

A- AND D-OPTIMAL ROW-COLUMN DESIGNS FOR TWO-COLOUR cDNA MICROARRAY EXPERIMENTS USING LINEAR MIXED EFFECTS MODELS

Dibaba Bayisa Gemechu

University of Pretoria

and

*Legesse Kassa Debusho*¹

University of Pretoria

e-mail: *legesse.debusho@up.ac.za*

and

Linda M. Haines

University of Cape Town

Key words: *A*-optimal, *D*-optimal, linear mixed effects model, microarray experiments, robust designs, row-column design.

Summary: Microarray experiments help scientists to study the expression level of thousands of genes simultaneously. These experiments have many design challenges, such as, for example, which mRNA samples should be co-hybridized together and which treatments should be labelled with which fluorescent dye. Therefore a carefully designed microarray experiment to obtain efficient and reliable data so as to ensure the precise estimate of comparisons of interest is required. The present paper is concerned with *A*- and *D*-optimal row-column designs for two-colour microarray experiments, with the array and dye effects treated as the column and row effects, respectively. Linear mixed effects models were used to describe experiments for which a comparison of all possible pairs of treatments is of particular interest by taking the arrays as random column effects. The results of this study show that the optimal row-column designs under the linear fixed effects model are not necessarily optimal under the linear mixed effects model setting.

1. Introduction

Microarray technology allows molecular biologists and geneticists to simultaneously measure the expression levels of thousands of genes in multiple biological samples, for example tissue from animals that are given different drugs. The gene expression level is the quantity of messenger ribonucleic acid (mRNA) produced by a gene. By understanding how expression levels change across

¹Corresponding author.

experimental conditions, researchers gain clues about gene function and learn how genes work together to carry out biological processes. The aim of microarray experiments is to identify genes for which the expression level distributions change in response to a treatment and such genes are called differentially expressed genes. The two-colour microarrays use two different dyes, Cyanine 3 (green) and Cyanine 5 (red), to measure gene expression. The use of two dyes allows two different nucleic acid samples, that is treatments, dyed with different dyes to be hybridized together on the same microarray slide or array. Hybridization of two target samples on the same microarray slide enables a comparison of two treatments on a single array. Since, in a two-colour microarray experiment, only two treatments labelled with two different dyes are accommodated on a single array, the designs for such an experiment constitute block designs with blocks of size two and with the array treated as the experimental block. According to Wit, Nobile and Khanin (2005), it is important to account for the dye effect when estimating the relative gene expression. If the dye effect is included in the model, that is if it is assumed that there is a gene specific dye effect, then there will be two blocking factors, array and dye (Kerr, 2003; Landgrebe, Bretz and Brunner, 2006). In such cases, the microarray experiment can be considered as a row-column design (Shah and Sinha, 1989), with dyes as rows and arrays as columns. Like any other experimental design a microarray experiment has different design challenges, such as, for example, which mRNA samples should be co-hybridized together and which treatment should be labelled with which fluorescent dye. To compound matters, microarrays are expensive. Therefore a carefully designed microarray experiment is required in order to obtain efficient and reliable data so as to ensure the precise estimation of parameters which are of interest.

Optimal or efficient row-column designs for two-colour microarray experiments for the estimation of pairwise contrasts of treatment effects have attracted enormous interest in the design literature (Kerr and Churchill, 2001a; Kerr and Churchill, 2001b; Churchill, 2002; Yang and Speed, 2002; Kerr, 2003; Wit et al., 2005; Landgrebe et al., 2006; Bailey, 2007; Grossman and Schwabe, 2008; Sarkar, Prasad, Guptha, Kashinath and Rathore, 2010; Bailey, Schiff and Hilgers, 2013; Gemechu, Debusho and Haines, 2014). Furthermore, row-column designs are used extensively in agricultural field trials. The design construction approaches adopted in the microarray studies are based almost exclusively on fixed effects models. However, Kerr and Churchill (2001a), Wolfinger et al. (2001) and Lee (2004) have suggested that the array effect should be taken as random. The present paper is broadly concerned with optimal or efficient row-column designs for two-colour microarray experiments for which comparisons of all possible pairs of treatments are of particular interest under the linear mixed effects model setting with the array effects assumed to be random. The essential problem is that of choosing the best row-column design layout for a given number of treatments, dyes, arrays and a given value of a parameter θ which lies in the interval $(0, 1)$ and which is a function of the error and the random column variances. This article is a follow-up to the paper by Gemechu et al. (2014) in which the A -optimal block designs for two-colour cDNA microarray experiments using the linear mixed effects model were discussed. The design problem relating to the linear mixed effects model was recognized by Sarkar et al. (2010) who generated A - and D -optimal or efficient row-column designs for two-colour microarray experiments using a linear fixed effects model. These authors did not construct efficient row-column designs for the linear mixed effects models directly but rather investigated their robustness to random effects. In particular, Sarkar et al. (2010) used the percent coefficient of variation of the lower bounds on the relative A - and D -efficiency for different

values of the parameter θ in the interval $(0, 1)$. However not all A- and D-optimal designs for the linear fixed effects model are necessarily optimal for the linear mixed effects model.

The aim of this paper therefore is to construct A- and D-optimal or near-optimal row-column designs under both linear fixed and mixed effects models separately and to compare these with the best available designs in the literature using A- and D-efficiency. The robustness of optimal row-column designs against different values of the parameter θ is also discussed. The model, information matrix and optimality criteria are introduced in Section 2 and the interchange-exchange algorithm used for numerical construction of A- and D-optimal or near-optimal designs is outlined in Section 3. Results are presented and discussed in Section 4 and some concluding remarks are given in Section 5. The notation and terminology introduced in Gemechu et al. (2014) will be used throughout the paper.

2. Preliminaries

Consider a two-colour microarray experiment comprising b arrays and v treatments replicated r_1, \dots, r_v times. The experiment can be modelled as a row-column design with dyes corresponding to rows and arrays to columns of size 2. Specifically, suppose that the treatment effects and the dye effects are fixed and the array effects are random. Then, following Shah and Sinha (1989), the appropriate linear mixed effects model can be expressed in matrix form as

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{X}_r\boldsymbol{\alpha} + \mathbf{X}_t\boldsymbol{\tau} + \mathbf{X}_c\mathbf{b} + \mathbf{e} \quad (1)$$

where \mathbf{y} is the $n \times 1$ response vector with $n = bv$, μ is the overall mean, $\mathbf{1}$ denotes the $n \times 1$ vector of ones, $\boldsymbol{\alpha}$ is a $k \times 1$ vector of fixed dye effects with attendant $n \times k$ design matrix \mathbf{X}_r , $\boldsymbol{\tau}$ is a $v \times 1$ vector of fixed treatment effects with attendant $n \times v$ design matrix \mathbf{X}_t , \mathbf{b} is a $b \times 1$ vector of random array effects with attendant $n \times b$ design matrix \mathbf{X}_c and \mathbf{e} is an $n \times 1$ vector of error terms. The random effect vector \mathbf{b} is taken to be distributed as $\mathcal{N}(\mathbf{0}, \sigma_b^2 \mathbf{I}_b)$ and the error term \mathbf{e} as $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$, where \mathbf{I}_b and \mathbf{I}_n denote the identity matrices of order b and n , respectively, and \mathbf{b} and \mathbf{e} are assumed to be independent. Thus $\text{Var}(\mathbf{y}) = \mathbf{V} = \sigma_b^2 \mathbf{X}_c \mathbf{X}_c' + \sigma^2 \mathbf{I}_n$ and its inverse is given by

$$\mathbf{V}^{-1} = \frac{1}{\sigma^2} (\mathbf{I}_n - \mathbf{X}_c \mathbf{W} \mathbf{X}_c')$$

where

$$\mathbf{W} = \text{diag}(\sigma_b^2/(\sigma^2 + k_1 \sigma_b^2), \dots, \sigma_b^2/(\sigma^2 + k_b \sigma_b^2))$$

and k_1, \dots, k_b are block sizes not necessarily equal to 2. Observe that $\mathbf{X}_t' \mathbf{1} = \mathbf{r} = (r_1, \dots, r_v)'$, $\mathbf{X}_t' \mathbf{X}_r = \mathbf{M}$ and $\mathbf{X}_t' \mathbf{X}_c = \mathbf{N}$ are the vector of treatment replications, treatments-rows incidence matrix and treatments-columns incidence matrix, respectively, $\mathbf{X}_t' \mathbf{X}_t = \mathbf{R} = \text{diag}(r_1, \dots, r_v)$ and $\mathbf{X}_r' \mathbf{X}_c = \mathbf{S}$ is a $k \times b$ row-column incidence matrix. Ignoring the constant $1/\sigma^2$, the treatment information matrix, also called the C-matrix, is given by

$$\mathbf{C} = \mathbf{R} - \mathbf{N} \mathbf{W} \mathbf{N}' - (\mathbf{M} - \mathbf{N} \mathbf{W} \mathbf{S}') (\mathbf{B} - \mathbf{S} \mathbf{W} \mathbf{S}')^{-1} (\mathbf{M}' - \mathbf{S} \mathbf{W} \mathbf{N}') \quad (2)$$

where $\mathbf{B} = \mathbf{X}_r' \mathbf{X}_r$ is a $k \times k$ diagonal matrix of rows size. Details of the derivation of \mathbf{C} are given in Appendix A.

Assume now that the rows and columns are orthogonal. Note that in a two-colour microarray experiment no treatment occurs more than once in an array and thus the design layout is that of a binary row-column design. Thus for column sizes $k_j = k, j = 1, \dots, b$, $\mathbf{S} = \mathbf{1}_k \mathbf{1}'_b$, $\mathbf{M} \mathbf{1}_k = \mathbf{r}$, and $\mathbf{B} = b \mathbf{I}_k$. It then follows immediately from the results pertaining to row-column designs with random array effects (Shah and Sinha, 1989) that the \mathbf{C} matrix can be expressed as

$$\mathbf{C} = \mathbf{R} - \frac{1}{k} \mathbf{N} \mathbf{N}' - \frac{1}{b} \mathbf{M} \mathbf{M}' + \frac{1}{bk} \mathbf{r} \mathbf{r}' + \theta \left(\frac{1}{k} \mathbf{N} \mathbf{N}' - \frac{1}{bk} \mathbf{r} \mathbf{r}' \right)$$

where $\theta = \sigma^2 / (\sigma^2 + k\sigma_b^2)$, $0 \leq \theta \leq 1$. Note that $k = 2$ for a two-colour microarray experiment. Note also that the case of $\theta = 0$ corresponds to a fixed effects model with

$$\mathbf{C} = \mathbf{R} - \frac{1}{2} \mathbf{N} \mathbf{N}' - \frac{1}{b} \mathbf{M} \mathbf{M}' + \frac{1}{2b} \mathbf{r} \mathbf{r}'$$

and that as θ approaches 1, so the model tends to a treatments only model with

$$\mathbf{C} = \mathbf{R} - \frac{1}{b} \mathbf{M} \mathbf{M}'.$$

The matrix \mathbf{C} is symmetric and positive semi-definite and, since $\mathbf{C} \mathbf{1}_v = \mathbf{0}$, has $\text{rank}(\mathbf{C}) \leq v - 1$. In the present study attention is restricted to designs for which the treatment information matrix \mathbf{C} has rank $v - 1$. It thus follows that the designs are connected and that all treatment contrasts, and in particular pairwise differences, are estimable (Haines, 2015). In this paper, attention is restricted to connected row-column designs.

One of the main aims of a microarray experiment is to identify differentially expressed genes. Thus it is natural to focus on design criteria which are based on the variances of the estimated pairwise treatment differences and to select designs from a set of competing designs comprising b arrays and v treatments for which the criterion of interest is in some sense optimal. Consider therefore all possible pairwise treatment differences expressed as $\mathbf{T} \boldsymbol{\tau}$, where \mathbf{T} is a $\binom{v}{2} \times v$ matrix with each row comprising an element 1, an element -1 and all other elements 0, and is such that $\mathbf{T}' \mathbf{T} = v \mathbf{I}_v - \mathbf{J}_v$ (Dey, 2010, p. 39). Then the variance-covariance matrix of the best linear unbiased estimator of $\mathbf{T} \boldsymbol{\tau}$ is given by $\mathbf{T} \mathbf{C}^- \mathbf{T}'$, where \mathbf{C}^- is an arbitrary g -inverse of \mathbf{C} .

Suppose that $\mathcal{D} = \mathcal{D}(v, k, b, \theta)$ represents the class of connected row-column designs with $k = 2$ rows, b columns and v treatments and that $\lambda_i, i = 1, \dots, v - 1$, denote the strictly positive eigenvalues of the treatment information matrix \mathbf{C} for a design in \mathcal{D} . Then the design $d_A^* \in \mathcal{D}$ is said to be A -optimal over all competing designs, that is over \mathcal{D} , if it minimizes the sum of the variances of the estimates of the pairwise treatment differences, that is the criterion $\Psi_A(d) = \text{trace}(\mathbf{T} \mathbf{C}^- \mathbf{T}') = \text{trace}(\mathbf{C}^- (v \mathbf{I}_v - \mathbf{J}_v)) = \text{trace}(\mathbf{C}^+)$, where \mathbf{C}^+ is the Moore-Penrose inverse of \mathbf{C} (Bailey, 2009). Note that the A -optimality criterion can be expressed in terms of the eigenvalues of \mathbf{C} as $\Psi_A(d) = \sum_{i=1}^{v-1} \lambda_i^{-1}$. In addition a design d_D^* is said to be D -optimal if it minimizes the determinant of $\mathbf{T} \mathbf{C}^- \mathbf{T}'$, that is the criterion $\Psi_D(d) = |\mathbf{T} \mathbf{C}^- \mathbf{T}'|$ or, equivalently, using the non-zero eigenvalues of \mathbf{C} , the D -criterion is given $\Psi_D(d) = \prod_{i=1}^{v-1} \lambda_i^{-1}$, over the set of competing designs \mathcal{D} . The optimal designs can be compared on the basis of their efficiencies (Atkinson, Donev and Tobias, 2007). Specifically, the A -efficiency of an arbitrary design $d \in \mathcal{D}$ with respect to the corresponding optimal design is defined as

$$A_{\text{eff}}(d) = \frac{\Psi_A(d_A^*)}{\Psi_A(d)},$$

where $\Psi_A(d_A^*)$ and $\Psi_A(d)$ are the A -scores at the designs d_A^* and d , that is the values of the A -optimality criterion, respectively. Similarly, the D -efficiency of an arbitrary design $d \in \mathcal{D}$ with respect to the corresponding optimal design is defined as

$$D_{\text{eff}}(d) = \left\{ \frac{\Psi_D(d_D^*)}{\Psi_D(d)} \right\}^{1/(v-1)},$$

where $\Psi_D(d_D^*)$ and $\Psi_D(d)$ are the D -scores at the designs d_D^* and d , respectively.

3. Candidate designs and an algorithm for constructing optimal row-column designs

3.1. Candidate designs

Consider $\mathcal{D} = \mathcal{D}(v, k, b, \theta)$, the class of candidate row-column designs with v treatments, $k = 2$ dyes, b arrays each of size 2 and assume that the model in Equation (1) can be used to model the observations generated from a design in \mathcal{D} . Then there is a set of $2 \times \binom{v}{2}$ possible arrays from which the arrays in a design d can be taken. In total therefore there are

$$N_d = \binom{2 \times \binom{v}{2}}{b}$$

candidate designs in \mathcal{D} . Certain of the arrays do not comprise all v treatments and certain of the remaining arrays are not connected. There would seem to be no general formula for calculating these numbers. Thus designs can be enumerated for small values of v and b computationally and a computer program was written in GAUSS (2013) to count the total number of designs N_d and the number of connected designs, N_c . As an illustration, the values of N_d and N_c for selected values of v and b are presented in Table 1. It is clear from the table that the number of designs in \mathcal{D} , N_d , and the number of connected designs, N_c , increase as v or b or both increase.

Table 1: Number of candidate designs N_d and the number connected N_c for selected values of v and b .

(v, b)	(3, 3)	(4, 4)	(4, 5)	(4, 6)	(5, 5)	(5, 6)	(5, 10)	(6, 6)
N_d	20	495	792	924	15,504	38,760	184,756	593,775
N_c	20	414	768	920	10,384	33,780	184,426	310,800

3.2. An algorithm for searching for optimal or near-optimal row-column designs

In this paper the following algorithm for constructing A - and D -optimal or near-optimal row-column designs under both linear fixed and mixed effects model settings is used. Note that, under the mixed model setting, the algorithm computes locally A - and D -optimal or near-optimal designs for fixed values of θ directly (Chernoff, 1953). The steps used in the algorithm are presented below:

Step 0: Generate the set of all $v(v-1)$ possible candidate arrays, denoted S .

Step 1: Initial Design

- (a) Construct an initial design by taking a design, denoted d_n , with b arrays at random and without replacement from the set S .
- (b) If the design does not comprise all v treatments or if the design is not connected, that is $\text{rank}(C) < v-1$, go to Step 1(a).

Step 2: Exchange Procedure

Let $\Psi_A(d_n)$ be the A -score of d_n .

- (a) Delete each treatment in turn from the design d_n and calculate the A -score for the resultant design, d_i using Equation (2) to calculate the appropriate C matrix. That is, compute the treatment information matrix C_{d_i} , for $i = 1, \dots, n$, where $n = bk$, by deleting the i th observation from the design d_n one at a time and then calculate the corresponding A -score $\Psi_A(d_i)$, $i = 1, \dots, n$.
- (b) Compute the difference, $\Delta_i = \Psi_A(d_i) - \Psi_A(d_n)$, and sort the results in increasing order. This step is required to identify the least important observation in the initial design. Thus, if Δ_i is a minimum for a given i compared to other differences, then the i th observation will be considered as the least important observation. Since the number of observations in the design is bk , which will be large for large sizes of b , record the first four observations whose deletion will result in the minimum differences Δ_i or record the four deletions for which the A -scores are largest and place these in descending order of scores as d_1, \dots, d_4 . Set $r = 1$.
- (c) For deletion d_r , replace the deleted treatment with that treatment, other than those in the same array as the deleted treatment, for which the A -score is a minimum.
- (d) If $r < 4$ or if the resultant design does not comprise all v treatments or if the design is not connected, set $r = r + 1$ and go to Step 2(b). Otherwise, if there is no improvement in the A -score, proceed to Step 3.

Step 3: Interchange Procedure:

- (a) Perform a dye-flip for each array and adopt that design for which the A -score is a minimum.
- (b) Repeat Step 3(a) until there is no further improvement in the A -score.

Step 4: Run Steps 1, 2 and 3 100 times and select the design or designs for which the A -score is a minimum.

Note that the exchange procedure follows closely that of Eccleston and Jones (1980) and Sarkar et al. (2010) but that the process for generating an initial design and the interchange procedure are new. Note also that the same algorithm can be invoked in order to construct D -optimal or near-optimal designs. The algorithm for both A - and D -optimality has been coded and implemented in the GAUSS (2013) programming language.

4. Results and Discussion

4.1. Optimal row-column designs for the fixed effects model

It is interesting to compare the results obtained using the proposed algorithm with some related results which have been reported in the literature. Sarkar et al. (2010) implemented a comprehensive review of D - and A -optimal or near-optimal designs for two-colour microarray experiments and, in particular, prepared a catalogue of design layouts for 139 optimal or near-optimal row-column designs. For comparison purposes, the parametric combinations used in Sarkar et al. (2010), that is $3 \leq v \leq 10$, $v \leq b \leq \binom{v}{2}$, $11 \leq v = b \leq 25$, and also settings with $v = b$, and $(v, b) = (11, 13), (12, 14), (13, 14), (13, 15)$, were considered here and the attendant D - and A -optimal or near-optimal designs constructed. Out of these 139 parametric combinations, 41 row-column designs which are more A -efficient than those of Sarkar et al. (2010) were found. Thus, for example, for $v = 8$, $b = 13$ and $\theta = 0$, Sarkar et al. (2010) report the row-column design presented in Table 2(a) as being the most A -efficient design, with an A -score 4.4436, whereas in the present study the row-column design in Table 2(b) with an A -score of 4.4238 was constructed. The parametric combinations of these 41 designs are presented in Table B.1 in Appendix B. The design layouts of these designs are available from the authors on request.

Table 2: (a) Row-column A -optimal designs of Sarkar *et al.* (2010) for $v = 8$ and $b = 13$.

Dye 1	6	8	4	3	2	5	4	7	6	2	1	3	8
Dye 2	1	2	5	8	5	6	1	3	8	7	7	5	4

(b) Obtained row-column A -optimal designs for $v = 8$ and $b = 13$.

Dye 1	6	2	2	8	3	1	5	3	7	4	8	4	1
Dye 2	8	7	6	4	6	8	2	5	1	7	5	3	3

4.2. Equireplicate row-column designs

Equireplicate designs for the parametric combinations used in Sarkar et al. (2010), together with those reported in Bailey (2007), are now considered. Among the 139 parametric combinations used, 29 and 45 equireplicate A - and D -optimal row-column designs, respectively, were constructed under the linear fixed effects model, that is for $\theta = 0$. Fifteen of these designs are similar to the equireplicate row-column designs recommended by Bailey (2007). However, more efficient equireplicate A -optimal row-column designs for the four parametric combinations, $(v, b) = (8, 12), (10, 15), (9, 18)$ and $(8, 20)$, were obtained. The design layouts, the numbers of replications and the corresponding A -scores for these four designs are presented in Table 3 and a comparison of the results with those of Bailey (2007) is summarised in Table 4.

Wit et al. (2005) investigated equireplicate designs with replications 2 and 4. However, the designs obtained here for the parametric combinations $(v, b) = (9, 9), (7, 14), (16, 16)$ and $(9, 18)$, and presented in Table 5, are better than those given in Wit et al. (2005) in terms of their A -efficiencies.

Table 5: Summaries of the comparison of the results of the present study with Wit et al. (2005) designs.

<i>v</i>	<i>b</i>	<i>r</i>	Obtained designs																	A-score		
9	9	2	Dye 1	4	6	7	2	5	9	3	8	1									13.3333	
			Dye 2	6	1	9	7	3	8	4	5	2										
7	14	4	Dye 1	2	6	6	1	5	7	7	5	1	4	4	2	3	3				2.7524	
			Dye 2	7	2	1	5	4	6	3	2	7	3	1	4	6	5					
16	16	2	Dye 1	14	9	4	16	15	7	6	12	10	11	5	13	2	1	8	3		28.3750	
			Dye 2	3	3	3	10	3	3	3	3	5	3	3	3	3	3	3	16			
9	18	4	Dye 1	7	9	8	3	2	9	5	6	2	6	7	4	4	1	1	3	8	5	3.9128
			Dye 2	6	5	2	8	4	1	3	2	5	1	9	3	7	8	4	6	9	7	

<i>v</i>	<i>b</i>	<i>r</i>	Wit, Nobile and Khanin (2005) designs																	A-score		
9	9	2	Dye 1	1	9	8	2	2	3	3	3	3										25.7778
			Dye 2	3	3	3	8	5	7	6	5	4										
7	14	4	Dye 1	1	1	1	4	4	4	4	6	6	3	2	5	5	7					3.4571
			Dye 2	7	5	6	1	3	2	6	3	2	7	7	2	3	5					
16	16	2	Dye 1	5	5	5	5	5	5	5	5	5	6	7	8	1	2	3	4			31.7500
			Dye 2	16	15	14	13	12	11	10	9	6	7	8	1	2	3	4	5			
9	18	4	Dye 1	9	9	9	8	7	7	7	5	5	4	4	4	3	3	3	2	2	2	4.5562
			Dye 2	8	6	1	2	8	1	6	6	8	5	7	9	5	4	1	1	6	3	

Under the linear mixed effects model 45 equireplicate optimal row-column designs which are A- and D-optimal or near-optimal for θ values of 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8 and 0.9 have been obtained. The design layouts and numbers of replications for selected v , b and θ values are presented in Table 6. A complete list of parametric combinations for these 45 designs are presented in Table B.2 in the Appendix B. Note that among the designs for which parametric combinations are listed, only those highlighted in bold are A-optimal row-column designs for the fixed effects model with numbers of replications equal to 2, that is for $\theta = 0$ and $r = 2$. Note also that the optimal designs with an odd number of replications of treatments are only given for even values of treatment numbers v (Bailey, 2007).

Table 6: Selected equireplicate row-column designs for linear mixed effects models.

<i>v</i>	<i>b</i>	<i>r</i>	θ	Obtained designs																				
7	14	4	0.1	Dye 1	2	7	1	7	6	5	2	4	4	3	6	5	3	1						
				Dye 2	4	5	3	3	1	2	6	1	7	2	7	4	6	5						
10	15	3	0.9	Dye 1	7	10	8	7	2	4	4	9	3	5	6	5	8	1	3					
				Dye 2	10	3	1	2	4	9	6	8	2	1	5	10	6	7	9					
15	15	2	0.6	Dye 1	8	7	6	13	9	1	15	10	3	14	5	12	2	11	4					
				Dye 2	12	9	5	2	13	14	7	4	6	11	1	3	8	10	15					
8	20	5	0.3	Dye 1	6	7	2	4	3	6	1	3	5	4	6	1	8	7	4	8	8	5	2	1
				Dye 2	5	8	7	7	5	3	2	1	8	3	4	5	2	6	2	3	1	4	6	7

4.3. Robustness aspects of A - and D -optimal row-column designs

To study the robustness properties of A - and D -optimal row-column designs for different values of θ , A - and D -optimal row-column designs under both fixed and mixed effects models for different values of θ were constructed. The designs so generated were then compared in order to ascertain whether or not the A - and D -optimal row-column designs under a fixed effects model are also A - and D -optimal under a mixed effects model. Among the parametric combinations considered, 83 and 58 have at least two different A - and D -optimal row-column designs, respectively, for various values of θ . This indicates that, for the respective parametric combinations, A - and D -optimal row-column designs under the linear fixed effects model are not necessarily A - and D -optimal under the linear mixed effects model. To illustrate this, consider a two-colour microarray experiment with the same number of arrays as treatments, that is $v = b$. The A - and D -optimal row-column designs results for $3 \leq v = b \leq 25$ are summarized below:

- (i) When $v \leq 9$, the A -optimal designs for $0 \leq \theta < 1$ are the loop designs $C_v^*(v)$, where $C_s^*(v)$ represents a design for v treatments in v arrays for which the associated graphs contain a circuit of length s (Bailey, 2007).
- (ii) When $v = 10$, the designs $C_4^*(10)$ and $C_5^*(10)$ are both A -optimal for $\theta = 0$, the design $C_5^*(10)$ is A -optimal for $0 < \theta < 0.00825$ and the loop design $C_{10}^*(10)$ is A -optimal for $0.00825 < \theta < 1$. These designs are shown in Figure 1.
- (iii) When $v = 11$, the design $C_4^*(11)$ is A -optimal for $0 \leq \theta < 0.00776$, the design $C_5^*(11)$ is A -optimal for $0.00776 < \theta < 0.02091$ and the loop design $C_{11}^*(11)$ is A -optimal for $0.02091 < \theta < 1$.
- (iv) When $v = 12$, the design $C_4^*(12)$ is A -optimal for $0 \leq \theta < 0.01335$, the design $C_5^*(12)$ is A -optimal for $0.01335 < \theta < 0.02833$, the design $C_6^*(12)$ is A -optimal for $0.02833 < \theta < 0.02983$ and the loop design $C_{12}^*(12)$ is A -optimal for $0.02983 < \theta < 1$.
- (v) When $v = 17$, the design $C_4^*(17)$ is A -optimal for $0 \leq \theta < 0.02744$, the design $C_5^*(17)$ is A -optimal for $0.02744 < \theta < 0.04174$, $C_6^*(17)$ is the A -optimal for $0.04174 < \theta < 0.04959$, the design $C_7^*(17)$ is A -optimal for $0.04959 < \theta < 0.05025$ and the loop design $C_{17}^*(17)$ is A -optimal design for $0.05025 < \theta < 1$.
- (vi) When $v = 18$, the designs $C_3^*(18)$ and $C_4^*(18)$ are both A -optimal designs for $\theta = 0$, the design $C_4^*(18)$ is A -optimal for $0 < \theta < 0.02894$, the design $C_5^*(18)$ is A -optimal for $0.02894 < \theta < 0.04307$, the design $C_6^*(18)$ is A -optimal for $0.04307 < \theta < 0.05082$, the design $C_7^*(18)$ is A -optimal for $0.05082 < \theta < 0.05222$ and the loop design $C_{18}^*(18)$ is A -optimal design for $0.05222 < \theta < 1$.
- (vii) When $v = 25$, the design $C_3^*(25)$ is A -optimal for $0 \leq \theta < 0.00773$, the design $C_4^*(25)$ is A -optimal for $0.00773 < \theta < 0.03531$, the design $C_5^*(25)$ is A -optimal for $0.03531 < \theta < 0.04844$, the design $C_6^*(25)$ is A -optimal for for $0.04844 < \theta < 0.05549$, the design $C_7^*(25)$ is A -optimal for $0.05549 < \theta < 0.05982$ and the loop design $C_{25}^*(25)$ is A -optimal design for $0.05982 < \theta < 1$.
- (viii) When $3 \leq v \leq 25$, the D -optimal design is the loop design $C_v^*(v)$ for $0 \leq \theta < 1$.

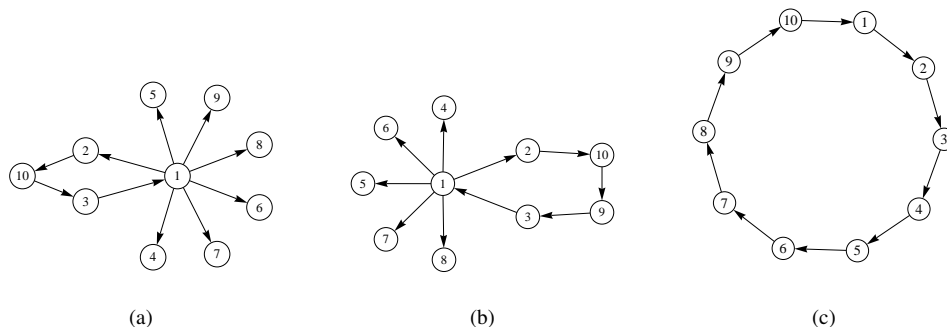


Figure 1: Three designs for ten treatments in ten arrays: $C_4^*(10)$ and $C_5^*(10)$ in (a) and (b) both are A -optimal designs for $\theta = 0$; $C_5^*(10)$ in (b) is A -optimal for $0 < \theta < 0.00825$; the loop design $C_{10}^*(10)$ in (c) is A -optimal for $0.00825 < \theta < 1$ and D -optimal for all θ .

The results for the A - and D -optimal designs when $\theta = 0$ agree with those of Bailey (2007). Plots of the A -scores against θ are shown in Figure 2 for $v = 10$ and $v = 18$ and indicate that the loop design is A -optimal for large values of θ . The plots also demonstrate that the A -optimal designs for $\theta = 0$ are not necessarily A -optimal for $\theta > 0$. The cut-off values given for θ in the list above are approximate and are obtained graphically using the "get coordinate" option in Mathematica (2013). Note that there will be two different designs which are A -optimal at each cut-off value for θ . For example, the two designs $C_{10}^*(10)$ and $C_5^*(10)$ will be optimal at the point where the curves in Figure 2(a) for $C_{10}^*(10)$ and $C_5^*(10)$ intersect. Finally, Tables B.3 and B.4 in Appendix B summarize the parametric combinations for A - and D -optimal designs, respectively, which are non-robust against the models under consideration. Recall that a design is non-robust against the linear mixed effects model if the design does not remain A - or D -efficient for every possible value of θ .

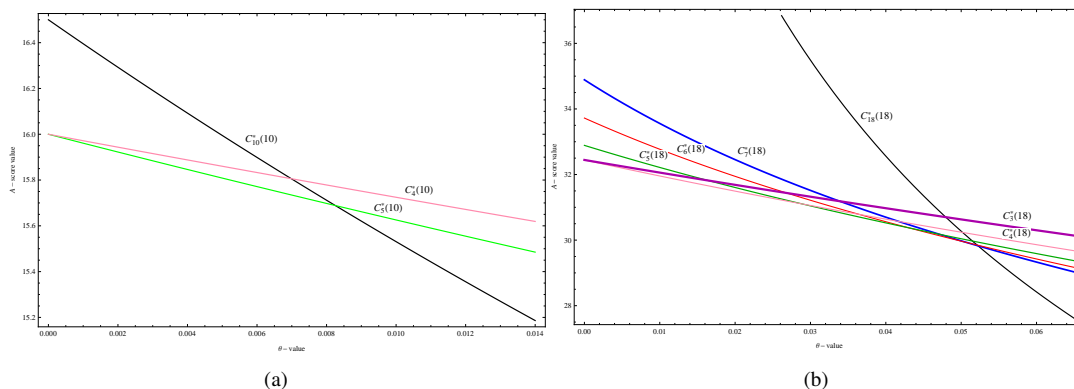


Figure 2: Plot of A -score against θ : for $v = 10$ in (a) and for $v = 18$ in (b).

5. Conclusion

The aim of this paper has been to compute A - and D -optimal designs for microarray experiments numerically for given parametric combinations (v, k, b, θ) where $k = 2$ and for $\theta = \sigma_b^2 / (\sigma^2 + k\sigma_b^2) \in (0, 1)$ where σ^2 and σ_b^2 are the error variance and the variance of random array effects, respectively. One of the key features of this study is the application of the exchange-interchange algorithm to search for microarray designs for which interest focuses on every pairwise comparison of treatment effects. A second important feature is the construction of A - and D -optimal or near-optimal designs for the linear mixed effects model with random array effects for a microarray experiment. In other words, since the random array effects introduce an additional term into the \mathbf{C} matrix of the linear fixed effects model, the A - and D -optimal designs for estimation of all pairwise contrasts of treatment effects found under the linear fixed effects model may not be optimal under the linear mixed effects model. Thus, while A -optimal row-column designs for two-colour microarray experiments under the linear fixed effects model are available in the literature, as detailed in particular by Bailey (2007) and by Sarkar et al. (2010), the designs discussed in the present paper represent viable options for application in microarray experiments for which $\theta > 0$. Note that the results discussed here focus mainly on two-colour microarray experiments but that the exchange-interchange algorithm can be easily modified to generate optimal designs for, for example, three-colour microarrays and also to accommodate other optimality criteria, such as E - and MV -optimality. Full details will be reported elsewhere. Finally note that, since the parameter θ in the \mathbf{C} matrix is usually unknown, the A - and D -optimal designs must necessarily be computed using a best guess for θ . Thus the A - and D -optimal designs presented here must be considered as locally optimal in the sense of Chernoff (1953). Work is currently in progress to address this issue by applying a Bayesian approach to the construction of optimal designs for row-column designs under the linear mixed effects model setting.

Acknowledgements

This research is supported by grants from the National Research Foundation (NRF) through the Competitive Support for Unrated Researchers Programme (Reference: SUR20110629000019843 and Grant No: 80407) and the Incentive Funding for Rated Researchers Programme (Grant No: UID 85456), the University of Pretoria and the University of Cape Town. The authors would like to thank Prof Ananta Sarkar for providing us with the catalogue of design layouts of row-column designs. Any opinion, finding and conclusion or recommendation expressed in this material is that of the authors and the NRF does not accept liability in this regard.

References

- ATKINSON, A. C., DONEV, A. N., AND TOBIAS, R. D. (2007). *Optimum Experimental Designs, with SAS*. Oxford University Press: Oxford, New York.
- BAILEY, R. A. (2007). Designs for two-colour microarray experiments. *Journal of the Royal Statistical Society*, **56**, 365–394.
- BAILEY, R. A. (2009). Designs for dose-escalation trials with quantitative responses. *Statistics in Medicine*, **28**, 3721–3738.

- BAILEY, R. A., SCHIFFL, K., AND HILGERS, R. D. (2013). A note on robustness of D -optimal block designs for two-colour microarray experiments. *Journal of Statistical Planning and Inference*, **143**, 1195–1202.
- CHERNOFF, H. (1953). Locally optimal designs for estimating parameters. *The Annals of Mathematical Statistics*, **24**, 586–602.
- CHURCHILL, G. A. (2002). Fundamentals of experimental design for cDNA microarrays. *Nature Genetics*, **32**, 490–495.
- DEY, A. (2010). *Incomplete Block Designs*. World Scientific: London.
- ECCLESTON, J. A. AND JONES, B. (1980). Exchange and interchange procedures to search for optimal row-column designs. *Journal of the Royal Statistical Society, Series B*, **42**, 372–376.
- GAUSS PROGRAMMING LANGUAGE (2013). Aptech Systems, Inc.: Black Diamond, WA.
- GEMECHU, D. B., DEBUSHO, L. K., AND HAINES, L. M. (2014). A -optimal designs for two-colour cDNA microarray experiments using the linear mixed effects model. *In Proceedings of the Annual Conference of the South African Statistical Association (SASA 2014)*. pp. 33–40. ISBN: 978-1-86822-659-7.
- GROSSMAN, H. AND SCHWABE, R. (2008). The relationship between optimal designs for microarray and paired comparison experiments. Accessed on 20 June 2013.
URL: <http://www.math.uni-magdeburg.de/~schwabe/>
- HAINES, L. M. (2015). An introduction to linear models. *In* DEAN, A., MORRIS, M., STUFKEN, J., AND BINGHAM, D. (Editors) *The Handbook of Design and Analysis of Experiments*. Chapman and Hall/CRC, Series: Chapman & Hall/CRC Handbooks of Modern Statistical Methods, Chapman and Hall/CRC: New York.
- KERR, M. K. (2003). Design considerations for efficient and effective microarray studies. *Biometrics*, **59**, 822–828.
- KERR, M. K. AND CHURCHILL, G. A. (2001a). Experimental design for gene expression microarrays. *Biostatistics*, **2** (2), 183–201.
- KERR, M. K. AND CHURCHILL, G. A. (2001b). Statistical design and the analysis of gene expression. *Genetics*, **77**, 123–128.
- LANDGREBE, J., BRETZ, F., AND BRUNNER, R. (2006). Efficient design and analysis of two-color factorial microarray design. *Computational Statistical Data Analysis*, **50**, 499–517.
- LEE, M. L. T. (2004). *Analysis of Microarray Gene Expression Data*. Springer: New York.
- MATHEMATICA, VERSION 9.0.1.0 (2013). Wolfram Research, Inc.: Champaign, IL.
- SARKAR, A., PRASAD, R., GUPTHA, V. K., KASHINATH, C., AND RATHORE, A. (2010). Efficient row-column designs for microarray experiments. *Journal of the Indian Society of the Agriculture Statistics*, **64** (1), 89–117.
- SHAH, K. AND SINHA, B. K. (1989). *Theory of Optimal Designs*. Springer: New York.
- WIT, E., NOBILE, A., AND KHANIN, R. (2005). Near-optimal designs for dual channel microarray studies. *Applied Statistics*, **54**, 817–830.
- WOLFINGER, R. D., GIBSON, G., WOLFINGER, E. D., BENNETT, L., HAMADEH, H., BUSHEL, P., AFSHARI, C., AND PAULES, R. S. (2001). Assessing gene significance from cDNA microarray expression data via mixed models. *Journal of Computational Biology*, **8**, 625–637.
- YANG, Y. H. AND SPEED, T. (2002). Design issues for cDNA microarray experiments. *Nature Reviews Genetics*, **3**, 579–588.

Appendix A

The model in Equation (1) can also be expressed in a general linear mixed model form as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon} \quad (3)$$

where $\mathbf{X} = [\mathbf{1} \ \mathbf{X}_r \ \mathbf{X}_t]$, $\boldsymbol{\beta}' = (\mu \ \boldsymbol{\alpha}' \ \boldsymbol{\tau}')$ and $\mathbf{Z} = \mathbf{X}_c$. Under the assumptions $\mathbf{b} \sim \mathcal{N}(\mathbf{0}, \sigma_b^2 \mathbf{I}_b)$, $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ and \mathbf{b} and \mathbf{e} are independent, the mean and variance of \mathbf{y} are given by $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ and $\mathbf{V} = \sigma_b^2 \mathbf{X}_c \mathbf{X}_c' + \sigma^2 \mathbf{I}_n$, respectively. The normal equations related to vector of fixed effects, $\boldsymbol{\beta}$, in Equation (3) are readily derived and can be expressed succinctly as

$$(\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}) \hat{\boldsymbol{\beta}} = \mathbf{X}' \mathbf{V}^{-1} \mathbf{y},$$

where $\hat{\boldsymbol{\beta}}$ is the maximum likelihood estimator of $\boldsymbol{\beta}$ and this can be re-expressed as

$$\begin{pmatrix} \mathbf{1}'_n \mathbf{V}^{-1} \mathbf{1}_n & \mathbf{1}'_n \mathbf{V}^{-1} \mathbf{X}_r & \mathbf{1}'_n \mathbf{V}^{-1} \mathbf{X}_t \\ \mathbf{X}'_r \mathbf{V}^{-1} \mathbf{1}_n & \mathbf{X}'_r \mathbf{V}^{-1} \mathbf{X}_r & \mathbf{X}'_r \mathbf{V}^{-1} \mathbf{X}_t \\ \mathbf{X}'_t \mathbf{V}^{-1} \mathbf{1}_n & \mathbf{X}'_t \mathbf{V}^{-1} \mathbf{X}_r & \mathbf{X}'_t \mathbf{V}^{-1} \mathbf{X}_t \end{pmatrix} \begin{pmatrix} \hat{\mu} \\ \hat{\boldsymbol{\alpha}} \\ \hat{\boldsymbol{\tau}} \end{pmatrix} = \begin{pmatrix} \mathbf{1}'_n \mathbf{V}^{-1} \mathbf{y} \\ \mathbf{X}'_r \mathbf{V}^{-1} \mathbf{y} \\ \mathbf{X}'_t \mathbf{V}^{-1} \mathbf{y} \end{pmatrix}$$

Furthermore, performing some straightforward matrix algebra, the reduced normal equations for treatment effects, $\boldsymbol{\tau}$, eliminating the fixed mean and row effects can be given by

$$\begin{aligned} & \left[\mathbf{X}'_t \mathbf{V}^{-1} \mathbf{X}_t - (\mathbf{X}'_t \mathbf{V}^{-1} \mathbf{X}_r) (\mathbf{X}'_r \mathbf{V}^{-1} \mathbf{X}_r)^{-1} \mathbf{X}'_r \mathbf{V}^{-1} \mathbf{X}_t \right] \hat{\boldsymbol{\tau}} = \\ & \left[\mathbf{X}'_t \mathbf{V}^{-1} - (\mathbf{X}'_t \mathbf{V}^{-1} \mathbf{X}_r) (\mathbf{X}'_r \mathbf{V}^{-1} \mathbf{X}_r)^{-1} \mathbf{X}'_r \mathbf{V}^{-1} \right] \mathbf{y} \end{aligned} \quad (4)$$

which we write in abbreviated form as

$$\mathbf{C} \hat{\boldsymbol{\tau}} = \mathbf{q}$$

with $\mathbf{C} = \mathbf{X}'_t \mathbf{V}^{-1} \mathbf{X}_t - (\mathbf{X}'_t \mathbf{V}^{-1} \mathbf{X}_r) (\mathbf{X}'_r \mathbf{V}^{-1} \mathbf{X}_r)^{-1} \mathbf{X}'_r \mathbf{V}^{-1} \mathbf{X}_t$ and \mathbf{q} is the term in the right hand side of Equation (4).

Using the inverse of the variance matrix, $\mathbf{V} = \sigma_c^2 \mathbf{X}_c \mathbf{X}_c' + \sigma^2 \mathbf{I}_n$, which is given by

$$\mathbf{V}^{-1} = \frac{1}{\sigma^2} (\mathbf{I}_n - \mathbf{X}_c \mathbf{W} \mathbf{X}_c')$$

where $\mathbf{W} = \text{diag}(\sigma_b^2/(\sigma^2 + k_1 \sigma_b^2), \dots, \sigma_b^2/(\sigma^2 + k_b \sigma_b^2))$, the terms in the \mathbf{C} matrix can be expressed as

1. $\mathbf{X}'_r \mathbf{V}^{-1} \mathbf{X}_r = \mathbf{B} - \mathbf{SWS}'$,
2. $\mathbf{X}'_t \mathbf{V}^{-1} \mathbf{X}_r = \mathbf{M} - \mathbf{NWS}'$, and
3. $\mathbf{X}'_t \mathbf{V}^{-1} \mathbf{X}_t = \mathbf{R} - \mathbf{NWN}'$.

Thus, substituting these results in Equation (4) it is readily follows, ignoring the constant $\frac{1}{\sigma^2}$, that

$$\mathbf{C} = \mathbf{R} - \mathbf{NWN}' - (\mathbf{M} - \mathbf{NWS}') (\mathbf{B} - \mathbf{SWS}')^{-1} (\mathbf{M}' - \mathbf{SWN}')$$

Appendix B

Table B.1: Lists parametric combinations of A-optimal row-column designs which are more efficient than Sarkar et al. (2010).

v	b	v	b	v	b	v	b
7	14	9	25	10	23	12	14
8	13	9	27	10	25	13	13
8	16	9	28	10	26	13	14
8	20	9	30	10	29	13	15
9	10	9	33	10	30	14	14
9	14	9	35	10	31	15	15
9	16	10	10	10	34	16	16
9	17	10	11	10	35	17	17
9	22	10	13	10	40		
9	23	10	14	11	11		
9	24	10	20	12	12		

Table B.2: Lists of parametric combinations of equireplicate A-and D-optimal row-column designs for all possible values of θ under consideration.

v	b	r	v	b	r	v	b	r	v	b	r
3	3	2	15	15	2	6	9	3	7	21	6
4	4	2	16	16	2	8	12	3	8	24	6
5	5	2	17	17	2	10	15	3	9	27	6
6	6	2	18	18	2	5	10	4	10	30	6
7	7	2	19	19	2	6	12	4	8	28	7
8	8	2	20	20	2	7	14	4	10	35	7
9	9	2	21	21	2	8	16	4	9	36	8
10	10	2	22	22	2	9	18	4	10	40	8
11	11	2	23	23	2	10	20	4	10	45	9
12	12	2	24	24	2	6	15	5			
13	13	2	25	25	2	8	20	5			
14	14	2	4	6	3	10	25	5			

Table B.3: Lists of parametric combinations of *A*-optimal row-column designs which are non-robust against the model.

v	b	v	b	v	b	v	b	v	b	v	b	v	b
4	5	7	8	8	10	9	11	9	26	10	17	12	14
4	6	7	9	8	11	9	12	9	29	10	18	13	14
5	6	7	11	8	12	9	13	9	30	10	19	13	15
5	7	7	12	8	13	9	14	9	33	10	25	11	11
5	9	7	13	8	14	9	15	10	10	10	26	12	12
6	7	7	15	8	15	9	16	10	11	10	28	13	13
6	8	7	17	8	18	9	19	10	12	10	31	14	14
6	10	7	18	8	19	9	20	10	13	10	35	15	15
6	13	7	18	8	20	9	21	10	14	10	36	16	16
6	14	7	20	8	22	9	24	10	15	10	38	17	17
6	15	8	9	9	10	9	25	10	16	11	13	18	18

Table B.4: Lists of parametric combinations of *D*-optimal row-column designs which are non-robust against the model.

v	b	v	b	v	b	v	b	v	b
4	5	7	11	8	15	9	21	10	19
4	6	7	12	8	18	9	25	10	26
5	6	7	13	8	20	9	26	10	28
5	7	7	15	8	22	9	29	10	31
5	9	7	16	9	10	9	30	10	36
6	7	7	17	9	11	10	11	10	38
6	8	7	18	9	12	10	12	11	13
6	10	8	9	9	14	10	13	12	14
6	14	8	10	9	15	10	14	13	14
6	15	8	11	9	16	10	15	13	15
7	8	8	13	9	19	10	17		
7	9	8	14	9	20	10	18		