# Non-parametric optimal design in dose finding studies

JOHN O'QUIGLEY*

*Department of Mathematics, University of California, San Diego, CA 92093, USA*
joquigle@ucsd.edu

XAVIER PAOLETTI, JEAN MACCARIO

*Unité 472 de l'INSERM, 16 Avenue P.V. Couturier, 94807 Villejuif, France*

SUMMARY

We describe a non-parametric optimal design as a theoretical gold standard for dose finding studies. Its purpose is analogous to the Cramer–Rao bound for unbiased estimators, i.e. it provides a bound beyond which improvements are not generally possible. The bound applies to the class of non-parametric designs where the data are not assumed to be generated by any known parametric model. Whenever parametric assumptions really hold it may be possible to do better than the optimal non-parametric design. The goal is to be able to compare any potential dose finding scheme with the optimal non-parametric benchmark. This paper makes precise what is meant by optimal in this context and also why the procedure is described as non-parametric.

*Keywords*: Clinical trials; Dose finding; Efficiency; Non-parametric; Optimality; Phase I.

## 1. INTRODUCTION

In oncology the anti-cancer effect of a new treatment generally increases with dose. However, the probability of toxicity also increases with dose. In practice, a certain amount of severe toxicity is accepted in order to have some anti-tumour effect. This means that for a given dose a proportion of patients will suffer toxic side effects. The maximum tolerated dose (MTD) is defined for a population of patients as the highest acceptable dose level among several levels. For an individual, toxic side effects are often reported following NCI grades and take values between 0 (no toxicity) and 4 (severe, possibly life threatening toxicity). In most phase I dose finding studies, the outcome variable of interest is a severe toxicity and is coded as a dichotomy (0, 1) i.e. above a certain grade threshold, the toxicity observed for an individual is reported as a dose limiting toxicity (DLT) and is coded by a 1. Otherwise the toxic response is coded by a 0. The distribution of the DLT among the patient population generates the dose–toxicity relation. We make the fundamental assumption that this relation is monotonic increasing. Eventually, the problem of identifying the MTD translates into one of targeting some percentile of this relation.

A key assumption here is that, for the toxic side effects of interest, the hypothesis of the toxicity increasing with dose results from the fact that a subject who experiences a severe toxicity at a given dose level would have had a severe toxicity at every higher level. In the same way we are leaning on the implicit assumption that a subject tolerating the treatment at some given dose level would have tolerated

*To whom correspondence should be addressed

the treatment at every lower level. This relatively uncontroversial, and easily accepted, albeit debatable, restriction, turns out to be central to the main development.

A number of approaches have been developed to identify the MTD (Anbar, 1978, 1984; O'Quigley *et al.*, 1990; Gatsonis and Greenhouse, 1992; Whitehead, 1997; Whitehead and Williamson, 1998; Storer, 1989, 1993, 1998; Goodman *et al.*, 1995; Faries, 1994; Møller, 1995; Piantadosi and Liu, 1996; O'Quigley and Shen, 1996; O'Quigley and Reiner, 1998). How well does any approach do? This question begs the further question: how well is it possible to do? Of course, the answer depends upon the objectives which may sometimes vary among studies. Here we restrict our attention to cases where the purpose is to identify a level whose associated probability of toxicity is as close as possible to a specified target. As we carry out design improvements we may ask ourselves, is there some limit beyond which further improvements can not be achieved without making strong parametric assumptions about the shape of the dose–toxicity curve? The purpose of this paper is to identify this limit.

We shall adopt the following notation. We consider dose–toxicity relations completely described by the $m$ pairs $(d_1, R_1), \ldots, (d_m, R_m)$ where $d_k$ indicates the dose at level $k$ and $R_k$ the associated probability of toxicity. Let $Y_{jk}$ ($k = 1, \ldots, m$; $j = 1, \ldots, n$) denote the random variable taking the value 1 when a DLT is observed for subject $j$ at level $d_k$. We restrict our attention to monotonic curves for which the toxicity increases with the dose, $R_1 < R_2 < \cdots < R_m$. We target the dose level $d_k$ where the probability of toxicity $R_k$ is closest to a given percentile denoted by $\theta$. More formally, the goal is to identify the level $d_i$ such that $|R_i - \theta| < |R_l - \theta|\ l = 1, \ldots, m; l \neq i$.

## 2. AN OPTIMAL DESIGN

We can define an optimal design using the concept of complete and incomplete information. It is the ability of the optimal method to exploit compete information, typically unavailable in practice, which confers the optimality property. We provide the following definitions to the terms *complete* and *incomplete information.*

### *Incomplete information*

Suppose we have subjects available for experimentation at six dose levels. In a real experiment each patient provides partial or incomplete information. The assumption of monotonic toxicity implies that if a subject had a toxic reaction at level $d_k$ ($k \leqslant 6$) then he/she would necessarily have a toxic reaction at $d_\ell$ ($k \leqslant \ell \leqslant 6$). As for his/her response at levels below $d_k$, we have no information on this. On the other hand, should the subject tolerate the treatment at level $d_k$ ($1 \leqslant k \leqslant 6$) then he/she would necessarily tolerate the treatment at all levels $d_\ell$ ($1 \leqslant \ell \leqslant k$). As an illustration let us suppose that subject $h$ experiences a toxicity at $d_5$ and subject $j$ a non-toxicity at level $d_3$. This is summarized in the table below where a * indicates missing or incomplete information.

|          | Dose  |       |       |       |       |       |
|----------|-------|-------|-------|-------|-------|-------|
|          | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ |
| $Y_{hk}$ | *     | *     | *     | *     | 1     | 1     |
| $Y_{jk}$ | 0     | 0     | 0     | *     | *     | *     |

### *Complete information*

If all the information were available, for each patient, we would know the response at every dose level. In other words, the highest tolerated dose level would be known. For instance, instead of the above

table, we could imagine a table for a subject for whom DLT appears from dose level 3. Conceptually, we might think of the experiment, for each subject $j$, being carried out on six clones of this same subject. Complete information is then as shown in the table below. Of course in a real trial, such information is not available. However in the framework of simulations or probabilistic calculations, complete information can be obtained.

| | Dose | | | | | |
|---|---|---|---|---|---|---|
| | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ |
| $Y_{jk}$ | 0 | 0 | 1 | 1 | 1 | 1 |

*Implementation*

Were we to have complete information from an experiment then we could use it to provide an efficient non-parametric estimate of the dose–toxicity relation. Once again these estimates will not be available in practice but they can be evaluated for the purposes of theoretical comparisons. To each subject $j$ we can associate a numerical value $v_j$ indicating his/her toxicity tolerance threshold to the treatment under study. Without loss of generality, the $v_j$ can be generated from a uniform distribution. Using these outcomes and the given values of $R_k$, it is then determined whether or not we see a toxicity at any dose level. On the basis of a sample of size $n$ the probabilities $R_k$ ($k = 1, \ldots, m$) at each dose level can be estimated by $\hat{R}_k = \sum_{j=1}^n Y_{jk}$, the empirical frequencies. We know that $R_1 \leqslant \cdots \leqslant R_k \leqslant \cdots \leqslant R_m$ and it is easily seen that $\hat{R}_1 \leqslant \cdots \leqslant \hat{R}_m$. The estimated MTD is the level whose estimate of the probability of toxicity minimizes $|\hat{R}_k - \theta|$. It can be readily shown that $\hat{R}_k$ is an unbiased estimator of $R_k$ and that $\text{Var}\,\hat{R}_k$ achieves the Cramer–Rao lower bound.

For every $k$ we claim our estimate of $R_k$ to be optimal in the sense of being unbiased and attaining the Cramer–Rao lower bound. Unless we are prepared to make parametric assumptions whereby information at any level provides some information at all levels then we can clearly not improve upon this estimator. This is the sense in which it is also non-parametric. Of course, in any practical setting, the observations are not complete, we are unable to independently replicate experimentation for the same subject at all levels, and so such an optimal method is not available to us. It can nonetheless be a useful tool in theoretical investigations, providing insights into relative as well as absolute performance.

## 3. ILLUSTRATION

As an illustration we compare a two-stage CRM (Møller, 1995; O'Quigley and Shen, 1996), based on initial 3 by 3 escalation, with the optimal method. We carried out 5000 simulations of the two procedures. Table 1 shows the recommendation distribution when the target is 0.2. We denote by $q_k(16)$ the proportion of times that the optimal method recommends level $k$ based on 16 patients and $p_k(16)$ the analogous quantity for the CRM. The information used in the table can also be expressed in terms of a cumulative distribution of errors. The error is simply the distance of the probability of toxicity at the recommended level from the probability of toxicity at the target level. An illustration is given in Figure 1 for the above situation (left panel). The right panel results from simulations over many hundred dose–toxicity relations, and highlights the poorer performance of a two-parameter CRM as compared to a one-parameter CRM (Shen and O'Quigley, 1996), even though the data are generated by this same two-parameter model. This poorer performance has been discussed in the CRM literature and the graph simply adds an angle to this fact.

Table 1. *Frequency of recommending dose level $d_k$*

|          | \multicolumn{6}{c}{$d_k$} | | | | | |
|----------|------|------|------|------|------|------|
|          | 1    | 2    | 3    | 4    | 5    | 6    |
| $R_k$    | 0.05 | 0.11 | 0.22 | 0.35 | 0.45 | 0.60 |
| $p_k(16)$| 0.05 | 0.26 | 0.42 | 0.21 | 0.06 | 0.0  |
| $q_k(16)$| 0.04 | 0.27 | 0.48 | 0.17 | 0.04 | 0.0  |



Fig. 1. Cumulative probability of error for the optimal and CRM designs.

## 4. FURTHER POINTS

We can associate a single measure of efficiency with the graphs in Figure 1. In order to summarize efficiency measures at different sample sizes it is common to consider asymptotic efficiency if it exists. For phase I studies it may be more sensible to average some measure of finite sample efficiency over a range of sample sizes likely to be of interest in practice. Based on experience we suggest

$$e(n) = \frac{\sum_k p_k(n)\, q_k(n)}{\max\left(\sum_k p_k(n)^2, \sum_k q_k(n)^2\right)}$$

to be the efficiency of the reference method for a sample of size $n$. Many other candidates could also be considered, the only requirement being the use of an appropriate distance measure between the distribution $p_k$ and $q_k$. Often the numerical differences between different measures will be slight and, since we are only interested in relative differences, it makes little difference which distance measure we choose. For the results shown in Table 1 we find $e(16) = 0.93$.

High efficiency is one desirable property of a design but not the only one. Objective evaluation is important but, as for any single summary measure, other aspects may be obscured. For instance, we might anticipate higher efficiency if we put patients at higher risk of toxicity. Such efficiency improvements may be deemed too costly and we would most often only want to choose a design on the basis of superior efficiency if the risks to patients in the study were comparable. In addition, the pattern of deviation of the optimal from any reference design is altogether lost in the efficiency index. Such information may be of real interest.

Note also that efficiency is always defined with respect to some particular class of situations. We may want to evaluate competing designs in cases where the target level is most likely to be one of the higher

levels under investigation, or we might like to focus attention on dose–toxicity relationships that are very shallow.

The optimal method will help in evaluation. It is also helpful as a conceptual tool since it indicates just how our data are incomplete. This alone can be instructive to those wishing to carry out dose finding studies with very few subjects and yet hoping to maintain some reasonable degree of accuracy. Designs with stopping rules and different escalation/de-escalation algorithms may cloud the fact that the very best we can ever do is limited by elementary binomial variation.

All of the results in this paper hinge upon the monotonicity assumption: i.e. a patient presenting a DLT at some level would necessarily present a DLT at all higher levels. Other plausible models might be considered; for example, a stochastic one in which the probability of toxicity for any given individual increases with dose but not as a (0, 1) step function. Our findings would not directly apply here to this case.

The example shows the CRM in a favourable light. Other situations may be less favourable to the CRM, although this is not our experience. Our purpose here is not to study the efficiency of the CRM, already shown to be fully efficient for large samples (Shen and O'Quigley, 1996). The example is illustrative. More extensive examples, using much broader families of dose–toxicity curves, and many competing designs can be found in Paoletti (2001).

## REFERENCES

AHN, C. (1998). An evaluation of phase I cancer clinical trial designs. *Statistics in Medicine* **17**, 1537–1549.

ANBAR, D. (1977). The application of stochastic methods to the bioassay problem. *Journal of Statistical Planning and Inference* **1**, 191–206.

ANBAR, D. (1978). A stochastic Newton–Raphson method. *Journal of Statistical Planning and Inference* **2**, 153–163.

ANBAR, D. (1984). Stochastic approximation methods and their use in bioassay and phase I clinical trials. *Communications in Statistics* **13**, 2451–2467.

BABB, J., ROGATKO, A. AND ZACKS, S. (1998). Cancer phase I clinical trials: efficient dose escalation with overdose control. *Statistics in Medicine* **17**, 1103–1120.

FARIES, D. (1994). Practical modifications of the CRM for phase I cancer study. *Journal of Biopharmaceutical Statistics* **4**, 147–164.

GATSONIS, C. AND GREENHOUSE, J. B. (1992). Bayesian methods for phase I clinical trials. *Statistics in Medicine* **11**, 1377–1389.

GOODMAN, S. N., ZAHURAK, M. L. AND PIANTADOSI, S. (1995). Some practical improvements in the continual reassessment method for phase I studies. *Statistics in Medicine* **14**, 1149–1161.

GOOLEY, T. A., MARTIN, P. J., FISHER, L. D. AND PETTINGER, M. (1994). Simulation as a design tool for phase I/II clinical trials: an example from bone marrow transplantation. *Controlled Clinical Trials* **15**, 450–462.

MØLLER, S. (1995). An extension of the CRM using a preliminary up and down design in a dose finding study in cancer patients, in order to investigate a greater range of doses. *Statistics in Medicine* **14**, 991–1022.

O'QUIGLEY, J., PEPE, M. AND FISHER, L. (1990). Continual reassessment method: a practical design for phase I clinical studies in cancer. *Biometrics* **46**, 33–48.

O'QUIGLEY, J. AND SHEN, Z. L. (1996). Continual reassessment method: a likelihood approach. *Biometrics* **52**, 163–174.

O'QUIGLEY, J. AND REINER, E. (1998). A stopping rule for the continual reassessment method. *Biometrika* **85**, 741–748.

O'QUIGLEY, J. (1992). Estimating the probability of toxicity at the recommended dose following a phase I clinical trial in cancer. *Biometrics* **48**, 853–862.

O'QUIGLEY, J., SHEN, L. AND GAMST, A. (1999). Two sample continual reassessment method. *Journal of Biopharmaceutical Statistics* **9**, 17–44.

PAOLETTI, X. (2001). Comparative evaluation of phase I trial designs, Ph.D. Thesis, University of Paris VII Jussieu.

PIANTADOSI, S. AND LIU, G. (1996). Improved designs for dose escalation studies using pharmacokinetics measurements. *Statistics in Medicine* **15**, 1605–1618.

SHEN, L. Z. AND O'QUIGLEY, J. (1996). Consistency of continual reassessment method in dose finding studies. *Biometrika* **83**, 395–406.

STORER, B. E. (1989). Design and analysis of phase I clinical trials. *Biometrics* **45**, 925–937.

STORER, B. E. (1993). Small-sample confidence sets for the MTD in a phase I clinical trial. *Biometrics* **49**, 1117–1125.

STORER, B. E. (1998). Phase I clinical trials, *Encylopedia of Biostatistics*. New York: Wiley.

WHITEHEAD, J. (1997). Bayesian decision procedures with application to dose-finding studies. *International Journal of Pharmaceutical Medicine* **11**, 201–208.

WHITEHEAD, J. AND WILLIAMSON, D. (1998). Bayesian decision procedures based on logistic regression models for dose-finding studies. *Journal of Biopharmaceutical Statistics* **8**, 445–467.