already, and  $\mathcal{X}_{d',Q}$  consisting of the first  $2^{14}$  points of a Sobol' sequence in  $[-1, 1]^{d'}$ , complemented by a  $2^{d'}$  full-factorial design when  $d' \leq 10$ . We consider designs of size  $k = 500$ . Equation (1.17) shows the importance of having  $\rho_H(\mathbf{S}_{2:k})$  as large as possible to obtain good projection properties in terms of packing radius for dimensions  $d'$  as small as possible. The  $2^{50-35}$  fractional factorial design  $\mathbf{X}_n$  of Example 2 has  $\rho_H(\mathbf{X}_n) = 13$ ;  $\mathbf{S}_{2:k}$  generated with Algorithm 3 satisfies  $\rho_H(\mathbf{S}_{2:500}) = 17$ , with  $\text{PR}(\mathbf{S}_{2:500}) = \sqrt{17} \simeq 4.1231$ , see Fig. 1.3. Therefore,  $\text{PR}_{d'}(\mathbf{S}_{500}) = 0$  for  $d' \leq 33$  from Prop. 6, and, due to the random choice of 100 projections only among  $\binom{50}{d'}$ ,  $\widehat{\text{PR}}_{d'}(\mathbf{S}_{500})$  is equal to zero with positive probability when  $d' \leq 33$ . From Remark 4,  $\widehat{\text{PR}}_{d'}(\mathbf{S}_{500}) = 0$  for  $d' \leq 8$ .  $\widehat{\text{PR}}_{d'}(\mathbf{S}_{500}) = \text{PR}_{d'}(\mathbf{S}_{500}) = \sqrt{d'}/2$  for  $d' \geq 44$ , in agreement with (1.17). Figure 1.5 presents the lower and upper bounds (1.18) on  $\text{CR}_{d'}(\mathbf{S}_{500})$  (black dotted lines) and the approximation  $\widehat{\text{CR}}_{d'}(\mathbf{S}_{500})$  based on  $q$  random projections (stars), together with  $\text{PR}_{d'}(\mathbf{S}_{500})$  given by (1.17) (black dotted line) and its approximation  $\widehat{\text{PR}}_{d'}(\mathbf{S}_{500})$  (stars).

Linear rescaling of  $\mathbf{s}_k$  by (1.13) with  $\gamma = 1$  affects the values of  $\widehat{\text{PR}}_{d'}(\widetilde{\mathbf{S}}_{500})$ , see Fig. 1.6. Although we have now  $\|P_{d',r}(\widetilde{\mathbf{x}}_j) - P_{d',r}(\widetilde{\mathbf{x}}_\ell)\| \neq 0$  for all projections  $P_{d',r}$  and all  $\widetilde{\mathbf{x}}_j \neq \widetilde{\mathbf{x}}_\ell$  in  $\widetilde{\mathbf{S}}_{500}$ , when  $d'$  is small this value remains very close to zero for some pairs of points and projections, therefore  $\widehat{\text{PR}}_{d'}(\widetilde{\mathbf{X}}_{500})$  is still very small: the difference is hardly visible on the figure for small  $d'$ ; compare the red stars on the plots of  $\widehat{\text{PR}}_{d'}$  in Figs. 1.5 and 1.6 for  $d' \lesssim 16$ . For larger  $d'$ , rescaling decreases  $\widehat{\text{PR}}_{d'}$ , but has a small positive effect on  $\widehat{\text{CR}}_{d'}$  which is slightly decreased for  $d' \lesssim 30$ . The values of  $\widehat{\text{CR}}_{d'}$  for a design  $\mathbf{X}_{500}^S$  given by the first 500 points of Sobol' sequence (black circles), or for a (non-incremental) Lh design  $\mathbf{X}_{500}^{Lh}$  optimised for the PR criterion (blue diamonds), are marginally better than  $\widehat{\text{CR}}_{d'}(\widetilde{\mathbf{S}}_{500})$ , but  $\widetilde{\mathbf{S}}_{500}$  is significantly better in terms of  $\widehat{\text{PR}}_{d'}$  for large  $d'$ . The construction of  $\mathbf{X}_{500}^{Lh}$  uses the ESE algorithm of [10] with the default tuning parameters suggested in that paper; 100 cycles are performed, requiring 500 000 evaluations of PR.

Rescaling with  $\gamma = d$  in (1.13) yields results intermediate between no rescaling (Fig. 1.5) and linear rescaling (Fig. 1.6); see Fig. 1.7. Performances with linear periodic rescaling using (1.14) with  $K = 50$  (Fig. 1.8) are close to those on Fig. 1.6, with some small improvement in terms of  $\widehat{\text{CR}}_{d'}$ . For large enough  $d' < d$ , performances in terms of  $\widehat{\text{PR}}_{d'}$  are significantly better than those obtained for a design  $\mathbf{S}_{500}^S$  generated by Algorithm 1 with the first  $n = 2^{19}$  points of Sobol' sequence as candidate set (note that Algorithm 3 only uses 32 768 candidate points), whereas the performances of  $\widetilde{\mathbf{S}}_{500}$  and  $\mathbf{S}_{500}^S$  in terms of  $\widehat{\text{CR}}_{d'}$  are fairly close. The value  $\text{PR}_d(\mathbf{S}_{500}^S)$  corresponds to the last point on the black dashed line in Fig. 1.3;  $\text{PR}_d(\mathbf{S}_{500})$  is smaller than  $\text{PR}_d(\widetilde{\mathbf{S}}_{2:500})$  on Fig. 1.4-Right due to the addition of the central point  $\mathbf{0}$ .

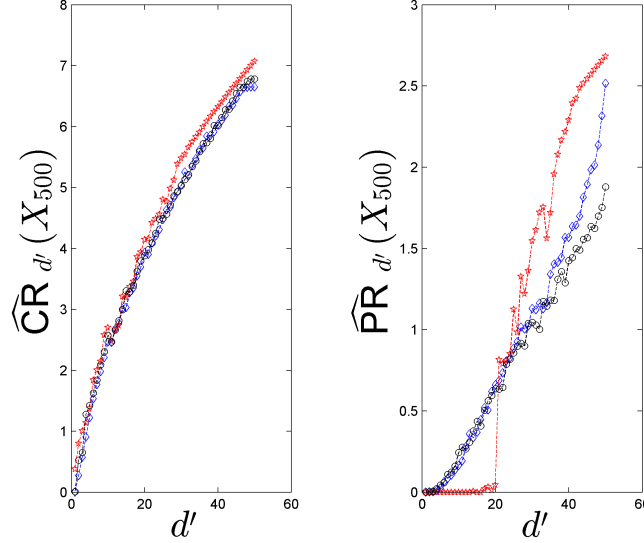


**Fig. 1.5** Lower and upper bounds (1.18) on  $CR_{d'}(\mathbf{S}_{500})$  and  $PR_{d'}(\mathbf{S}_{500})$  given by (1.17) (black dotted lines);  $\widehat{CR}_{d'}(\mathbf{S}_{500})$  and  $\widehat{PR}_{d'}(\mathbf{S}_{500})$  (red stars);  $\mathbf{S}_{500}$  is generated by Algorithm 3.

Above, we have used a  $2^{d-m}$  fractional factorial design with  $m = 35$  as candidate set in Algorithm 3, which allows the construction of incremental designs of size up to 32 769. If we are only interested in shorter sequences, we may increase  $m$  and get a design with larger  $\rho_H$  value. For instance, when taking  $m = 41$  instead of  $m = 35$ , the design  $\mathbf{X}_n$  obtained with the algorithm of Sect. 1.5.2 has 512 points, resolution III and  $\rho_H(\mathbf{X}_n) = 20$  ( $CR_H(\mathbf{X}_n) \geq 17$ ). Algorithm 3 then yields a  $\mathbf{S}_k$  such that  $\rho_H(\mathbf{S}_{2:500}) = 20$  (and  $PR_{d'}(\mathbf{S}_{500}) = 0$  for  $d' \leq 30$ , instead of  $d' \leq 33$  when  $m = 35$ , see (1.17)).

### 1.8 Summary and future work

In situations where the number  $d$  of factors is too large to inspect all vertices of the hypercube  $\mathcal{C}_d = [-1, 1]^d$  to construct a design, we suggest to use a fractional factorial design  $\mathbf{X}_n$  to thin the search space. When  $\mathbf{X}_n$  has minimum Hamming distance at least  $d/4$ , the coffee-house rule permits to construct a sequence of nested designs, with flexible size up

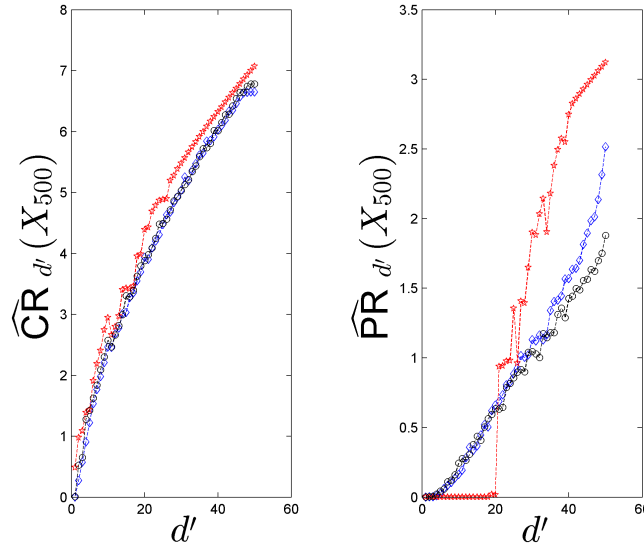


**Fig. 1.6**  $\widehat{CR}_{d'}(\widetilde{\mathbf{S}}_{500})$  and  $\widehat{PR}_{d'}(\widetilde{\mathbf{S}}_{500})$  after linear rescaling of  $\mathbf{S}_{500}$  using (1.13) (red stars;  $r = 10^{-6}$ ,  $K = 500$  and  $\gamma = 1$ );  $\widehat{CR}_{d'}$  and  $\widehat{PR}_{d'}$  for  $\mathbf{X}_{500}^S$  given by the first 500 points of Sobol' sequence (black circles) and a 500 point Lh design  $\mathbf{X}_{500}^{Lh}$  optimised for the PR criterion with the ESE algorithm of [10] (blue diamonds).

to  $n + 1$ , each design along the sequence having at least 50% packing (maximin) and covering (minimax) efficiency.

The packing and covering properties of designs projected in lower dimensional subspaces have been investigated. The covering performances are slightly worse than those obtained for more classical space-filling designs, but their packing performance is significantly better when projecting on a subspace with large enough dimension.

A natural drawback of the construction is that all design points (except the first one, taken at the centre) are vertices of the hypercube. A rescaling rule has been proposed to populate the interior of  $\mathcal{C}_d$ , but, like for the multi-layer designs of [1], all rescaled design points lie along the diagonals of  $\mathcal{C}_d$ . Other rules could be considered that deserve further investigations. For instance, the compromise between placing points on vertices, which is favourable for packing in the full dimensional space, and in the interior of  $\mathcal{C}_d$ , which is favourable to the performance of projected designs, may rely on interlacing the sequence proposed in the paper with a low discrepancy se-



**Fig. 1.7** Same as Fig. 1.6, but with nonlinear rescaling ((1.13) with  $\gamma = d$ ) of  $\mathbf{S}_{500}$  generated by Algorithm 3 (red stars).

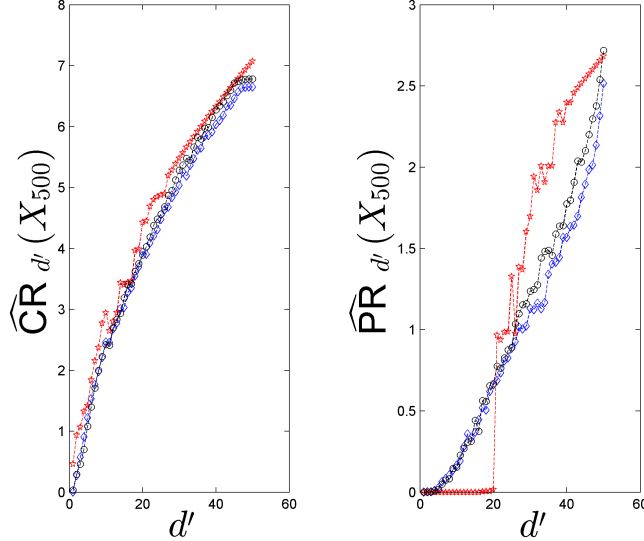
quence. Combination with other space-filling sequences could be considered as well, see, e.g., [27, 28]. We leave such developments for further work.

**Acknowledgements** This work benefited from the support of the project INDEX (Incremental Design of EXperiments) ANR-18-CE91-0007 of the French National Research Agency (ANR).

## Appendix

As shown in [7], a design for which a basic factor is not used in generators cannot have maximum resolution; see Sect. 1.3. It also has poor performance in terms of Hamming distance. Indeed, suppose without any loss of generality that the first factor is not used, and consider  $\mathbf{x}_i = (1, \mathbf{x}_{\setminus 1}) \in \mathbf{X}_n$ , where  $\mathbf{x}_{\setminus 1}$  is the vector obtained omitting the first coordinate of  $\mathbf{x}$ . The point  $\mathbf{x}' = (-1, \mathbf{x}_{\setminus 1})$  also belongs to  $\mathbf{X}_n$ , and  $\rho_H(\mathbf{X}_n) \leq d_H(\mathbf{x}_i, \mathbf{x}') = 1$ , implying that  $B_1(\mathbf{X}_n) \geq 1$ . As shown below, the reverse property holds true.

**Proposition 7** *If each basic factor is used in the construction of generators for  $\mathbf{X}_n$ , then  $\rho_H(\mathbf{X}_n) \geq 2$  and  $B_1(\mathbf{X}_n) = 0$ .*



**Fig. 1.8**  $\widehat{CR}_{d'}(\widetilde{S}_{500})$  and  $\widehat{PR}_{d'}(\widetilde{S}_{500})$  after linear periodic rescaling of  $S_{500}$  (red stars;  $K = 50$ ,  $r = 10^{-6}$ ,  $\gamma = 1$ ); the black circles correspond to the design obtained with Algorithm 1 using the first  $n = 2^{19}$  points of Sobol' sequence as candidate set (the same design is used for the black dashed line in Fig. 1.3); the blue diamonds correspond to a 500 point Lh design  $\mathbf{X}_{500}^{Lh}$  optimised for the PR criterion with the ESE algorithm of [10].

*Proof.* We use the notation of Sect. 1.5.1. Consider any pair of points  $\mathbf{x}_i = (\underline{\mathbf{x}}_i, \mathbf{g}(\underline{\mathbf{x}}_i))$  and  $\mathbf{x}_j = (\underline{\mathbf{x}}_j, \mathbf{g}(\underline{\mathbf{x}}_j))$  of  $\mathbf{X}_n$ . If  $d_H(\underline{\mathbf{x}}_i, \underline{\mathbf{x}}_j) \geq 2$ , then  $d_H(\mathbf{x}_i, \mathbf{x}_j) \geq 2$ . Otherwise,  $\underline{\mathbf{x}}_i$  and  $\underline{\mathbf{x}}_j$  differ by one coordinate only, say the  $k$ th. Since the  $k$ th basic factor is used within generators,  $d_H(\mathbf{g}(\underline{\mathbf{x}}_i), \mathbf{g}(\underline{\mathbf{x}}_j)) \geq 1$ , implying  $d_H(\mathbf{x}_i, \mathbf{x}_j) \geq 2$ .  $\square$

In the rest of the appendix we show how a similar reasoning can be used to construct  $2^{d-m}$  designs with a larger minimum Hamming distance  $\rho_H$ .

**Proposition 8**  $\rho_H(\mathbf{X}_n) \geq 3$  if and only if in the construction of generators

(i) each basic factor is used at least twice

and

(ii) for each pair of basic factors, one of the factors appears at least once separately.

*Proof.* Consider the point  $\mathbf{x} = (\mathbf{1}_{d-m}, \mathbf{g}(\mathbf{1}_{d-m})) \in \mathbf{X}_n$ . Suppose that (i) is not satisfied, with the first basic factor appearing only once among generators. Then  $\mathbf{x}'_1 = (-1, \mathbf{1}_{d-m-1}, \mathbf{g}(-1, \mathbf{1}_{d-m-1}))$  belongs to  $\mathbf{X}_n$  and  $d_H(\mathbf{g}(\mathbf{1}_{d-m}), \mathbf{g}(-1, \mathbf{1}_{d-m-1})) = 1$ , implying that  $d_H(\mathbf{x}, \mathbf{x}'_1) = 2$ . Also, when the first two factors only appear as a pair, then  $\mathbf{x}'_2 = (-1, -1, \mathbf{1}_{d-m-2},$

$\mathbf{g}(-1, -1, \mathbf{1}_{d-m-2}) \in \mathbf{X}_n$  and  $d_H(\mathbf{g}(\mathbf{1}_{d-m}), \mathbf{g}(-1, -1, \mathbf{1}_{d-m-2})) = 0$ , implying that  $d_H(\mathbf{x}, \mathbf{x}'_2) = 2$ . This shows that (i) and (ii) are necessary to have  $\rho_H(\mathbf{X}_n) \geq 3$ .

We show that the condition is sufficient. From Prop. 2, we only need to consider the nearest neighbour to the point  $\mathbf{x} = (\mathbf{1}_{d-m}, \mathbf{g}(\mathbf{1}_{d-m}))$ , which, up to a reordering of basic factors, is given by  $\mathbf{x}'_1$  or  $\mathbf{x}'_2$ . Now, (i) implies that  $d_H(\mathbf{g}(\mathbf{1}_{d-m}), \mathbf{g}(-1, \mathbf{1}_{d-m-1})) \geq 2$  and (ii) implies that  $d_H(\mathbf{g}(\mathbf{1}_{d-m}), \mathbf{g}(-1, -1, \mathbf{1}_{d-m-2})) \geq 1$ , showing that  $d_H(\mathbf{x}, \mathbf{x}'_1) \geq 3$  and  $d_H(\mathbf{x}, \mathbf{x}'_2) \geq 3$ .  $\square$

*Example 3.* Consider the design given by the half fraction  $2^{d-1}$  with the product of all basic factors as generators:  $g_1 = \prod_{j=1}^{d-1} x_j$  ( $m = 1$ ). The condition of Prop. 7 is satisfied, but none of the conditions (i) and (ii) of Prop. 8 is; therefore,  $\rho_H(\mathbf{X}_n) = 2$ . Direct calculation gives  $A_0(\mathbf{X}_n) = 1$  and  $A_{2q}(\mathbf{X}_n) = \binom{d-1}{2q-1} + \binom{d-1}{2q}$  for  $1 \leq q < d/2$ , with  $A_d(\mathbf{X}_n) = 1$  when  $d$  is even, all  $A_i$  with  $i$  odd being equal to zero.

**Proposition 9**  $\rho_H(\mathbf{X}_n) \geq 4$  if and only if in the construction of generators

(i) each basic factor is used at least three times

and

(ii-a) for each pair of basic factors, one of the factors appears at least twice separately

or

(ii-b) for each pair of basic factors, each one of the factors in the pair appears at least once separately

and

(iii-a) each triple of basic factors appears at least once

or

(iii-b) within each triple of basic factors, each factor appears at least once without the other two.

The proof uses arguments similar to that used for Prop. 8.

*Example 4.* Consider the case  $d = 9$  and  $m = 5$ , with basic factors  $a, b, c, d$  and generators  $abc, abd, acd, bcd$  and  $abcd$ . Conditions (i), (ii-b) and (iii-a) are satisfied and we get  $\mathcal{A}(\mathbf{X}_n) = [1 \ 0 \ 0 \ 4 \ 14 \ 8 \ 0 \ 4 \ 1 \ 0]$  and  $\mathcal{B}(\mathbf{X}_n) = [1 \ 0 \ 0 \ 0 \ 6 \ 8 \ 0 \ 0 \ 1 \ 0]$ .

When the generators are  $ab, abd, acd, bc$  and  $cd$ , then (iii-b) is satisfied instead of (iii-a) and we get  $\mathcal{A}(\mathbf{X}_n) = [1 \ 0 \ 0 \ 6 \ 9 \ 9 \ 6 \ 0 \ 0 \ 1]$  and  $\mathcal{B}(\mathbf{X}_n) = [1 \ 0 \ 0 \ 0 \ 9 \ 0 \ 6 \ 0 \ 0 \ 0]$ . Note that the first design is preferable both in terms of aberration and maximin distance.

*Example 5.* Consider the case where  $d = 2m$ ,  $m \geq 4$ , and where each of the  $m$  generators is the product of all basic factors but one; that is, with obvious notation,  $g_i = \prod_{j=1, j \neq i}^m x_j$ , for  $i = 1, \dots, m$ . Conditions (i), (ii-b) and (iii-a) of Prop. 9 are satisfied, and direct calculation shows that the word length pattern of  $\mathbf{X}_n$  (with  $n = 2^{d-m} = 2^m$ ) satisfies

$$\begin{aligned}
A_0(\mathbf{X}_n) &= 1, \\
A_m(\mathbf{X}_n) &= \sum_{p=0}^{\lfloor m/2 \rfloor} \binom{m}{2p+1} \text{ and } A_{4p}(\mathbf{X}_n) = \binom{m}{2p} \text{ for } p = 1, \dots, \lfloor m/2 \rfloor \\
&\text{if } m \bmod 4 \neq 0, \\
A_m(\mathbf{X}_n) &= \binom{m}{m/2} + \sum_{p=0}^{\lfloor m/2 \rfloor} \binom{m}{2p+1} \text{ and } A_{4p}(\mathbf{X}_n) = \binom{m}{2p} \text{ for } p \neq m/4 \\
&\text{if } m \bmod 4 = 0,
\end{aligned}$$

all other  $A_i$  being equal to zero (and the design has resolution  $IV$ ). Direct calculation also indicates that  $\mathcal{B}(\mathbf{X}_n) = \mathcal{A}(\mathbf{X}_n)$ , with therefore  $\rho_H(\mathbf{x}_n) = 4$ .

One can check that the sphere packing bound (1.6) implies that a  $2^{d-m}$  design with  $d = 2m$  cannot reach  $\rho_H(\mathbf{x}_n) \geq 5$  for  $m < 7$ . However, designs with better maximin properties can be obtained for larger  $m$ . For instance, when  $m = 8$  ( $d = 16$ ,  $n = 256$ ), the construction above yields a design  $\mathbf{X}_n$  with  $\mathcal{B}(\mathbf{X}_n) = \mathcal{A}(\mathbf{X}_n) = [1 \ 0 \ 0 \ 0 \ 28 \ 0 \ 0 \ 0 \ 198 \ 0 \ 0 \ 0 \ 28 \ 0 \ 0 \ 0 \ 1]$  and  $\rho_H(\mathbf{X}_n) = 4$ , whereas the design with generators  $abcdefgh$ ,  $defgh$ ,  $befgh$ ,  $acegh$ ,  $bdgh$ ,  $cefh$ ,  $adfh$ , and  $abeh$  (and basic factors  $1, b, c, d, e, f, g, h$ ) has distance distribution  $\mathcal{B}(\mathbf{X}_n) = [1 \ 0 \ 0 \ 0 \ 0 \ 24 \ 44 \ 40 \ 45 \ 40 \ 28 \ 24 \ 10 \ 0 \ 0 \ 0 \ 0]$ , with  $\rho_H(\mathbf{x}_n) = 5$  (again,  $\mathcal{A}(\mathbf{X}_n) = \mathcal{B}(\mathbf{X}_n)$ , with  $\mathbf{X}_n$  thus having resolution  $V$ ).  $\text{CR}_H(\mathbf{X}_n) = 4$  for both designs.

## References

1. S. Ba and V.R. Jospheh. Multi-layer designs for computer experiments. *Journal of the American Statistical Association*, 106(495):1139–1149, 2011.
2. G.E.P. Box and J.S. Hunter. The  $2^{k-p}$  fractional factorial designs. part I. *Technometrics*, 3(3):311–351, 1961.
3. G.E.P. Box and J.S. Hunter. The  $2^{k-p}$  fractional factorial designs. part II. *Technometrics*, 3(4):449–458, 1961.
4. C.-S. Cheng and B. Tang. A general theory of minimum aberration and its applications. *Annals of Statistics*, 33(2):944–958, 2005.
5. K.-T. Fang and R. Mukerjee. Uniform design: theory and application. *Biometrika*, 87(1):193–198, 2000.
6. M.F. Franklin and B.A. Bailey. Selection of defining contrasts and confounded effects in two-level experiments. *Applied Statistics*, 26(3):321–326, 1977.
7. A. Fries and W.G. Hunter. Minimum aberration  $2^{k-p}$  designs. *Technometrics*, 22(4):601–608, 1980.
8. T.F. Gonzalez. Clustering to minimize the maximum intercluster distance. *Theoretical Computer Science*, 38:293–306, 1985.
9. F.J. Hickernell. A generalized discrepancy and quadrature error bound. *Mathematics of Computation*, 67(221):299–322, 1998.
10. R. Jin, W. Chen, and A. Sudjianto. An efficient algorithm for constructing optimal design of computer experiments. *Journal of Statistical Planning and Inference*, 134:268–287, 2005.

11. P.W.M. John, M.E. Johnson, L.M. Moore, and D. Ylvisaker. Minimax distance designs in two-level factorial experiments. *Journal of Statistical Planning and Inference*, 44:249–263, 1995.
12. M.E. Johnson, L.M. Moore, and D. Ylvisaker. Minimax and maximin distance designs. *Journal of Statistical Planning and Inference*, 26:131–148, 1990.
13. V.R. Joseph, E. Gul, and S. Ba. Maximum projection designs for computer experiments. *Biometrika*, 102(2):371–380, 2015.
14. R.W. Kennard and L.A. Stone. Computer aided design of experiments. *Technometrics*, 11(1):137–148, 1969.
15. M.K. Kerr. Bayesian optimal fractional factorials. *Statistica Sinica*, 11:605–630, 2001.
16. S. Lin and C. Xing. *Coding Theory. A First Course*. Cambridge University Press, Cambridge, 2004.
17. J.L. Loewy, J. Sacks, and W.J. Welch. Choosing the sample size of a computer experiment: a practical guide. *Journal of the American Statistical Association*, 51(4):366–376, 2009.
18. T.J. Mitchell, M.D. Morris, and D. Ylvisaker. Two-level fractional factorials and Bayesian prediction. *Statistica Sinica*, 5:559–573, 1995.
19. M.D. Morris and T.J. Mitchell. Exploratory designs for computational experiments. *Journal of Statistical Planning and Inference*, 43:381–402, 1995.
20. W.G. Müller. Coffee-house designs. In A. Atkinson, B. Bogacka, and A. Zhigljavsky, editors, *Optimum Design 2000*, chapter 21, pages 241–248. Kluwer, Dordrecht, 2001.
21. W.G. Müller. *Collecting Spatial Data*. Springer, Berlin, 2007. [3rd ed.].
22. H. Niederreiter. *Random Number Generation and Quasi-Monte Carlo Methods*. SIAM, Philadelphia, 1992.
23. R.L. Plackett and J.P. Burman. The design of optimum multifactorial experiments. *Biometrika*, 33(4):305–325, 1946.
24. M. Plotkin. Binary codes with specified minimum distance. *IRE Transactions on Information Theory*, 6(4):445–450, 1960.
25. L. Pronzato. Minimax and maximin space-filling designs: some properties and methods for construction. *Journal de la Société Française de Statistique*, 158(1):7–36, 2017.
26. L. Pronzato and W.G. Müller. Design of computer experiments: space filling and beyond. *Statistics and Computing*, 22:681–701, 2012.
27. L. Pronzato and A.A. Zhigljavsky. Bayesian quadrature and energy minimization for space-filling design, 2018. hal-01864076, arXiv:1808.10722v1.
28. L. Pronzato and A.A. Zhigljavsky. Measures minimizing regularized dispersion. *J. Scientific Computing*, 78(3):1550–1570, 2019.
29. R.C. Singleton. Maximum distance  $q$ -nary codes. *IEEE Transactions on Information Theory*, 10(2):116–118, 1964.
30. F. Sun, Y. Wang, and H. Xu. Uniform projection designs. *Annals of Statistics*, 47(1):641–661, 2019.
31. Y. Tang, H. Xu, and D.K.J. Lin. Uniform fractional factorial designs. *Annals of Statistics*, 40(2):891–907, 2012.
32. J.H. van Lint and R.M. Wilson. *A Course in Combinatorics*. Cambridge University Press, Cambridge, 1992.
33. C.F.J. Wu and M.S. Hamada. *Experiments: Planning Analysis, and Optimization*. Wiley, New York, 2009. [2nd ed.].
34. Q. Xiao and H. Xu. Construction of maximin distance designs via level permutation and expansion. *Statistica Sinica*, 28:1395–1414, 2018.
35. H. Xu and C.F.J. Wu. Generalized minimum aberration for asymmetrical fractional factorial designs. *Annals of Statistics*, 29(4):1066–1077, 2001.
36. Y.-D. Zhou and H. Xu. Space-filling fractional factorial designs. *Journal of the American Statistical Association*, 109(507):1134–1144, 2014.