

## ADDITIONAL REFERENCES

- Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. *2nd Internat. Symp. on Information Theory* (Petrov, B.N. and Csaki, F., eds.) 267–281, Akademiai Kiado, Budapest. (Reproduced in *Breakthroughs in Statistics*, 1, 1992 (S. Kotz and N.L. Johnson, eds.) 610–624, Springer-Verlag, New York.)
- Kass, R. E., Tierney, L. and Kadane, J. B. (1990). The validity of posterior asymptotic expansions based on Laplace's method. In *Bayesian and Likelihood Methods in Statistics and Econometrics*, (S. Geisser, J. S. Hodges, S. J. Press and A. Zellner, eds.), North-Holland, New York.
- Tierney, L. and Kadane, J. B. (1986). Accurate approximations for posterior moments and marginal densities. *J. Amer. Statist. Assoc.* **81**, 82 - 86.

Rahul Mukerjee

Indian Institute of Management

I begin by congratulating the authors, Professors Rao and Wu, on this very illuminating and scholarly piece of work which will inspire future researchers in this area. They have done an enormous job of which we are the beneficiaries.

Considerable attention has been given in this paper on the important problem of selecting an appropriate sub-model starting from the linear model (2.1). I, therefore, find it relevant to briefly discuss some related issues in design of experiments. The discussion will be focussed primarily on discrete designs. Incidentally, experimental design problems under model uncertainty have been of substantial interest in recent years (Dey and Mukerjee, 1999; Wu and Hamada, 2000).

To motivate the ideas, consider the setup of a  $2^n$  factorial experiment, a situation where there are  $n$  factors each at two levels. Suppose interest lies in identifying the active factors, i.e., the ones with nonzero main effects, under the absence of all interactions. A *factor screening experiment* is one which can achieve this. Interpreting the factors as regressors, the problem here is the same as that initiated by (2.1) and (2.2). The model (2.1) now consists of the general mean and the main effects of the two-level factors, each main effect being represented by a single parameter. Clearly, then at least  $n + 1$  observations are needed to examine (2.1) and all possible sub-models thereof.

Unfortunately, in many practical situations, especially in exploratory studies,  $n$  can be quite large and even  $n + 1$  observations are not affordable. This poses additional concerns. A way out is possible via consideration of the phenomenon of effect sparsity. Often it is known that among the  $n$  factors at most  $k$  are active where the known quantity  $k$  is small compared to  $n$ . However, there is no knowledge about which factors are active. The problem is then to identify the active factors, which are at most  $k$  in number, and to estimate the corresponding main effects as well as the general mean.

The notion of search designs, pioneered by Srivastava (1975), plays a crucial role in handling problems of this kind. To highlight the underlying ideas, we consider a linear model  $E(\mathbf{Y}) = X_1\beta_1 + X_2\beta_2$ ,  $\text{Cov}(\mathbf{Y}) = \sigma^2 I_N$ , where  $\mathbf{Y}$  is the  $N \times 1$  observational vector,  $X = [X_1, X_2]$  is the known design matrix,  $\beta_1$  is an unknown parametric vector,  $\sigma^2$  is the common variance of the observations and  $I_N$  is the  $N \times N$  identity matrix. The parametric vector  $\beta_2$  is partially known. It is known that at most  $k$  its elements are nonzero, where  $k$  is small compared to the dimension of  $\beta_2$ . No knowledge is, however, available about which elements of  $\beta_2$  are possibly nonzero and what their values are. Interest lies in estimating  $\beta_1$  and searching and estimating the possibly nonzero elements of  $\beta_2$ . Following Srivastava (1975), an experimental design  $d$  (associated with a choice of  $X$ ), that enables one to achieve this, is called a *search design* with resolving power  $(\beta_1, \beta_2, k)$ . Observe that if  $\beta_1$  is taken as the scalar representing the general mean and  $\beta_2$  is taken as the vector of the main effects of the factors then the search design problem reduces to the factor screening problem introduced above. Of course, the idea of search designs applies to many other situations. For example, one can consider more complex factorials including asymmetric factorials and entertain interactions too in addition to the general mean and main effects.

The following result, due to Srivastava (1975), is a fundamental tool in the study of search designs.

**Theorem 1.** *For a design to have resolving power  $(\beta_1, \beta_2, k)$  in the above setup, it is necessary that for every  $N \times 2k$  submatrix  $X_{20}$  of  $X_2$ , the matrix  $[X_1, X_{20}]$  has full column rank. Furthermore, in the noiseless case  $\sigma^2 = 0$ , this condition is also sufficient for to have resolving power  $(\beta_1, \beta_2, k)$ .*

For  $\sigma^2 = 0$ , if the rank condition of Theorem 1 is satisfied then the true model can be identified with probability unity. Further discussion on the actual identification of the true model is available in Srivastava (1975). The design problem here consists of finding a design such that the rank condition holds. The construction of such experimental plans and related combinatorics have received much attention in the literature and we refer to Gupta (1990) and Ghosh (1996) for informative reviews. The latter article also discusses the important issue of sequential experimentation in this and related contexts.

For  $\sigma^2 > 0$ , even under the rank condition, detection of the true model is not possible with probability unity. Discussion on the search procedure in this case is available in Srivastava (1975, 1996). The design problem here consists of ensuring the rank condition as well as attaining a high probability of correct search. Significantly new grounds have been broken in this direction by Shirakura et al. (1996). The actual construction of optimal designs for the case  $\sigma^2 > 0$  deserves further attention.

With reference to the problem of factor screening, *supersaturated designs* have also been of much interest. With a  $2^n$  factorial, under the absence of all interactions, let  $X$  be the  $N \times (n+1)$  design matrix where  $N$  is the number of observations and the columns of  $X$  correspond to the general mean and the  $n$  main effects. As before, let  $N < n+1$ , and suppose the objective is to identify the active factors under effect sparsity. Since  $N < n+1$ , one cannot achieve orthogonality in the sense of making  $X^T X$  a diagonal matrix, where  $X^T$  is the transpose of  $X$ . The idea of supersaturated designs, which dates back to Booth and Cox (1962) and has witnessed a revival of interest in recent years, essentially aims at choosing  $X$  such that the sum of squares of the off-diagonal elements of  $X^T X$  is minimized. Li and Wu (1997) describe several other criteria for choosing such designs. Data analytic techniques for the use of supersaturated designs in the identification of the correct model, via detection of the nonzero main effects, have been discussed in Lin (1993, 1995); see Dey and Mukerjee (1999) and Wu and Hamada (2000) for more references on supersaturated designs.

Cheng et al. (1999) consider another approach towards the design problem for the study of model robustness and model selection. In connection with regular fractions of symmetric factorials, i.e., the ones specified by a system of linear equations over a finite field, they introduce the notion of *estimation capacity* which is a criterion of model robustness. The objective here is to retain full information on the main effects and as much information as possible on the two-factor interactions in the sense of entertaining the maximum possible model diversity, under the absence of interactions involving three or more factors. Cheng and Mukerjee (1998) report further theoretical results on designs with maximum estimation capacity.

Turning to continuous experimental designs, an important design problem in model selection is the one where an objective is the identification of the appropriate degree of the polynomial in a polynomial regression model. Innovative results on this problem, via the use of canonical moments, have been reported by Dette (1995) and Franke (2000), where further references are available.

To summarize, the experimental design problem for model selection has already been an active area of research and even greater activity, catering to both frequentist and Bayesian inference, is anticipated in this area in the near future. The elegant exposi-

tion given by Professors Rao and Wu in the present work will definitely act as a great stimulator for future research in this direction.

#### ADDITIONAL REFERENCES

- Booth, K.H.V. and Cox, D.R. (1962). Some systematic supersaturated designs. *Technometrics* **4**, 489-495.
- Cheng, C.S. and Mukerjee, R. (1998). Regular fractional factorial designs with minimum aberration and maximum estimation capacity. *Ann. Statist.* **26**, 2289-2300.
- Cheng, C.S., Steinberg, D.M. and Sun, D.X. (1999). Minimum aberration and model robustness for two-level factorial designs. *J. Roy. Statist. Soc. B* **61**, 85-93.
- Dette, H. (1995) Optimal designs for identifying the degree of a polynomial regression. *Ann. Statist.* **23**, 1248-1266.
- Dey, A. and Mukerjee, R. (1999). *Fractional Factorial Plans*. John Wiley, New York.
- Franke, T. (2000). *D- und  $D_1$ -optimale Versuchsplane unter Nebenbedingungen und gewichtete Maximin-Versuchsplane bei polynomialer regression*. Dissertation, Ruhr-Universität Bochum.
- Ghosh, S. (1996). Sequential assembly of fractions in factorial experiments. In *Handbook of Statistics*, **13** (S. Ghosh and C.R. Rao, eds.), 407-435. North-Holland, Amsterdam.
- Gupta, B.C. (1990). A survey of search designs for factorial experiments. In *Probability, Statistics and Design of Experiments* (R.R. Bahadur, ed.), 329-345. Wiley Eastern, New Delhi.
- Li, W.W. and Wu, C.F.J. (1997). Columnwise-pairwise algorithms with applications to the construction of supersaturated designs. *Technometrics* **39**, 171-179.
- Lin, D.K.J. (1993). A new class of supersaturated designs. *Technometrics* **35**, 28-31.
- Lin, D.K.J. (1995). Generating systematic supersaturated designs. *Technometrics* **37**, 213-225.
- Shirakura, T., Takahashi, T. and Srivastava, J.N. (1996). Searching probabilities for nonzero effects in search designs for the noisy case. *Ann. Statist.* **24**, 2560-2568.
- Srivastava, J.N. (1975). Designs for searching non-negligible effects. In *A Survey of Statistical Design and Linear Models* (J.N. Srivastava, ed.), 507-519. North-Holland, Amsterdam.

- Srivastava, J.N. (1996). A critique of some aspects of experimental design. In *Handbook of Statistics*, **13** (S. Ghosh and C.R. Rao, eds.), 309-341. North-Holland, Amsterdam.
- Wu, C.F.J. and Hamada, M. (2000). *Experiments: Planning, Analysis and Parameter Design Optimization*. John Wiley, New York.

## REJOINDER

C. R. Rao and Y. Wu

The authors would like to thank S. Konishi and R. Mukerjee for their valuable comments. Konishi suggests an extension of the GIC criterion using a penalized maximum likelihood estimator of the unknown parameters. This new method may provide some robustness to the choice of a model. To what extent is the selection of the model affected by the particular choice of the prior distribution of parameters and models suggested by Konishi needs some investigation. Mukerjee raised the problem of design of experiments to provide the minimum number of observations needed for model selection ensuring some robustness. This is, indeed, a new area of research, but much depends on the accuracy of apriori information regarding the unknown parameters. For instance, in the example mentioned by Mukerjee, the number of active factors out of a large number  $n$  of factors is known to be a given number  $k < n$ , and the problem is that of generating a minimum number of observations to determine which subset of  $k$  factors is active. It would be interesting, perhaps more relevant in practice, to know whether supersaturated designs suggested for this purpose can also be used to select a subset of factors which are more active than the others. The problem of model selection needs more discussion in terms of objectives, the use of prior information, appropriate methodology and robustness. We hope our review with the additional material contributed by Konishi and Mukerjee will stimulate further research in statistical model selection.