

## Report 2001-001

***Citation for published version (APA):***

Eurandom (2001). *Report 2001-001*. (Report Eurandom; Vol. 2001001). Eurandom.

***Document status and date:***

Published: 01/01/2001

***Document Version:***

Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

***Please check the document version of this publication:***

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

***General rights***

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

[www.tue.nl/taverne](http://www.tue.nl/taverne)

***Take down policy***

If you believe that this document breaches copyright please contact us at:

[openaccess@tue.nl](mailto:openaccess@tue.nl)

providing details and we will investigate your claim.

Report 2001-001

**Understanding aliasing using Gröbner bases**

G. Pistone, E. Riccomagno and H. P. Wynn

**Replications with Gröbner bases**

A. M. Cohen, A. Di Bucchianico and E. Riccomagno

ISSN: 1389-2355



# Understanding aliasing using Gröbner bases

Giovanni Pistone  
Eva Riccomagno  
Henry P. Wynn

**ABSTRACT:** The now well-established Gröbner basis method in experimental design (see the authors' monograph "Algebraic Statistics") had the understanding of aliasing as a key motivation. The basic method asks: given an experimental design, what is estimable, or more generally what is the alias structure? The paper addresses the following related question: given a set of conditions which the design is known to satisfy, what can we say about the alias structure? Some classical and non-classical construction methods are included.

**KEYWORDS:** Aliasing; Gröbner basis; polynomial conditions

## 1 The Gröbner basis method

We summarise the method briefly in a number of steps.

Step 1 Define a design  $D \subset \mathbf{R}^d$  as a set of  $n$  distinct points:  $D = [a^{(i)}]_{i=1}^n$ .

Step 2 Set up a series of polynomial equations whose solutions give precisely  $D$ . This, mathematically, amounts to representing the design as a zero dimensional algebraic variety. The design ideal,  $\text{Ideal}(D)$ , is the set of all polynomials whose zeros include the design points.

Step 3 Select a so-called monomial ordering  $\tau$ . This is a total well-ordering on the monomials such that  $x^\alpha \prec_\tau x^\beta$  implies  $x^\alpha x^\gamma \prec_\tau x^\beta x^\gamma$  for all  $\gamma \neq 0$ .

Step 4 Generate a Gröbner basis for  $\text{Ideal}(D)$ , given  $\tau$ , namely a special representation of  $D$  as the solutions of polynomial equations

$$\{g_j(x) = 0 : j = 1, \dots, k\}$$

The full details are omitted.

Step 5 List the leading terms  $l_j(x) = \text{LT}(g_j(x))$ ,  $j = 1, \dots, k$  with respect to the monomial ordering  $\tau$ .

Step 6 List all monomials not divisible by any leading term  $l_j(x)$  ( $j = 1, \dots, k$ ). Call this list  $\text{Est}_\tau(D)$  and note that

- (i)  $\#\text{Est}_\tau(D) = n$ , that is the sample size of the design.
- (ii) If  $\text{Est}_\tau(D) = \{x^\alpha : \alpha \in L\}$  then the monomial terms are the basis of a saturated and estimable regression model

$$\sum_{\alpha \in L} \theta_\alpha x^\alpha$$

with non singular  $X$ -matrix

$$X = \{x^\alpha\}_{x \in D, \alpha \in L}$$

- (iii)  $\text{Est}_\tau(D)$  is an order ideal, that is if  $x^\alpha \in \text{Est}_\tau(D)$  then  $x^\beta \in \text{Est}_\tau(D)$  for any  $\beta$  component wise smaller than  $\alpha$ .

A lot can be said about this process with regard to appropriate computer algebra. For example methods are available for directly computing the Gröbner basis and  $\text{Est}_\tau(D)$  from  $D$ . See Pistone, Riccomagno and Wynn (2000) for details. For the present paper we note simply that the equations  $\{g_j(x) = 0 : j = 1, \dots, k\}$  defining the design essentially also give some alias structure. For example each leading term can be written

$$l_j(x) = \sum_{\alpha \in L} \theta_\alpha^{(j)} x^\alpha$$

so that  $g_j(x) = l_j(x) - \sum_{\alpha \in L} \theta_\alpha^{(j)} x^\alpha$ . That is to say “higher order” terms with respect to  $\tau$  can be written in terms of polynomials constructed from monomials in  $\text{Est}_\tau(D)$ . Apart from the fact that the equations are dependent on  $\tau$  all the alias structure can be captured from such equations. A generic member of the ideal  $\text{Ideal}(D)$  is written  $\sum_{j=1}^k s_j(x)g_j(x)$  where the  $s_j(x)$ ’s are generic polynomials. Setting this to zero for arbitrary  $s_j(x)$  gives all possible alias relations.

## 2 A theorem on aliasing

Historically there have been many combinatorial constructions of experimental design of which the standard Abelian group construction of symmetric and asymmetric factorial design is perhaps the most celebrated. In such constructions one exhibits a set of conditions which the designs must satisfy. For example to construct a  $2^{3-1}$  the equations are

$$x_1^2 = 1, \quad x_2^2 = 1, \quad x_3^2 = 1, \quad x_1 x_2 x_3 = 1$$

In the previous section this could be considered as Step 2.

In this section we discuss some simple properties that can be predicted for  $\text{Est}_\tau(D)$  given the construction equations but in advance of computing  $\text{Est}_\tau(D)$  itself.

The equation  $x_1x_2x_3 = 1$  for the  $2^{3-1}$  above implies that the interaction  $x_1x_2x_3$  and the constant term are aliased, in particular the vector obtained by evaluating  $x_1x_2x_3$  at the design points is the unit vector,  $(1, \dots, 1)$ . Theorem 1 generalises this observation.

**Theorem 1** *Let the design  $D$  be known to satisfy the polynomial equation*

$$h(x) = 0$$

*in  $\mathbf{R}^d$  and let  $\tau$  be a monomial ordering. Let  $M(\neq \emptyset)$  be the set of monomials with non zero coefficients in  $h$ . Then*

$$M \not\subseteq \text{Est}_\tau(D)$$

*Proof.* This is by contradiction. Suppose  $M \subseteq \text{Est}_\tau(D)$ . Then

$$h(x) = \sum_{\alpha \in M \subseteq L} \phi_\alpha x^\alpha$$

where  $\text{Est}_\tau(D) = \{x^\alpha : \alpha \in L\}$  and all  $\phi_\alpha \neq 0$  for  $\alpha$  such that  $x^\alpha \in M$ . But this is false since all  $x^\alpha$  are linearly independent over  $D$  and  $h(x) = 0$  for all  $x \in D$ .

A proof relying on more classical arguments from matrix theory is as follows. Let  $h(x) = \sum_{\alpha \in M} \phi_\alpha x^\alpha$  and consider the matrix  $X = \{x^\alpha\}_{x \in D, \alpha \in M}$  with columns  $X(\alpha) = \{x^\alpha\}_{x \in D}$ . The matrix  $X$  is singular because the condition  $h(x) = 0$  implies that the linear combination of the columns of  $X$ ,  $\{\sum_{\alpha \in M} \phi_\alpha X(\alpha)\}$  is the zero vector. Thus the monomials  $x^\alpha$ ,  $\alpha \in M$  are linearly dependent over  $D$  and cannot all be included in a model identifiable by  $D$ .  $\square$

To repeat the result of the theorem: any  $M$  must have at least one “non zero” term not in  $\text{Est}_\tau(D)$ . Note also that if  $x^\beta \notin \text{Est}_\tau(D)$  it also follows from Gröbner basis theory that  $x^\gamma \notin \text{Est}_\tau(D)$  for all  $\gamma > \beta$  component wise.

**Corollary 1** *If  $h(x) = x^\alpha - c$  for some index  $\alpha$  and constant  $c$  then if  $D \neq \emptyset$  then  $x^\alpha$  cannot be in  $\text{Est}_\tau(D)$ .*

*Proof.* Since  $D$  is not empty the constant must be in  $\text{Est}_\tau(D)$ . This follows from Step 6 (iii). Thus by Theorem 1,  $x^\alpha \notin \text{Est}_\tau(D)$ .  $\square$

Typically in construction we may know that  $h_j(x) = 0$  on  $D$  for  $j = 1, \dots, r$ . Then Theorem 1 applies to each corresponding  $M_j$ .

Sometimes we may construct designs as exactly all solutions of  $h_j(x) = 0$ ,  $j = 1, \dots, r$

$$D = \{x : h_j(x) = 0, \quad j = 1, \dots, r\}$$

However it is advisable to replace this by the Gröbner basis representation to obtain a more tractable description of aliasing.

**Example 1** Factorial design. The corollary makes a strong connection to the fractional factorial construction mentioned above since those consist typically of solving sets of equations of the form

$$\{x^{\alpha(j)} - c_j : j = 1, \dots, r\}$$

### 3 Further examples

**Example 2** Mixture. Here a basic equation is  $\sum x_i - 1 = 0$ . Theorem 1 simply says that not all of 1 and  $x_i$ ,  $i = 1, \dots, d$  can be in  $\text{Est}_\tau(D)$  confirming the standard redundancy in this case. See Giglio, Riccomagno and Wynn (2000).

**Example 3** Other groups. Any design  $D$  invariant under a group  $G$  on  $\mathbf{R}^d$  will preserve the maximal invariants,  $\pi_j(x)$ , under  $G$ . Thus candidates for  $h_j(x)$  are

$$h_j(x) = \pi_j(x) - c_j$$

As a very simple example consider designs on a circle in  $\mathbf{R}^2$  satisfying

$$x_1^2 + x_2^2 = 1$$

Then we can conclude that both of  $x_1^2$  and  $x_2^2$  cannot be in  $\text{Est}_\tau(D)$ . Since maximal invariants are constant on orbits any design constructed as an orbit will be invariant

$$D = \{x = G(x_0) : \text{for a point } x_0 \in D\}$$

In the above example one can easily construct arbitrary large designs in this way and still not have  $x_1^2$  and  $x_2^2$  in  $\text{Est}_\tau(D)$ . This can easily be extended to rotations in  $\mathbf{R}^d$ .

An important class of groups in design theory are reflection groups. Indeed the conditions above  $x_1^2 = x_2^2 = x_3^2 = 1$  and  $x_1 x_2 x_3 = 1$  are precisely a set of invariants for the subgroup generated by the reflections

$$(x_1, x_2, x_3) \longrightarrow \begin{cases} (-x_1, -x_2, x_3) \\ (-x_1, x_2, -x_3) \end{cases}$$

and the design  $D = \{(1, 1, 1), (1, -1, -1), (-1, 1, -1), (-1, -1, 1)\}$  is an orbit.

**Example 4** Lattices. One generator lattice designs are equally spaced designs on the integer grid defined as

$$D = \{gk \pmod{n} : k = 0, \dots, n-1\}$$

where  $g$  is a vector of integers. If the components of  $g$  and  $n$  have greatest common divisor equal to one, then  $D$  has exactly  $n$  distinct points. For  $g = (1, 2)$  and  $n = 5$ ,  $D = \{(0, 0), (1, 2), (2, 4), (3, 1), (4, 3)\}$ . The Gröbner basis computed modulo 5 and with respect to any term-ordering for which  $x_2$  is smaller than  $x_1$ , includes the polynomial  $x_1 + 2x_2$ . Thus every point in  $D$  has to satisfy the equation  $x_1 + 2x_2 = 0 \pmod{5}$ . The full Gröbner basis is

$$\begin{aligned} &x_1 + 2x_2, \\ &x_2^5 - x_2 \end{aligned}$$

and  $\text{Est}_\tau(D)$  is  $\{1, x_2, x_2^2, x_2^3, x_2^4\}$ . This shows algebraically that modulo 5 the design  $D$  is a one dimensional object.

Over the real numbers and with respect to a lexicographic term ordering (see Cox, Little and O'Shea, 1996) with again  $x_2$  smaller than  $x_1$  the Gröbner basis is

$$\begin{aligned} g_1(x) &= x_2^5 - 10x_2^4 + 35x_2^3 - 50x_2^2 + 24x_2, \\ g_2(x) &= x_1 + 5/6x_2^4 - 20/3x_2^3 + 50/3x_2^2 - 83/6x_2 \end{aligned}$$

with the same  $\text{Est}_\tau(D)$ . The condition  $\pi(x) = x_1 + 2x_2$  is rewritten over  $\text{Est}_\tau(D)$  as

$$-5/6x_2^4 + 20/3x_2^3 - 50/3x_2^2 + 95/6x_2 = \pi(x) - g_2(x)$$

With respect to an ordering that does not favour either  $x_1$  or  $x_2$  so strongly, namely  $\text{tdeg}$  (see Char, Geddes, Gonnet, Leong, and Monogan, 1991) the set  $\text{Est}_\tau(D)$  is  $\{1, x_1, x_2, x_1x_2, x_2^2\}$ . This example shows that term orderings can be chosen to determine the structure of  $\text{Est}_\tau(D)$  as far as the design allows.

## 4 Conclusion

The theory and examples in this paper are relatively simple but, we hope, show the power of the method. The challenge is to revisit many of the classical and some of the more recent constructions in design, such as lattices, to relate the special algebra used in each case to the wider Gröbner basis theory. The list should include notions such as blocking, dummyming, trend resistance, cross-over which are of considerable practical importance, but where aliasing is not yet fully understood.

## References

- Char, B., Geddes, K., Gonnet, G., Leong, B. and Monogan, M. (1991). *MAPLE V Library Reference Manual*. Springer-Verlag, New York.
- Cox, D., Little, J. and O'Shea, D. (1996). *Ideals, Varieties, and Algorithms*. Springer, New York (second edition).
- Giglio, B., Riccomagno E. and Wynn, H.P. (2000). Gröbner basis methods in mixture experiments and generalisations. In: Atkinson, A., Bogacka, B. and Zhigljavsky, A. (eds). *Optimum Design 2000. Improvements of "Optimum experimental designs: Prospects for the new millennium?"*.
- Pistone, G., Riccomagno E. and Wynn, H.P. (2000). *Algebraic Statistics*. Chapman and Hall, New York.



# Replications with Gröbner Bases

A.M. Cohen  
A. Di Bucchianico  
E. Riccomagno

**ABSTRACT:** We present an extension of the Gröbner basis method for experimental design introduced in Pistone and Wynn (1996) to designs with replicates. This extension is presented in an abstract regression analysis framework, based on direct computations with functions and inner products. Explicit examples are presented to illustrate our approach.

**KEYWORDS:** Gröbner basis; replicates; orthonormalisation; projection

## 1 Introduction

Recently tools from algebraic geometry have been introduced in experimental design. See Pistone and Wynn (1996), Pistone, Riccomagno and Wynn (2000) and Riccomagno (1997). They are particularly useful in the analysis of complex experiments where there is a large number of factors and runs and the structure of the design is not regular, for example there are missing observations from a standard full factorial experiment. Confounding relations among factors and interactions are encoded in the Gröbner bases associated with a design allowing us to interpret confounding relations of the kind  $I = AB$  (where  $A$  and  $B$  are factors and  $I$  is the constant term) for a large class of designs and models.

A major requirement for the application of this technique is that the design has no replicates. However, there are several practical situations where replicates are useful. In the present work we extend the algebraic methods to designs with replicates.

The main idea is to introduce a new variable that counts how many times a point appears in the design. For example, the one-dimensional design with five points  $\mathcal{D}^* = \{0, 0, 1, 1, 2\}$  becomes the two-dimensional object  $\mathcal{D} = \{(0, 1), (0, 2), (1, 1), (1, 2), (2, 1)\}$ . This is encoded in the following set of polynomials in two indeterminates

$$x^3 - 3x^2 + 2x, \quad x^2t - xt - x^2 + x, \quad t^2 - 3t + 2,$$

where  $x$  represents the design factor and  $t$  counts the number of replicates.

The polynomials above can be used to construct a polynomial system of equations whose zeros are the points in  $\mathcal{D}$ . The zeros of the first polynomial, involving only the  $x$  indeterminates, are the distinct points in  $\mathcal{D}^*$ .

Least squares models are fitted to the data with replications as polynomial interpolators using the Gröbner basis method. In order to accommodate this process, we present a vector space setting for regression analysis in terms of functions on the design points. We suggest to perform estimation after orthonormalisation of the model terms (see also Giglio et al., 2000). The traditional sums of squares appear naturally as the lengths of the terms in the orthonormalised model. The coefficients from the non-orthonormalised model are obtained simply by comparing coefficients. A pleasant feature from the computational point of view is that to compute regression coefficients, we do not need to perform matrix inversion as in the standard matrix way of computing regression coefficients. We present several explicit examples to illustrate our method.

## 2 Basic setup

We start by fixing notation. A design without replicates is a finite subset of  $\mathbf{R}^d$ . The main idea behind the algebraic geometry approach to experimental design is to view a design as a variety, i.e. the set of common zeroes of a finite set of polynomials. Statistical analysis of data starts with finding a polynomial that interpolates the data at the design points. If points of a design are replicated, then strictly speaking we are dealing with multi-sets rather than ordinary sets. This causes problems for the algebraic geometric approach. Namely there is no polynomial (function) that takes different values at the same point. We overcome this difficulty by introducing an extra variable that counts how many times a point appears in the design

**Definition 2.1** *A design  $\mathcal{D}$  with replicates is a finite set of points in  $\mathbf{R}^d \times \mathcal{L}$ , where  $\mathcal{L}$  is a finite ordered set (the label set). The associated unreplicated design  $\mathcal{D}^*$  is defined by  $\mathcal{D}^* = \{a^* \in \mathbf{R}^d \mid \exists \ell \in \mathcal{L} \text{ such that } (a^*, \ell) \in \mathcal{D}\}$ . Each element  $a$  of  $\mathcal{D}$  is of the form  $a = (a^*, \ell)$ . Thus we may alternatively define  $\mathcal{D}^*$  as  $\mathcal{D}^* := \{a^* \mid a \in \mathcal{D}\}$ .*

Designs without replicates are special designs such that for each  $a^* \in \mathcal{D}^*$  there is exactly one  $\ell \in \mathcal{L}$  such that  $(a^*, \ell) \in \mathcal{D}$ . Two designs are isomorphic if their associated unreplicated designs coincide and there is a bijection between the two designs. In general, the unreplicated design  $\mathcal{D}^*$  is obtained by projecting  $\mathcal{D}$  onto the first  $d$  factors. The operation of projection does not take into account the number of replicates. It has a nice algebraic counterpart (see Theorem 4.3 below).

**Notation 2.2** *Let  $\mathcal{D} \subset \mathbf{R}^d \times \mathcal{L}$  be a design. The set of real-valued functions on  $\mathcal{D}$  is denoted by  $\mathcal{L}(\mathcal{D})$ .*

The inner product on  $\mathcal{L}(\mathcal{D})$  given in Definition 2.3 below is directly related to least squares estimation.

**Definition 2.3** *If  $\mathcal{D}$  is a design, then for all  $f, g \in \mathcal{L}(\mathcal{D})$  we define an inner product by  $\langle f, g \rangle_{\mathcal{D}} := \sum_{a \in \mathcal{D}} f(a)g(a)$ . A norm is defined on  $\mathcal{L}(\mathcal{D})$  by  $\|f\|_{\mathcal{D}} = \sqrt{\langle f, f \rangle_{\mathcal{D}}}$ .*

Note that weighted least squares is easily incorporated in this setup by slightly changing the definition of the inner product  $\langle \cdot, \cdot \rangle_{\mathcal{D}}$ .

Let  $\mathcal{D}$  be a design. Suppose our statistical model is

$$Y(x) = f(x, \theta) + \varepsilon(x), \quad (1.1)$$

where  $\theta \in \mathbf{R}^p$  and  $\varepsilon(x)$  is a real-valued random variable for all  $x \in \mathcal{D}^*$  with  $\mathbf{E}\varepsilon(x) = 0$  and  $\mathbf{V}\varepsilon(x) = \sigma^2$ . Suppose we have observations  $Y_1, \dots, Y_N$  from this model, where  $Y_i$  is  $Y(a_i^*)$  for  $a_i^* \in \mathcal{D}^*$  and replications are allowed, i.e.  $a_i^*$  may be equal to  $a_j^*$  for  $i \neq j$ . Then the **least squares estimator** for the parameter vector  $\theta$  is given by

$$\hat{\theta} = \min_{\theta \in \Theta} \sum_{i=1}^N |Y_i - f(a_i^*, \theta)|^2. \quad (1.2)$$

Let  $g \in \mathcal{L}(\mathcal{D})$  be the unique function in  $\mathcal{L}(\mathcal{D})$  such that  $g(a_i) = Y_i$  for all  $i = 1, \dots, N$ . Since

$$\hat{\theta} = \min_{\theta \in \Theta} \sum_{i=1}^N |Y_i - f(a_i^*, \theta)|^2 = \min_{\theta \in \Theta} \|g - f(\cdot, \theta)\|_{\mathcal{D}}^2, \quad (1.3)$$

we see that least squares estimation corresponds to a minimum distance problem in  $\mathcal{L}(\mathcal{D})$  with the inner product in Definition (2.3). Note that a function  $f(x)$  for  $x = (x_1, \dots, x_d) \in \mathcal{D}^*$  can be naturally extended for  $x = (x_1, \dots, x_d, x_{d+1}) \in \mathcal{D}$  by  $(x_1, \dots, x_d, x_{d+1}) \mapsto f(x_1, \dots, x_d)$ .

### 3 Identifiability of linear models

In the sequel we restrict ourselves to linear models, i.e. models such that  $f(x, \theta)$  is a linear function of the components of the parameter vector  $\theta$

$$Y(x) = \sum_{\alpha \in \mathcal{M}} \theta_{\alpha} p_{\alpha}(x) + \varepsilon(x), \quad (1.4)$$

where  $p_{\alpha}$  ( $\alpha \in \mathcal{M}$ ) is an element of  $\mathcal{L}(\mathcal{D}^*)$ . Clearly the  $p_{\alpha}$ 's can be viewed as elements of  $\mathcal{L}(\mathcal{D})$ . Recall that  $\mathcal{L}(\mathcal{D})$  is a vector space over the real numbers and  $\mathcal{L}(\mathcal{D})$  is isomorphic to  $\mathbf{R}^N$ , where  $N$  is the number of points in  $\mathcal{D}$ .

**Definition 3.1** *A linear model (1.4) is identifiable by a design  $\mathcal{D}$  if the functions  $p_\alpha$  ( $\alpha \in \mathcal{M}$ ) are linearly independent elements of  $\mathcal{L}(\mathcal{D})$ .*

The classical notion of identifiability is equivalent to our definition. Indeed, let  $Y = X\theta + \varepsilon$  be a linear model where  $X$  is a matrix with  $p$  columns. If the design matrix  $X$  has rank less than  $p$ , then  $\theta$  is not identifiable since different values of  $\theta$  yield the same value of  $X\theta$ . This actually means that the model coincides for different parameter values when restricted to the design points. In other words, the functions on  $\mathcal{D}$  that take as values the components of the columns of  $X$  are linearly dependent.

For linear models, least squares estimation is the orthogonal projection of  $g$  onto  $\text{span}\{p_\alpha \mid \alpha \in \mathcal{M}\}$ . Note that if  $\{p_\alpha \mid \alpha \in \mathcal{M}\}$  is an orthogonal subset of  $\mathcal{L}(\mathcal{D})$ , then elementary linear algebra arguments yield that

$$\hat{\theta}_\alpha = \frac{\langle g, p_\alpha \rangle_{\mathcal{D}}}{\langle p_\alpha, p_\alpha \rangle_{\mathcal{D}}}. \quad (1.5)$$

The functional description of least squares estimation has some advantages over the usual vector space description. It is more natural in our opinion since the model description is also at a functional level. A numerical advantage is that we do not need matrix inversion to compute the coefficient estimates. Indeed, orthogonalisation by the Gram-Schmidt procedure becomes a simple recursive procedure. Note that contrary to the classical use of Gram-Schmidt in the case of  $\mathbf{R}^N$ , we use Gram-Schmidt in a symbolic way in the space of polynomials. In this polynomial setting rewriting the estimated orthogonalised model in terms of the original model corresponds simply to collect coefficients.

The functional description given here differs from the abstract setting to linear models initiated by Kruskal (1961). See Drygas (1970) for a self-contained treatment. Specifically we extensively use computations with polynomials in the next sections. A paper which is closer in spirit to our paper is Neumaier and Seidel (1992), where a design is seen as a normalized measure and optimal designs are derived using arguments in  $\mathcal{L}(\mathcal{D})$ .

## 4 A polynomial algebraic representation of $\mathcal{L}(D)$

The set of real-valued functions over a finite set of distinct points can be described using particular classes of polynomials. More precisely let  $\mathcal{D}$  be a design in  $\mathbf{R}^d \times \mathcal{L}$ , let  $\mathbf{R}[x_1, \dots, x_{d+1}]$  be the polynomial ring in  $d + 1$  indeterminates with real coefficients and let  $\text{Ideal}(\mathcal{D}) \subset \mathbf{R}[x_1, \dots, x_{d+1}]$  be the set of all polynomials whose zeros include the design points. Then the quotient space  $\mathbf{R}[x_1, \dots, x_{d+1}]/\text{Ideal}(\mathcal{D})$  is a description or representation of  $\mathcal{L}(D)$ . Moreover vector space bases of  $\mathbf{R}[x_1, \dots, x_{d+1}]/\text{Ideal}(\mathcal{D})$  made of monomials can be determined with Gröbner basis methods. We require the definition of a term ordering.

**Definition 4.1** A term ordering  $\tau$  on the monomials of  $\mathbf{R}[x_1, \dots, x_d]$  is a total well-ordering such that  $x^\alpha \prec_\tau x^\beta$  implies  $x^\alpha x^\gamma \prec_\tau x^\beta x^\gamma$  for all  $\gamma \neq 0$ .

**Theorem 4.2** Given a design  $\mathcal{D} \subset \mathbf{R}^d \times \mathcal{L}$ , a term ordering  $\tau$  and a Gröbner basis  $G \subset \mathbf{R}[x_1, \dots, x_{d+1}]$  for  $\mathcal{D}$  with respect to  $\tau$ , then a vector space basis of  $\mathbf{R}[x_1, \dots, x_{d+1}]/\text{Ideal}(\mathcal{D})$  is given by

$$\begin{aligned} \text{Est}_{\mathcal{D},\tau} &:= \{x^\alpha \mid x^\alpha \text{ is not divisible} \\ &\quad \text{by any of the leading terms of the elements of } G\} \\ &= \{x^\alpha \mid \alpha \in L_{\mathcal{D},\tau}\}. \end{aligned}$$

Moreover, if the set  $\{p_\alpha \mid \alpha \in \mathcal{M}\}$  in Model (1.4) is a subset of  $\text{Est}_{\mathcal{D},\tau}$ , then Model (1.4) is identifiable. The set  $\text{Est}_{\mathcal{D},\tau}$  has exactly  $N$  elements where  $N$  is the cardinality of  $\mathcal{D}$ .

**Proof.** For the first part see for example Cox et al. (1996) and for the second and third parts see Pistone, Riccomagno and Wynn (2000). ■

Note that Theorem 4.2 applies to any design with no replicates, namely to a set of distinct points. Designs defined according to Definition 2.1 are particular examples of sets of distinct points where the “label indeterminate”,  $x_{d+1}$  distinguishes replicated points. For designs with replicates the trick here is to consider in Model (1.4) only terms of  $\text{Est}$  not involving  $x_{d+1}$ .

For statistical inference we need a design, a model, and observations. In a classical screening setup a model is chosen first. However, we may also choose the model after seeing the design (for example the planned design was not completed and there are missing points, see Holliday et al., 1999). In this case, Theorem 4.2 provides a powerful tool in the choice of a regression vector for a linear model of the type in (1.4).

In general different term orderings give different  $\text{Est}$  sets and also typically  $\text{Est}$  includes monomials involving the label indeterminate  $x_{d+1}$  which clearly should not be included in Model (1.4). This suggests to partition  $\text{Est}$ , equivalently  $L$ , in three disjoint parts

$$L_{\mathcal{D},\tau} = L_x^* \cup L_{x_{d+1}} \cup L_{x,x_{d+1}}$$

where  $L_x^*$  includes all the elements of  $L_{\mathcal{D},\tau}$  that do not involve  $x_{d+1}$ ,  $L_{x_{d+1}}$  includes all the monomials in  $\text{Est}_{\mathcal{D},\tau}$  which involve only  $x_{d+1}$  and  $L_{x,x_{d+1}}$  includes the remaining terms. The set  $\mathcal{M}$  in Model (1.4) can be chosen to be a subset of  $L_x^*$ . The combination of the choice of the term ordering and of the structure of the design determines these three parts.

A reasonable choice for the term ordering is one that eliminates the  $x_{d+1}$  variable. For elimination term ordering we refer to Cox et al. (1996) and here simply observe that an effect of eliminating  $x_{d+1}$  is that the number of monomials in  $L_{x_{d+1}}$  is as small as possible.

We conclude this section by showing with polynomial algebra techniques that identifiability is not affected by replications. The elimination of  $x_{d+1}$

from  $\text{Ideal}(\mathcal{D})$  corresponds to projecting  $\mathcal{D} \subset \mathbf{R}^d \times \mathcal{L}$  onto  $\mathbf{R}^d$ . For some term orderings the Gröbner basis of  $\text{Ideal}(\mathcal{D}^*)$  and  $\text{Est}_{\mathcal{D}^*}$  can be easily deduced from the Gröbner basis of  $\text{Ideal}(\mathcal{D})$  and  $\text{Est}_{\mathcal{D}}$ .

**Theorem 4.3** 1)  $\text{Ideal}(\mathcal{D}) \cap \mathbf{R}[x_1, \dots, x_d] = \text{Ideal}(\mathcal{D}^*)$ . 2) Let  $G$  be the Gröbner basis of  $\text{Ideal}(\mathcal{D})$  with respect to a term ordering eliminating  $x_{d+1}$ . The Gröbner basis of  $\text{Ideal}(\mathcal{D}^*)$  is  $G \cap \mathbf{R}[x_1, \dots, x_d]$ .

**Proof.** 1) Assume  $f \in \text{Ideal}(\mathcal{D}) \cap \mathbf{R}[x_1, \dots, x_d]$ . Then for all  $a = (a^*, \ell) \in \mathcal{D}$ , we have that  $0 = f(a) = f(a^*, \ell) = f(a^*)$  as  $f \in \mathbf{R}[x_1, \dots, x_d]$ . This implies  $f \in \text{Ideal}(\mathcal{D}^*)$ . The converse is obvious. 2) See Cox et al. (1996). ■

Clearly Theorem 4.3 applies when instead of  $x_{d+1}$  we need to eliminate some other variable. The projection is now on the remaining variables and replications may appear. For example the projection of the  $2^2$  design at levels  $\pm 1$  over the first factor gives  $\pm 1$  replicated twice.

## 5 Examples

The analysis of observations suggested in the paper proceeds as follows. Given a design  $\mathcal{D}$ , compute  $\text{Est}_{\mathcal{D}, \tau}$  where  $\tau$  is a term ordering that eliminates the extra variable  $t$ . Orthonormalise the terms of  $\text{Est}_{\mathcal{D}, \tau}$  that do not involve  $t$ . Collect coefficients to determine the parameters of the wanted model from the estimated coefficients in the orthonormalised model.

### Example 5.1 ( $2^2$ full factorial design with centre points)

Consider the  $2^2$  design at levels  $\pm 1$ . The standard model associated with it is

$$Y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_{12} x_1 x_2.$$

Clearly  $1$ ,  $x_1^2$  and  $x_2^2$  are confounded on  $\mathcal{D}$ . Suppose that we want to test linearity by adding quadratic terms to the model. The simplest way to extend this design such that quadratic terms become identifiable, is to add centre points. We add four observations at  $(0, 0)$  and the design becomes

$$\begin{aligned} \mathcal{D} &= \{(-1, -1, 1), (-1, 1, 1), (1, -1, 1), (1, 1, 1), \\ &\quad (0, 0, 1), (0, 0, 2), (0, 0, 3), (0, 0, 4)\} \\ &= \{a_i \mid i = 1, \dots, 8\}. \end{aligned}$$

We use a term ordering  $\sigma$  that eliminates the variable  $t$  and is a degree reverse lexicographic ordering on  $x_1$  and  $x_2$  (Cox et al., 1996). We obtain

$$\text{Est}_{\mathcal{D}} = \{1, x_1, x_2, x_1 x_2, x_2^2, t, t^2, t^3\}.$$

An orthonormal basis for the linear span of the terms not involving  $t$  is

$$\left\{ \frac{1}{\sqrt{8}}, \frac{x_1}{2}, \frac{x_2}{2}, \frac{x_1 x_2}{2}, \frac{x_2^2 - \frac{1}{2}}{\sqrt{2}} \right\}.$$

Let  $g$  be the polynomial that interpolates the observations  $Y_1, \dots, Y_8$  at the design points,  $g(a_i) = Y_i$ ,  $i = 1, \dots, 8$ . The traditional sums of squares correspond to the squares of the inner products. In particular, using Equation (1.5) the average over the centre points minus the average over the full factorial is computed as

$$SS_{\text{pure quadratic}} = \frac{1}{8} (Y_1 + \dots + Y_4 - (Y_5 + \dots + Y_8))^2 = \left\langle g, \frac{x_2^2 - \frac{1}{2}}{\sqrt{2}} \right\rangle_{\mathcal{D}}^2.$$

By simple comparison of terms, we read off the coefficients of the wanted model from the coefficients of the orthonormalised model.

**Example 5.2 (Star composite design with centre points)**

If the analysis of the  $2^2$  full factorial design with centre points indicates that there is curvature, then it is practice to study the quadratic terms. We choose to break the aliasing by augmenting the design with four axial points at  $(0, \pm 2)$  and  $(\pm 2, 0)$ . The new design  $\mathcal{D}$  is given below

$$\begin{aligned} \mathcal{D} = \{ & (-1, -1, 1), (-1, 1, 1), (1, -1, 1), (1, 1, 1), (-2, 0, 1), (2, 0, 1), \\ & (0, 2, 1), (0, -2, 1), (0, 0, 1), (0, 0, 2), (0, 0, 3), (0, 0, 4) \}. \end{aligned}$$

We use again the term ordering  $\sigma$  used in Example 5.1. The monomials in  $Est_{\mathcal{D}}$  not involving the counting variable are

$$\{1, x_1, x_2, x_1^2, x_2^2, x_1 x_2, x_2^3, x_1 x_2^2, x_2^4\}$$

An orthonormal basis for the linear span of these terms by applying the Gram-Schmidt procedure to  $Est_{\mathcal{D}}$  in the order above, yields

$$\left\{ \frac{\sqrt{3}}{6}, \frac{\sqrt{3}}{6} x_1, \frac{\sqrt{3}}{6} x_2, \frac{\sqrt{6}}{12} (x_1^2 - 1), \frac{\sqrt{3}}{24} (x_1^2 + 3x_2^2 - 4), \right. \\ \left. \frac{x_1 x_2}{2}, \frac{\sqrt{6}}{12} x_1 (3x_2^2 - 1), \frac{\sqrt{6}}{12} x_2 (x_2^2 - 3) \right\}.$$

**Example 5.3 ( $2^{3-1}$  fractional factorial design with centre points)**

Consider the standard  $2^{3-1}$  design with generator  $I = ABC$  and four additional centre points. The design is

$$\begin{aligned} \mathcal{D} = \{ & (1, -1, -1, 1), (-1, 1, -1, 1), (-1, -1, 1, 1), (1, 1, 1, 1), \\ & (0, 0, 0, 1), (0, 0, 0, 2), (0, 0, 0, 3), (0, 0, 0, 4) \}. \end{aligned}$$

Again using the term ordering  $\sigma$ , we obtain  $Est_{\mathcal{D}} = \{1, x_1, x_2, x_3, x_3^2, t, t^2, t^3\}$ . An orthonormal basis for the linear span of the terms not involving  $t$  is

$$\left\{ \frac{1}{\sqrt{8}}, \frac{x_1}{2}, \frac{x_2}{2}, \frac{x_3}{2}, \frac{x_3^2 - \frac{1}{2}}{\sqrt{2}} \right\}.$$

## References

- Caboara, M. and Robbiano, L. (1997). Families of ideals in statistics. In: Küchlin, W. (ed). *ISSAC'97, Proceedings of the International Symposium on Symbolic and Algebraic Computation, Hawaii*, pp. 404–417. ACM Press, New York.
- Cox, D., Little, J. and O'Shea, D. (1996). *Ideals, Varieties, and Algorithms*. Springer, New York. Second edition.
- Drygas, H. (1970). *The Coordinate-Free approach to Gauss-Markov Estimation*. Springer-Verlag, Berlin.
- Giglio, B., Riccomagno, E. and Wynn, H. P. (2000). Gröbner bases in regression. *Journal of Applied Statistics*, **27**, 923-928.
- Holliday, T., Pistone, G., Riccomagno, E. and Wynn, H. P. (1999). The application of computational algebraic geometry to the analysis of designed experiments: a case study. *Computational Statistics*, **14**, 213-231.
- Kruskal, W. (1961). The coordinate-free approach to Gauss-Markov estimation, and its application to missing and extra observations. In: Neyman, J. (ed). *Proceedings of the 4th Berkeley Symposium on Mathematical Statistics and Probability, Vol. I*, pp. 435–451. University of California Press, Berkeley.
- Neumaier, A. and Seidel, J.J. (1992). Measures of strength  $2e$  and optimal designs of degree  $e$ . *Sankhyā*, **54**, 299-309.
- Pistone, G., Riccomagno, E. and Wynn, H.P. (2000). *Algebraic Statistics: Computational Commutative Algebra in Statistics*. Chapman & Hall / CRC Press, London.
- Pistone, G. and Wynn, H.P. (1996). Generalised confounding with Gröbner bases. *Biometrika*, **83**, 653-666.
- Riccomagno, E. (1997). Algebraic identifiability in experimental design and related topics. Ph.D. thesis, University of Warwick.