

Adaptive sample size modification in clinical trials: start small then ask for more?

Christopher Jennison^{a*†} and Bruce W. Turnbull^b

We consider sample size re-estimation in a clinical trial, in particular when there is a significant delay before the measurement of patient response. Mehta and Pocock have proposed methods in which sample size is increased when interim results fall in a 'promising zone' where it is deemed worthwhile to increase conditional power by adding more subjects. Our analysis reveals potential pitfalls in applying this approach. Mehta and Pocock use results of Chen, DeMets and Lan to identify when increasing sample size, but applying a conventional level α significance test at the end of the trial does not inflate the type I error rate: we have found the greatest gains in power per additional observation are liable to lie outside the region defined by this method. Mehta and Pocock increase sample size to achieve a particular conditional power, calculated under the current estimate of treatment effect: this leads to high increases in sample size for a small range of interim outcomes, whereas we have found it more efficient to make moderate increases in sample size over a wider range of cases. If the aforementioned pitfalls are avoided, we believe the broad framework proposed by Mehta and Pocock is valuable for clinical trial design. Working in this framework, we propose sample size rules that apply explicitly the principle of adding observations when they are most beneficial. The resulting trial designs are closely related to efficient group sequential tests for a delayed response proposed by Hampson and Jennison. Copyright © 2015 John Wiley & Sons, Ltd.

Keywords: group sequential test; sample size re-estimation; adaptive design; clinical trial; optimal design; promising zone

1. Introduction

Adaptive strategies to extend a clinical trial have been proposed for a wide variety of situations: when there is uncertainty over the value of a nuisance parameter such as response variance; when there are co-primary endpoints with differing relative effect sizes; when the goals of superiority and non-inferiority are considered simultaneously; or in 'enrichment' designs where treatment efficacy may be assessed in both a general population and a specified sub-population. Nevertheless, considerable attention has been given to the relatively simple setting of a parallel design with a single primary outcome variable, where investigators re-consider the size of treatment effect they wish to detect with high probability and, typically, increase sample size to give adequate power at a smaller effect size. Such a change of objective is commonly based on an unblinded interim estimate of the treatment effect. It is this last form of sample size re-assessment that we shall examine in detail, prompted by recently published proposals.

The statistical considerations arising in clinical trials are many and complex and include: multiple primary and secondary endpoints, adverse safety events, compliance and treatment switching, quality of life, data quality, cost-benefit and risk analysis, intent-to-treat issues, and multiple sites and stratification. However, the section of the protocol that concerns the justification of the sample size, power, and interim analysis plan is typically based on a single primary endpoint (even though implications for other endpoints may be noted). If there is a plan for sample size modification at an interim analysis, specific instructions for this should be laid out in the protocol; this would then be in keeping with the US Food and Drug Administration's 'Draft Guidance on Adaptive Design' [1], which strongly endorses pre-specified,

^aDepartment of Mathematical Sciences, University of Bath, Bath, U.K.

^bSchool of Operations Research and Information Engineering, Cornell University, Ithaca, NY, U.S.A.

*Correspondence to: Christopher Jennison, Department of Mathematical Sciences, University of Bath, Bath, U.K.

†E-mail: C.Jennison@bath.ac.uk

non-flexible rules, stated in the study protocol. In a trial that is blinded to the sponsor, the Data Monitoring Committee (DMC) will be given flexibility to make decisions in some areas, such as responding to adverse safety events, but members of the DMC will have made a commitment to follow the rules for sample size modification and stopping boundaries by signing their agreement to the DMC Charter, which is a legal document.

Consider designing a trial where θ is the effect parameter of primary interest and the null hypothesis $H_0: \theta \leq 0$ is to be tested against the one-sided alternative $\theta > 0$. The type I error rate is to be protected at level α , and the goal is to achieve power of $1 - \beta$ at some positive treatment effect $\theta = \Delta$. Dispute may arise over the choice of Δ , for example, investigators may consider using a minimum effect of interest Δ_1 or a more optimistic anticipated effect size Δ_2 . As an example, suppose response variance and other features of a trial design are such that two fixed sample designs with sample sizes 500 and 1000 give rise to the power curves shown in Figure 1. If the treatment effect is equal to the minimum clinically significant effect, Δ_1 , a sample size of 1000 would give reasonable power, but a design with only 500 subjects would be under-powered. However, a sample size of 500 does provide good power under $\theta = \Delta_2$ and, if the true effect size were really this large, a sample size of 1000 would be unnecessarily high.

While regulatory guidances such as ICH E9 quite rightly emphasize the protection of the type I error rate, trial sponsors have a strong interest to design trials that achieve sufficient power in an efficient manner: reducing the sample size of a trial by just a few percent can easily translate into savings between \$100,000 and \$1million. From an ethical viewpoint, there is a desire to keep sample sizes small, producing results to support an effective new treatment as rapidly as possible and minimizing the numbers of subjects randomized to an inferior treatment. At the same time, it is important to avoid conducting under-powered studies.

Adaptive designs offer a mechanism to adjust sample size in response to updated treatment effect estimates. The proposals of Bauer and Köhne [2], Fisher [3], and Cui *et al.* [4] have the form:

- Start with a fixed sample size design;
- Examine interim data;
- Add observations to increase power when appropriate.

Group sequential designs (see, for example, [5]) follow a different route:

- Specify the desired power function, considering the range of possible effect sizes;
- Set the maximum sample size a little higher than the fixed sample size;
- Stop at an interim analysis if the data support an early conclusion.

When viewed holistically, these two approaches produce similar types of design. In each case, there is an overall maximum possible sample size but, depending on the observed data, the actual sample size

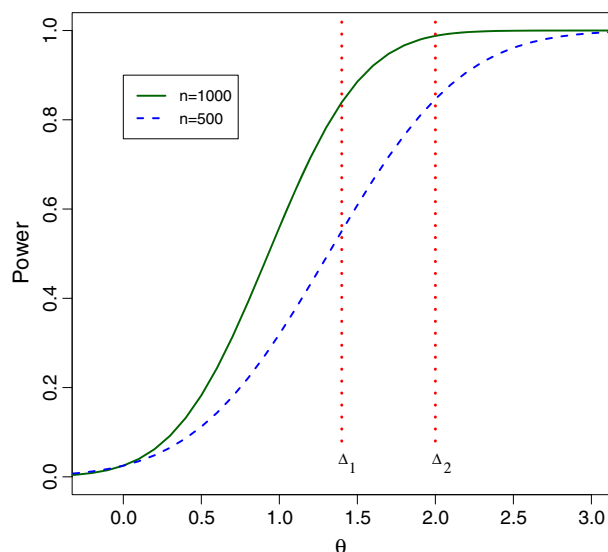


Figure 1. Choosing the sample size for a trial.

can be smaller than this. We are assuming here that an adaptive design has a pre-specified adaptation rule, as recommended in the US Food and Drug Administration Guidance [1]. Anderson and Liu used the phrase ‘start small and ask for more’ to describe the first form of design in a March 2004 presentation at the Conference on Adaptive Design for Clinical Trials in Philadelphia. However, with a pre-specified adaptation rule, there should be no doubt that sample size will be increased when this is required (this assumption does not exclude the case of a trial sponsor with limited resources who is able to rely on the promise of further investment if interim results are promising). With a full specification of the trial design in place, any adaptive or group sequential design has an overall power curve that can be computed at the start of the trial. Designs with similar power curves can be compared in terms of their average sample size functions, $E_{\theta}(N)$. Figure 2 compares two designs with essentially identical power curves: because the average sample size function of one is lower than that of the other for all values of θ , there can be little doubt that this design is the one to be preferred.

In previous work, we, and others, have compared group sequential designs (GSDs) and adaptive designs and concluded in favor of GSDs as they achieve given power with lower average sample size than published proposals for adaptive designs; see, for example, Jennison and Turnbull [6–8], Tsiatis and Mehta [9], and Fleming [10]. Suppose we specify a type I error rate α , power $1 - \beta$ at a designated effect size $\theta = \delta$, and a set of values of θ at which low $E_{\theta}(N)$ is desired. We impose the constraints that there are at most $K (\geq 2)$ analyses, and the maximum sample size is at most $R (> 1)$ times that which is required by a fixed sample size test. A GSD has a fixed sequence of group sizes and stopping boundaries that ensure the type I error rate and power requirements are met. An adaptive GSD (AGSD) is a generalization of a group sequential design in which, at any stage $k \in \{1, \dots, K - 1\}$, future group sizes are allowed to depend on the responses observed thus far. Because AGSDs form a larger class, they have the opportunity to be more efficient. However, the benefits of adapting group sizes are small: in examples presented by Jennison and Turnbull [7] the efficiency gain of optimal K -group AGSDs over optimal K -group GSDs is only about 2%; Lokhnygina and Tsiatis [11] compare adaptive and non-adaptive two-stage designs and report differences in efficiency of 1%; Banerjee and Tsiatis [12] investigate two-stage designs for a binary response and find adaptive designs to give decreases of 3% to 5% in expected sample size under the null hypothesis (although, because of the discreteness of the response, type I and type II error rates are not matched exactly in this comparison). It follows that, for *any* K -group adaptive design, there is a (simpler) K -group GSD of almost equal efficiency. In fact, in our investigations of published proposals for adaptive designs, we found these designs to use sub-optimal sample size rules, and this led to their being **significantly less efficient** than well-chosen GSDs.

Despite the aforementioned arguments in support of non-adaptive GSDs, Mehta and Pocock [13] have proposed a new adaptive procedure, termed the ‘promising zone’ approach. We shall refer to this paper hereafter as ‘MP’. Because the claims in MP are counter to the findings we have reported, it is appropriate to revisit the ‘GSD versus AGSD’ question. A significant feature of MP’s Example 1 is that response is measured some time after treatment, so many patients have been treated but are yet to produce a response at the interim analysis (we refer to these as *pipeline* subjects). Delayed response is common and not easily handled by standard GSDs, and the effects of pipeline subjects on average sample size need to

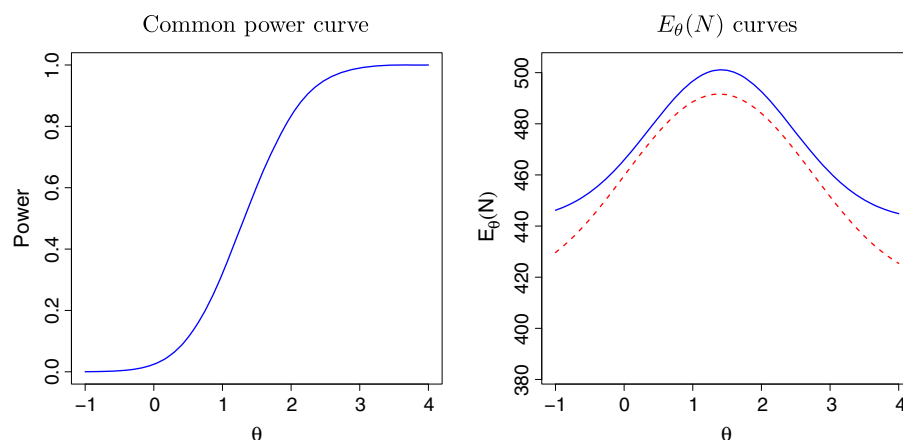


Figure 2. Power and average sample size curves for competing procedures.

be considered when comparing trial designs. Hampson and Jennison [14] have presented a new form of GSD, which handles delayed response and, in an investigation of two-stage versions of their designs, they found little benefit in adaptive choice of group sizes in their delayed response GSDs.

We note that Mehta and Pocock also recommend their methods for trials with a rapidly observed response—but then the conclusions of previous investigations [6–12] support the use of standard GSDs, as discussed earlier. In MP's second example of an acute coronary syndromes trial, the primary endpoint is measured within 48 hours of randomization. Meretoja *et al.* [15] have recently described a trial to be conducted using MP's promising zone approach, in which the primary endpoint is observed after 24 hours.

The organization of this paper is as follows. First, we introduce MP's Example 1 and examine the statistical properties of their proposed design for this problem. We present a two-group GSD, which can be applied to the same problem and which achieves essentially the same power curve. This rather naive GSD makes no use of the potential information from pipeline subjects at the interim analysis, but these subjects are counted in the expected sample size when the trial stops early: even so, we find that this naive GSD outperforms the MP design. We then take up the task of finding more competitive designs within MP's general framework. We show that alternative sample size rules can reduce average sample size, to an extent. However, for the best results, we find it is necessary to abandon the Chen *et al.* [16] approach for controlling the type I error rate. In Section 5, we use a combination test statistic [2] in the final hypothesis test. In his discussion of Mehta and Pocock [13], Glimm [17] suggests using the 'conditional error function' approach. A combination test is a special case of this method, and we see that, with a suitable sample size rule, a design employing a combination test can have very good efficiency; in particular, it can achieve the power of MP's design with a lower $E_\theta(N)$ function than the naive GSD that ignores pipeline data. In Sections 6 and 7, we return to tests using the usual Wald statistic in the final analysis. The design presented in Section 7 is an example of a 'delayed response group sequential design', as proposed by Hampson and Jennison [14], and we note the similarity of this design to adaptive designs created in our extension of MP's framework.

2. Mehta and Pocock's example

Mehta and Pocock's Example 1 concerns a Phase 3 trial of treatments for schizophrenia in which a new drug is to be tested against an active comparator. The efficacy endpoint is improvement in the Negative Symptoms Assessment score from baseline to week 26. Responses, denoted by Y_{Bi} , $i = 1, 2, \dots$, on the new treatment and Y_{Ai} , $i = 1, 2, \dots$, on the comparator treatment, are assumed to be normally distributed with known variance 7.5^2 , so each

$$Y_{Ai} \sim N(\mu_A, \sigma^2) \quad \text{and} \quad Y_{Bi} \sim N(\mu_B, \sigma^2),$$

where $\sigma^2 = 7.5^2$. The treatment effect is

$$\theta = \mu_B - \mu_A.$$

An initial plan is for a total of $n_2 = 442$ patients, 221 on each treatment, and the final analysis will reject $H_0: \theta \leq 0$, if $Z_2 > 1.96$, where

$$Z_2 = \frac{\bar{Y}_B(n_2) - \bar{Y}_A(n_2)}{\sqrt{\{4\sigma^2/n_2\}}} \quad (1)$$

and $\bar{Y}_A(n_2)$ and $\bar{Y}_B(n_2)$ are treatment means from a total of n_2 observations. This test has one-sided type I error rate 0.025 and power 0.8 at $\theta = 2$. Higher power, for example, power of 0.8 at $\theta = 1.6$, would be desirable, but the sponsors are only willing to increase sample size if interim results are promising.

An interim analysis is planned after observing $n_1 = 208$ responses. Because of staggered accrual and the 26-week delay in obtaining a response, another 208 pipeline subjects will have been treated but will not have completed 26 weeks of follow up at the time of the interim analysis. The purpose of the interim analysis is to revise the total sample size. The minimum value is the original figure of 442, which includes the pipeline subjects and an additional 26 new subjects; with 'promising data', an increase of up to 884 subjects is permitted. MP use conditional power to define their 'promising zone' and to determine the appropriate increase in sample size for particular interim results.

At the interim analysis, the estimated treatment effect is $\hat{\theta}_1 = \bar{Y}_B(n_1) - \bar{Y}_A(n_1)$ and the standardized test statistic is $Z_1 = \hat{\theta}_1 / \sqrt{(4\sigma^2/n_1)}$. The conditional power $CP_\theta(z_1)$ is defined to be the probability that the final test, with the original $n_2 = 442$ observations, rejects H_0 given $Z_1 = z_1$ if the effect size is θ . That is,

$$CP_\theta(z_1) = P_\theta\{Z_2 > 1.96 \mid Z_1 = z_1\}. \quad (2)$$

MP divide possible outcomes at the interim analysis into three regions, based on the conditional power under $\theta = \hat{\theta}_1$. These regions and the implications for sample size are:

<i>Favorable</i>	$CP_{\hat{\theta}_1}(z_1) \geq 0.8$	<i>Continue to $n_2 = 442$;</i>
<i>Promising</i>	$0.365 \leq CP_{\hat{\theta}_1}(z_1) < 0.8$	<i>Increase n_2;</i>
<i>Unfavorable</i>	$CP_{\hat{\theta}_1}(z_1) < 0.365$	<i>Continue to $n_2 = 442$.</i>

If the final decision is made by comparing the usual Wald statistic, Z_2 , to a standard normal distribution, the data dependent choice of n_2 may lead to inflation of the type I error rate; see [18]. It is, of course, crucial to protect the type I error rate when increasing sample size in the promising zone. MP do this by using a result of Chen, DeMets and Lan [16] and a subsequent extension of this result by Gao *et al.* [19]. Suppose at the interim analysis, the final sample size is increased to $n_2^* > n_2$ and a final test is carried out without adjustment for this adaptation, so H_0 is rejected if

$$Z_2(n_2^*) = \frac{\bar{Y}_B(n_2^*) - \bar{Y}_A(n_2^*)}{\sqrt{\{4\sigma^2/n_2^*\}}} > 1.96.$$

Chen, DeMets and Lan (CDL) proved the very neat result that the overall type I error probability will not increase if n_2 is only increased when

$$CP_{\hat{\theta}_1}(z_1) > 0.5. \quad (3)$$

Gao *et al.* [19] proved an extension of this result for values of $\hat{\theta}_1$ which are too low to satisfy (3). For an interval of $\hat{\theta}_1$ values, they showed the conditional type I error probability does not increase when the final sample size is increased to n_2^* , as long as n_2^* is higher than a lower bound, which depends on $\hat{\theta}_1$. We shall refer to an adaptive design using this result as following the CDL+Gao approach. In MP's Example 1, with an upper limit for n_2^* of 884, the final sample sizes permitted when using the CDL+Gao approach are shown in Figure 3. For values of $\hat{\theta}_1$ in the range (1.21, 1.40), n_2^* can be set equal to 426 or to a value

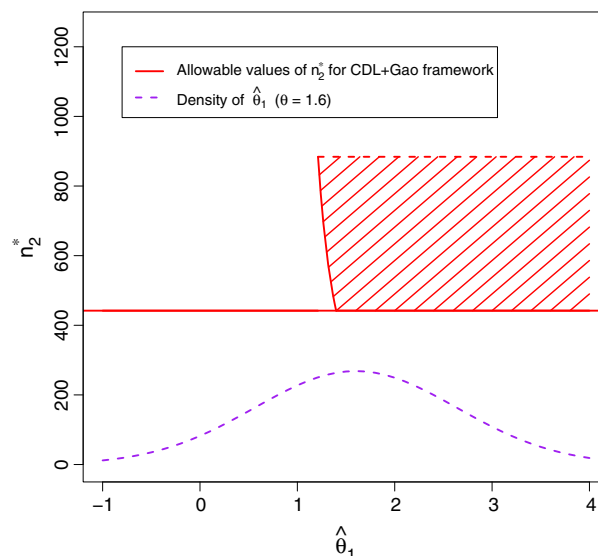


Figure 3. Sample size increases permitted by the method of Gao *et al.*

in the shaded region above the minimum value for this $\hat{\theta}_1$. At $\hat{\theta}_1 = 1.21$, the lowest value for which a sample size increase is permitted, the conditional power under $\theta = \hat{\theta}_1$ is 0.365, which explains the lower limit of the promising zone in MP's design.

In the lower panel of Figure 3, we display the density of $\hat{\theta}_1$ under $\theta = 1.6$ (on a different vertical scale) as an indication of the variability in $\hat{\theta}_1$. The distribution of $\hat{\theta}_1$ under other values of θ is shifted but has the same variance.

In their design, MP increase sample size to the value n_2^* that yields conditional power of 0.8 under $\theta = \hat{\theta}_1$, where this is possible. In the 'unfavorable' region, no increase in n_2 is permitted; in the 'favorable' region, $CP_{\hat{\theta}_1}(z_1)$ is already 0.8 or higher, and no increase is needed; in the 'promising' zone, n_2 is increased to the value that makes conditional power under $\theta = \hat{\theta}_1$ equal to 0.8, truncating this value to 884 if it is larger than that. This sample size rule is illustrated in Figure 4, where again we include the density of $\hat{\theta}_1$ under $\theta = 1.6$ for reference.

Because of the high variance of $\hat{\theta}_1$, increases in n_2 occur in a region of quite small probability, regardless of the true value of θ . The left-hand panel of Figure 5 shows that the resulting increase in power is quite small. Although it was stated that power 0.8 at $\theta = 1.6$ would be desirable, the power at this effect size has only risen from 0.61 to 0.658. The cost of this increase in power is the higher expected sample size function shown in the right-hand panel of Figure 5.

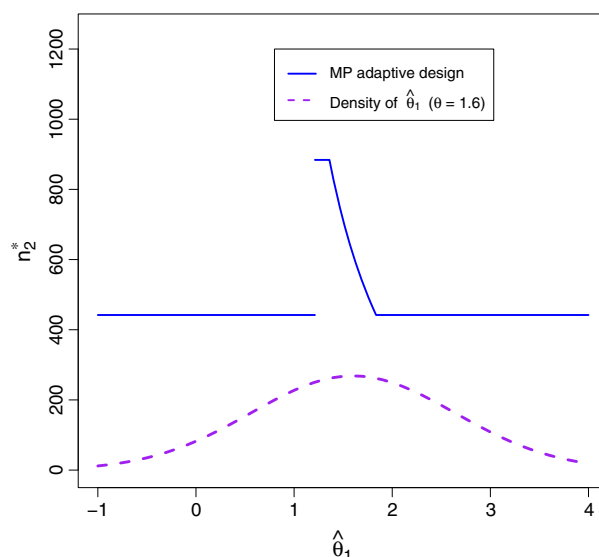


Figure 4. Sample size increase rule for Mehta and Pocock's design.

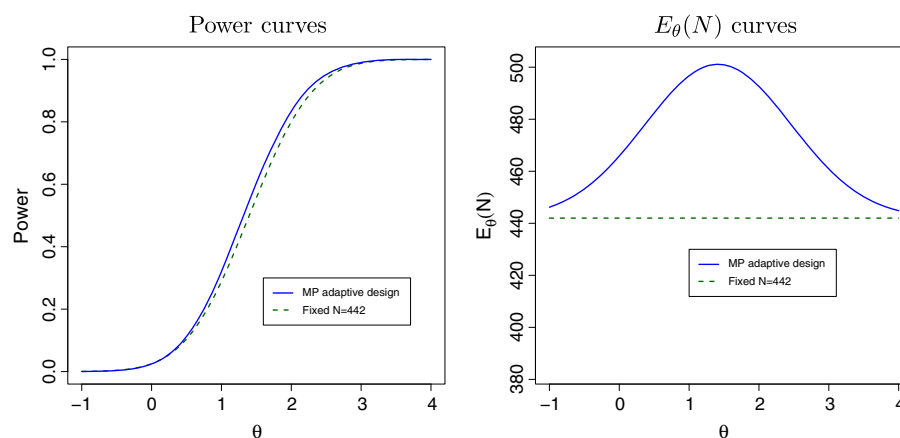


Figure 5. Power and average sample size curves for Mehta and Pocock's (MP) design and a fixed sample size design with $N = 442$.

Calculations underlying the plots in Figure 5 were carried out using the numerical integration methods described by Jennison and Turnbull in Chapter 19 of [5]. The same methods are used for all the results presented in this paper, and approximation errors in these results are negligible.

We note that the construction of the promising zone and the sample size increase function is very sensitive to the value of $\hat{\theta}_1$, and we have seen this is a highly variable estimate of θ . In fact, the value of $\hat{\theta}_1$ is used twice in determining the conditional power that underlies the sample size function: once through the value of z_1 in the conditional power (2) and again because this conditional power is evaluated at $\theta = \hat{\theta}_1$. This double role of $\hat{\theta}_1$ has been noted by Glimm [17] who describes this practice as ‘dangerous’ and recommends that the MP design ‘should be carefully inspected for its operating characteristics’.

It is permissible to modify the MP design by setting the sample size to achieve a higher conditional power under $\theta = \hat{\theta}_1$, or by raising the maximum for n_2 above 884. However, we have found that the resulting gains in power are small for the increases in $E_\theta(N)$. This leads us to consider alternatives to the MP design.

3. Alternatives to the Mehta and Pocock design

Suppose we are satisfied with the overall power function of the MP design. We shall present two more types of design which have the same power function and compare their properties with those of the MP design.

3.1. A fixed sample design

Emerson *et al.* [20] note that a fixed sample size study with 490 subjects has essentially the same power curve as the MP design. Comparison with Figure 5 shows that this figure of 490 is lower than the expected sample size of the MP design for effect sizes θ between 0.8 and 2.0. Although this fixed sample size is 11% more than the minimum of 442 required by the MP design, it is much lower than the MP design’s maximum sample size of 884.

3.2. A group sequential design

Despite the delayed response, we can still consider a group sequential design (GSD) with an interim analysis after 208 observed responses. However, if the trial stops to reject H_0 or accept H_0 at the first analysis, the sample size must be counted as 416 in order to include all the subjects admitted to the study and treated thus far.

We consider a two-stage error spending design for a one-sided test using a ρ -family error spending function with $\rho = 2$ that has type I error rate $\alpha = 0.025$ under H_0 and power 0.8 at $\theta = 1.9$; see Chapter 7 of [5]. This design has an interim analysis after 208 responses, with 208 pipeline subjects at that time, and a final analysis after a total of 514 subjects have been admitted and observed. The stopping rule and decision rule are:

At analysis 1

If $Z_1 \geq 2.54$	Stop, reject H_0
If $Z_1 \leq 0.12$	Stop, accept H_0
If $0.12 < Z_1 < 2.54$	Continue

At analysis 2

If $Z_2 \geq 2.00$	Reject H_0
If $Z_2 < 2.00$	Accept H_0

The sample size rules for the MP design, the fixed sample size trial, and our GSD are compared in Figure 6. For the GSD, the lower dot-dash line represents the 208 responses observed at the interim analysis, while the dashed line gives the sample size of 416 when the trial stops at this analysis and the final sample size of 514 when the trial continues on to the second group of observations. The maximum sample size of 514 is a factor $R = 1.05$ times the 490 needed for the same power in a fixed sample design.

The performance properties of the three designs are compared in Figure 7. By construction, all three designs have essentially the same power curve. It is evident that the GSD dominates the MP design

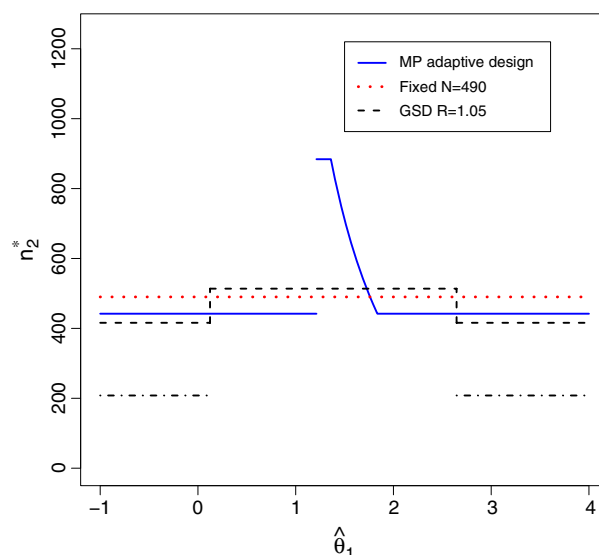


Figure 6. Sample size rules for Mehta and Pocock (MP), fixed ($N = 490$), and group sequential (GSD) designs.

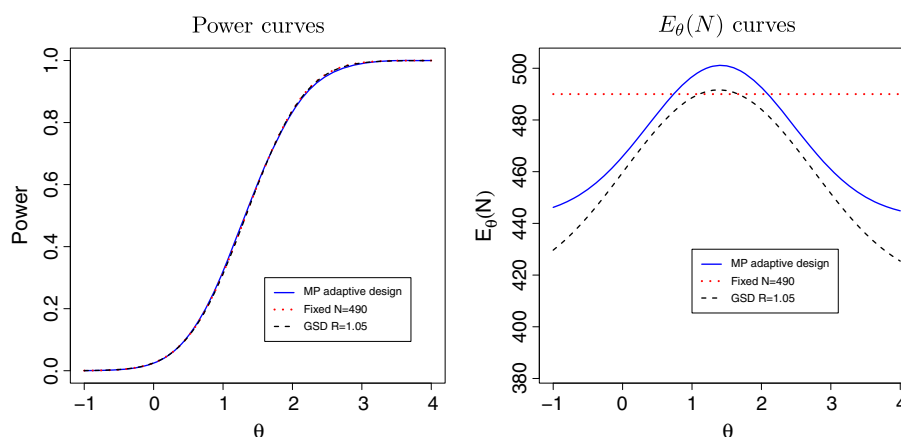


Figure 7. Power and average sample size curves for three designs.

everywhere with respect to average sample size. This is despite counting the 208 pipeline subjects when the trial stops at the interim analysis without using the information they could provide. Note that the two-stage GSD we have described is different from the two-stage GSD referred to as ‘Plan 3’ in [13]. That GSD has the higher power of 0.8 at $\theta = 1.6$, which, although desirable, is not achieved by the MP design. Because our design has the same overall power as the MP design, it is fairer to compare this design with the MP design with respect to average sample size.

One might simply decide at this point that conclusions about the relative merits of GSDs and AGSDs for the case of an immediate response carry over to trials with a delayed response. However, we recognize that our two-stage GSD has unsatisfactory features, and a design that made use of information from the pipeline subjects would be preferable. The framework of the MP design is appealing in that the information from the pipeline subjects at the interim analysis is eventually used, while interim data are used to decide how many additional subjects will be recruited for the final analysis; a similar strategy was previously proposed by Faldum and Hommel [21]. Also, the principle of adding observations when they will do the most good is attractive. The questions we shall address in the remainder of this paper concern: the form of sample size rule that gives the most additional power for the accompanying increase in sample size; and whether an alternative to the CDL+Gao approach for protecting the type I error rate is needed to create more efficient trial designs.

4. Deriving efficient sample size rules in the Mehta and Pocock framework

We continue to study MP's Example 1 and retain the basic elements of their design. The interim analysis takes place after $n_1 = 208$ observed responses. A final sample size n_2^* is to be chosen based on $\hat{\theta}_1$ or, equivalently, $Z_1 = \hat{\theta}_1 / \sqrt{\{4\sigma^2/n_1\}}$. We allow values of $n_2^* \in [442, 884]$ that in addition satisfy the CDL+Gao conditions, as displayed in Figure 3. At the final analysis, H_0 is rejected if $Z_2 > 1.96$, where Z_2 is the standard Wald statistic with no adjustment for the adaptive choice of sample size.

In our search for efficient sample size rules in the aforementioned framework, we specify a rule that makes a tradeoff between the competing goals of high conditional power and low sample size. Suppose we observe $Z_1 = z_1$ and are considering a final sample size n_2^* . Let

$$Z_2(n_2^*) = \frac{\bar{Y}_B(n_2^*) - \bar{Y}_A(n_2^*)}{\sqrt{\{4\sigma^2/n_2^*\}}} = \frac{\hat{\theta}(n_2^*)}{\sqrt{\{4\sigma^2/n_2^*\}}}.$$

For a given $\tilde{\theta}$, denote the conditional power under $\theta = \tilde{\theta}$, given $Z_1 = z_1$ and sample size n_2^* , by

$$CP_{\tilde{\theta}}(z_1, n_2^*) = P_{\tilde{\theta}}\{Z_2(n_2^*) > 1.96 \mid Z_1 = z_1\}.$$

In order to implement the idea of using additional sample size where it is most effective in increasing power, we specify a parameter γ , which represents the acceptable 'rate of exchange' between sample size and conditional power, evaluated at an effect size of interest $\theta = \tilde{\theta}$ (such as $\tilde{\theta} = 1.6$ in MP's example). For each value of z_1 , or equivalently $\hat{\theta}_1$, we choose n_2^* to optimize the combined objective

$$CP_{\tilde{\theta}}(z_1, n_2^*) - \gamma(n_2^* - 442). \quad (4)$$

At high values of $\hat{\theta}_1$, the gradient of $CP_{\tilde{\theta}}(z_1, n_2^*)$ decreases as n_2^* increases above 442, in which case our rule is to increase sample size up to the point where the improvement in conditional power on adding one more observation is less than γ . For lower values of $\hat{\theta}_1$, the gradient of $CP_{\tilde{\theta}}(z_1, n_2^*)$ increases and then decreases with n_2^* : in some cases, the maximum of (4) is at a value $n_2^* > 442$ where the gradient of $CP_{\tilde{\theta}}(z_1, n_2^*)$ falls back below γ ; in others, the maximum is at $n_2^* = 442$. We have subtracted 442 from n_2^* in (4) to emphasize that we are concerned with the cost of additional observations although, of course, this does not affect where the maximum occurs. To start with, we shall set $\tilde{\theta} = 1.6$, an effect size where power is lower than desired and we wish to 'buy' additional power.

Before using the criterion (4) to derive sample size rules, it is important to comment on the role of this objective function. Although one could base the value of γ on the cost of treating each subject in the trial and an estimate of the financial return from a positive outcome, this is certainly not required. Instead, γ can be regarded as a tuning parameter that controls the degree to which sample size may be increased when interim data are promising but not overwhelming. With such a criterion in place, interim data sets with different values of Z_1 will be treated consistently, with the benefits of any increase in sample size being measured on a common scale. Because the value of γ determines the operating characteristics of the optimized design, we can choose a value that gives a design with a specific overall power: we shall do precisely this to obtain designs with power curves matching that of the MP design seen in Figure 7.

Although the criterion (4) concerns *conditional* probabilities given the interim data, choosing a sample size rule to optimize this objective function also yields a design with an overall optimality property expressed in terms of *unconditional* power. Let the function $n_2^*(z_1)$ specify a sample size rule for choosing the total sample size n_2^* when $Z_1 = z_1$ is observed at the interim analysis. Then, we can write

$$P_{\tilde{\theta}}(\text{Reject } H_0) - \gamma E_{\tilde{\theta}}(N) = \int \{CP_{\tilde{\theta}}(z_1, n_2^*(z_1)) - \gamma n_2^*(z_1)\} f_{\tilde{\theta}}(z_1) dz_1, \quad (5)$$

where $f_{\tilde{\theta}}(z_1)$ denotes the density of Z_1 under $\theta = \tilde{\theta}$. Because our sample size rule maximizes $CP_{\tilde{\theta}}(z_1, n_2^*(z_1)) - \gamma n_2^*(z_1)$ for every z_1 , it also maximizes the right-hand side of (5). It follows that such a rule has the minimum $E_{\tilde{\theta}}(N)$ among all rules that achieve the same power under $\theta = \tilde{\theta}$. We shall see that working with this simple optimality property involving $E_{\tilde{\theta}}(N)$ at a single value of θ is sufficient to

explore aspects of the MP design and to give rules that improve on its performance over a range of θ values. More general criteria involving a weighted average of $E_{\theta_i}(N)$ taken over several effect sizes θ_i or a weighted integral of $E_{\theta}(N)$ can also be defined, and we shall discuss these in Section 6.

Having specified $\tilde{\theta}$ and γ , finding the sample size n_2^* that maximizes the combined objective function (4) is straightforward. For a given value of $\hat{\theta}_1$, and hence z_1 , we simply compute $CP_{\tilde{\theta}}(z_1, n_2^*)$ over the range of values of n_2^* and choose the n_2^* that gives the maximum value of $CP_{\tilde{\theta}}(z_1, n_2^*) - \gamma(n_2^*(z_1) - 442)$. We have chosen to use $\gamma = 0.140/(4\sigma^2) = 0.140/(4 \times 7.5^2)$ with $\tilde{\theta} = 1.6$, as this leads to a design with a power curve almost identical to that of the MP design.

The left panel of Figure 8 shows plots of $CP_{\tilde{\theta}}(z_1, n_2^*)$ and $CP_{\tilde{\theta}}(z_1, n_2^*) - \gamma(n_2^*(z_1) - 442)$ for the case $\tilde{\theta} = 1.6$, $\gamma = 0.140/(4\sigma^2)$ and $\hat{\theta}_1 = 1.5$, so $z_1 = 1.44$. The combined objective function has its maximum at $n_2^* = 654$, a somewhat lower value than the 712 for the MP design. Note that the slope of the function $CP_{\tilde{\theta}}(z_1, n_2^*)$ is γ at the optimum n_2^* because the higher slope at lower values of n_2^* implies that (4) is increasing, and the lower slope at higher values of n_2^* means (4) is decreasing.

The right panel of Figure 8 plots the same pair of functions for the less promising interim estimate $\hat{\theta}_1 = 1.3$. In this case, the slope of the conditional power curve is initially higher, and there is greater benefit from taking additional observations. The optimum n_2^* , where the derivative of $CP_{\tilde{\theta}}(z_1, n_2^*)$ has fallen to γ , occurs later, and (4) is maximized at $n_2^* = 707$. This sample size is substantially lower than the final sample size of 884 in the MP design.

The left panel of Figure 9 shows our optimized sample size n_2^* as a function of the interim estimate $\hat{\theta}_1$, labeled 'CDL+Gao Min $E(N)$ at $\theta = 1.6$ ' in the legend and compares this with the sample size rule of the MP design. The right panel shows $E_{\theta}(N)$ for these two designs and for the two-stage GSD described in

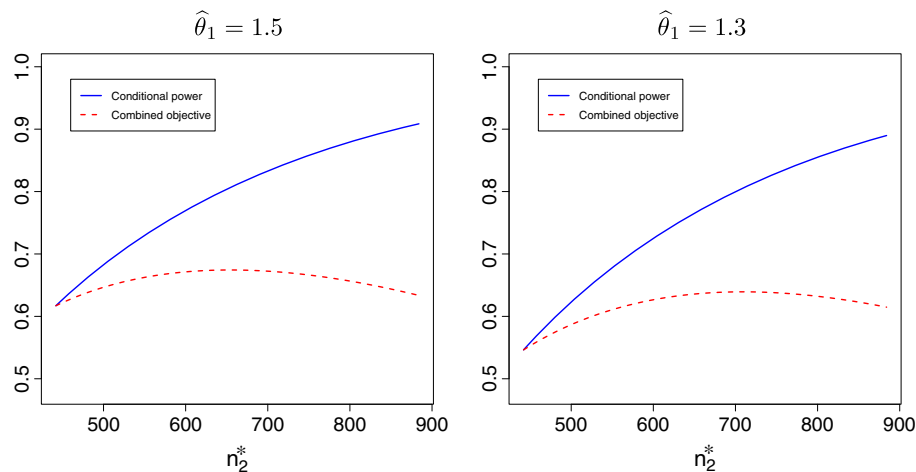


Figure 8. Plots of $CP_{\tilde{\theta}}(z_1, n_2^*)$ and $CP_{\tilde{\theta}}(z_1, n_2^*) - \gamma(n_2^*(z_1) - 442)$ for $\tilde{\theta} = 1.6$ and $\gamma = 0.140/(4\sigma^2)$.

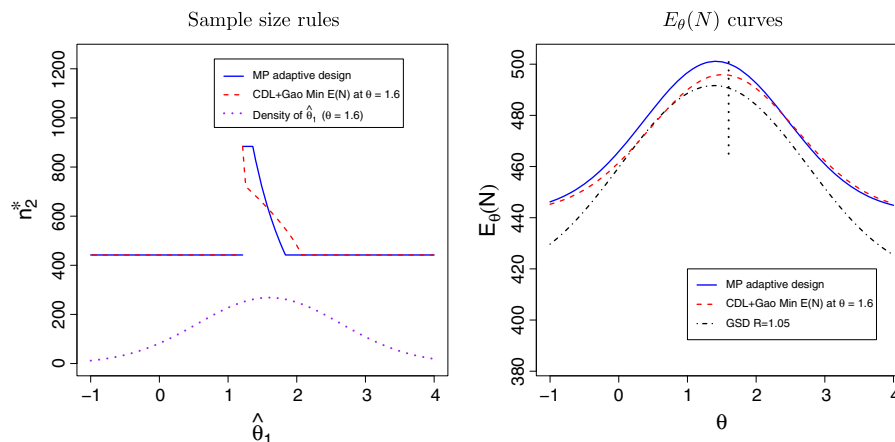


Figure 9. Efficient rules within the MP framework.

Section 3. The design using our optimal rule for $\gamma = 0.140/(4\sigma^2)$ has power 0.658 at $\theta = 1.6$, the same as the MP design. Because of the sigmoid shape of the power curves, matching the type I error probability of 0.025 at $\theta = 0$ and power 0.658 at $\theta = 1.6$ implies that the two designs have essentially identical power curves over the whole range of θ , as we saw in the left panel of Figure 7. The vertical dotted line in the right panel of Figure 9 is included to aid comparison with values of $E_\theta(N)$ at $\theta = 1.6$. Our optimal rule has the lowest possible $E_\theta(N)$ at $\theta = 1.6$ among all rules following the CDL+Gao framework that achieve power 0.658 at $\theta = 1.6$. Although this only guarantees an improvement over the MP design for $\theta = 1.6$, the right panel of Figure 9 shows we have achieved lower $E_\theta(N)$ for θ values up to 2.4, and there are only slight differences at higher values of θ .

Despite our new design's optimality property, we see from Figure 9 that it still has higher $E_\theta(N)$ than the two-stage GSD proposed in Section 3, which does not make use of data from pipeline subjects (although these are still counted in $E_\theta(N)$). We conclude that the constraints implicit in the CDL+Gao framework do not allow the most efficient designs to be realized. A first feature of the CDL+Gao construction is its conservatism: the actual one-sided type I error rate is less than the permitted $\alpha = 0.025$ and, as a consequence, the power is reduced. A second aspect of the CDL+Gao construction is the constraint on values of $\hat{\theta}_1$ for which the sample size can be increased. It is clear from Figure 6 that the two-stage GSD of Section 3 increases sample size over a wider range of $\hat{\theta}_1$ values than MP's promising zone, and the shape of our optimized sample size rule in the left panel of Figure 9 suggests it would help to increase n_2^* at lower values of $\hat{\theta}_1$. A third feature of MP's application of the CDL+Gao method is that they only allow sample size to be increased. In practical terms, it is possible to terminate recruitment at the interim analysis, in which case the pipeline subjects will bring the final numbers up to 416, a lower figure than the original sample size of 442. The two-stage GSD of Section 3 had 416 as its sample size when stopping at the interim analysis—but then it failed to use data from the pipeline subjects at all. Allowing a choice of final sample size with a minimum of 416, the value determined by the number of pipeline subjects at the interim analysis, is a natural way to extend the group sequential approach to the case of a delayed response.

In order to relax the aforementioned constraints that underlie MP's implementation of the CDL+Gao approach, we need an alternative method to protect the type I error rate in an adaptive design. We shall do this by using combination test statistics, as proposed by Bauer and Köhne [2], which enable more general adaptations.

5. Using combination test statistics

In creating a 'combination test' framework, we first revisit the initial plan, with a fixed total of $n_2 = 442$ subjects and an interim analysis with $n_1 = 208$ observations, and we express the final test in terms of the two sets of data observed before and after the interim analysis. We define the standardized statistics based on these two sets of observations as

$$V_1 = \frac{\sum_{i=1}^{n_1/2} Y_{Bi} - \sum_{i=1}^{n_1/2} Y_{Ai}}{\sqrt{(n_1\sigma^2)}}$$

and

$$V_2 = \frac{\sum_{i=(n_1/2)+1}^{n_2/2} Y_{Bi} - \sum_{i=(n_1/2)+1}^{n_2/2} Y_{Ai}}{\sqrt{\{(n_2 - n_1)\sigma^2\}}}$$

and note that the definition of the overall test statistic in (1) is equivalent to setting

$$Z_2 = w_1 V_1 + w_2 V_2$$

where $w_1 = \sqrt{(n_1/n_2)}$ and $w_2 = \sqrt{\{(n_2 - n_1)/n_2\}}$. (There is deliberate redundancy of notation here as V_1 is the same as the previously defined Z_1 . This is because we wished to reflect the similarity between terms V_1 and V_2 and were unable to use the name Z_2 , which is already defined as the standardized statistic for the combined data.) With n_2 fixed, $Z_2 \sim N(0, 1)$ under $\theta = 0$ and the decision rule is to reject H_0 : $\theta \leq 0$ when $Z_2 > 1.96$.

Now suppose we adapt the second stage sample size based on first stage data, increasing the total sample size from n_2 to n_2^* . The standardized statistic for the second stage data alone is now

$$V_2^* = \frac{\sum_{i=(n_1/2)+1}^{n_2^*/2} Y_{Bi} - \sum_{i=(n_1/2)+1}^{n_2^*/2} Y_{Ai}}{\sqrt{\{(n_2^* - n_1)\sigma^2\}}}. \quad (6)$$

The weighted inverse normal combination test statistic is defined as

$$Z^* = w_1 V_1 + w_2 V_2^*$$

where, importantly, w_1 and w_2 are the *original* weights, defined in terms of n_1 and n_2 . The combination test based on this test statistic rejects $H_0: \theta \leq 0$ when $Z^* > 1.96$.

To see that this combination test has one-sided type I error rate 0.025, consider first the case $\theta = 0$. In this case, V_2^* has a conditional $N(0, 1)$ distribution given any set of interim data; hence, V_2^* has an unconditional $N(0, 1)$ distribution and is, therefore, statistically independent of V_1 . Because $V_1 \sim N(0, 1)$, $V_2^* \sim N(0, 1)$, and $w_1^2 + w_2^2 = 1$, we have $Z^* \sim N(0, 1)$, and the combination test has type I error rate 0.025. It remains to consider cases where $\theta < 0$. Here, V_1 is normal with variance 1 and a fixed negative mean. The conditional distribution of V_2^* given the interim data is normal with variance 1 and a mean which depends on the choice of n_2^* , but is always negative. It follows that the overall distribution of Z^* is stochastically smaller than a $N(0, 1)$ variate when $\theta < 0$, and so the combination test's type I error rate is less than 0.025.

Some comments on the combination test approach are appropriate before we apply this method to MP's Example 1. Bauer and Köhne [2] proposed the use of combination tests in adaptive clinical trial design, focusing on a combination rule based on the product of P -values introduced by R. A. Fisher [22]. There is a relationship between combination tests and procedures that preserve the conditional type I error probability when adaptation occurs, as proposed by, for example, Proschan and Hunsberger [18] and Müller and Schäfer [23]. Jennison and Turnbull [6] note that combination tests preserve the conditional type I error probability given the interim data at the time of adaptation, and they go on to show that any flexible design, which allows a choice of whether or not to adapt, must preserve this conditional error rate in order for the overall type I error rate to be protected. Glimm [17] argues that the MP designs can be viewed as procedures that protect the conditional type I error probability, but with some conservatism. It follows from the equivalence referred to earlier that the MP designs can also be interpreted as a conservative form of combination test that applies a higher threshold than 1.96 to the statistic Z^* defined in (6) when the CDL+Gao method is used to ensure the type I error rate is protected.

We now return to MP's Example 1. We shall apply the weighted inverse normal combination test with weights $w_1 = \sqrt{(n_1/n_2)} = \sqrt{(208/442)}$ and $w_2 = \sqrt{\{(n_2 - n_1)/n_2\}} = \sqrt{(234/442)}$. Applying this test after the total sample size has been changed to n_2^* , the conditional power given $V_1 = Z_1 = z_1$ under treatment effect θ is

$$CP_\theta(z_1, n_2^*) = P_\theta\{Z^* > 1.96 \mid Z_1 = z_1\} = P_\theta\{V_2^* > (1.96 - w_1 z_1)/w_2 \mid Z_1 = z_1\}, \quad (7)$$

where the conditional distribution of V_2^* given $Z_1 = z_1$ is $N(\theta\sqrt{(n_2^* - n_1)/(2\sigma)}, 1)$. We are at liberty to modify the sample size at all values of z_1 , or equivalently $\hat{\theta}_1$, in the knowledge that the type I error rate will be protected at level 0.025. We use this freedom: firstly, to avoid the conservatism when increasing sample size for values of $\hat{\theta}_1$ covered by the CDL+Gao approach; secondly, to increase sample size for values of $\hat{\theta}_1$ where this is not permitted in the CDL+Gao approach; and thirdly, to decrease sample size for some values of $\hat{\theta}_1$, noting that the pipeline subjects imply a minimum sample size of 416.

In order to create a design with low $E_{\hat{\theta}}(N)$ under $\theta = \hat{\theta}$, we follow our previous strategy and choose n_2^* to maximize the combined objective

$$CP_{\hat{\theta}}(z_1, n_2^*) - \gamma(n_2^* - 442), \quad (8)$$

where now $CP_{\hat{\theta}}(z_1, n_2^*)$ is defined by (7). The argument presented in Section 4 can be applied to show that the resulting design has the minimum value of $E_{\hat{\theta}}(N)$ among all designs in this larger class that achieve the same power under $\theta = \hat{\theta}$.

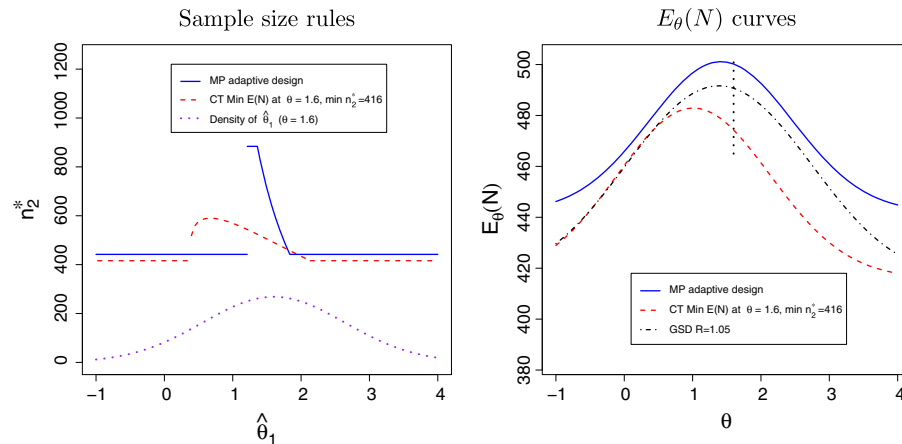


Figure 10. Efficient combination test (CT) designs.

As before, we consider combination test (CT) designs optimized for $\tilde{\theta} = 1.6$. The optimal CT design with $\gamma = 0.250/(4\sigma^2)$ matches the MP test's power of 0.658 at $\theta = 1.6$. The sample size rule for this design is shown in the left panel of Figure 10. The form of this optimized sample size rule is quite different from that of the MP design, with the greatest increases in sample size occurring at values of $\hat{\theta}_1$ below the range for which the CDL+Gao approach allows any change of sample size. Plots of average sample size in the right panel of Figure 10 show the CT design has lower $E_\theta(N)$ than the MP design across the range of θ values; this new design also improves on the simple group sequential design, which ignores (but is charged for) data from pipeline subjects when stopping at the interim analysis. Examination of intermediate designs that exploit only some of the three freedoms listed earlier show that the improved performance depends in roughly equal parts on all three aspects of the CT design: absence of conservatism; permitting sample size modification for all values of $\hat{\theta}_1$; and allowing reduction in sample size in some cases.

There is clearly scope to increase sample sizes beyond those shown for the optimized CT design in the left panel of Figure 10. Thus, it would be quite feasible to reduce the parameter γ in (8) and produce efficient designs with higher power curves. However, we shall continue to calibrate designs by matching the MP design's power curve in order to draw comparisons with this design.

One conclusion from examining the optimized CT design is that the best opportunities for investing additional resource are *not* in MP's promising zone. The left panel of Figure 10 shows this directly when interest is primarily in reducing $E_\theta(N)$ for $\theta = 1.6$. However, the fact that the CT design has lower $E_\theta(N)$ for all values of θ implies that the MP design cannot be ideal for *any* value of θ .

We have presented a CT design which minimizes $E_{\tilde{\theta}}(N)$ for $\tilde{\theta} = 1.6$. However, we see from the $E_\theta(N)$ curve in the right hand panel of Figure 10 that this design has robust efficiency as it performs well under a wide range of θ values. One could, instead, optimize a CT design for a different value of $\tilde{\theta}$, choosing γ in (4) to meet a specified power requirement. We have constructed such designs with power curves matching the MP design and examined their properties: their sample size rules vary gradually with $\tilde{\theta}$, as do the resulting $E_\theta(N)$ curves.

Several authors have pointed out that the combination test approach leads to the use of a non-sufficient statistic in the final decision rule. In order to give credibility to the final decision, Burman and Sonesson [24] suggest use of a 'dual test' that rejects H_0 overall only if it is rejected by both the combination test and the naive test based on the sufficient statistic and ignoring the adaptive sampling. The MP design can be viewed as applying the dual test when, following the CDL+Gao approach, it applies the naive test rather than the more permissive combination test. Interestingly, in using the combination test to facilitate sample size increases not permitted by the CDL+Gao approach, we are in step with the dual test as here the combination test is *stricter* than the naive test. See also the discussion in Section 4 of [25]. Clearly, issues of credibility must be considered when specifying a trial design and the method of analysis. However, we believe that combination tests are well understood, and we would expect the CT design illustrated in Figure 10 to be accepted as a valid method.

6. Further extensions of the Mehta and Pocock framework

We have seen how the weighted inverse normal combination test can be used to create a trial design that has low average sample size under a specified effect size. The same methodology can be extended in two further ways.

6.1. Minimizing a weighted average sample size

Rather than aim to minimize $E_\theta(N)$ at a single value $\theta = \tilde{\theta}$, we could consider a weighted average of terms $E_{\theta_i}(N)$ over a number of effect sizes θ_i , or an integral

$$\int h(\theta) E_\theta(N) d\theta. \quad (9)$$

Here, the function $h(\theta)$ should reflect both the likelihood of each θ value and the importance placed on reducing sample size at that effect size. After scaling h so that its integral is 1, the expression (9) can be regarded as the average sample size when θ is drawn from a prior distribution with density $h(\theta)$. Assuming we still wish to use the weighted inverse normal combination test, we can derive the sample size function that maximizes

$$\int h(\theta) \{P_\theta(\text{Reject } H_0) - \gamma E_\theta(N)\} d\theta \quad (10)$$

for a specified value of γ . A little algebra shows that, on observing $Z_1 = z_1$, n_2^* should be chosen to maximize

$$\int g(\theta|z_1) CP_\theta(z_1, n_2^*) d\theta - \gamma n_2^*, \quad (11)$$

where $g(\theta|z_1)$ is the posterior distribution of θ given prior $h(\theta)$ and $Z_1 = z_1$, and $CP_\theta(z_1, n_2^*)$ is the conditional power function defined by (7). The calculations simplify when $h(\theta)$ is a normal density as the conditional distribution of V_2^* , as defined in (6), is also normal and the integral in (11) is a single normal tail probability. Use of the combination test implies that the resulting procedure will have the required type I error probability. One may then search over values of γ to find a design with a suitably high power function.

6.2. Generalizing the final decision rule

Another option is to replace the inverse normal combination test by a more general final decision rule. Suppose we wish to find the design in this general class, which maximizes the optimality criterion (10) with a specific choice of $h(\theta)$, subject to a type I error constraint. This problem is very close to a two-stage version of the K -stage adaptive design problem solved by Jennison and Turnbull [7], but now, the presence of pipeline subjects implies a minimum value for n_2^* , and the criterion (10) contains an integral of the power function rather than power at a single value of θ . Because power curves follow, almost exactly, a one-parameter family, it is just as easy to calibrate a design through this integral of power.

The optimal design can be found by following the procedure described in Appendix 2 of [7]. To summarize briefly, we formulate a Bayes sequential decision problem with a mixture prior for θ comprising a point mass at $\theta = 0$ and the remaining probability distributed with a density proportional to $h(\theta)$ and add a cost for a type I error when $\theta = 0$ to the criterion (10). Given an observed value $Z_1 = z_1$ and a candidate sample size n_2^* , the optimal final decision on observing V_2^* can be derived as a Bayes rule and the posterior risk given $Z_1 = z_1$, using sample size n_2^* and the Bayes optimal decision rule can be evaluated. The value of n_2^* with the smallest posterior risk is the optimal choice for the case $Z_1 = z_1$. Properties of the resulting design are found by integrating over values z_1 . Hence, for a given value of γ , the cost of a type I error can be chosen so that the optimal design has type I error probability 0.025. A second level of searching over values of γ can then be conducted to find the optimal design with the desired power characteristics.

One benefit of working through this optimization process is that it provides a benchmark of the best possible sampling and decision rules for a given objective. The final decision rule has the attractive feature that, because it is a Bayes rule, it is expressed as a function of the sufficient statistic for θ [7]. Working in this general class, we have found the design for MP's Example 1 that minimizes $E_\theta(N)$ for $\theta = 1.6$ among all designs with $n_1 = 208$, n_2^* in the range 442 to 884, type I error probability 0.025 and power 0.658 at

$\theta = 1.6$. We found this design to be almost indistinguishable from the optimized CT design of Section 5 whose properties are shown in Figure 10. The fact that there is little room to improve on the weighted inverse normal combination test in this case can be attributed to the small range of n_2^* values used in the optimal design. One should expect greater benefit from more general final decision rules when the range of n_2^* values is greater, and we have seen that this is in fact the case in examples with smaller numbers of pipeline subjects, and hence, a lower minimum value that n_2^* can take. In the next Section, we shall explain how this optimized design is closely related to a form of group sequential design developed to accommodate delayed response and the resulting pipeline subjects at interim analyses.

7. Using a delayed response group sequential design

We have noted that MP's Example 1 has the distinctive feature of a large number of subjects in the pipeline at the time of the interim analysis. Until recently, little attention had been paid to the effect of a delayed response on standard GSDs. This is rather surprising, given that almost all endpoints are measured some time after treatment, and in many trials, this delay is substantial. Hampson and Jennison [14] review proposals for incorporating observations from pipeline subjects that accrue after the decision has been made to stop a trial at an interim analysis and note that it has proved difficult to use this information effectively; they then propose a new form of GSD, which anticipates data from pipeline subjects from the outset. These delayed response group sequential designs (DRGSDs) combine a rule that stipulates when patient recruitment ceases and a final decision rule, applied after observations have been obtained from pipeline subjects, which determines whether or not the null hypothesis is rejected.

We have constructed a DRGSD with just two analyses for MP's Example 1. This design matches the MP design in having type I error probability 0.025 and power 0.658 when $\theta = 1.6$. As in other designs, the first analysis takes place after $n_1 = 208$ observed responses. If termination is halted at this point, the study continues until responses are obtained from the 208 pipeline patients, and the final sample size is 416. If recruitment continues beyond the first analysis, a further 102 patients are admitted to the trial, and the final sample size is 518. This maximum sample size of 518 was obtained by multiplying the fixed sample size needed to achieve power 0.658 under $\theta = 1.6$ by an inflation factor of $R = 1.05$. Given these sample sizes, four constants are needed to complete the specification of the DRGSD. In the stopping rule and decision rule stated below, we have chosen the critical values to minimize $E_\theta(N)$ at $\theta = 1.6$ among designs with type I error probability 0.025 and power 0.658 when $\theta = 1.6$. Computation of these values is achieved by solving a related Bayes decision problem, in a similar way to the optimizations described in Section 6.

At analysis 1 (208 responses)

If $Z_1 \leq 0.088$ or $Z_1 \geq 1.999$

Halt recruitment, continue
to final sample of 416 (Case A)

If $0.088 < Z_1 < 1.999$

Recruit a further 102 subjects,
continue to final sample of 518 (Case B)

Case A: At analysis 2 with 416 responses

If $Z_2 \geq 1.948$

Reject H_0

If $Z_2 < 1.948$

Accept H_0

Case B: At analysis 2 with 518 responses

If $Z_2 \geq 1.984$

Reject H_0

If $Z_2 < 1.984$

Accept H_0

The sample size rule for this DRGSD is shown in the left panel of Figure 11. In contrast to the CT design shown in the same display, the DRGSD has only two possible total sample sizes. Nevertheless, the average sample size curves in the right panel of Figure 11 demonstrate that the DRGSD is almost as efficient as the CT design. Having just two scenarios to plan for after the interim analysis may have logistical benefits for trial organizers. There is certainly less of a problem of 'information leakage', whereby parties who are blinded in order to protect the integrity of the trial may be able to make deductions about the interim data from knowledge of the target sample size after the interim analysis.

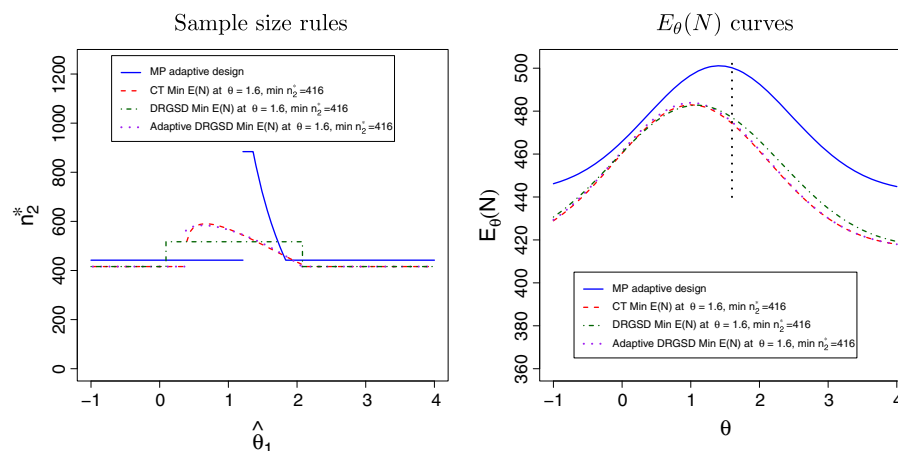


Figure 11. Comparison of Mehta and Pocock (MP) and combination test (CT) designs with delayed response group sequential designs (DRGSD).

Hampson and Jennison [14] considered adaptive versions of their DRGSD in which the final sample size could take any value, subject to the minimum value determined by the number of pipeline subjects who would automatically be observed. The adaptive DRGSD that minimizes $E_{\theta}(N)$ at $\theta = 1.6$ is also included in Figure 11: the left panel shows that its sample size rule is almost identical to that of the CT design and the average sample size curves in the right panel are so close that they look like a single curve. In fact, this adaptive DRGSD is exactly the same as the optimal design with a general form of final decision rule discussed in Section 6 because both are constructed by minimizing $E_{\theta}(N)$ at $\theta = 1.6$ within the class of two-stage designs with $n_1 = 208$, $n_2^* \geq 416$, type I error probability 0.025 and power 0.658 at $\theta = 1.6$.

From the aforementioned comparisons, we deduce that, in this example, there is minimal benefit to be gained from fine-tuning the final sample size of a DRGSD in response to interim data. Hampson and Jennison [14] report similar findings in other examples of two-stage designs with different numbers of pipeline subjects. Thus, we have strong evidence that Jennison and Turnbull's [7] argument that adaptive choice of group sizes in a GSD offers at best modest gains can be extended from the case of an immediate response to the case of delayed responses.

If planning for a data-dependent total sample size does not pose logistical problems, a trial design with a fully adaptive sample size may still be considered as an attractive proposition. All we would note is that the sample size rule and final decision rule need to be chosen carefully. We have seen that an optimized CT design works well in MP's Example 1; however, when the number of pipeline subjects is smaller, the optimal CT design is not so efficient, and more general adaptive designs should be considered. On the other hand, it is known that DRGSDs with fixed group sizes based on error spending functions from the class described in Section 4.1 of [14] provide highly efficient designs for a variety of pipeline sizes. Moreover, the error spending versions of Hampson and Jennison's DRGSDs are able to adapt to an unpredictable number of pipeline patients.

8. Conclusions

We have studied in depth Mehta and Pocock's [13] Example 1 and evaluated their proposed trial design for this problem. The nature of the delayed response in this example means that standard GSDs are not readily applicable. In the discussion at the 2004 workshop 'Adaptive Clinical Trial Designs: Ready for Prime Time?' [26], Mehta stated 'I do not agree with Jay Siegel's assertion that this overrun problem will be exactly the same in the adaptive situation as it is in the group sequential situation. The two situations are not analogous.' The need for a new type of trial design that allows a data-dependent choice of sample size while making proper use of the information anticipated from pipeline (or overrun) subjects motivates the MP design.

Mehta and Pocock [13] note it is desirable to spend additional resource when this will have the greatest benefit. We fully agree with this objective and have pursued it in the example. Our investigations have led us to some clear conclusions. The first is that designs that rely on the result of Chen *et al.* [16] to support

a simple form of final analysis are at a disadvantage because this result does not allow sample size to be increased in situations where the greatest benefits might accrue. Secondly, we have found that setting sample size to attain a specific conditional power under $\theta = \hat{\theta}_1$ is sub-optimal. Instead, we have derived efficient designs by using a weighted inverse normal combination test to control the type I error rate and by choosing the final sample size to optimize a criterion that balances the gain in conditional power under a fixed effect size θ against the extra sample size, imposing the constraint that the final sample size must include current pipeline subjects. The calculations required for this optimization are quite simple and involve only conditional probabilities and expectations, even though the resulting design also has optimal unconditional properties. Although we have focused on the particular case of MP's Example 1, we have found the aforementioned conclusions to apply quite generally for different delays in observing the final endpoint and, hence, different numbers of pipeline patients.

Another option is to follow the conventional group sequential framework more closely and limit the choice at the interim analysis to halting recruitment and waiting to observe responses for the pipeline subjects, or adding a fixed number of additional subjects. This gives a DRGSD, as defined in [14], and one can follow the recommendations of Hampson and Jennison for using error spending designs to create designs with robust efficiency over a range of effect sizes. The error spending approach is useful in dealing with departures of observed sample sizes from the numbers originally planned. Additional features of Hampson and Jennison's methods described in [14] are: they extend easily to designs with three or more analyses; if measurements can be taken on a short-term endpoint, which is correlated with the primary endpoint (for example, one might consider the Negative Symptoms Assessment score measured at 10 weeks in MP's Example 1), this information can be used at an interim analysis to gain greater benefit from the pipeline subjects.

References

1. US Food and Drug Administration Guidance for Industry: *Adaptive Design Clinical Trials for Drugs and Biologics (draft)*. FDA: Silver Spring, MD, 2010. (Available at <http://www.fda.gov/downloads/Drugs/GuidanceComplianceRegulatoryInformation/Guidances/UCM201790.pdf>), [Accessed 1 May 2015].
2. Bauer P, Kohne K. Evaluation of experiments with adaptive interim analyses. *Biometrics* 1994; **50**:1029–1041. Correction 52:380.
3. Fisher LD. Self-designing clinical trials. *Statistics in Medicine* 1998; **17**:1551–1562.
4. Cui L, Hung HMJ, Wang S-J. Modification of sample size in group sequential clinical trials. *Biometrics* 1999; **55**:853–857.
5. Jennison C, Turnbull BW. *Group Sequential Methods with Applications to Clinical Trials*. Chapman & Hall/CRC: Boca Raton, 2000.
6. Jennison C, Turnbull BW. Mid-course sample size modification in clinical trials based on the observed treatment effect. *Statistics in Medicine* 2003; **22**:971–993.
7. Jennison C, Turnbull BW. Adaptive and nonadaptive group sequential tests. *Biometrika* 2006; **93**:1–21.
8. Jennison C, Turnbull BW. Efficient group sequential designs when there are several effect sizes under consideration. *Statistics in Medicine* 2006; **25**:917–932.
9. Tsiatis AA, Mehta C. On the inefficiency of the adaptive design for monitoring clinical trials. *Biometrika* 2003; **90**:367–378.
10. Fleming TR. Standard versus adaptive monitoring procedures: a commentary. *Statistics in Medicine* 2006; **25**:3305–3312.
11. Lokhnygina Y, Tsiatis AA. Optimal two-stage group-sequential designs. *Journal of Statistical Planning and Inference* 2008; **138**:489–499.
12. Banerjee A, Tsiatis AA. Adaptive two-stage designs in phase II clinical trials. *Statistics in Medicine* 2006; **25**:3382–3395.
13. Mehta CR, Pocock SJ. Adaptive increase in sample size when interim results are promising: a practical guide with examples. *Statistics in Medicine* 2011; **30**:3267–3284.
14. Hampson LV, Jennison C. Group sequential tests for delayed responses (with discussion). *Journal of the Royal Statistical Society: Series B* 2013; **75**:3–54.
15. Meretoja A, Churilov L, Campbell BC, Aviv RI, Yassi N, Barras C, Mitchell P, Yan B, Nandurkar H, Bladin C, Wijeratne T, Spratt NJ, Jannes J, Sturm J, Rupasinghe J, Zavala J, Lee A, Kleinig T, Markus R, Delcourt C, Mahant N, Parsons MW, Levi C, Anderson CS, Donnan GA, Davis SM. The Spot sign and Tranexamic acid On Preventing ICH growth — AUSTALASIA Trial (STOP-AUST): Protocol of a phase II randomized, placebo-controlled, double-blind, multicenter trial. *International Journal of Stroke* 2014; **9**:519–524.
16. Chen YHJ, DeMets DL, Lan GKK. Increasing the sample size when the unblinded interim result is promising. *Statistics in Medicine* 2004; **23**:1023–1038.
17. Glimm E. Comments on 'Adaptive increase in sample size when interim results are promising: a practical guide with examples'. *Statistics in Medicine* 2012; **31**:98–99.
18. Proschan MA, Hunsberger SA. Designed extension of studies based on conditional power. *Biometrics* 1995; **51**:1315–1324.
19. Gao P, Ware JH, Mehta C. Sample size re-estimation for adaptive sequential design in clinical trials. *Journal of Biopharmaceutical Statistics* 2008; **18**:1184–1196.
20. Emerson SS, Levin GP, Emerson SC. Comments on 'Adaptive increase in sample size when interim results are promising: a practical guide with examples'. *Statistics in Medicine* 2011; **30**:3285–3301.

21. Faldum A, Hommel G. Strategies for including patients recruited during interim analysis of clinical trials. *Journal of Biopharmaceutical Statistics* 2007; **17**:1211–1225.
22. Fisher RA. *Statistical Methods for Research Workers* 5th edition. Oliver and Boyd: Edinburgh, 1934.
23. Müller HH, Schäfer H. Adaptive group sequential designs for clinical trials: combining the advantages of adaptive and of classical group sequential approaches. *Biometrics* 2001; **57**:886–891.
24. Burman CF, Sonesson C. Are flexible designs sound? *Biometrics* 2006; **62**:664–669.
25. Wang SJ, Brannath W, Brückner M, Hung JHM, Koch A. Unblinded adaptive statistical information design based on clinical endpoint or biomarker. *Statistics in Biopharmaceutical Research* 2013; **5**:293–310.
26. Golub H, Mehta C, Tsiatis B, Temple R, D’Agostino R, Fleming T, Siegel J, Ellenberg S, Hung J, Chuang-Stein C, *et al.* Papers from the 2004 Harvard-MIT Division of Health Science and Technology Workshop, ‘Adaptive Clinical Trial Designs: Ready for Prime Time?’ group discussion. *Statistics in Medicine* 2006; **25**:3326–3347.