

Construction and analysis of Es^2 efficient supersaturated designs

Yufeng Liu^a, Shiling Ruan^b, Angela M. Dean^{b,*}

^a*Department of Statistics and Operations Research, Carolina Center for Genome Sciences, University of North Carolina, Chapel Hill, NC 27599, USA*

^b*Department of Statistics, The Ohio State University, 1958 Neil Avenue, Columbus, OH 43210, USA*

Available online 4 October 2006

Abstract

In this paper, we construct supersaturated designs for large numbers of two-level factors and $10 \leq n \leq 22$ runs by augmenting k -circulant designs [Liu, Y., Dean, A.M., 2004. k -circulant supersaturated designs. *Technometrics* 46, 32–43] with interaction columns or by deleting columns from k -circulant designs. Most of the designs presented have Es^2 efficiencies above 0.90 and they extend the range of efficient supersaturated designs available in the literature.

Difficulties encountered in the use of supersaturated designs in detecting active factors are addressed. We show that, when only one factor is active, the regression technique of forward selection is guaranteed to select the correct factor as active under the idealized conditions that non-active factors have negligible effects and the errors are small. Under similar conditions, we derive bounds on the maximum allowable correlation between the columns of the model matrix that guarantee the correct selection of the “most active” factor when two or more factors are non-negligible. Further, we obtain conditions for the correct selection of the two most active factors using subset selection in regression. A number of designs that satisfy these conditions are identified.

© 2006 Elsevier B.V. All rights reserved.

MSC: Primary 62K15; Secondary 62K05

Keywords: Augmented design; Correlation; Cyclic generator; Efficient design; Interaction; k -circulant design

1. Introduction

Liu and Dean (2004) introduced a class of supersaturated designs called *k-circulant supersaturated designs* which can be obtained very simply from row generators, by cycling k elements at a time. The designs, which are generalizations of the saturated Plackett–Burman designs (Plackett and Burman, 1946) have $n = 2t$ runs and $m = k(n - 1)$ two-level factors. The k -circulant class includes the cyclic balanced incomplete block (cyclic BIBD) based designs of Nguyen (1996), Liu and Zhang (2000) and Eskridge et al. (2004) as special cases.

We represent a supersaturated design d , having m factors and n treatment combinations (runs), by its $n \times m$ design matrix T_d , where the (i, j) th element of T_d denotes the level of the j th factor in the i th treatment combination. The treatment combinations need to be randomly ordered before the experiment is run.

* Corresponding author. Tel.: +1 614 292 0292; fax: +1 614 292 2096.

E-mail address: amd@stat.ohio-state.edu (A.M. Dean).

Throughout this paper, we use the main effects model

$$Y = \mu \mathbf{1} + X\beta + \epsilon, \quad \epsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n), \quad (1.1)$$

where $\mathbf{1}$ is a vector of 1's corresponding to the mean parameter μ , and $X = \{x_{ij}\}$ is the model matrix without the intercept column. The j th column of X is denoted by x_j and represents the main effect contrast between the levels of factor j corresponding to the j th element of the parameter vector β ($j = 1, 2, \dots, m$).

We consider only designs that are “mean-orthogonal”; that is, every factor is to be observed t times at each of its high and low levels (This property is often called “balanced” in the literature). The high and low levels of each factor are coded as +1 and -1, respectively. With this particular coding in the design matrix, T_d , and in the absence of interactions in the model, the model matrix X without the intercept column of 1's has the same form as the design matrix. However, since the two matrices play different roles (one for the design and the other for the model), we will keep the distinction in this paper.

Construction of k -circulant designs, in common with other construction methods for efficient supersaturated designs described in the literature, is restricted to limited values of m and n . A few authors have considered adding extra columns to, or deleting columns from, “base” supersaturated designs. For example, Cheng (1997) obtained Es^2 optimal designs for $n = 8$ runs and $m = 8, \dots, 35$ factors by adding columns to, or deleting columns from, cyclic BIBD-based supersaturated designs, where Es^2 is defined as $\sum_{i < j} s_{ij}^2 / \binom{m}{2}$, and s_{ij} is the $\{i, j\}$ th element of $T_d' T_d$. Wu (1993) suggested an alternative approach of appending “interaction” columns to the design matrices of Plackett–Burman designs. The (i, j) th interaction column, $c_{i,j}$, is formed by the product of corresponding elements in the i th and j th columns, c_i and c_j , of T_d . If $c_{i,j}$ is distinct from the columns already present in T_d , it can be appended to T_d and the augmented design can be used to measure the main effect of an additional factor.

Following this idea of Wu (1993), Liu and Dean (2004) discussed conditions under which interaction columns could be appended to any base k -circulant design and gave one example. In Section 2, we explore this idea further and give efficient sequential orderings for addition of interaction columns to k -circulant designs with 12, 16 and 20 runs. In Section 3, we present efficient sequential orderings of deletion of columns from k -circulant designs to obtain a wider range of supersaturated designs with 10, 14, 18 and 22 runs. As in Liu and Dean (2004), we use Es^2 optimality as our major criterion. Most of the designs presented have Es^2 efficiencies of at least 0.9 and many are close to 1.0. The efficiencies are calculated relative to the bounds of Butler et al. (2001) or Bulutoglu and Cheng (2004) as appropriate. If two designs have the same value of Es^2 , we use $r_{\max} = S_{\max}/n$ as a secondary criterion to distinguish designs, where $S_{\max} = \max |s_{ij}|$. Our range of designs extends that previously available in the literature and, where overlap exists, we identify previously available Es^2 optimal designs.

Srivastava (1975) showed that, for any given design, any set of p contrasts can be estimated if every subset of $2p$ columns in the model matrix X contains linearly independent columns. Chen and Lin (1998) showed further that, to be able to estimate any set of p contrasts in a supersaturated design with factors at two levels, the absolute value of the correlation between any two columns in the model matrix can be at most $1/(p-1)$ (and less than $1/(p-1)$ for p even). Even when these conditions are satisfied, it is well known that the main drawback in the use of designs with such small numbers of observations is the difficulty in identifying correctly which factors are active (that is, have a substantial effect on the response). Examples highlighting these difficulties have been given, for example, by Abraham et al. (1999) using forward selection and best subset selection.

In Section 4, we discuss the effect of correlations between the columns of the model matrix, X , on the correct selection of active factors using the methods of forward selection and best subset selection which are based on the same underlying formulae. Chen and Lin (1998, Section 3) explored the effect of the correlation and error variability on the correct identification of the most active factors, but did not investigate relationship between the maximum allowable correlation and the relative magnitude of the main effects of the active factors. Here, we assume that the definition of “active” is that the main effect of the active factor is considerably larger than the error standard deviation so that the effect of error is minor. Under the idealized conditions that non-active factors have negligible effects, we show that, when only one factor is active, forward selection is guaranteed to select the correct factor as active. Under similar conditions, we derive bounds on the maximum allowable correlation between the columns of the model matrix that guarantee the correct selection of the “most active” factor, when two or more factors are non-negligible. Further, we obtain conditions for the correct selection of the two most active factors using subset selection in regression. The various conditions are satisfied by a large number of the k -circulant designs of Liu and Dean (2004), many of the

designs constructed in this paper, as well as other supersaturated designs in the literature. Some of these are listed in Section 4.2.

2. Construction of Es^2 efficient designs via augmentation

Suppose there are m^I candidate interaction columns for possible augmentation of a base k -circulant supersaturated design d with m_1 factors and design matrix T_d . To obtain an “augmented” design d_a with $m = m_1 + m_2$ factors, we select m_2 of the interaction columns from the m^I candidates to append to the m_1 columns of T_d .

We append interaction columns to a mean-orthogonal design matrix only if the design has run size $n = 0 \pmod 4$. Otherwise, the added factors would be observed different numbers of times at their high and low levels (and $X \neq T_d$). Theorem 5.1 of Liu and Dean (2004) gives an efficient method of identifying the interaction columns that are candidates for appending to a base k -circulant design matrix so that no candidate interaction column is identical to any column in the base design and all candidate interaction columns are distinct. They give an example for $n = 8$ runs in which Es^2 optimal designs for $m = 14, 15, \dots, 35$ factors are constructed by appending the optimal selection of up to seven interaction columns to k -circulant design matrices for $k = 2, 3, 4, 5$.

When m^I is large, the number of possible selections of m_2 columns out of m^I is extensive and an exhaustive search becomes computationally intensive and even prohibitive to use. To solve this computational difficulty, we follow the method of Wu (1993) and select the interaction columns sequentially. The computation for this approach requires a search over only $m_2(2m^I - m_2 + 1)/2$ augmented designs. Although in the sequential search we only consider a small proportion of all possible designs, our results show that this method works well, resulting in highly Es^2 efficient supersaturated designs.

Wu (1993) suggested a specific sequential order of addition of interaction columns to base Plackett–Burman designs with $n = 12$ and 20 runs. Plackett–Burman designs are 1-circulant designs with $m_1 = n - 1$ factors and, thus, his suggested order is one possible solution of our sequential search. Using random starts, we have been able to improve upon his order for $n = 12$ (see Section 2.1). Plackett–Burman designs with $n = 0 \pmod 8$ are fractional factorial designs in which all interaction columns are completely aliased with main effect columns and so the base designs cannot be augmented in this manner. However, k -circulant designs with $k > 1$ can be used for obtaining designs with the number of runs equal to $n = 0 \pmod 8$. In the following subsections, we examine the augmentation of designs having 12, 16, 20 runs. In our tables, we use the notation (k, n, d) to refer to the designs in the tables of Liu and Dean (2004). For example, $(k3n12d2)$ corresponds to the second recommended design with $k = 3$ and $n = 12$ in their Table 3.

2.1. Designs with $n = 12$ runs

The Plackett–Burman design with $n = 12$ runs is a 1-circulant design with row generator $(-1 -1 1 -1 -1 -1 1 1 1 -1 1 1)$. Wu (1993) pointed out that all 55 interaction columns can be appended to this base design to obtain Es^2 efficient designs with up to 66 factors and he suggested the following order

$$c_{1,2}, c_{1,3}, \dots, c_{1,12}, c_{2,3}, \dots, c_{2,12}, \dots, c_{11,12}. \quad (2.1)$$

Under the restriction of $r_{\max} = 0.33$, our search has revealed an alternative order which, although it cannot be expressed as neatly as (2.1), does improve upon the Es^2 efficiency for $m = 39, \dots, 62$ factors, where the efficiencies are calculated relative to the bound of Butler et al. (2001, Theorem 1) or, equivalently, that of Bulutoglu and Cheng (2004, Theorem 3.1). Our suggested order of addition of interaction columns $c_{i,j}$ is listed in the upper part of Table 1. All of these designs have Es^2 efficiency of at least 0.833. For $m \geq 23$, the efficiencies of our augmented k -circulant designs are always above 0.90 and, for $m \geq 38$ factors, are plotted as the dotted curve in Fig. 1. We note that Butler et al. (2001) list Es^2 optimal designs for the range $m = 13, \dots, 22$.

It is also possible to perform a sequential search for 12-run designs using the k -circulant designs with $k = 2, \dots, 6$ as the base designs. For example, using the 6-circulant design $(k6n12d2)$, we are able to obtain 12-run designs for $67 \leq m \leq 231$ factors with $r_{\max} = 0.67$. We have tabulated recommended sequential orderings for appending interaction columns to design $(k6n12d2)$ in the lower part of Table 1 and have shown the Es^2 efficiencies as the dashed curve in Fig. 1. Note that, although these designs have very high Es^2 efficiencies, they may not be optimal since our search was not exhaustive (even among the sequential orderings).

Table 1
Sequential order of augmentation by interaction columns of designs ($k1n12d1$) and ($k6n12d2$) of Liu and Dean (2004)

| Design | Recommended interaction order |
|--|--|
| $k1n12d1$ $r_{\max} = 0.33$ $13 \leq m \leq 66$ | $C_{1-2} C_{2-3} C_{1-3} C_{3-4} C_{2-4} C_{1-4} C_{4-5} C_{3-5} C_{2-5} C_{1-5} C_{5-6} C_{4-6} C_{3-6} C_{2-6} C_{1-6} C_{6-7} C_{5-7} C_{4-7}$ $C_{3-7} C_{2-7} C_{7-1} C_{7-8} C_{6-8} C_{5-8} C_{4-8} C_{8-1} C_{3-8} C_{8-2} C_{8-9} C_{7-9} C_{6-9} C_{9-1} C_{5-9} C_{9-2} C_{4-9} C_{9-3}$ $C_{9-10} C_{8-10} C_{10-1} C_{7-10} C_{10-2} C_{6-10} C_{10-3} C_{5-10} C_{10-4} C_{10-11} C_{11-1} C_{9-11} C_{11-2} C_{8-11} C_{11-3}$ $C_{7-11} C_{11-4} C_{6-11} C_{11-5}$ |
| $k6n12d2$ $r_{\max} = 0.67$ $67 \leq m \leq 231$ | $C_{1-2} C_{7-8} C_{13-14} C_{19-20} C_{19-23} C_{31-35} C_{19-24} C_{61-1} C_{7-15} C_{31-32} C_{31-36} C_{1-19} C_{43-1} C_{49-1}$ $C_{1-6} C_{49-54} C_{13-18} C_{61-66} C_{43-48} C_{13-58} C_{43-49} C_{19-25} C_{19-37} C_{7-52} C_{13-17} C_{61-5} C_{43-44}$ $C_{43-22} C_{1-5} C_{31-38} C_{49-62} C_{55-65} C_{13-19} C_{13-23} C_{25-32} C_{55-62} C_{55-61} C_{61-62} C_{7-20} C_{61-65}$ $C_{25-43} C_{19-32} C_{31-44} C_{49-53} C_{43-61} C_{25-40} C_{55-63} C_{37-45} C_{49-50} C_{25-35} C_{55-11} C_{25-31} C_{7-31}$ $C_{37-61} C_{55-60} C_{25-30} C_{7-25} C_{37-55} C_{37-59} C_{25-49} C_{49-56} C_{55-56} C_{61-8} C_{1-7} C_{37-41} C_{37-38}$ $C_{31-39} C_{43-47} C_{1-11} C_{13-20} C_{13-21} C_{25-26} C_{19-27} C_{7-13} C_{13-31} C_{31-37} C_{61-3} C_{61-18} C_{25-4}$ $C_{1-24} C_{61-2} C_{25-33} C_{55-2} C_{49-64} C_{37-43} C_{1-8} C_{43-50} C_{7-14} C_{31-41} C_{13-35} C_{55-7} C_{49-7} C_{43-53}$ $C_{1-16} C_{25-38} C_{49-55} C_{55-59} C_{7-17} C_{7-12} C_{19-29} C_{19-64} C_{25-29} C_{37-42} C_{49-57} C_{31-49} C_{1-46}$ $C_{49-59} C_{13-50} C_{13-26} C_{19-26} C_{37-50} C_{13-28} C_{25-48} C_{1-38} C_{1-14} C_{61-19} C_{61-13} C_{37-47} C_{37-60}$ $C_{43-51} C_{49-6} C_{43-56} C_{43-65} C_{1-9} C_{25-47} C_{55-26} C_{55-4} C_{61-17} C_{37-52} C_{37-44} C_{19-34} C_{49-5}$ $C_{7-22} C_{31-55} C_{61-10} C_{19-42} C_{7-29} C_{19-41} C_{31-53} C_{13-37} C_{37-8} C_{43-58} C_{31-54} C_{25-62} C_{31-46}$ $C_{13-36} C_{43-66} C_{1-23} C_{7-30} C_{19-56} C_{19-43} C_{61-32} C_{61-40} C_{55-34} C_{55-12} C_{1-25} C_{7-44} C_{43-14}$ $C_{55-13} C_{31-2} C_{49-20} C_{49-28} C_{7-11} C_{31-10} C_{37-16}$ |

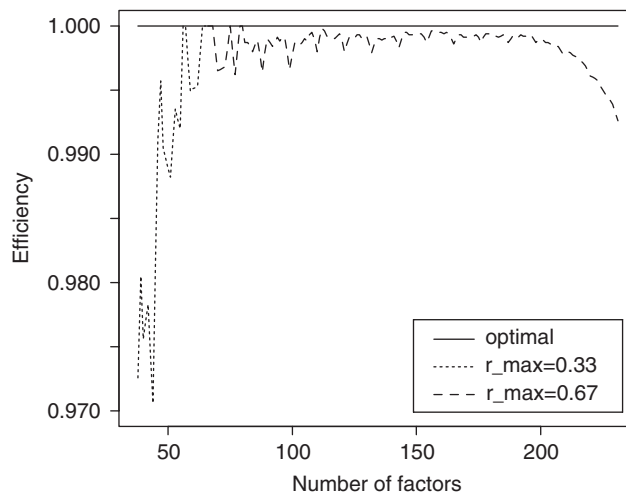


Fig. 1. The Es^2 efficiencies of the augmented designs listed in Table 1. The dotted and dashed curves correspond to base designs ($k1n12d1$) and ($k6n12d2$), with $r_{\max} = 0.33$ and 0.67 , respectively.

2.2. Designs with $n = 16$ runs

Since Es^2 optimal 1-circulant designs with $n = 16$ do not have interaction columns distinct from those of the base design, we examine augmentation of 16-run k -circulant designs with $k = 2, \dots, 6$. As listed in Table 2, we can obtain augmented designs with 31 to 90 factors and $r_{\max} = 0.5$ by augmenting ($k2n16d2$). However, in the range $45 \leq m \leq 66$ factors, more efficient 16-run designs with $r_{\max} = 0.5$ are obtained by augmenting ($k3n16d3$). We plot the Es^2 efficiencies for $30 \leq m \leq 44$ and $67 \leq m \leq 90$ of ($k2n16d2$) and $45 \leq m \leq 66$ of ($k3n16d3$) as the dotted curve in Fig. 2. These Es^2 efficiencies are all above 0.92. For $r_{\max} = 0.75$, design ($k3n16d4$) can be augmented from 50 to 255 factors and the Es^2 efficiencies (which are all above 0.956) are represented by the dashed curve in Fig. 2.

Table 2
 Sequential order of augmentation by interaction columns of designs ($k2n16d2$), ($k3n16d3$), and ($k3n16d4$) of Liu and Dean (2004)

| Design | Recommended interaction order |
|--|---|
| $k2n16d2$ $r_{\max} = 0.5$ $31 \leq m \leq 90$ | C1-3 C5-7 C17-19 C9-15 C29-5 C21-23 C7-9 C25-1 C3-5 C19-21 C23-25 C27-29 C9-11 C13-15 C5-11 C15-21 C13-19 C11-13 C19-25 C15-17 C25-27 C3-9 C23-29 C29-1 C1-7 C11-17 C21-27 C7-13 C17-23 C27-3 C1-9 C5-13 C17-25 C13-21 C9-17 C21-29 C25-3 C27-5 C23-1 C7-15 C11-19 C19-27 C3-11 C15-23 C29-7 C1-14 C13-26 C15-28 C29-12 C7-20 C21-4 C23-6 C5-18 C17-30 C3-16 C19-2 C9-22 C25-8 C11-24 C27-10 |
| $k3n16d3$ $r_{\max} = 0.5$ $46 \leq m \leq 66$ | C1-10 C7-16 C13-22 C19-28 C1-20 C22-41 C28-2 C34-8 C40-14 C19-2 C40-4 C16-35 C31-40 C37-1 C25-34 C19-9 C43-32 C25-15 C22-3 C37-26 C4-29 |
| $k3n16d4^a$ $r_{\max} = 0.75$ $50 \leq m \leq 255$ | C1-10 C7-16 C13-22 C19-28 C22-42 C40-7 C37-1 C19-38 C5-9 C10-19 C22-31 C43-7 C16-25 C1-19 C28-37 C25-34 C34-43 C31-27 C14-17 C10-28 C43-16 C22-18 C16-34 C7-25 C31-9 C5-8 C40-18 C2-5 C37-4 C16-35 C37-21 C43-17 C40-36 C11-14 C4-45 C25-9 C7-26 C13-9 C22-36 C4-13 C28-41 C19-32 C38-42 C7-5 C34-8 C37-33 C28-26 C31-29 C43-41 C25-44 C13-26 C31-15 C40-8 C28-45 C16-14 C34-32 C22-35 C4-17 C31-44 C37-35 C1-44 C40-4 C16-29 C28-2 C25-38 C13-30 C13-11 C20-23 C1-42 C29-32 C44-2 C28-24 C35-38 C17-20 C43-39 C31-40 C22-39 C37-9 C7-3 C10-6 C19-37 C19-15 C25-21 C19-36 C4-21 C7-20 C31-4 C7-36 C4-2 C10-8 C34-18 C34-2 C38-41 C10-29 C16-12 C10-39 C22-20 C34-30 C43-15 C10-23 C37-5 C1-14 C1-30 C43-11 C7-24 C16-33 C4-23 C13-42 C43-10 C7-30 C40-12 C23-26 C16-28 C40-14 C26-29 C25-23 C22-6 C32-35 C31-3 C19-17 C34-6 C22-41 C19-3 C40-38 C25-42 C13-32 C41-44 C10-27 C8-11 C19-31 C37-11 C28-40 C31-5 C7-19 C40-24 C43-27 C26-30 C4-16 C1-18 C16-45 C25-37 C31-43 C13-25 C1-13 C4-33 C10-22 C34-1 C25-3 C22-45 C22-34 C25-43 C13-36 C34-7 C1-20 C28-12 C11-15 C37-10 C41-45 C4-22 C43-21 C37-15 C13-31 C8-12 C32-36 C22-40 C1-24 C28-1 C40-13 C20-24 C7-21 C17-21 C29-33 C16-30 C37-6 C10-33 C16-39 C23-27 C28-42 C31-45 C40-9 C1-15 C34-12 C25-39 C43-12 C4-27 C10-24 C13-27 C19-42 C19-33 C34-3 C2-6 C44-3 C14-18 C28-6 C35-39 C4-18 C1-21 C7-27 C28-3 C37-12 C13-33 C31-6 C16-36 C43-18 C4-24 C25-45 C10-30 C34-9 C19-39 C40-15 |

^aThe augmented designs have $r_{\max} = 0.75$ for $m \geq 50$, i.e., augmenting $m_1 \geq 5$ interactions.

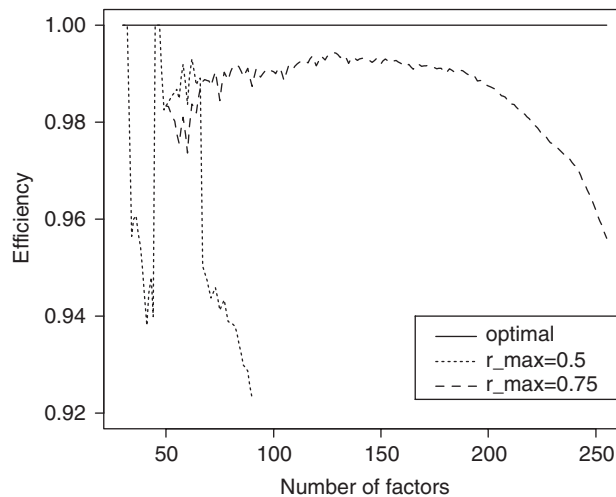


Fig. 2. The Es^2 efficiencies of the augmented designs listed in Table 2. The dotted curve corresponds to the augmented designs with $30 \leq m \leq 44$ and $67 \leq m \leq 90$ for ($k2n16d2$) and $45 \leq m \leq 66$ for ($k3n16d3$), all with $r_{\max} = 0.5$. The dashed curve corresponds to ($k3n16d4$) with $50 \leq m \leq 255$ and $r_{\max} = 0.75$.

2.3. Designs with $n = 20$ runs

For $n=20$ runs, Wu (1993) augmented a Plackett–Burman 1-circulant design having generator $(-1\ 1\ -1\ -1\ 1\ 1\ 1\ 1\ -1\ 1\ -1\ 1\ -1\ -1\ -1\ -1\ 1\ 1\ 1\ -1)$ to obtain designs with up to $m = 124$ factors with $r_{\max} = 0.6$. In fact, all two-factor interactions can be appended, so this design can be augmented up to $m = 190$ factors. Similar to $n = 12$, Wu’s order can

Table 3
Sequential order of augmentation by interaction columns of designs ($k2n20d2$) and ($k2n20d4$) of Liu and Dean (2004)

| Design | Recommended interaction order |
|---|---|
| $k2n20d2$ $r_{\max} = 0.4$ $39 \leq m \leq 60$ | $C_{1-6} C_{3-8} C_{5-10} C_{7-12} C_{17-22} C_{29-34} C_{31-5} C_{15-20} C_{27-32} C_{5-17} C_{11-16} C_{23-28} C_{13-18} C_{35-2}$ $C_{19-24} C_{33-38} C_{37-4} C_{9-14} C_{31-36} C_{21-26} C_{25-30} C_{17-24}$ |
| $k2n20d4$ $r_{\max} = 0.6$ $40 \leq m \leq 266$ | $C_{1-9} C_{5-13} C_{17-25} C_{21-29} C_{9-21} C_{5-17} C_{13-25} C_{27-1} C_{35-5} C_{29-6} C_{5-26} C_{1-36} C_{13-34} C_{31-28} C_{2-10}$ $C_{27-35} C_{10-28} C_{3-15} C_{37-14} C_{35-18} C_{19-27} C_{31-5} C_{9-17} C_{37-7} C_{21-4} C_{13-30} C_{3-24} C_{15-2} C_{21-36}$ $C_{29-37} C_{11-23} C_{21-38} C_{19-31} C_{23-20} C_{19-28} C_{7-15} C_{1-22} C_{23-10} C_{1-16} C_{25-5} C_{7-19} C_{1-13} C_{7-32}$ $C_{7-22} C_{31-8} C_{7-25} C_{17-4} C_{31-14} C_{35-9} C_{31-10} C_{15-33} C_{13-31} C_{24-32} C_{9-24} C_{29-8} C_{23-3} C_{11-29}$ $C_{11-26} C_{27-14} C_{35-15} C_{27-7} C_{23-6} C_{21-33} C_{33-13} C_{5-23} C_{37-24} C_{3-21} C_{25-12} C_{11-19} C_{15-27} C_{38-18}$ $C_{11-36} C_{9-26} C_{23-35} C_{31-1} C_{29-9} C_{1-26} C_{23-38} C_{15-36} C_{24-36} C_{29-3} C_{17-32} C_{17-38} C_{29-16} C_{23-31}$ $C_{13-21} C_{15-30} C_{19-2} C_{33-7} C_{37-11} C_{13-28} C_{13-38} C_{21-1} C_{37-20} C_{25-33} C_{27-10} C_{33-30} C_{25-2} C_{33-4}$ $C_{29-12} C_{19-6} C_{9-18} C_{35-6} C_{37-17} C_{22-2} C_{3-28} C_{9-27} C_{3-11} C_{31-18} C_{9-34} C_{17-35} C_{7-28} C_{21-8} C_{18-26}$ $C_{25-8} C_{33-3} C_{33-20} C_{15-23} C_{11-32} C_{19-34} C_{33-10} C_{23-2} C_{34-8} C_{27-4} C_{8-16} C_{20-32} C_{35-22} C_{17-29}$ $C_{25-37} C_{5-20} C_{3-18} C_{19-37} C_{36-10} C_{31-11} C_{1-19} C_{3-12} C_{33-16} C_{2-20} C_{16-34} C_{5-30} C_{6-24} C_{18-30}$ $C_{29-26} C_{27-6} C_{9-30} C_{1-10} C_{35-12} C_{21-18} C_{11-8} C_{5-2} C_{13-10} C_{1-18} C_{11-28} C_{8-26} C_{37-8} C_{26-38}$ $C_{19-16} C_{5-14} C_{17-34} C_{15-12} C_{30-4} C_{36-6} C_{21-30} C_{30-10} C_{27-36} C_{7-24} C_{14-26} C_{31-2} C_{3-20} C_{9-6}$ $C_{22-30} C_{35-14} C_{13-22} C_{7-4} C_{27-24} C_{25-34} C_{25-22} C_{25-4} C_{15-32} C_{15-24} C_{17-14} C_{35-32} C_{12-24} C_{5-22}$ $C_{11-20} C_{3-38} C_{29-38} C_{12-30} C_{8-20} C_{19-36} C_{12-20} C_{23-32} C_{32-12} C_{33-12} C_{26-6} C_{28-8} C_{7-16} C_{37-16}$ $C_{28-36} C_{2-14} C_{37-34} C_{34-14} C_{34-4} C_{17-26} C_{6-14} C_{30-38} C_{28-2} C_{22-34} C_{36-16} C_{38-8} C_{14-22} C_{32-6}$ $C_{4-12} C_{32-2} C_{10-18} C_{4-22} C_{16-28} C_{20-28} C_{4-16} C_{10-22} C_{14-32} C_{16-24} C_{24-4} C_{26-34} C_{6-18}$ $C_{38-12} C_{18-36} C_{20-38}$ |

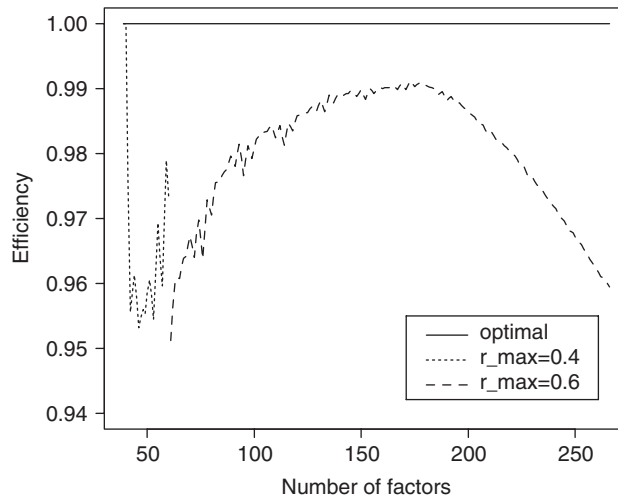


Fig. 3. The Es^2 efficiencies of the augmented designs listed in Table 3. The dotted and dashed curves correspond to base designs ($k2n20d2$) and ($k2n20d4$), with $r_{\max} = 0.4$ and 0.6 , respectively.

be improved in terms of Es^2 efficiencies. However, instead of augmenting the Plackett–Burman design, we may reduce the maximum correlation to $r_{\max} = 0.4$ by augmenting design ($k2n20d2$) for $39 \leq m \leq 60$ factors. Design ($k2n20d4$) can be augmented up to 266 factors with $r_{\max} = 0.6$ with Es^2 efficiencies similar to those obtained from the 1-circulant design (some slightly higher and some slightly lower). The recommended sequential order of addition of interaction columns is given in Table 3. The corresponding Es^2 efficiencies are plotted in Fig. 3 and are all above 0.95. Es^2 optimal designs with $n = 20$ are available in the tables of Butler et al. (2001) for $m = 39, \dots, 43, 45, 46, 47$ factors.

3. Construction of Es^2 efficient designs via column deletion

As mentioned in Section 2, when $n = 2 \pmod 4$, interaction columns formed from the elementwise product of main effect columns are not mean-orthogonal. In this section, we discuss the alternative method of construction of designs with $n = 10, 14, 18, 22$ runs by deleting columns from k -circulant designs.

For 10-run designs, we can construct designs for $m = 10, \dots, 17$ factors by sequential deletion of columns from the 18-factor 2-circulant design ($k2n10d1$) and for $m = 19, \dots, 26$ by sequential deletion of columns of 3-circulant 27-factor design ($k3n10d1$). Similarly, we can construct designs by deleting columns from ($k4n10d1$), ($k5n10d1$), and ($k6n10d1$) to obtain designs for $m = 28, \dots, 53$ factors. Recommended deletion orderings are given in Table 4, together with the Es^2 efficiencies of the resulting designs. All listed designs with $m \geq 15$ factors have $r_{\max} \leq 0.6$ and Es^2 efficiency at least 0.95 and several have Es^2 values that attain the bound of Bulutoglu and Cheng (2004, Theorem 3.1).

In Tables 5–7, recommended orderings are given of deletion of columns of k -circulant designs for $n = 14, 18$ and 22 runs. As shown in the tables, except when the number of factors is close to the number of runs, the designs created by deletion have high efficiencies. For the sizes $(n = 10, m = 15)$, $(n = 10, m = 14)$, $(n = 14, m = 19)$ and $(n = 14, m = 18)$, better designs exist in the literature. Es^2 optimal designs have been found by Bulutoglu and Cheng (2004) using the design search algorithm Gendex of Nguyen (1996) whereas, for these sizes, our designs have efficiencies 0.95, 0.88, 0.82, 0.80.

Table 4
Sequential order of deletion of columns from k -circulant designs ($k2n10d1$), ($k3n10d1$), ($k4n10d1$), ($k5n10d1$), and ($k6n10d1$) for $10 \leq m \leq 53$ factors

| Design | r_{\max} | Deletion order; corresponding Es^2 efficiencies | | | | | | | |
|-----------|------------|---|------------------|------------------|------------------|------------------|------------------|------------------|------------------|
| $k2n10d1$ | 0.6 | c_{18} 1.0 | c_{16} 1.0 | c_{14} 0.95 | c_{12} 0.88 | c_{10} 0.82 | c_8 0.73 | c_6 0.78 | c_4 0.85 |
| $k3n10d1$ | 0.6 | c_{25} 0.98 | c_{22} 0.96 | c_{19} 0.96 | c_{16} 0.96 | c_{13} 0.96 | c_{10} 0.98 | c_7 1.0 | c_4 1.0 |
| $k4n10d1$ | 0.6 | c_{36} 1.0 | c_{33} 1.0 | c_{32} 0.99 | c_{27} 0.99 | c_{24} 0.98 | c_{21} 0.97 | c_{19} 0.97 | c_{15} 0.98 |
| $k5n10d1$ | 0.6 | c_{41} 0.99 | c_{36} 0.99 | c_{31} 0.99 | c_{26} 0.99 | c_{21} 0.99 | c_{16} 1.0 | c_{11} 1.0 | c_6 1.0 |
| $k6n10d1$ | 0.6 | c_{54} 1.0 | c_{53} 1.0 | c_{47} 1.0 | c_{43} 0.99 | c_{37} 0.99 | c_{36} 0.99 | c_{29} 0.99 | c_{26} 0.99 |

Table 5
Sequential order of deletion of columns from k -circulant designs ($k2n14d2$), ($k3n14d3$), and ($k4n14d5$) for $14 \leq m \leq 51$ factors

| Design | r_{\max} | Deletion order; corresponding Es^2 efficiencies | | | | | | | | | | | |
|-----------|------------|---|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|
| $k2n14d2$ | 0.43 | c_{26} 1.0 | c_{24} 1.0 | c_{22} 0.95 | c_{20} 0.89 | c_{18} 0.87 | c_{16} 0.84 | c_{14} 0.82 | c_{12} 0.80 | c_{10} 0.72 | c_8 0.63 | c_6 0.69 | c_4 0.79 |
| $k3n14d3$ | 0.43 | c_{39} 0.97 | c_{36} 0.95 | c_{33} 0.94 | c_{30} 0.94 | c_{27} 0.94 | c_{24} 0.93 | c_{21} 0.91 | c_{18} 0.90 | c_{15} 0.89 | c_{12} 0.90 | c_9 0.91 | c_6 0.88 |
| $k4n14d5$ | 0.43 | c_{52} 1.0 | c_{50} 1.0 | c_{48} 0.99 | c_{37} 0.99 | c_{36} 0.98 | c_{26} 0.98 | c_{24} 0.98 | c_{22} 0.98 | c_{13} 0.98 | c_{12} 0.97 | c_{47} 0.97 | c_{41} 0.97 |

Table 6
Sequential order of deletion of columns from k -circulant designs ($k2n18d2$) and ($k3n18d5$) for $18 \leq m \leq 50$ factors

| Design | r_{\max} | Deletion order; corresponding Es^2 efficiencies | | | | | | | | | | | | | |
|-----------|------------|---|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| $k2n18d2$ | 0.33 | c_{34} | c_{33} | c_{32} | c_{28} | c_{26} | c_{24} | c_{20} | c_{16} | c_{14} | c_{12} | c_{10} | c_4 | c_{30} | c_{22} |
| | | 1.0 | 1.0 | 0.96 | 0.91 | 0.90 | 0.89 | 0.86 | 0.82 | 0.79 | 0.75 | 0.72 | 0.71 | 0.66 | 0.54 |
| $k3n18d5$ | 0.33 | c_{18} | c_8 | | | | | | | | | | | | |
| | | 0.59 | 0.66 | | | | | | | | | | | | |
| | | c_{50} | c_{47} | c_{44} | c_{41} | c_{26} | c_{23} | c_{29} | c_{16} | c_9 | c_{33} | c_3 | c_{17} | c_{11} | c_{48} |
| | | | | | | | | | | | | | | | |
| | | 0.98 | 0.96 | 0.96 | 0.96 | 0.96 | 0.95 | 0.94 | 0.93 | 0.93 | 0.92 | 0.91 | 0.89 | 0.87 | 0.87 |
| | | c_{37} | c_{34} | | | | | | | | | | | | |
| | | 0.87 | 0.84 | | | | | | | | | | | | |

Table 7
Sequential order of deletion of columns from k -circulant designs ($k2n22d4$) for $22 \leq m \leq 41$ factors

| Design | r_{\max} | Deletion order; corresponding Es^2 efficiencies | | | | | | | | | | | | | |
|-----------|------------|---|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| $k2n22d4$ | 0.27 | c_{42} | c_{40} | c_{38} | c_{36} | c_{34} | c_{22} | c_{20} | c_{18} | c_{16} | c_{14} | c_{12} | c_{32} | c_{30} | c_{28} |
| | | 1.0 | 1.0 | 0.96 | 0.92 | 0.92 | 0.91 | 0.88 | 0.85 | 0.83 | 0.81 | 0.79 | 0.76 | 0.71 | 0.65 |
| | | c_{26} | c_{24} | c_{10} | c_8 | c_6 | c_4 | | | | | | | | |
| | | 0.60 | 0.59 | 0.55 | 0.43 | 0.46 | 0.50 | | | | | | | | |

4. Searching for active factors

Abraham et al. (1999) illustrated the problems encountered by the methods of forward selection in regression and best subset selection in correctly identifying active factors using supersaturated designs. Based on the Plackett–Burman design used by Williams (1968) in analyzing an experiment on rubber making process, Abraham et al. (1999) used the “branching column method” of Lin (1993) to select eight different half-fractions from the complete design. Their results showed that the different half-fractions lead to different factors being selected as active and none agrees entirely with the analysis from the full experiment. They also showed that the number and magnitude of active effects, as well as the particular assignment of design matrix columns to the active factors, affect which factors are selected as active. These problems arise from the non-orthogonality of the columns in the model matrix. In this section, we investigate how large the correlation between columns can be while still being able to identify the active factors correctly.

Throughout we assume without loss of generality that $\beta_t = a_t \beta_1$ ($t=2, \dots, m$) with $\beta_1 \geq 0$ and $1 \geq |a_2| \geq \dots \geq |a_m| \geq 0$, so that $|\beta_1| \geq |\beta_2| \geq \dots \geq |\beta_m| \geq 0$. Further, we make the assumption that the number of active effects is much smaller than the number of factors (“factor sparsity”, see Box and Meyer, 1986). We consider only “mean-orthogonal” designs where every factor is observed the same number of times at the high and low levels and assume that no two columns in the design are identical. This means that the absolute value of the inner product of any two columns is at most $n - 4$.

4.1. Selection of the most active factor

We use model (1.1) and, for $i = 1, \dots, m$, define the “hat” matrix H_i as

$$H_i = x_i^* (x_i^{*'} x_i^*)^{-1} x_i^{*'} \quad \text{with } x_i^* = [\mathbf{1} \ x_i], \tag{4.1}$$

and

$$F_i^* = \frac{Y' [H_i - n^{-1} J] Y}{(n - 2)^{-1} Y' [I - H_i] Y}, \tag{4.2}$$

where x_i and $\mathbf{1}$ are defined in Section 1 and where $J = \mathbf{1}'\mathbf{1}$.

In the first step of forward selection, the factor with the largest F^* value will be added to the regression model when the F statistic is greater than a critical value (see, for example, [Seber, 1977](#), Section 4.5). Therefore, on comparing F_i^* with F_j^* , factor i will be selected over factor j in the first step of forward selection if

$$Y'H_iY > Y'H_jY. \tag{4.3}$$

For any mean-orthogonal design and main effects model (1.1), $\mathbf{1}'\mathbf{x}_j = 0$, $\mathbf{J}\mathbf{X} = \mathbf{0}$ and $\mathbf{x}'_j\mathbf{x}_j = n$, so $\mathbf{H}_j = n^{-1}\mathbf{x}_j^*\mathbf{x}_j^{*'} = n^{-1}(\mathbf{J} + \mathbf{x}_j\mathbf{x}'_j)$. Then, using a well known result concerning expectation of quadratic forms (see, for example, Theorem 1.7, [Seber, 1977](#)) and the fact that \mathbf{H}_j in (4.1) is symmetric and idempotent, we have

$$\begin{aligned} E[Y'H_jY] &= \text{trace}(\mathbf{H}_j\sigma^2) + (\mu\mathbf{1} + \mathbf{X}\boldsymbol{\beta})'\mathbf{H}_j(\mu\mathbf{1} + \mathbf{X}\boldsymbol{\beta}) \\ &= 2\sigma^2 + n^{-1}(\mu\mathbf{1} + \mathbf{X}\boldsymbol{\beta})'(\mathbf{J} + \mathbf{x}_j\mathbf{x}'_j)(\mu\mathbf{1} + \mathbf{X}\boldsymbol{\beta}) \\ &= C + n^{-1}\left[\sum_{t=1}^m \beta_t \mathbf{x}'_t \mathbf{x}_j\right]^2 \\ &= C + n^{-1}\beta_1^2\left[a_j n + \sum_{t \neq j=1}^m a_t \mathbf{x}'_t \mathbf{x}_j\right]^2, \end{aligned} \tag{4.4}$$

where $C = 2\sigma^2 + n\mu^2$. Then factor 1 (the factor with the largest absolute effect) is expected to be selected correctly at step 1 of forward selection if $E[Y'H_1Y] - E[Y'H_jY] > 0$ for all $j = 2, \dots, m$, where

$$\begin{aligned} E[Y'H_1Y] - E[Y'H_jY] &= n^{-1}\beta_1^2\left[n + \sum_{t=2}^m a_t \mathbf{x}'_t \mathbf{x}_1\right]^2 - n^{-1}\beta_1^2\left[na_j + \mathbf{x}'_1 \mathbf{x}_j + \sum_{t \neq j=2}^m a_t \mathbf{x}'_t \mathbf{x}_1\right]^2. \end{aligned} \tag{4.5}$$

We now examine some special cases.

(i) *One non-negligible factor*: Suppose that factor 1 is the only non-negligible factor (with $\beta_1 \neq 0$). All other factors have negligible effects with the idealized situation of $\beta_t = 0, t > 1$, then for $j \neq 1$, Eq. (4.5) gives

$$E[Y'H_1Y] - E[Y'H_jY] = n^{-1}\beta_1^2[n^2 - (\mathbf{x}'_1 \mathbf{x}_j)^2].$$

Since $(\mathbf{x}'_1 \mathbf{x}_j)^2 \leq (n - 4)^2$, it follows that $E[Y'H_1Y] > E[Y'H_jY]$ and the correct factor is expected to be selected. If the error variance is small compared with the magnitude of the single active effect, then the active effect will *always* be selected correctly, as follows. Under model (1.1), and using (4.1),

$$Y'H_jY = n^{-1}(n\mu + \boldsymbol{\epsilon}'\mathbf{1})^2 + n^{-1}(\beta_1 \mathbf{x}'_1 \mathbf{x}_j + \boldsymbol{\epsilon}'\mathbf{x}_j)^2.$$

Now $\mathbf{x}'_1 \mathbf{x}_1 = n$ and $|\mathbf{x}_j \mathbf{x}_1| \leq n - 4$ so, if $|\mathbf{x}'_1 \boldsymbol{\epsilon}| < |\beta_1|$ for all $i = 1, 2, \dots, m$, then

$$Y'H_1Y - Y'H_jY > n^{-1}\beta_1^2[(n - 1)^2 - (n - 3)^2] > 0$$

and the correct factor will be selected.

(ii) *Two non-negligible factors*: Suppose that factors 1 and 2 are the only factors with non-negligible effects and $\beta_2 = a_2\beta_1$, for $0 < |a_2| \leq 1$ and all other factors are inactive with the idealized situation of $\beta_t = 0, t > 2$. Then, from (4.5)

$$E[Y'H_1Y] - E[Y'H_jY] = n^{-1}\beta_1^2[n + a_2 \mathbf{x}'_1 \mathbf{x}_2]^2 - n^{-1}\beta_1^2[\mathbf{x}'_1 \mathbf{x}_j + a_2 \mathbf{x}'_2 \mathbf{x}_j]^2. \tag{4.6}$$

Then, for $j = 2$,

$$E[Y'H_1Y] - E[Y'H_2Y] = n^{-1}\beta_1^2(n^2 - (\mathbf{x}'_1 \mathbf{x}_2)^2)(1 - a_2^2) > 0. \tag{4.7}$$

This is positive since $|\mathbf{x}'_1 \mathbf{x}_2| \leq (n - 4)$. Thus, factor 1 is expected to be selected over factor 2 at the first step of the forward selection procedure.

Now, let r_{\max} be the maximum absolute value of the correlation between any pair of columns of X and let $c = nr_{\max}$ ($\leq n - 4$). Set the right-hand side of (4.6) to its minimum value, then

$$E[Y'H_1Y] - E[Y'H_jY] \geq n^{-1} \beta_1^2 [(n - |a_2|c)^2 - (|a_2|c + c)^2] = n^{-1} \beta_1^2 [n^2 - 2|a_2|cn - 2c^2|a_2| - c^2]. \tag{4.8}$$

Factor 1 will be selected on average if expression (4.8) is positive, that is, if

$$c < n(2|a_2| + 1)^{-1}, \tag{4.9}$$

or, equivalently, if the correlation between the columns measuring the active factors and the inactive factors satisfies

$$r_{\max} < (2|a_2| + 1)^{-1}. \tag{4.10}$$

Again, if the error variance is small compared with β_1 and β_2 , factor 1 will always be correctly selected if (4.10) is satisfied.

If $|a_2| = 0.5$, for example, then we require, $c < n/2$ or $r_{\max} < \frac{1}{2}$. As $|a_2|$ becomes smaller, the absolute difference between the regression coefficient β_1 for the most active factor and the regression coefficient β_2 for the second active factor increases, and it is easier to select the correct most active factor at step 1. The most difficult case for selecting factor 1 correctly at step 1 occurs when $|a_2|$ is close to 1.0. In this case, we require $c < n/3$ or $r_{\max} < 0.33$ for the correct selection of factor 1 at step 1. Some supersaturated designs whose column correlations have small values of r_{\max} are identified later in this section.

We conjecture that the right-hand side of (4.8) cannot be attained in practice (at least for moderate sized n) and so the bound (4.9) is lower than needed. In the idealized situation being considered and $|a_2| < 1$, we conjecture that factor 1 will always be selected over factor j ($j \geq 2$) although, as yet, we have no formal proof of this. However, if there are sizeable errors present, then the incorrect factor can be selected as shown in the case study of Abraham et al. (1999, p. 138) which had $r_{\max} = 0.43$, $c = 0.43n$. Thus, it is still advisable to select a design satisfying $r_{\max} < 0.33$ if one exists.

(iii) *Three or more non-negligible factors*: For three active factors, with $\beta_j = a_j \beta_1 > 0$, $j = 2, 3$, for $0 < |a_3| \leq |a_2| \leq 1$ and $\beta_t = 0$, $t > 3$, Eq. (4.5) becomes

$$E[Y'H_1Y] - E[Y'H_jY] = n^{-1} \beta_1^2 [n + a_2 x'_1 x_2 + a_3 x'_1 x_3]^2 - n^{-1} \beta_1^2 [x'_1 x_j + a_2 x'_2 x_j + a_3 x'_3 x_j]^2 \tag{4.11}$$

for $j = 1, \dots, m$. Using (4.11) with $j > 3$, it can be shown that, provided the errors are small compared with the magnitude of the active effects, factor 1 is expected to be selected correctly over factors 4, 5, ..., m at step 1 if the maximum correlation between pairs of columns in the design satisfies $c < n/(1 + 2|a_2| + 2|a_3|)$. Similarly, using (4.11) with $j = 2, 3$, factor 1 is expected to be selected correctly over factors 2 and 3 if $c < n(1 - |a_2|)/(1 + 2|a_3| - |a_2|)$. Now, if $|a_2| = 1$, either factor 1 or 2 would be a correct selection at step 1. Using (4.11), it can be shown that factor 1 or 2 can always be correctly selected over factors 4, 5, ..., m and factor 1 or 2 can be selected over factor 3 if $c < n(1 - |a_3|)/(3 - |a_3|)$. Finally, if $|a_2| = |a_3| = 1$, then any of factors 1, 2 or 3 would be a correct selection at step 1 and would be selected if the errors are small and if $c < n/5$. Comparing these various bounds leads to the following set of bounds for correct selection of the most active factor at step 1 of forward regression when there are three non-negligible effects:

$$c < \frac{n(1 - |a_2|)}{1 + 2|a_3| - |a_2|} \quad \text{for } |a_2| + |a_3| \geq 1, \quad 0 < |a_3| \leq |a_2| < 1, \tag{4.12}$$

$$c < \frac{n}{1 + 2(|a_2| + |a_3|)} \quad \text{for } |a_2| + |a_3| \leq 1, \quad 0 < |a_3| \leq |a_2| < 1, \tag{4.13}$$

$$c < \frac{n(1 - |a_3|)}{3 - |a_3|} \quad \text{for } 0 < |a_3| < |a_2| = 1, \tag{4.14}$$

$$c < \frac{n}{5} \quad \text{for } |a_3| = |a_2| = 1. \tag{4.15}$$

Note that (4.12) and (4.13) not only provide somewhat different bounds for r_{\max} from those given by Chen and Lin (1998) for identifiability, but they also depend on the relative sizes of the active factors. For example, when

$|a_2| + |a_3| = 1$, both conditions become $c < n/3$, but when $|a_2| = 0.6$, $|a_3| = 0.5$, condition (4.12) becomes $c < 0.28n$. When $|a_2| + |a_3| < 0.45$, condition (4.13) becomes $c < 0.526n$ or $r_{\max} < 0.526$. The cases of $\beta_2 = \beta_1$ and $\beta_3 = \beta_2 = \beta_1$ are unlikely to arise in practice, so the bounds (4.14) and (4.15) are less useful. If $|a_2|$ and $|a_3|$ are both close to, but not equal to, 1.0, then (4.12) shows that a supersaturated design with sufficiently small column correlations cannot exist, which helps to explain the difficulty in selecting the correct most active effect in the presence of several large competing effects. However, it can be argued that, in this case, it would not be a great error to select the second or third most active effect.

The case of 3 active factors can be extended to $p > 3$ active factors and it can be shown that factor 1 (the most active factor) is expected to be correctly selected at step 1 if the following bounds hold. Similar comments to those above apply concerning the use of these bounds:

$$c < \frac{n(1 - |a_2|)}{1 + 2(\sum_{i=3}^p |a_i|) - |a_2|} \quad \text{for } \sum_{i=2}^p |a_i| \geq 1, \quad 0 < |a_i| < 1, \quad i = 2, \dots, p, \tag{4.16}$$

$$c < \frac{n}{1 + 2(\sum_{i=2}^p |a_i|)} \quad \text{for } \sum_{i=2}^p |a_i| \leq 1, \quad 0 < |a_i| < 1, \quad i = 2, \dots, p, \tag{4.17}$$

$$c < \frac{n(1 - |a_{k+1}|)}{1 + 2(k - 1) + \sum_{i=k+2}^p |a_i| - |a_{k+1}|} \quad \text{for } 0 < |a_p| < |a_{p-1}| = \dots = |a_2| = 1, \tag{4.18}$$

$$c < \frac{n}{1 + 2(p - 1)} \quad \text{for } |a_2| = \dots = |a_p| = 1, \tag{4.19}$$

4.2. *Supersaturated designs with small r_{\max}*

Many of the 2-circulant supersaturated designs of Liu and Dean (2004, Table 2) have $r_{\max} < \frac{1}{3}$. For example, design $(k2n20d4)$ has $r_{\max} = 0.2$ and $(k2n16d2)$ has $r_{\max} = 0.25$. The k -circulant design $(k2n22d4)$ and the designs obtained by deleting columns from it (Table 7 of this paper) all have $r_{\max} = 0.27$. The designs obtained by deletion from $(k2n18d2)$ and $(k3n18d5)$ in Table 6 and the augmented designs of $(k1n12d1)$ in Table 1 all have $r_{\max} = 0.333$, as do the following k -circulant designs of Tables 2–6 of Liu and Dean (2004): $(k2n6d1)$, $(k2n12d2)$, $(k2n18d2)$, $(k3n12d3)$, $(k3n18d5)$, $(k4n12d5)$, $(k5n12d5)$, $(k6n12d5)$.

Many other designs have $r_{\max} < 0.5$. For example, designs with r_{\max} between 0.4 and 0.46 include the 20-run augmented designs of $(k2n20d2)$ in Table 3, the 14-run designs obtained by deletion of columns in Table 5, and the following k -circulant designs of Tables 2–6 of Liu and Dean (2004): $(k2n14d2)$, $(k2n20d1)$, $(k2n20d2)$, $(k2n20d3)$, $(k2n22d2)$, $(k2n22d3)$, $(k3n14d3)$, $(k4n14n5)$.

The following designs listed in Tables 3–6 of Liu and Dean (2004) all have $r_{\max} = 0.5$: $(k2n8d1)$, $(k2n16d1)$, $(k3n8d1)$, $(k3n16d1)$, $(k3n16d2)$, $(k3n16d3)$, $(k3n16d4)$, $(k4n8d1)$, $(k4n14d5)$, $(k5n8d1)$. Also, the augmented designs of $(k2n16d2)$ and $(k3n16d3)$ listed in Table 2 of this paper have $r_{\max} = 0.5$.

4.3. *Selection of the two most active factors*

For selection of the two most active factors via forward selection in regression, we may either select factor 1 followed by factor 2 or select factor 2 followed by factor 1. We consider here, instead, the equivalent method of “best subset” or “all subset” selection, which results in the same bounds on the column correlations of a supersaturated design as those given by the two steps of forward regression.

Consider the subset of factors $S_p = \{i_1, i_2, \dots, i_p\}$, with $p \leq m$, selected from factors $1, \dots, m$. Define

$$\mathbf{x}_{S_p}^* = [1, \mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \dots, \mathbf{x}_{i_p}],$$

where \mathbf{x}_{i_j} is column i_j of the model matrix \mathbf{X} . Also define the matrix \mathbf{H}_{S_p} to be

$$\mathbf{H}_{S_p} = \mathbf{x}_{S_p}^* (\mathbf{x}_{S_p}^{*'} \mathbf{x}_{S_p}^*)^{-1} \mathbf{x}_{S_p}^{*'} \tag{4.20}$$

In “best subset” or “all subset” selection, the subset S_p of factors is selected if it yields maximum $R_{S_p}^2 = 1 - \text{SSE}_{S_p} / \text{SSTO}$ (see, for example, **Seber, 1977**, Theorem 4.3), where $\text{SSE}_{S_p} = \mathbf{Y}'[\mathbf{I} - \mathbf{H}_{S_p}]\mathbf{Y}$ and $\text{SSTO} = \mathbf{Y}'(\mathbf{I} - n^{-1}\mathbf{J})\mathbf{Y}$, over all subsets S_p of p factors. Now, $R_{S_p^{(1)}}^2 > R_{S_p^{(2)}}^2$ if and only if $\mathbf{Y}'\mathbf{H}_{S_p^{(1)}}\mathbf{Y} > \mathbf{Y}'\mathbf{H}_{S_p^{(2)}}\mathbf{Y}$, where $S_p^{(1)}$ and $S_p^{(2)}$ are two different subsets of p factors. (When $p = 1$, this is the same comparison as for forward selection).

Define $k_{ij} = \mathbf{x}'_i\mathbf{x}_j$ for $i \neq j$ and consider the subset $S_2 = \{i, j\}$ of size 2. Then,

$$\begin{aligned} E[\mathbf{Y}'\mathbf{H}_{ij}\mathbf{Y}] &= \text{trace}(\mathbf{H}_{ij}\sigma^2) + (\mu\mathbf{1} + \mathbf{X}\boldsymbol{\beta})'\mathbf{H}_{ij}(\mu\mathbf{1} + \mathbf{X}\boldsymbol{\beta}) \\ &= 3\sigma^2 + n\mu^2 + (\mathbf{X}\boldsymbol{\beta})'[\mathbf{x}_i \ \mathbf{x}_j] \begin{bmatrix} n & k_{ij} \\ k_{ij} & n \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{x}'_i \\ \mathbf{x}'_j \end{bmatrix} (\mathbf{X}\boldsymbol{\beta}) \\ &= C + (n^3 - nk_{ij}^2)^{-1}\beta_1^2 \left[n^2 \left(\sum_{t=1}^m a_t k_{it} \right)^2 + n^2 \left(\sum_{t=1}^m a_t k_{jt} \right)^2 \right. \\ &\quad \left. - 2nk_{ij} \left(\sum_{t=1}^m a_t k_{it} \right) \left(\sum_{t=1}^m a_t k_{jt} \right) \right], \end{aligned} \tag{4.21}$$

where $C = 3\sigma^2 + n\mu^2$. The subset of size $p = 2$ containing the two most active factors (factors 1 and 2) is expected to be selected correctly if

$$E[\mathbf{Y}'\mathbf{H}_{12}\mathbf{Y}] > E[\mathbf{Y}'\mathbf{H}_{ij}\mathbf{Y}], \quad i, j = 1, \dots, m \text{ and } (i, j) \neq (1, 2),$$

where \mathbf{H}_{12} and \mathbf{H}_{ij} are defined as in (4.20) with subsets $\{1, 2\}$ and $\{i, j\}$, respectively.

Suppose that factors 1 and 2 are the only factors with non-negligible effects and that $\beta_2 = a_2\beta_1$ ($0 < |a_2| \leq 1$) and all other factors are inactive with the idealized situation of $\beta_t = 0, t > 2$. Then, from (4.21),

$$\begin{aligned} E[\mathbf{Y}'\mathbf{H}_{ij}\mathbf{Y}] &= C + (n^2 - k_{ij}^2)^{-1}\beta_1^2 [n(k_{1i} + a_2k_{2i})^2 + n(k_{1j} + a_2k_{2j})^2 \\ &\quad - 2k_{ij}(k_{1i} + a_2k_{2i})(k_{1j} + a_2k_{2j})] \end{aligned} \tag{4.22}$$

and with $i = 1$,

$$E[\mathbf{Y}'\mathbf{H}_{1j}\mathbf{Y}] = C + \beta_1^2(n + 2a_2k_{12}) + (n^2 - k_{1j}^2)^{-1}\beta_1^2 a_2^2 (nk_{12}^2 + nk_{2j}^2 - 2k_{12}k_{1j}k_{2j}). \tag{4.23}$$

For $(i, j) = (1, 2)$, (4.21) reduces to

$$E[\mathbf{Y}'\mathbf{H}_{12}\mathbf{Y}] = C + \beta_1^2(n + 2a_2k_{12} + na_2^2), \tag{4.24}$$

and so

$$\begin{aligned} E[\mathbf{Y}'\mathbf{H}_{12}\mathbf{Y}] - E[\mathbf{Y}'\mathbf{H}_{1j}\mathbf{Y}] &= n\beta_1^2 a_2^2 - (n^2 - k_{1j}^2)^{-1}\beta_1^2 a_2^2 (nk_{12}^2 + nk_{2j}^2 - 2k_{12}k_{1j}k_{2j}) \\ &= (n^2 - k_{1j}^2)^{-1}\beta_1^2 a_2^2 [n^3 - n(k_{12}^2 + k_{1j}^2 + k_{2j}^2) + 2k_{12}k_{1j}k_{2j}]. \end{aligned} \tag{4.25}$$

The right-hand side of (4.25) is non-negative, since the expression in square brackets is the determinant of the non-negative definite matrix $\mathbf{M}'\mathbf{M}$ where $\mathbf{M} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \mathbf{x}_j]$. So the subset of factors 1 and 2 is expected to be selected over the subset of factors $(1, j)$ ($j > 2$). Similarly,

$$E[\mathbf{Y}'\mathbf{H}_{12}\mathbf{Y}] - E[\mathbf{Y}'\mathbf{H}_{2j}\mathbf{Y}] = \beta_1^2(n^2 - k_{2j}^2)^{-1}[n^3 - n(k_{12}^2 + k_{1j}^2 + k_{2j}^2) + 2k_{12}k_{1j}k_{2j}],$$

which is similar to (4.25) apart from the constant a_2^2 . Therefore, if the errors are small, the subset of factors 1 and 2 will always be selected over the subset of factors 1 and j and over the subset of factors 2 and j , ($j > 2$). Furthermore, from (4.22) and (4.24), for $(i, j) \neq (1, 2)$, and $0 < |a_2| \leq 1.0$,

$$E[\mathbf{Y}'\mathbf{H}_{12}\mathbf{Y}] - E[\mathbf{Y}'\mathbf{H}_{ij}\mathbf{Y}] \geq \frac{\beta_1^2}{n - c} [-2(1 + |a_2| + |a_2|^2)c^2 - nc(1 + |a_2|)^2 + n^2(1 + |a_2|^2)],$$

which is greater than zero if

$$c < (n/4) \left[\frac{\sqrt{(1 + |a_2|)^4 + 8(1 + |a|^2)(1 + |a_2| + |a_2|^2)} - (1 + |a_2|)^2}{1 + |a_2|^2 - |a_2|} \right]. \quad (4.26)$$

The term in square brackets in (4.26) can be shown to be larger than $\frac{4}{3}$. Consequently, best subset selection will select the correct two active factors if the errors are small and the design satisfies $c < n/3$, that is, $r_{\max} < 0.33$.

Obtaining a theoretical formula for $E[\mathbf{Y}'\mathbf{H}_{S_p^{(1)}}\mathbf{Y}] - E[\mathbf{Y}'\mathbf{H}_{S_p^{(2)}}\mathbf{Y}]$ for $S_p^{(1)} = \{1, \dots, p\}$ and $S_p^{(2)} = \{i_1, \dots, i_p\}$ becomes progressively more involved as p increases and the details are omitted, but the same general principles apply.

5. Conclusions

In Sections 2 and 3, we presented several series of efficient supersaturated designs obtained by augmenting k -circulant designs with interaction columns or by deleting columns from k -circulant designs. Where better designs exist in the literature, these were identified.

For n a multiple of 4, we listed augmented k -circulant supersaturated designs for $n = 12$ runs with $m = 13, \dots, 231$ factors; for $n = 16$ runs with $m = 31, \dots, 255$ factors; and for $n = 20$ runs with $m = 39, \dots, 266$ factors. The Es^2 efficiencies of these designs, as compared with the bounds of Butler et al. (2001) or Bulutoglu and Cheng (2004), are mostly above 0.95 as shown in Figs. 1–3.

For $n = 2 \pmod{4}$, designs were obtained by deletion of columns from k -circulant supersaturated designs for $n = 10$ with $m = 10, \dots, 53$; for $n = 14$ with $m = 14, \dots, 51$; for $n = 18$ with $m = 18, \dots, 50$; and for $n = 22$ with $m = 22, \dots, 41$. Except when m is close to n , the Es^2 efficiencies of these designs are above 0.83 and many have efficiency above 0.90. (See Tables 4–7).

In Section 4, we investigated the potential difficulty of using a supersaturated design to select active factors. Bounds were derived for the maximum correlation between any two columns of the design matrix that will allow the “most active” factor to be selected correctly using forward regression when the errors are small in comparison with the size of the active effect. In addition, we have obtained conditions for the correct selection of two active factors using subset selection, again when the errors are small. The results assume that all other factors have negligible effects and that the errors are small in comparison with the effects of the few active factors. We have listed a number of k -circulant supersaturated designs that achieve the bounds.

Acknowledgments

The work of Shing Ruan and Angela Dean was partly supported by Grant SES-0437251 from the National Science Foundation and that of Yufeng Liu by Grant DMS-0606577 from the National Science Foundation and by a UNC Junior Faculty Development Award.

References

- Abraham, B., Chipman, H., Vijayan, K., 1999. Some risks in the construction and analysis of supersaturated designs. *Technometrics* 41 (2), 135–141.
- Box, G.E.P., Meyer, R.D., 1986. An analysis for unreplicated fractional factorials. *Technometrics* 28, 11–18.
- Bulutoglu, D.A., Cheng, C.S., 2004. Construction of $E(s^2)$ -optimal supersaturated designs. *Ann. Statist.* 32 (4), 1662–1678.
- Butler, N.A., Mead, R., Eskridge, K.M., Gilmour, S.G., 2001. A general method of constructing $E(s^2)$ -optimal supersaturated designs. *J. R. Statist. Soc. B* 63, 621–632.
- Chen, J., Lin, D.K.J., 1998. On the identifiability of a supersaturated design. *J. Statist. Plann. Inference* 72, 99–107.
- Cheng, C.-S., 1997. $E(s^2)$ -optimality of supersaturated designs. *Statistica Sinica* 7, 929–939.
- Eskridge, K.M., Gilmour, S.G., Mead, R., Butler, N., Travnicek, D.A., 2004. Large supersaturated designs. *J. Statist. Comput. Simulation* 74, 525–542.
- Lin, D.K.J., 1993. A new class of supersaturated designs. *Technometrics* 35, 28–31.
- Liu, M., Zhang, R., 2000. Construction of $E(s^2)$ optimal supersaturated designs using cyclic BIBDs. *J. Statist. Plann. Inference* 91, 139–150.
- Liu, Y., Dean, A.M., 2004. k -circulant supersaturated designs. *Technometrics* 46, 32–43.
- Nguyen, N.-K., 1996. An algorithmic approach to constructing supersaturated designs. *Technometrics* 38, 69–73.
- Plackett, R.L., Burman, J.P., 1946. The design of optimum multifactorial experiments. *Biometrika* 33, 305–325.

Seber, G.A.F., 1977. *Linear Regression Analysis*. Wiley, New York.

Srivastava, J.N., 1975. Designs for searching for non-negligible effects. In: Srivastava, J.N. (Ed.), *A Survey of Statistical Design and Linear Models*. North-Holland, Amsterdam, pp. 507–519.

Williams, K.R., 1968. Designed experiments. *Rubber Age* 100, 65–71.

Wu, C.F.J., 1993. Construction of supersaturated designs through partially aliased interactions. *Biometrika* 80, 661–669.