# Continual Reassessment Method: A Practical Design for Phase 1 Clinical Trials in Cancer

John O'Quigley,* Margaret Pepe, and Lloyd Fisher

Fred Hutchinson Cancer Research Center and Department of Biostatistics,
University of Washington, Seattle, Washington 98104, U.S.A.

## SUMMARY

This paper looks at a new approach to the design and analysis of Phase 1 clinical trials in cancer. The basic idea and motivation behind the approach stem from an attempt to reconcile the needs of dose-finding experimentation with the ethical demands of established medical practice. It is argued that for these trials the particular shape of the dose toxicity curve is of little interest. Attention focuses rather on identifying a dose with a given targeted toxicity level and on concentrating experimentation at that which all current available evidence indicates to be the best estimate of this level. Such an approach not only makes an explicit attempt to meet ethical requirements but also enables the use of models whose only requirements are that locally (i.e., around the dose corresponding to the targeted toxicity level) they reasonably well approximate the true probability of toxic response. Although a large number of models could be contemplated, we look at a particularly simple one. Extensive simulations show the model to have real promise.

## 1. Introduction

This work is motivated by a particular type of trial in cancer patients. Consider a trial where (1) the patients are at very high risk of death in the short term under all standard therapies (some of which may have already failed); (2) the proposed new therapy at high doses will have severe, usually fatal toxicity; (3) at low doses little or no efficacy is expected from the drug; and (4) relatively little is known about the appropriate dose range for possible efficacy with tolerable toxicity. Current examples at the Fred Hutchinson Cancer Research Center include a potential new therapy based on radiolabeled tumor-specific monoclonal antibodies as well as the more traditional trials of new combination chemotherapies.

In new drug development, Phase 1 studies are the studies where a drug is initially given to humans. These studies with more benign drugs often use healthy volunteers. Drugs are then initiated at very low doses and slowly escalated to show safety at a level where some biologic activity takes place (an example would be an antihypertensive drug). Later when the pharmacologic and safety information is available the drug is introduced to the patient population, again with an emphasis on safety. Observations can then be made on efficacy (Phase 2 studies). This classic paradigm does not fit the motivating examples where, if the treatments do not work, the patient will likely die before any other therapy can be attempted. Thus, because of the large potential gain to the patients, a slow dose escalation, as in more classical designs, may put the patient at increased risk. Indeed, to give the patient an optimal chance of a favorable response, one may be willing to accept a nonnegligible probability of severe toxic reaction. This occurs because, generally, and, in particular, with

the therapies considered here, there is more prospect for cure, or prolongation of life, as doses associated with a new therapy are increased. Unfortunately, in parallel with such improving prospects goes an increasing risk of toxic reaction.

In evaluating the goals of such a Phase 1 study we have to consider that in many circumstances the benefits of a new therapy may not be known for months, perhaps years, after enrollment. There is therefore no immediate possibility of assessing the relative advantages of the treatment against its toxic disadvantages where, in these examples, such toxicities manifest themselves in terms of days, weeks, or months rather than years. In this context of life-threatening illness then, given the hoped-for benefit, we decided to aim for an "acceptable" level of toxic response. This level could be set fairly high when toxic response corresponds to transitory symptoms, and obviously lower when such response, as is frequent in Phase 1 cancer trials, corresponds to death.

There are some rather stringent constraints on the conduct of a Phase 1 trial, which need to be considered in its design. Not only are we concerned with efficient evaluation of a level for use in Phase 2 or Phase 3 trials, but also, those patients in the Phase 1 trial must be treated in accordance with standard ethical considerations. These dictate that no patient be given a treatment which the attending physician knows to be inferior. Given the high mortality without new treatment, potential gain may be better served by knowingly accepting the dose best estimated to give some targeted toxicity level. On the other hand, neither can the patient be administered a dose with a believed unacceptably high probability of toxic response, as seen in many animal studies. Subjects enter the trial in a sequential fashion, as they become available. Once again, for ethical reasons as well as those of statistical efficiency, information already available from the toxicity responses of previously entered patients ought to be used in determining the dose assigned to a patient. In general, a continuum of doses is not available. For practical reasons, we usually have a fixed number of doses (about six in our experience) from which to choose. Finally, the number of patients available to be entered in the Phase 1 trials considered here is usually rather small. In order that the trials be completed in a timely manner, coupled with the understandable desire of the physicians involved to proceed as rapidly as possible to a Phase 2 trial, we have found that a maximum of about 25 patients can be accrued. This figure might be increased if, after 25 patients, we had not arrived at a sufficiently accurate estimate of an acceptable dose level, and sufficient interest in the trial remained to warrant continuation.  .

Apart from pharmacologically based studies using animal data (Collins et al., 1986), a number of somewhat ad hoc protocols for Phase 1 trials are currently in common use. For one widely used variant, described by Storer (1989), a fixed number of patients (very often three) is treated at an initial dose level, the lowest possible at the beginning of the experiment. If no severe toxicity is observed then the dose is escalated to the next highest level. Otherwise, an additional three are treated at the same dose. If fewer than two toxicities are observed amongst the six, then the dose is escalated to the next highest level. Otherwise the trial is terminated. The dose in use at trial termination is recommended as the maximum tolerated dose (MTD) for use in a Phase 2 trial. Some ad hoc modification of the dose is recommended if there are many toxic responses in the last group of six patients.

Such sequential designs, whilst appearing to address the essential requirements of a Phase 1 trial, albeit with an initial dose level generally far below the MTD, can have very poor operating characteristics. Not only do we expect patients entered early in the trial to be treated suboptimally, if the drug is as effective as hoped, but the apparent caution of such designs often masks real anticonservative properties. We come back to this point in Section 6. From a statistical viewpoint, perhaps the major criticism of such protocols, however, is that in general the recommended dose has no interpretation as an estimate of the dose which yields a specified rate of severe toxicity. If the trial is such that there are many dose

levels below the MTD then the standard protocol will choose a dose too low with greater probability than if there were fewer dose levels below the MTD. Indeed the nature of the experiment does not yield an estimate of the probability of toxicity at the recommended dose level, though this might be regarded as one minimal objective of a Phase 1 trial. The poor performance of the standard protocol in practice, together with its inability to comply with what, in our view and that of our medical colleagues, are compelling ethical demands, has prompted us to consider an alternative design.

We call this the continual reassessment method (CRM) since we are continually updating our ideas and treating at that level which current available evidence indicates to be the target level. This method, described in the next section, is essentially Bayesian, the Bayesian setting being particularly well adapted to decision making problems. The key aspect of a Phase 1 trial is one of decision although, subsequent to the trial, we may wish to investigate inferential questions. Given the Bayesian framework, in which we set the problem, Bayesian inference is straightforward. It turns out though that classical inference, based on maximum likelihood theory, is also possible and we outline how this can be done in Section 6. In Section 3 some possible modifications of the basic method are considered. Section 4 presents an illustration of how the method would work in practice and Section 5 looks at large-scale simulations under a variety of circumstances.

## 2. Continual Reassessment Method

### 2.1 *Motivation*

We can reasonably assume that the probability of toxic response increases monotonically with increasing dose. Other assumptions, such as the existence of some lower dose at which the probability of toxic response is well approximated by zero, or some idea as to an appropriate family of parametric shapes for the dose response curve, may help guide us in the choice of a suitable model. However, such considerations are not essential to the methods we develop here and can, unless the evidence is strong to do otherwise, be left aside. In most trials, and in particular in cancer chemotherapy trials, it is rare that we have no idea at all about the dose–response relationship. Such information is implicitly used in the choice of dose levels available for use in the trial, and part of our approach here is to make an attempt to quantify such information.

The procedure we propose is to update our notion of the dose–response relationship as observations on severe toxicity become available. In addition, patients are always treated at the dose whose response probability, according to our current knowledge, is closest to the desired level. This achieves two things. First, we meet, as best we can, a compelling ethical criterion. Second, by concentrating experimentation around that dose corresponding to the anticipated target toxicity level, we would expect reasonable estimates from any relatively flexible model. This would be the case even if, in terms of dose–response, it performs poorly away from the target level. Thus, a one-parameter model, at least for the sample sizes considered here, will perform as well as a two-parameter model for instance, as long as interest focuses only on local estimation of dose–response. We come back to this point in the discussion of Section 6.

### 2.2 *The Method*

Suppose that, from the dose range $x$, the dose levels $x_i$ ($i = 1, \ldots, k$) are chosen for experimentation. Let $Y_j$ be a binary random variable (0, 1), where 1 denotes severe toxic response for the $j$th patient ($j = 1, \ldots, n$) entered into the trial. Suppose the probability of toxic response at $x^*$ (not necessarily one of those levels selected) is equal to $\theta$, where $\theta$ is

the probability of response corresponding to the aimed-for target level. Consider some simple dose–response function for $E(Y_j)$ and denote this by $\psi(x_i, a)$. All we assume is that this function is monotonic in $x_i$ and $a$ and that for some $a$, say $a_0$, from the set $\mathscr{A}$ of possible values of $a$, we have $\psi(x^*, a_0) = \theta$. Thus $a_0$ will be the true state of nature and $x^*$ the dose giving the target toxicity level. We want a model rich enough so that for any dose, say $\tilde{x}$, and probability of response $\tilde{\theta}$, there is a parameter, say $\tilde{a}$, such that $\psi(\tilde{x}, \tilde{a}) = \tilde{\theta}$. We will also require that $\tilde{a}$ be unique. Thus our one-parameter model is sufficiently flexible to reproduce the probability of toxic response at the target level, provided such a level is included in our range. In fact, if our range does not include the target level this will be reflected in our estimates of $a_0$.

Let $\Omega_j = \{y_1, \ldots, y_{j-1}\}$ and let $f(a, \Omega_j)$ be a nonnegative function summarizing all available information about the parameter $a_0$. This is our current prior before experimentation on the $j$th subject. In the sequel we will take $\mathscr{A}$ to be $(0, \infty)$ so that

$$\int_0^\infty f(a, \Omega_j) \, da = 1 \quad (j = 1, \ldots, n).$$

Defined in this way, and given the response of the $j$th patient, we can exploit a formulation of Bayes' theorem to obtain $f(a, \Omega_{j+1})$ from $f(a, \Omega_j)$, thus updating our information about the parameter $a_0$ as observations become available. Next, we will need some estimate of the probability of toxic response at level $i$ given the accumulated information on the first $j - 1$ patient responses. Denote this as $\theta_{ij}$ where

$$\theta_{ij} = \int_0^\infty \psi(x_i, a) f(a, \Omega_j) \, da \quad (i = 1, \ldots, k). \tag{2.1}$$

Rather than work with the expected values of the probabilities over $\mathscr{A}$, we could work directly with the expected value of $a$ over $\mathscr{A}$. This leads to the alternative estimate

$$\theta'_{ij} = \psi\{x_i, \mu(j)\} \quad (i = 1, \ldots, k), \qquad \mu(j) = \int_0^\infty a f(a, \Omega_j) \, da.$$

Finally, let $\Delta(v, w)$ denote some measure of distance between $v$ and $w$, for example, $\Delta(v, w) = (v - w)^2$. Then for the $j$th entered patient in the trial, choose dose level $x_i$ such that $\Delta(\theta_{ij}, \theta)$, $\Delta(\theta'_{ij}, \theta)$, or $\Delta(x_i, \psi^{-1}_{a=\mu(j)}(\theta))$, depending on which criterion we choose to work with, is a minimum. The latter two of these criteria have the advantage of being computationally economic, reducing the need to perform $k$ infinite integrals to that of performing a single integral.

Having chosen the treatment level for the $j$th patient, denoting it by $x(j)$, and having observed whether the $j$th patient experiences a toxic response, we can now evaluate the function $f(a, \Omega_{j+1})$ which, in the light of the most recent observations, updates our knowledge about $a_0$. Suppose that

$$\phi(x(j), y_j, a) = \psi^{y_j}(x(j), a)\{1 - \psi(x(j), a)\}^{(1-y_j)}$$

and that

$$g(a) = f(a, \Omega_1).$$

Thus, $g(a)$ is our prior distribution for $a_0$ and reflects whatever knowledge we have of the dose toxicity relationship before experimentation begins. We examine the question of

determining this function in the next section. We then have

$$f(a, \Omega_{j+1}) = \frac{\phi(x(j), y_j, a)f(a, \Omega_j)}{\int_0^\infty \phi(x(j), y_j, u)f(u, \Omega_j)\ du} \tag{2.2}$$

$$= \frac{g(a)\ \prod_{l=1}^{j} \phi\{x(l), y_l, a\}}{\int_0^\infty g(u)\ \prod_{l=1}^{j} \phi\{x(l), y_l, u\}\ du}. \tag{2.3}$$

Now use equation (2.1) and calculate $\theta_{i,j+1}$ or, alternatively, calculate $\mu(j + 1)$ and the approximation $\theta'_{i,j+1}$ and treat the next patient at dose level $x_i$ such that $\Delta(\theta_{i,j+1}, \theta)$, or one of the other criteria mentioned above, is minimized. Continue in this way until the results of the last patient entered are available. The recommended dose will be that dose $x_i$ ($i = 1$, ..., $k$) such that $\Delta(\theta_{i,n+1}, \theta)$, $\Delta(\theta'_{i,n+1}, \theta)$, or $\Delta(x_i, \psi^{-1}_{a=\mu(n+1)}(\theta))$ is minimized, depending on which criterion we choose to work with.

The function $g(a)$ should reflect in part our prior information as to the nature of the dose at which it is believed the probability of toxic response is $\theta$. If experimentation is to be begun at level $x_l$ then we should choose $g(a)$ such that $\theta_{l1} = \theta$. In practice it may be simpler to consider $g(a)$ such that $\theta'_{l1} = \theta$, even if subsequently we choose to work with $\theta_{ij}$ rather than $\theta'_{ij}$. A second consideration in selecting $g(a)$ is the amount of uncertainty in our prior feeling that $\theta_{l1}$ (or $\theta'_{l1}$) $= \theta$. Discussions with the physicians and pharmacologists involved in the study should lead to identifying $\theta^L$ and $\theta^U$, the lower and upper bounds that $\theta_{l1}$ (or $\theta'_{l1}$) may plausibly take. These bounds will reflect the investigators' prior beliefs as to just how toxic or nontoxic the new treatment may turn out to be, these beliefs being transformed into bounds for $\theta_{l1}$. Such bounds have no frequentist interpretation and it seems helpful to describe them such that those setting up the study are, in some sense, $100(1 - \alpha)\%$ confident that $\theta_{l1}$ lies between $\theta^L$ and $\theta^U$. Some practical discussion on establishing priors and how to view them is given by Martz and Waller (1982). Jaynes (1986) discusses some of the deeper conceptual issues raised by the Bayesian approach.

Given the above formulation, we will be able to find constants $a_1$ and $a_2$ such that $\psi(x_l, a_1) = \theta^L$ and $\psi(x_l, a_2) = \theta^U$, and we choose $g(\cdot)$ to be a density having $\alpha/2$ and $1 - \alpha/2$ quantiles equal to $a_1$ and $a_2$, respectively. In general then, three constraints are imposed on the function $g(a)$ and to meet these would require a three-parameter function. However, having found $\theta^L$ and $\theta^U$, we can adjust these in a conservative fashion such that a two-parameter density will suffice. In the examples of Sections 4 and 5 the choice $g(a) = \exp(-a)$ corresponds to a simple vague prior and seems to work well. Our feeling is that for many cases this choice is quite adequate. Otherwise, an outline of how to go about a more general specification is provided in the next section.

### 2.3 *Fitting Parameters in the Function* $g(a)$

Given that $\mathscr{A} = (0, \infty)$, we suggest

$$g(a) = \lambda^c a^{c-1} \exp\{-(\lambda a)\}/\Gamma(c), \quad \Gamma(c) = \int_0^\infty \exp(-u)u^{c-1}\ du,$$

the gamma density with scale parameter $\lambda$ and shape parameter $c$. Other choices for $g(a)$, such as a truncated normal or log-normal, are, of course, possible. The gamma is suggested here because of its simplicity and its positivity constraint. Unless we really have strong knowledge about the probabilities of toxicity before experimentation is undertaken, we

would not like our methods to depend very strongly on the choice of $g(a)$, within a family of equally plausible alternatives.

Martz and Waller (1982) detail the necessary steps in fitting a gamma prior on the basis of the upper and lower points of our prior confidence region. Here, the problem is slightly more involved since, not only do we wish to choose a prior incorporating these upper and lower points, but also we would like the mean (possibly median) to correspond to our targeted toxicity rate at the starting dose. Note that since $\psi$ depends monotonically on $a$ as well as on $x$ and, since we are free to use whatever units we wish for $x$, we can, without loss of generality, fix $\mu(1) = 1$. Thus $\lambda = c$, hence $\lambda$ can be found by solving

$$\Gamma(\lambda)^{-1} \int_{a_2}^{\infty} \lambda^{\lambda} a^{\lambda-1} \exp\{-(\lambda a)\} \, da = \Gamma(\lambda)^{-1} \int_{0}^{a_1} \lambda^{\lambda} a^{\lambda-1} \exp\{-(\lambda a)\} \, da = \alpha/2.$$

No exact solution is likely to exist. However, suppose that for the first expression, $\lambda$ lies between the integers $m_2$ and $m_2 + 1$; then, exploiting the fact that the gamma density reduces to the Erlang for $c$ an integer (Hastings and Peacock, 1975), and defining

$$w(m_2, a_2) = \exp(-m_2 a_2) \left\{ \sum_{i=0}^{m_2-1} (m_2 a_2)^i / i! \right\}, \quad w(0, \cdot) = 1,$$

we find that $w(m_2 + 1, a_2) < \alpha/2 < w(m_2, a_2)$ and a good approximate solution to the first integral expression is given by

$$\lambda_2 = m_2 + \{w(m_2, a_2) - \alpha/2\}/\{w(m_2, a_2) - w(m_2 + 1, a_2)\}.$$

It is fairly easy to locate the particular value of $m_2$ that satisfies $w(m_2 + 1, a_2) < \alpha/2 < w(m_2, a_2)$. Next we need to find the value of $m_1$ such that $1 - w(m_1, a_1) < \alpha/2 < 1 - w(m_1 + 1, a_1)$ and then we use as a good approximate solution to the second integral expression,

$$\lambda_1 = m_1 + \{w(m_1, a_1) + \alpha/2 - 1\}/\{w(m_1, a_1) - w(m_1 + 1, a_1)\}.$$

We would not anticipate $\lambda_1$ and $\lambda_2$ to give the same values. Taking $\lambda = \min(\lambda_1, \lambda_2)$ gives us that value which corresponds to the greater dispersion and consequently the more conservative prior.

### 2.4 *Establishing Dose Levels*

Practical considerations, such as preparation and packaging, will generally result in the number of available dose levels, $k$, for testing, being small. In our experience this has generally been six although there is no reason not to try to include more levels if feasible. It makes sense to start experimentation at the $l$th level, where $l$ is the largest integer smaller than $k/2$. The preparation at this dose should correspond to the experimenter's best prior estimate at that amount which will result in severe toxic response in $\theta\%$ of patients.

Having calculated $a_1$ and $a_2$ as outlined in the previous section, and considering these as corresponding in some sense to worst and best cases, we would like our dose range to include somewhere a level at which the probability of toxic response is close to $\theta$. For instance, should $a_0$ turn out in truth to be close to $a_2$, and the treatment levels in consequence far more toxic than initially estimated, we would still hope that the lowest dose available would have an associated probability of toxic response not much greater than $\theta$. $\theta_1$ should then be chosen such that those setting up the experiment are 95% confident that an upper limit for $\theta_1$ will still be less than $\theta$. Considerations of an appropriate lower limit would then define $\theta_k$, the remaining $\theta_i$ being defined by interpolation and practical considerations of drug packaging. The investigators choose doses which they

believe, a priori, correspond to the desirable range of toxicity probabilities for testing, i.e., which correspond to the probabilities $\theta_1, \ldots, \theta_k$.

The next point is that our simple model for dose–response should, in order to mirror the investigators' prior assumptions, provide an exact fit over our a priori estimates of toxicity at each dose level. To do this we will need to transform the original units in which the dose levels are expressed, to new ones defined by the equations

$$\theta_i = \int_0^\infty \psi(x_i, a) g(a) \, da \quad (i = 1, \ldots, k),$$

where $\theta_1, \ldots, \theta_k$ ($\theta_l = \theta$) are our a priori estimates of the probabilities of toxicity at levels $x_1, \ldots, x_k$. Rather than solve these $k$ integral equations to establish the $x_i$, an easier computational approach would be to transform via

$$x_i = \psi_{a=1}^{-1}(\theta_i)$$

so that our exact fit is now assured with respect to expectation in $\mathscr{A}$-space rather than the $(0, 1)$-space associated with our a priori probabilities.

## 3. Delayed Response

Often we will wish to include a new patient in the study before we have observed the response to treatment of the previously entered patient. In general, there may be a number of new patients to include and responses themselves may be grouped.

Suppose we have observed responses on $j - 1$ patients and since this time an additional $m^*$ patients have been recruited. In terms of allocation, one way to proceed would be to calculate the appropriate experimentation level, $x_i$, for the $j$th patient, and then consider that, since we have no additional information, it makes sense to treat the remaining $m^* - 1$ patients at this same level. From this set of $m^*$ patients at some future time point, we will be able to observe the presence or absence of responses on some subset, say $m$ ($\leq m^*$) patients. Then calculate $\theta_{ij}$ or $\theta'_{ij}$ as in Section 2.2, and defining

$$\phi_m(x_i, y_j, \ldots, y_{j+m-1}; a) = \psi^{\bar{y}_m}(x_i, a)\{1 - \psi(x_i, a)\}^{1-\bar{y}_m},$$

where $\bar{y}_m = \sum y_l$, summation running from $j$ to $j + m - 1$, the right-hand side of equation (2.2) can now be replaced by

$$\frac{\phi_m(x_i, y_j, \ldots, y_{j+m-1}; a) f(a, \Omega_j)}{\int_0^\infty \phi_m(x_i, y_j, \ldots, y_{j+m-1}; u) f(u, \Omega_j) \, du}.$$

The current uncertainty, after $j - 1$ patients have been treated, in the estimated target level may be such that it can be ethically justified to simultaneously experiment at levels, say, just above and below the currently estimated target. A wider spread could be envisaged although, since in practice $m^*$ is not likely to greatly exceed three or four, this is probably unnecessary.

In order to justify experimentation at the higher or lower level though we need reasonable evidence that either of these levels may have a probability of toxic response close to $\theta$. The simplest way to do this would be to calculate confidence intervals for the probability of response at the various dose levels and to use these to subjectively choose the next level for experimentation. Alternatively, closeness could be defined relative to the level for which the point estimate is nearest to $\theta$. To do this we suggest defining a range $R$, comprising $\theta$ and values near enough to $\theta$ to be considered acceptable, and to calculate $q_R(h)$ where

$$q_R(h) = \frac{\int_{A_R(x_{i+h})} \psi(x_{i+h}, a) f(a, \Omega_j) \, da}{\int_{A_R(x_i)} \psi(x_i, a) f(a, \Omega_j) \, da}, \quad A_R(x_i) = \{\psi_{x_i}^{-1}(\theta^*): \theta^* \in R\}, \quad h = -1, 1.$$

If $q_R(h) > q$ for some fixed $q$, then we might feel justified in experimenting at level $x_{i+h}$. A value of .5 for $q$ seems reasonable, although of course this is an arbitrary figure and would require agreement on the part of the investigators participating in the trial.

## 4. Examples

Two simulated studies are used to illustrate the method. The outcome patterns in these two studies are, in our experience, fairly typical and help indicate how we might anticipate the scheme to work in practice. A more detailed examination of this, along with estimated operating characteristics, is presented in the next section.

Table 1 shows the results of a simulation in which six levels $(x_1, \ldots, x_6)$ are chosen for experimentation, $x_3$ being, on the basis of available knowledge, that level which would produce the target response of 20% (i.e., $\theta = .2$). First we need some simple one-parameter dose–response model and one plausible candidate would be

$$\psi(x_i, a) = \{(\tanh x_i + 1)/2\}^a.$$

Figure 1 shows what this function looks like for various values of $a$ and clearly it meets the requirements of Section 2.2. Experimentation will begin at $x_3$, subsequent allocation being

**Table 1**
*Sequential trial of 25 patients*; $\Pr(Y = 1) = \{(\tanh x_i + 1)/2\}^{1/2}$

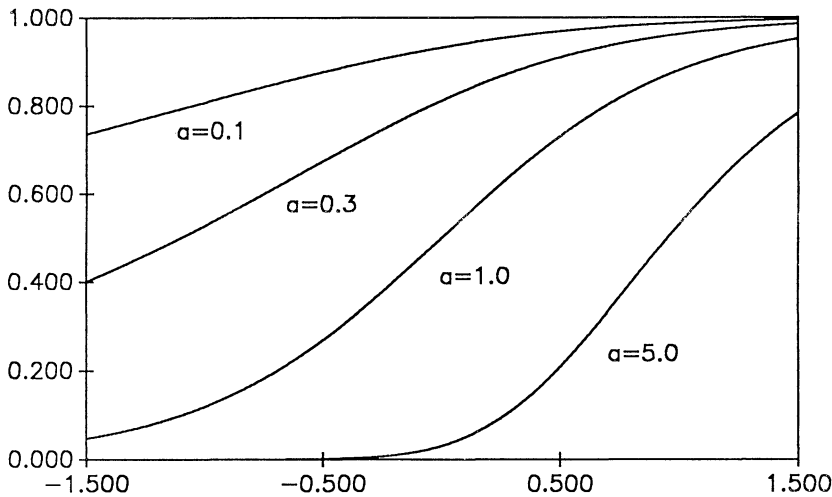| $j$ | $\mu(j)$ | $x_i$ | $y_j$ | $j$ | $\mu(j)$ | $x_i$ | $y_j$ |
|---|---|---|---|---|---|---|---|
| 1 | 1.00 | $x_3$ | 0 | 14 | .64 | $x_2$ | 0 |
| 2 | 1.38 | $x_4$ | 0 | 15 | .67 | $x_2$ | 0 |
| 3 | 1.68 | $x_4$ | 1 | 16 | .69 | $x_2$ | 0 |
| 4 | .92 | $x_3$ | 0 | 17 | .72 | $x_2$ | 0 |
| 5 | 1.07 | $x_3$ | 1 | 18 | .74 | $x_2$ | 0 |
| 6 | .71 | $x_2$ | 1 | 19 | .76 | $x_2$ | 1 |
| 7 | .49 | $x_1$ | 0 | 20 | .67 | $x_2$ | 0 |
| 8 | .55 | $x_1$ | 0 | 21 | .69 | $x_2$ | 0 |
| 9 | .60 | $x_1$ | 0 | 22 | .71 | $x_2$ | 1 |
| 10 | .65 | $x_2$ | 0 | 23 | .64 | $x_2$ | 0 |
| 11 | .69 | $x_2$ | 0 | 24 | .65 | $x_2$ | 1 |
| 12 | .73 | $x_2$ | 0 | 25 | .61 | $x_1$ | 1 |
| 13 | .77 | $x_2$ | 1 | (25 + 1) | .56 | $x_1$ | |



**Figure 1.** $y = \{(\tanh x + 1)/2\}^a$ for some values of $a$.

based, for computational ease, on the third of the criteria given in Section 2.2. We rescale such that

$$(\tanh x_3 + 1)/2 = \theta = .2,$$

and we can then assume that $g(a)$ satisfies $\int_0^\infty ug(u) \, du = 1$. The considerations of Sections 2.3 and 2.4 suggest that $g(a) = \exp(-a)$ would suffice if our prior notions of the probability of toxic response at $x_3$ are fairly vague, i.e., the true probability could, with, in some sense, 95% confidence, be between .003 and .96. Rather that stringently apply Section 2.4 we will use it as a guide, establishing dose levels on the basis of the probabilities of toxic response under our assumed model and our assumed uncertainty. We obtain the following probabilities of toxicity for $a_0 = 1$:

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| $x_i$ | −1.47 | −1.1 | −.69 | −.42 | 0.0 | .42 |
| $(\tanh x_i + 1)/2$ | .05 | .1 | .2 | .3 | .5 | .7 |

In reality let us suppose the data are generated by $\Pr(Y = 1) = \{(\tanh x_i + 1)/2\}^{1/2}$ so that we underestimate the treatment's toxic potential which is truly as below for $a_0 = .5$:

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| $x_i$ | −1.47 | −1.1 | −.69 | −.42 | 0.0 | .42 |
| $\{(\tanh x_i + 1)/2\}^{1/2}$ | .22 | .32 | .45 | .54 | .69 | .80 |

Experimentation commences at $x_3$, where we have a probability of toxic response more than twice that we believed to be the case. Despite this, no toxic response is observed for the first entered patient. Applying the methods of Section 2.2 indicates that we should try the second patient at $x_4$ and once again no response is observed. Even so the method indicates that this level should be maintained for the third patient, and this time a toxicity is observed. Level $x_3$ is used for the following two patients; the one toxicity leads to experimentation for the seventh entered patient at the lowest level $x_1$. For the remainder of the experiment only levels $x_1$ and $x_2$ are used. Indeed 15 consecutive patients are tried at $x_2$ and for the first 14 of these a total of three toxic responses are observed. Since $\frac{3}{14} = .21$ it is not surprising that the method hesitates before returning to the lowest level. The observed toxicity seen for the last of these 15 consecutive patients suggests that the final patient entered into the study be treated at level $x_1$. A toxic response is seen here leading to $\mu(25 + 1) = .55$, not far removed from the population value of .5. This leads to the estimate (slightly biased) of the mean response probability at $x_1$ of .19 and this is the level we would recommend for use in future experimentation.

A number of other points are highlighted by a study of Table 1. Here the value of $\mu(j)$ always increases if no toxicity is observed and always decreases when we see a toxicity. However, the amount by which it increases or decreases will depend on the current accumulated information. For instance, many nontoxicities are seen at $x_2$. The first toxicity at $x_2$ occurs before any of these and produces a change of $-.23$ in $\mu(j)$. The second toxicity at $x_2$ produces a change of $-.13$ in $\mu(j)$ and the third a change of only $-.09$ in $\mu(j)$.

Table 2 presents a simulation where $\Pr(Y = 1) = \{(\tanh x_i + 1)/2\}^2$ so that we are being overconservative. This is seen in the following table of probabilities of toxicity for $a_0 = 2$:

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| $x_i$ | −1.47 | −1.1 | −.69 | −.42 | 0.0 | .42 |
| $\{(\tanh x_i + 1)/2\}^2$ | .002 | .01 | .04 | .09 | .24 | .49 |

The appropriate level to choose is $x_5$ and, as can be seen from Table 2, this is the one selected despite the fact that the last entered patient, treated at this same level, experiences

**Table 2**
*Sequential trial of 25 patients;* $\Pr(Y = 1) = \{(\tanh x_i + 1)/2\}^2$

| $j$ | $\mu(j)$ | $x_i$ | $y_j$ | $j$ | $\mu(j)$ | $x_i$ | $y_j$ |
|---|---|---|---|---|---|---|---|
| 1 | 1.00 | $x_3$ | 0 | 14 | 1.43 | $x_4$ | 0 |
| 2 | 1.38 | $x_4$ | 0 | 15 | 1.49 | $x_4$ | 0 |
| 3 | 1.68 | $x_4$ | 0 | 16 | 1.55 | $x_4$ | 0 |
| 4 | 1.90 | $x_5$ | 0 | 17 | 1.60 | $x_4$ | 0 |
| 5 | 2.19 | $x_5$ | 0 | 18 | 1.64 | $x_4$ | 0 |
| 6 | 2.44 | $x_5$ | 1 | 19 | 1.69 | $x_4$ | 0 |
| 7 | 1.75 | $x_4$ | 0 | 20 | 1.73 | $x_4$ | 0 |
| 8 | 1.86 | $x_5$ | 1 | 21 | 1.77 | $x_5$ | 0 |
| 9 | 1.50 | $x_4$ | 0 | 22 | 1.84 | $x_5$ | 0 |
| 10 | 1.59 | $x_4$ | 1 | 23 | 1.91 | $x_5$ | 0 |
| 11 | 1.22 | $x_4$ | 0 | 24 | 1.97 | $x_5$ | 0 |
| 12 | 1.30 | $x_4$ | 0 | 25 | 2.03 | $x_5$ | 1 |
| 13 | 1.37 | $x_4$ | 0 | (25 + 1) | 1.86 | $x_5$ | |

a toxic response. Only the first patient is treated at $x_3$ and thereafter the method is trying to decide between $x_4$ and $x_5$, finally settling for $x_5$ after 20 patients have been included.

These two examples are not too extreme and give some feeling as to what we might expect in practice. Here the model generating the data is from the same family as the initial working model. It turns out, and this is examined in Section 5, that for misspecified models, CRM behaves in much the same way and will still give good estimates of the probability of toxic response at the finally selected dose, its performance at other doses worsening with increasing degrees of misspecification. Some more extreme cases, both when the model form is correct and when it is incorrect, are among those considered in the next section.

## 5. Simulations: Operating Characteristics and Model Misspecification

In this section we consider a variety of situations. We suppose there are six ordered dose levels, $x_1, \ldots, x_6$, and that the probability of toxic response at each level is generated in one of three ways:

(i) $\Pr(Y = 1) = \{\tanh x_i + 1)/2\}^{a'}$

(ii) a logistic model

(iii) $p(x_i) \le p(x_i'), \quad i < i'$

where $p(x_i)$ is the probability of toxic response at $x_i$. In all cases the prior is $g(a) = \exp(-a)$ and the targeted toxicity level, $\theta$, is equal to .2, as in the previous section.

Each of the following tables is based on estimates from 200 simulations. The first row gives the level, the second the frequency at which experimentation was performed at that level, the third the frequency at which that level was the one finally recommended, and the final row gives the true probabilities generating the data.

In the first entry of Table 3 the model generating the data is the same as that of our a priori estimates and this results in our recommending the correct level 44% of the time, either this or that level immediately above or below 96% of the time. The remainder of this table gives results of the simulations when the form of the model is correct but the initial value assumed for $a_0$ is incorrect. In general the results look very encouraging.

In Table 4 the data are generated by the more classical logistic model. Many interesting observations are to be made. The impression given is that, in many situations, what matters

**Table 3**
*Probability of toxic response generated by* $\Pr(Y = 1) = \{(\tanh x_i + 1)/2\}^a$

| i | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| Observed toxicities = 22% | | | | | | |
| % experimentation at $x_i$ | 5.34 | 17.58 | 36.56 | 33.82 | 6.70 | .00 |
| % recommendation for $x_i$ | 1.00 | 19.50 | 44.50 | 32.00 | 3.00 | .00 |
| **Probability of toxicity at $x_i$** | .05 | .10 | .20 | .30 | .50 | .70 |
| Observed toxicities = 23% | | | | | | |
| % experimentation at $x_i$ | 18.58 | 29.26 | 31.32 | 18.68 | 2.14 | .02 |
| % recommendation for $x_i$ | 14.00 | 39.00 | 36.50 | 10.50 | .00 | .00 |
| **Probability of toxicity at $x_i$** | .09 | .16 | .27 | .38 | .57 | .75 |
| Observed toxicities = 34% | | | | | | |
| % experimentation at $x_i$ | 73.88 | 13.56 | 9.04 | 3.28 | .24 | .00 |
| % recommendation for $x_i$ | 92.50 | 6.50 | 1.00 | .00 | .00 | .00 |
| **Probability of toxicity at $x_i$** | .30 | .40 | .52 | .61 | .76 | .87 |
| Observed toxicities = 17% | | | | | | |
| % experimentation at $x_i$ | .50 | 1.02 | 9.06 | 35.08 | 52.46 | 1.88 |
| % recommendation for $x_i$ | .00 | .00 | 1.00 | 30.00 | 66.50 | 2.50 |
| **Probability of toxicity at $x_i$** | .00 | .00 | .04 | .09 | .25 | .49 |
| Observed toxicities = 18% | | | | | | |
| % experimentation at $x_i$ | 1.24 | 4.32 | 20.98 | 50.30 | 22.92 | .24 |
| % recommendation for $x_i$ | .00 | .00 | 16.00 | 64.00 | 20.00 | .00 |
| **Probability of toxicity at $x_i$** | .01 | .03 | .09 | .16 | .35 | .59 |

are the toxic probabilities and not so much the model generating these probabilities. Indeed we can compare the results for the logistic model generating probabilities (.07, .1, .2, .35, .76, .96) with the first entry of Table 3 when the assumed model is the true one. For the first four levels at least, and this is where about 95% of experimentation is performed in either case, the generated probabilities are almost the same. The results are comparable, although there is an indication we do slightly better for data generated under the logistic model. This is most likely because the difference between the generated toxic probabilities at $x_4$, i.e., between .3 and .35, is sufficient to enable the technique to detect that the toxic probability associated with level 4 is indeed slightly too high.

Another important point which emerges is that, in the range of situations considered here, almost regardless of the configuration, any level with a toxic response probability much in excess of 30%, provided some lower level exists, is likely to be recommended with small probability. Conversely, any level with a toxic response probability below about 8%, provided some greater level exists, is also likely to be recommended with very small probability.

In Table 5, we consider how the model performs under some general situations. In the first of these, one case looks at what happens when all toxic probabilities are extremely high (90%) apart from the very lowest one which is at the appropriate level. Not surprisingly, this level is recommended over 95% of the time but what is rather more encouraging is that only about 11% of the patients are ever tried at a level other than the lowest. One situation in which the method seems slightly disappointing is the last one considered. Here the uppermost level is close to the target one and the level before this is at 11%. The tendency would seem to be to select the more conservative level rather more frequently than the appropriate one. Note, though, that if this 11% falls only a slight amount to 6% the position is retrieved, the "correct level" then being chosen about 60% of the time.

**Table 4**
*Probabilities generated according to two-parameter logistic models*

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| Observed toxicities = 15% | | | | | | |
| % experimentation at $x_i$ | 1.66 | 2.28 | 10.70 | 31.20 | 48.84 | 5.32 |
| % recommendation for $x_i$ | .00 | 1.50 | 3.00 | 28.50 | 58.00 | 9.00 |
| **Probability of toxicity at $x_i$** | .05 | .06 | .08 | .11 | .19 | .34 |
| Observed toxicities = 21% | | | | | | |
| % experimentation at $x_i$ | 3.98 | 11.70 | 32.04 | 43.06 | 9.10 | .12 |
| % recommendation for $x_i$ | .50 | 11.50 | 30.50 | 54.00 | 3.50 | .00 |
| **Probability of toxicity at $x_i$** | .06 | .08 | .14 | .23 | .53 | .84 |
| Observed toxicities = 19% | | | | | | |
| % experimentation at $x_i$ | 2.82 | 7.18 | 24.76 | 47.26 | 17.88 | .10 |
| % recommendation for $x_i$ | .00 | 4.00 | 23.50 | 57.00 | 15.50 | .00 |
| **Probability of toxicity at $x_i$** | .06 | .08 | .12 | .18 | .40 | .71 |
| Observed toxicities = 22% | | | | | | |
| % experimentation at $x_i$ | 21.34 | 43.04 | 29.94 | 5.46 | .22 | .00 |
| % recommendation for $x_i$ | 13.50 | 61.00 | 25.00 | .50 | .00 | .00 |
| **Probability of toxicity at $x_i$** | .08 | .14 | .35 | .65 | .96 | 1.00 |
| Observed toxicities = 23% | | | | | | |
| % experimentation at $x_i$ | 8.28 | 26.44 | 42.16 | 21.38 | 1.74 | .00 |
| % recommendation for $x_i$ | 4.50 | 28.00 | 55.00 | 12.50 | .00 | .00 |
| **Probability of toxicity at $x_i$** | .07 | .11 | .23 | .43 | .84 | .98 |
| Observed toxicities = 22% | | | | | | |
| % experimentation at $x_i$ | 7.86 | 21.18 | 43.48 | 25.50 | 1.98 | .00 |
| % recommendation for $x_i$ | 4.50 | 21.00 | 55.00 | 19.50 | .00 | .00 |
| **Probability of toxicity at $x_i$** | .07 | .10 | .20 | .35 | .76 | .96 |
| Observed toxicities = 29% | | | | | | |
| % experimentation at $x_i$ | 35.45 | 14.30 | 18.49 | 19.47 | 11.47 | .82 |
| % recommendation for $x_i$ | 45.00 | 17.00 | 15.00 | 14.00 | 8.00 | 1.00 |
| **Probability of toxicity at $x_i$** | .27 | .28 | .29 | .30 | .32 | .35 |
| Observed toxicities = 21% | | | | | | |
| % experimentation at $x_i$ | 14.60 | 13.40 | 22.40 | 25.48 | 22.42 | 1.70 |
| % recommendation for $x_i$ | 17.00 | 12.00 | 19.00 | 25.00 | 23.50 | 3.50 |
| **Probability of toxicity at $x_i$** | .19 | .19 | .20 | .21 | .22 | .25 |

## 6. Discussion and Further Points

### 6.1 *Models and Sequential Experimentation*

The examples and simulations were rerun, after rescaling, using the even simpler model $(\sin x_i)^a$. Although this would appear somewhat restrictive in view of the probabilities of toxic response being 0 and 1 at $x_i = 0$ and $x_i = \pi/2$, respectively, the model performs identically to the one used in Sections 4 and 5 with respect to the first two distance criteria of Section 2.2 and very similarly with respect to the third criterion. This is because of the tendency of the method to home in fairly quickly on one or two levels over which the approximation is adequate. Intuitively we tend to think in terms of a dose–response model, valid for all values of $x$, and which monotonically increases from zero to one. Here, however, such restrictions, apart from the monotonicity one, are unnecessary and our only requirement is that the model reasonably approximates reality at those doses most used in the experiment.

We chose the number 25 in our simulations since this is a figure which seems to be used quite often in practice. Performance would clearly improve with increasing sample size and, for instance, the figure of 44% "correct allocation" in the first part of Table 3 becomes

**Table 5**
*Unmodelled probabilities satisfying monotonic constraint*

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| Observed toxicities = 28% | | | | | | |
| % experimentation at $x_i$ | 88.52 | 6.96 | 4.09 | .43 | .00 | .00 |
| % recommendation for $x_i$ | 95.50 | 4.50 | .00 | .00 | .00 | .00 |
| **Probability of toxicity at $x_i$** | .20 | .90 | .90 | .90 | .90 | .90 |
| Observed toxicities = 13% | | | | | | |
| % experimentation at $x_i$ | .26 | .56 | 6.22 | 18.12 | 64.48 | 10.36 |
| % recommendation for $x_i$ | .00 | .00 | .00 | 6.50 | 75.00 | 18.50 |
| **Probability of toxicity at $x_i$** | .00 | .00 | .03 | .05 | .10 | .50 |
| Observed toxicities = 13% | | | | | | |
| % experimentation at $x_i$ | .24 | .36 | 5.76 | 17.04 | 59.42 | 17.18 |
| % recommendation for $x_i$ | .00 | .00 | .00 | 4.00 | 61.50 | 34.50 |
| **Probability of toxicity at $x_i$** | .00 | .00 | .03 | .05 | .10 | .30 |
| Observed toxicities = 10% | | | | | | |
| % experimentation at $x_i$ | .14 | .46 | 5.50 | 14.04 | 52.28 | 27.58 |
| % recommendation for $x_i$ | .00 | .00 | .00 | 3.00 | 39.00 | 58.00 |
| **Probability of toxicity at $x_i$** | .00 | .00 | .03 | .05 | .06 | .22 |
| Observed toxicities = 12% | | | | | | |
| % experimentation at $x_i$ | .20 | .56 | 6.82 | 20.24 | 54.22 | 17.96 |
| % recommendation for $x_i$ | .00 | .00 | .00 | 10.00 | 53.00 | 37.00 |
| **Probability of toxicity at $x_i$** | .00 | .00 | .03 | .05 | .11 | .22 |

close to 70% when the number of patients is as many as 50. Indeed there is no difficulty in explicitly aiming for given levels of precision and continuing experimentation until such precision is achieved. Thus a fixed number of patients would not be determined in advance. We would experiment sequentially, solving at each step the equation

$$\int_{a_L}^{a_U} f(a, \Omega_j)\, da = 1 - \gamma$$

for $a_U$ and $a_L$. We can then calculate the corresponding intervals for the response probabilities of each $x_i$, given by $\{\psi(x_i, a_L), \psi(x_i, a_U)\}$. The quantity $1 - \gamma$ is the amount of confidence we deem necessary before terminating the experiment. As stated, this formulation will not uniquely determine $a_U$ and $a_L$ and we will need to invoke some additional requirement such as symmetry or perhaps that $a_L = 0$. All that then remains is for us to decide on a range $\mathscr{R}$, comprising $\theta$ and values near enough to $\theta$ to be considered acceptable, and to continue experimentation until for some $x_i$, $\{\psi(x_i, a_L), \psi(x_i, a_U)\}$ is wholly contained in $\mathscr{R}$.

### 6.2 *Maximum Likelihood Estimation*

At completion of the trial the likelihood, $e^{L(a)}$, of the data is $\Pr(Y_1, \ldots, Y_n, X_1, \ldots, X_n)$. Elementary manipulations show that

$\Pr(Y_1, \ldots, Y_n, X_1, \ldots, X_n)$

$\quad = \Pr(Y_n \mid Y_1, \ldots, Y_{n-1}, X_1, \ldots, X_n)$

$\qquad \times \Pr(X_n \mid Y_1, \ldots, Y_{n-1}, X_1, \ldots, X_{n-1})$

$\qquad \times \Pr(Y_1, \ldots, Y_{n-1}, X_1, \ldots, X_{n-1})$

$\quad = \Pr(Y_1, X_1) \prod_{i=2}^{n} \Pr(Y_i \mid Y_1, \ldots, Y_{i-1}, X_1, \ldots, X_i)\Pr(X_i \mid Y_1, \ldots, Y_{i-1}, X_1, \ldots, X_{i-1}).$

Note that $X_1$ is determined at the start of the experiment and so is not random. Furthermore, under any of the allocation rules considered in this paper, $X_i$, conditional on $(X_1, \ldots, X_{i-1}, Y_1, \ldots, Y_{i-1})$, is deterministic. In consequence,

$$e^{L(a)} = \Pr(Y_1, X_1) \prod_{i=2}^{n} \Pr(Y_i \mid Y_1, \ldots, Y_{i-1}, X_1, \ldots, X_i).$$

The outcome $Y_i$ is dependent only on the dose $X_i$ and not on previous outcomes or doses allocated. Therefore, we can further simplify:

$$e^{L(a)} = \prod_{i=1}^{n} \Pr(Y_i \mid X_i) = \phi_* \{x(1), y_1, \beta\} \prod_{l=2}^{n} \phi_* \{x(l), y_l, \beta\},$$

where

$$\phi_* \{x(l), y_l, \beta\} = \psi_*^{y_l}(x(l), \beta)\{1 - \psi_*(x(l), \beta)\}^{(1-y_l)}$$

and $\psi_*(x, \beta)$ is some chosen function, with a vector of parameters $\beta$, to model the probability of toxicity at dose $x$. When $\psi_*(x, \beta) = \psi(x, a)$ then inference is based on our simple one-parameter working model. Nothing prevents us fitting more complicated models, the logistic and probit models being obvious choices. Indeed we may begin the trial with a simple one-parameter model, at some intermediate stage fit a more complex model, and then continue the trial with this more complex model if, for instance, it seems to provide a noticeably improved fit. Our general impression, however, is that, for dose allocation, we are as well to stay with the original model, even if, at some point, for inferential purposes, we work with a more complex model. We give some discussion of this in the next section.

## 6.3 *One-Parameter Versus Two-Parameter Models*

Our intuition might in general cause us to be reluctant to use a one-parameter model, its lack of flexibility proving to be something of a handicap in trying to find a good fit to data. This lack of flexibility, though, can turn into an advantage in our situation since a better overall fit to the data may be coupled with a poorer point estimate of the level corresponding to our aimed-for "acceptable" toxicity rate. This is quite clearly seen in a simulation under the same conditions as the third set of conditions in Table 4. The level nearest the targeted one has a toxic response probability of 18% and the one-parameter model selects this level 57% of the time. A two-parameter logistic model with similarly defined vague priors ends up selecting this level only 48% of the time, and this despite the fact that the one-parameter model is of the wrong form. The overall probability estimates from the two-parameter model, however, give a better fit to the data. Bearing in mind though that overall fit is not the primary objective of a Phase 1 study, the one-parameter model seems preferable here, and in a large number of other cases we have looked at. Nonetheless, in situations where severe model inadequacy begins to tell, such as in the last two entries of Table 5, a two-parameter model can work better. The "correct allocation" rates of 58% and 37% in those tables become 66% and 53% in simulations based on a two-parameter model.

Another aspect of the two-parameter model, stemming from its greater flexibility, is the greater rapidity with which it changes dose levels early in the experiment. Under vague priors, starting out experimentation at level 3, three nontoxicities result in experimentation for the fourth subject at level 6. This seems a little incautious, could easily be rectified by modifying our priors, but seems already to be well catered for by the "dampening" effect of the one-parameter model. Our feeling, having looked at a wide variety of circumstances, is that the one-parameter model performs, in general, remarkably well and we have no

hesitations in recommending its use. Although there may be situations, which we are currently trying to identify, in which clear improvements can be made via more complex models, our general finding is that its performance is very much superior to many of those rules widely used in present-day Phase 1 trials. We finish this discussion by a look at how poorly some of these rules behave. In view of Storer's (1989) thorough study on such rules a detailed investigation is not appropriate.

### 6.4 *"Up and Down" Schemes*

There are many variants of these schemes, the four most widely used at present having been extensively studied by Storer (1989). He refers to these as Design A (traditional, groups of 3 patients, described in the introduction), Design B (single patients enter sequentially), Design C (a more conservative version of B), and Design D (a modified version of A).

   We generated data according to a six-level toxic response model with probabilities (.06, .08, .12, .18, .40, .71), as in the third entry of Table 4. Design A selected the extremely toxic level, level 6, nearly 35% of the time despite the fact that even level 5 would seem to be too toxic. Design B performed very similarly in this situation, recommending level 6 34% of the time, even when overriding the sequential rule and forcing the scheme to continue until all 25 patients had been entered. Design B is indeed so bad that it should never be used in studies of this type, even if we aim to sample around the median. For a model with associated toxicity probabilities (.08, .14, .35, .65, .96, 1.0), as in the fourth entry of Table 4, and again including 25 patients, the highly toxic level, level 4, is recommended 65% of the time! Even Design A, in this situation, ends up recommending this level, or levels yet more toxic, over 50% of the time. Storer indicates that Design B will be sampling around the 50th percentile and that the modified version of this (Design C) around the 33rd. Even so, simulations with Design C, on the basis of the first of the above two sets of probabilities, show it to recommend level 6 (71% chance of toxic response) in over 10% of runs, while for CRM and a targeted toxicity of .34 the corresponding percentage was less than one. At the other end of the toxicity spectrum for this set, CRM had a 2% chance of selecting any one of the first four levels (toxicity less than or equal to .18), in striking contrast to Design C's 25% chance.

   These operating characteristics are dismal and although some patching up can be done [Storer (1989), for instance, suggests some possibilities], we are convinced of CRM's superiority. The reason such schemes behave badly is quite simple. Not only do they not make efficient use of accumulated data, they make no use of such data at all, beyond say the previous three, or sometimes six, responses. One consequence of this is that, while for CRM the probability of a recommended change in dose level between the $n$th and $(n + k)$th entered patient ($k$ fixed) goes to zero as $n$ becomes large, no such property holds for any of the standard schemes.

### 6.5 *Further Work*

In Sections 4 and 5 only limited cases were considered, i.e., $\theta = .2$, the working model was taken to be a simple power function, the prior to be standard exponential, each trial included 25 patients and used 6 dose levels. What happens in other situations? We have carried out simulations at smaller sample sizes, for models having the form of one-parameter or two-parameter Weibull distributions, for models having the form of one-parameter or two-parameter logistic models, for priors of gamma and truncated normal form, and for values of $\theta$ between .1 and .5. Our overall conclusions remain unaltered although a deeper investigation could highlight the influence of these different factors.

RÉSUMÉ

Dans le cadre d'un modèle statistique simple, on propose une approche au déroulement et à l'analyse d'un essai clinique de phase 1 en cancérologie. L'aspect spécifique de cette approche est de permettre de réconcilier deux exigences conflictuelles de tels essai: la première concerne l'estimation efficace de la relation dose–effet (toxicité), la seconde le problème éthique qui conduit à rechercher le traitement le mieux adapté à chaque patient inclus dans l'essai. Nous avons simulé plusieurs cas de figure possibles dont certains très éloignés du modèle proposé. Les résultats de ces simulations sont très encourageants.

REFERENCES

Collins, J. M., Zaharko, D. S., Dedrick, R. L., and Chabner, B. A. (1986). Potential roles for preclinical pharmacology in Phase 1 clinical trials. *Cancer Treatment Reports* **70,** 73–80.
Hastings, N. A. J. and Peacock, J. B. (1975). *Statistical Distributions.* New York: Halsted Press.
Jaynes, E. T. (1986). Bayesian methods: General background. In *Maximum Entropy and Bayesian Methods in Applied Statistics,* J. H. Justice (ed.). Cambridge: Cambridge University Press.
Martz, H. F. and Waller, R. A. (1982). *Bayesian Reliability Analysis.* New York: Wiley.
Storer, B. E. (1989). Design and analysis of Phase I clinical trials. *Biometrics* **45,** 925–937.