



Randomized Play-the-Winner Clinical Trials: Review and Recommendations

William F. Rosenberger, PhD

*Department of Mathematics and Statistics, University of Maryland, Baltimore County,
and Department of Epidemiology and Preventive Medicine, University of Maryland
School of Medicine, Baltimore, Maryland*

ABSTRACT: The randomized play-the-winner rule is an adaptive randomized design, based on an urn model, that is used occasionally in clinical trials. This paper discusses practical and theoretical issues arising from its use, including stratification, delayed response, operating characteristics, selection of urn parameters, and inference. The paper also discusses recent experience with adaptive clinical trials within the pharmaceutical industry. The author concludes that the randomized play-the-winner rule is appropriate for some clinical trials, but intense and thoughtful planning must take place in the design phase. Such planning should incorporate considerations of variability, power, and appropriate techniques. *Control Clin Trials* 1999;20:328–342 © Elsevier Science Inc. 1999

KEY WORDS: *Adaptive designs, ethics in clinical trials, randomization, treatment allocation*

INTRODUCTION

The randomized play-the-winner (RPW) rule [1] employs a simple probability model to randomize subjects sequentially in a clinical trial. The rule applies specifically to clinical trials with binary outcomes (“success” or “failure,” loosely defined) and two treatments (“A” and “B”). The original formulation assumed that the previous subject’s outcome will be available before the next patient is randomized. At the start of the clinical trial, an urn contains α_A balls of type A and α_B balls of type B, where α_A and α_B are positive integers. When a subject is recruited, a ball is drawn and replaced. If it is a type A ball, the subject receives treatment A; if it is type B, the subject receives treatment B. When the subject’s outcome is available, the urn is updated. A success on treatment A or a failure on treatment B will generate an additional β type-A balls in the urn. A success on treatment B or a failure on treatment A will generate an additional β type-B balls in the urn, where β is a positive integer. In this way, the urn builds up more balls representing the more successful (or less unsuccessful) treatment. This scheme increases the probability that a volunteer will receive the treatment performing better thus far, an allocation

*Address reprint requests to: Professor W. F. Rosenberger, Department of Mathematics and Statistics,
University of Maryland, Baltimore County, 1000 Hilltop Circle, Baltimore, MD 21250.
Received June 10, 1998; accepted April 6, 1999.*

likelier to be popular with volunteers and their physicians than equal assignment.

The RPW rule is a particular example of an *adaptive design*, which uses sequentially accruing data to affect future selection of design points. An important subdiscipline of experimental design, adaptive designs are popular in engineering and the life sciences, to name two areas (see Flournoy and Rosenberger [2], for example). They have a long history in clinical trials methodology (see Rosenberger [3] for a brief history). However, very few adaptive clinical trials have occurred in practice, perhaps because a single clinical trial, commonly known as the extracorporeal membrane oxygenation (ECMO) trial [4], failed owing to its small sample size ($n = 12$) and to its inclusion of only one control patient. The controversy surrounding this trial is well documented (see Royall [5] and Rosenberger and Lachin [6] for details). Recently, however, several ongoing adaptive clinical trials have indicated something of a resurgence of interest in adaptive designs. Reports of two such trials have been published [7, 8].

This paper updates my previous paper [6], which discussed many logistical and statistical issues that needed resolution in order for RPW clinical trials to gain popularity. Having worked for several years on some of the statistical issues, I now have a firmer understanding of complex properties underlying the seemingly simple RPW rule. Also, Eli Lilly and Co. has implemented a clinical trial using the rule [7], giving us a stronger comprehension of the practical difficulties of such a task, as well as information on the participants' perceptions. This paper discusses underlying assumptions of the model, blinding, stratification, operating characteristics of the design (including its variability) and inference, with a view to practical implementation.

SOME EXAMPLES

Example 1: ECMO (1985)

The RPW rule was applied in a clinical trial of ECMO, a surgical procedure for newborns with respiratory failure used when infants are moribund and unresponsive to conventional treatment, which includes ventilation and pharmacologic therapy. Investigators chose the RPW rule because

- (1) the outcome of each case [was] known soon after randomization, making it possible to use;
- (2) [it was] anticipated that most ECMO patients would survive and most control patients would die, so significance could be reached with a modest number of patients; [and]
- (3) it was a reasonable approach to the scientific ethical dilemma [4, p. 480].

The trial employed the RPW rule with $\alpha_A = \alpha_B = \beta = 1$. The first patient was assigned to ECMO and survived. The second patient was assigned to conventional treatment and died. All subsequent patients were assigned to ECMO and survived. Investigators implemented a stopping rule after having randomized 12 patients. The RPW rule did assign more patients to the more successful therapy, but the trial had only one infant assigned to control.

Subsequent commentary labeled the trial a failure (see Royall [5] for a lively exchange). (Interestingly enough, the ECMO debate has not yet ended; a recent

clinical trial in the United Kingdom randomized 93 patients to ECMO and 92 to conventional therapy. A total of 30 infants died on the ECMO arm and 54 on the conventional therapy [9].)

Example 2: The Fluoxetine Trial (1994)

Eli Lilly employed the RPW rule recently in its multicenter clinical trial comparing fluoxetine to placebo for depressive disorder [7]. In addition to employing an adaptive structure, the design of this trial possessed a number of subtleties:

1. Staggered entry. Entry time was approximately uniformly distributed over a 270-day time interval [10].
2. The trial stratified patients into two groups: normal and shortened rapid eye movement latency (REML).
3. The trial assigned the first six patients in each stratum with permuted block randomization. Thereafter, the RPW rule with initial composition (1, 1) and $\beta = 1$ was employed by stratum (i.e., the trial used two independent urns).
4. The primary outcome, a reduction of 50% or greater in Hamilton Depression Scale (HAM-D₁₇) between baseline and final active visit after a minimum of 3 weeks of therapy, could only be ascertained after approximately 8 weeks. Determining that this was too long a period in which to run an adaptive trial, investigators used a surrogate marker to update the urn. They defined a surrogate responder as a patient exhibiting a reduction greater than 50% in HAM-D₁₇ in two consecutive visits after at least 3 weeks of therapy.
5. The trial was stopped after 61 patients had responded in accordance with the surrogate criterion. No further surrogate response was obtained for the remaining 89 patients.

The fluoxetine trial showed no significant difference between the treatment groups in either stratum, and the adaptive allocation came out very close to 50:50.

ASSUMPTIONS OF THE MODEL AND OTHER CONCERNS

Statement of the Obvious

Several assumptions underlie adaptive designs. First, we assume that the better treatment will actually perform better in the clinical trial. This seems like a statement of the obvious, and it is certainly an assumption that is in place every time we analyze any clinical trial. In adaptive designs the consequences of its violation are serious because if the assumption does not hold, more subjects will receive an inferior treatment. Second, we assume that the treatment with better outcome does not cause some serious toxicity. If it does, then most of the patients will receive a treatment with serious toxicity. There is nothing we can do about the first assumption. For the second assumption, it is important to use adaptive designs only when earlier experiments have thoroughly reviewed

drugs for toxicity. One can also incorporate combined endpoints such as toxicity and efficacy into an adaptive framework (e.g., Rosenberger [3]).

Delayed Response

The RPW rule does not require instantaneous outcomes, or even that they be available before randomization of the next subject. Investigators can update the urn when a subject's outcome is ascertained. The effect of this will be to "slow" the adaptation, and hence there will be less benefit to subjects, particularly those recruited early. As a secondary effect it will be virtually impossible to obtain theoretical results on operating characteristics of the RPW rule. In practice, we often use simulation as a tool to determine properties of the RPW rule under delayed response, either assuming a deterministic response time or some failure-time distribution. We shall discuss techniques to handle estimation under a delayed-response model later in the paper.

Yao and Wei [11] describe a multistage RPW rule that updates the urn only at certain time points. At the end of a prespecified time interval, the urn is updated. The number of type A balls added is the number of successes on A and failures on B during that time interval. Likewise, the number of type B balls added is the number of failures on A and successes on B during the time interval. Time intervals can be adjusted in accordance with the accrual rate in the trial.

Despite the theoretical feasibility of employing the RPW rule when outcomes are not immediate, many trials, such as long-term survival trials, have in practice a limited recruitment period and few outcomes are available before most volunteers are recruited. Investigators for these types of trials are not likely to be interested in a simple Bernoulli response. Several recent papers have described techniques for adaptive survival trials [11–13] that are more appropriate than the RPW rule. To ensure maximum time on study, a limited recruitment period is often essential when there is an administrative censoring date. Limited recruitment conflicts directly with the concept of adaptive designs, however, because subjects may be recruited long before any outcome data are available. Investigators may wish to consider lengthening the recruitment period to allow adaptation. Later in the paper, we shall discuss appropriate operating characteristics of the RPW rule with delayed response and also examine some inferential techniques that lend themselves to adaptive designs with delayed response.

Example 2 (continued)

In the fluoxetine trial, investigators used a surrogate outcome to update the urn when the primary outcome was not available quickly enough. It is unclear, however, how to analyze data on the primary outcome when investigators use a surrogate outcome in this way unless we ignore the design of the trial in its analysis. We discuss the use of permutation tests, as used by the Lilly investigators, later in this paper.

Blinding

Because the RPW rule is a fully randomized design it enjoys the same benefits as standard randomization procedures. In particular, for large samples, the

probability of a covariate imbalance is negligible. As with any randomization procedure, investigators should blind clinical trials using the RPW rule to maintain the integrity of the study by minimizing selection bias. RPW clinical trials should be double-blind, if feasible. As suggested in Rosenberger and Lachin [6], a subject needs only to know that he or she will receive one of two treatments and that the probability of treatment assignment will depend on the relative frequency of the favorable and unfavorable responses to treatments of previously treated volunteers. Investigators should not disclose the current allocation probability to the volunteer or physician. Adaptive designs may, however, lead to a unique type of bias, *accrual bias* [3], by which volunteers may wish to be recruited later in the study so as to benefit from the full impact of previous outcomes, because earlier subjects will have had higher probabilities of receiving the inferior treatment. Hence, volunteers should be blinded to their sequence numbers in the trial. Investigators have yet to explore the acceptability to patients and physicians of this form of blinding. Clearly, though, in trials dealing with emergency therapies, such as emergency surgical techniques, accrual bias is irrelevant.

Stratification

Most multicenter clinical trials that prestratify by clinic maintain balance on unknown covariates across clinics. Prestratification often occurs on major demographic variables, such as gender, race, or age. Similarly, clinical trials employing urn models can also prestratify by using separate urns within each clinic or other stratum. If one expects, a priori, heterogeneity across strata with respect to the underlying probability of outcome, one can run separate urns within strata. If this is not done, subjects could be randomized on the basis of responses from subjects who do not share the same underlying probability of outcome. If the trial employs prestratification, results can be stratum-specific, or combined across independent strata, using likelihood-based techniques. A number of problems arise in the context of adaptive designs, however. Stratification by clinic may lead to a relatively small sample size in each clinic, and some clinics may have very low recruitment. Often clinics with very small enrollment are pooled in some ad hoc way, but this is not possible if one wants to maintain the integrity of adaptive randomization. In the fluoxetine trial, investigators decided not to prestratify by clinic; however, they stratified by another characteristic, normal and shortened REML, and used two urns [7]. Using separate urns among a large number of strata reduces the adaptive nature of the design.

Despite the large-sample property of randomization that makes negligible the probability of a covariate imbalance, some imbalances always result in clinical trials. Post-stratification, or modeling techniques, can adjust for these imbalances. Unfortunately, regression modeling usually assumes independent observations, and this is not the case with adaptive designs, where subjects respond to treatments that are allocated on the basis of previous outcomes. Standard results such as asymptotic normality and consistency of regression estimators of covariate effects other than a treatment effect have not yet been rigorously proven under the RPW rule. This is *not* to say that regression estimators do not have those properties, but from a theoretical standpoint, this is

still an open problem. As of now, we do know that the maximum likelihood estimators of the success probabilities from each treatment (say p_A and p_B) have the usual properties [14, 15] and same variances, and so we can carry out inference on the simple difference or odds ratio. Coad explored prestratification and post hoc modeling with adaptive designs in his papers [16, 17], but did not study the RPW rule.

OPERATING CHARACTERISTICS

Let N_A/N be the proportion of subjects assigned to treatment A out of N subjects. Also, let $q_A = 1 - p_A$ and $q_B = 1 - p_B$ be the failure probabilities. As N tends to infinity, $E(N_A/N)$ converges to $q_B/(q_A + q_B)$ [18]. The most extreme proportion of subjects assigned to A is 0.9 (when $p_A = 0.9$ and $p_B = 0.1$). A table in Rosenberger and Sriram [15] shows how large N must be for N_A/N to be in a tolerance range of $q_B/(q_A + q_B)$. Convergence is fairly quick unless values of p_B are very large. We can use this table to determine how large the sample size should be before it is acceptable to rely on asymptotic results.

It is interesting to note that the number of subjects assigned to A, relative to the number assigned to B (denoted N_A/N_B), converges to q_B/q_A . Hence asymptotically, the rule allocates according to the relative risk of failure on B as compared to A, giving an interesting justification to the ad hoc RPW rule.

Asymptotic results may not be relevant for small trials, but we can compute the exact expected value of N_A/N for any value of N using the formula given in Rosenberg and Sriram [15] programmed in MATLAB (program available from author upon request). Although the asymptotic formula does not depend on the initial urn composition, the exact formula does. Table 1A gives the exact expected value for $\alpha_A = \alpha_B = 1$ and $N = 25$. Table 1B gives the same for $\alpha_A = \alpha_B = 5$. We set $\beta = 1$ in all the tables for simplicity of presentation, although the MATLAB program can incorporate different values of β . One would expect the allocation proportions to be less extreme when the initial urn composition is increased, as the urn will not favor the better treatment as highly. We can see this in comparing the two tables. In Table 1A, the maximum value of $E(N_A/N)$ is 0.81, compared with 0.69 in Table 1B. Note how far these values are from the maximum of the asymptotic results. Clearly, the asymptotic results should be used cautiously.

Variability in Allocation Proportions

Studies of adaptive designs often neglect variability. Because adaptive designs induce random processes, different sequences can arise from the same set of responses to treatment. One recent paper discusses variability in adaptive designs (but not the RPW rule) [19] and finds that mean squared errors are quite large. It is well-known that the sample path of a stochastic process will depend on its initial starting values, and therefore selection of α_A and α_B are of critical importance. Recently, we derived an expression for the variance [20] of the proportion of subjects assigned to treatment A, allowing us to explore rigorously the variability of the urn model for various initial compositions, underlying success rates, and sample sizes.

Table 1 Exact Expected Value (SD) of N_A/N

	A. $\alpha_A = \alpha_B = 1$					B. $\alpha_A = \alpha_B = 5$				
	0.1	0.3	0.5	0.7	0.9	0.1	0.3	0.5	0.7	0.9
0.1	0.500 (0.073)	0.444 (0.086)	0.377 (0.098)	0.294 (0.103)	0.191 (0.096)	0.500 (0.076)	0.460 (0.085)	0.415 (0.092)	0.364 (0.097)	0.306 (0.097)
0.3	0.555 (0.086)	0.500 (0.103)	0.431 (0.118)	0.345 (0.130)	0.237 (0.130)	0.539 (0.085)	0.500 (0.095)	0.454 (0.104)	0.403 (0.111)	0.344 (0.114)
0.5	0.622 (0.098)	0.568 (0.118)	0.500 (0.139)	0.412 (0.158)	0.299 (0.169)	0.584 (0.092)	0.545 (0.104)	0.500 (0.115)	0.448 (0.125)	0.388 (0.131)
0.7	0.705 (0.103)	0.654 (0.130)	0.587 (0.158)	0.500 (0.187)	0.384 (0.212)	0.635 (0.097)	0.596 (0.111)	0.552 (0.125)	0.500 (0.138)	0.439 (0.148)
0.9	0.808 (0.096)	0.762 (0.130)	0.700 (0.169)	0.615 (0.212)	0.500 (0.255)	0.693 (0.097)	0.655 (0.114)	0.611 (0.131)	0.560 (0.148)	0.500 (0.164)

Tables 1A and 1B also give the exact standard deviation of N_A/N for $N = 25$ with $\alpha_A = \alpha_B = 1$ and $\alpha_A = \alpha_B = 5$, respectively, using the formula in Matthews and Rosenberger [20]. One can see that the variability increases as p_A and p_B increase. Also, as expected, the variability decreases as the initial urn composition gets larger. For instance, in Table 1A, the maximum value is 0.26, as compared to Table 1B, where the maximum value is 0.16. Hence, the RPW rule is extremely variable with small sample sizes, and its use is questionable. As the sample size increases, however, the variability is much less. For example, when $N = 100$, $\alpha_A = \alpha_B = 5$, $p_A = 0.7$ and $p_B = 0.3$, the standard deviation is just 0.06.

Example 1 (continued)

As described, the ECMO trial used the RPW rule with $\alpha_A = \alpha_B = 1$ and a stopping rule that terminated the trial at $N = 12$ [4]. At that point, 11 patients had been allocated to the experimental therapy and only one to the control therapy. The extreme allocation is consistent with the variability of the allocation proportions. Earlier studies suggested a success rate of the new therapy of between 56% and 70%, and control patients were expected to have a 20% success rate [4]. If we use the most extreme prior information, that is, $p_A = 0.70$ and $p_B = 0.20$, we have $E(N_A/N) = 0.71$ and $SD(N_A/N) = 0.16$. The actual allocation was $N_A/N = 0.92$, which is about 1.3 standard deviations above the mean. Investigators suggested that the trial would have had a more reasonable allocation if $\alpha_A = \alpha_B = 5$. In this event, $E(N_A/N) = 0.58$ and $SD(N_A/N) = 0.14$. An actual allocation 1.3 standard deviations above the mean would have led to 9:3 allocation, which is more reasonable. The ECMO study demonstrated clearly that the variability in the urn model could lead to success runs and unfortunate results. We should remember, though, that any trial of size $N = 12$ is unlikely to be very convincing, regardless of the allocation. (Note that these computations ignore the stopping rule.)

Delayed Response

Exact computations preclude delayed response because of theoretical difficulties, so we use exact results for immediate response and then simulate the urn model under a simple time-to-response mechanism to explore the effects of delayed response. In this paper, we employ two simple delayed-response models:

1. Deterministic: One patient arrives per time unit, and each patient responds as success or failure after two time units.
2. Stochastic: One patient arrives per time unit. Response is randomly generated as either immediate, after one time unit, or after two time units, each with probability 1/3.

In both models, responses in a time unit are assumed to occur after patient arrival. It is easy to program a deterministic rule, such as rule 1. Stochastic rules, like rule 2, require more complex data structures, such as priority queues, to keep track of staggered entry, censoring, and delayed response, as was done in [13]. The stochastic rule in 2 is relatively simplistic. We used simulation to

Table 2 Simulated Mean (SD) of N_A/N Under Delayed Response Model 2

	0.1	0.3	0.5	0.7	0.9
0.1	0.500 (0.078)	0.460 (0.086)	0.418 (0.093)	0.371 (0.097)	0.316 (0.098)
0.3	0.539 (0.087)	0.500 (0.096)	0.456 (0.105)	0.409 (0.110)	0.354 (0.115)
0.5	0.581 (0.092)	0.542 (0.103)	0.499 (0.113)	0.449 (0.124)	0.394 (0.129)
0.7	0.629 (0.096)	0.593 (0.110)	0.548 (0.124)	0.499 (0.135)	0.444 (0.143)
0.9	0.684 (0.098)	0.647 (0.114)	0.604 (0.128)	0.555 (0.145)	0.502 (0.156)

$\alpha_A = \alpha_B = 5, \beta = 1, N = 25.$

explore rules 1 and 2 for $\alpha_A = \alpha_B = 1$ and 5 and found that delayed response made the allocation proportions slightly less skewed toward the better treatment, but not markedly so, and the variability remained the same as for immediate response. Each simulation had 10,000 replications. Table 2 gives the values for $\alpha_A = \alpha_B = 5$ and rule 2. Note that even the most skewed value (0.684) under delayed response is only slightly less than that for immediate response (0.693). We conclude that delays in response do not affect the variability of the RPW rule and that the allocation proportions are only slightly more conservative.

Proportion of Treatment Failures

Tables 3A and 3B examine the simulated expected proportion of failures under immediate response and delayed-response model 2, respectively, for $\alpha_A = \alpha_B = 5$ and $n = 25$. Again, delayed response does not affect variability, and the expected failures increase by a factor of approximately half a percent at the most.

ESTIMATION AND INFERENCE

Large Sample Inference

If we define S_A and S_B as the number of successes on treatments A and B in an RPW trial, respectively, then we give the likelihood of the data by the product of binomials:

$$L \propto p_A^{S_A} q_A^{N_A - S_A} p_B^{S_B} q_B^{N_B - S_B}.$$

The design, however, is not ancillary to the likelihood because, unlike in the equal allocation case where the sufficient statistics are S_A and S_B , there is an additional sufficient statistic, N_A . Hence, the allocation proportions are sufficient statistics, and provide information about p_A and p_B . Conditioning on N_A yields an identical likelihood as for equal allocation, but we would lose the information the design has to tell us. It is thus best to consider unconditional inference.

Let $\hat{p}_A = S_A/N_A$ and $\hat{p}_B = S_B/N_B$ be the maximum likelihood estimators. Note that these do not have binomial distributions because the RPW rule imposes

Table 3 Simulated Mean (SD) of the Proportion of Failures with $\alpha_A = \alpha_B = 5, \beta = 1, N = 25$

	A. Immediate Response						B. Under Delayed Response Model 2													
	0.1		0.3		0.5		0.7		0.9		0.1		0.3		0.5		0.7		0.9	
	0.1	(0.060)	0.793	(0.083)	0.667	(0.100)	0.518	(0.110)	0.345	(0.104)	0.900	(0.059)	0.791	(0.083)	0.667	(0.101)	0.523	(0.108)	0.353	(0.102)
0.1	0.900	(0.060)	0.793	(0.083)	0.667	(0.100)	0.518	(0.110)	0.345	(0.104)	0.900	(0.059)	0.791	(0.083)	0.667	(0.101)	0.523	(0.108)	0.353	(0.102)
0.3	0.793	(0.082)	0.700	(0.091)	0.590	(0.099)	0.461	(0.105)	0.307	(0.097)	0.793	(0.083)	0.700	(0.091)	0.590	(0.100)	0.463	(0.104)	0.312	(0.096)
0.5	0.666	(0.101)	0.592	(0.100)	0.501	(0.099)	0.388	(0.099)	0.255	(0.089)	0.667	(0.099)	0.590	(0.101)	0.500	(0.100)	0.391	(0.099)	0.259	(0.089)
0.7	0.518	(0.109)	0.460	(0.103)	0.389	(0.100)	0.299	(0.091)	0.188	(0.078)	0.522	(0.108)	0.463	(0.104)	0.389	(0.098)	0.298	(0.092)	0.190	(0.078)
0.9	0.346	(0.103)	0.305	(0.096)	0.256	(0.089)	0.189	(0.078)	0.102	(0.060)	0.354	(0.102)	0.313	(0.096)	0.258	(0.089)	0.189	(0.078)	0.100	(0.060)

dependencies in the response sequence. The usual large-sample approach considers a test statistic based on the simple difference (or the odds ratio), for example:

$$\frac{\hat{p}_A - \hat{p}_B - (p_A - p_B)}{\left\{ \frac{p_A q_A}{N_A} + \frac{p_B q_B}{N_B} \right\}^{1/2}} \quad (1)$$

(with the usual substitutions under H_0), which is asymptotically $N(0, 1)$ with RPW allocation [14]. Therefore use of asymptotic inference allows us to ignore the design in the analysis.

Power

If we use a test statistic such as (1), when $N_A \neq N_B$, power can be lost relative to equal allocation designs and hence, to maintain power, we would have to increase the sample size of the trial. In that event, more patients will receive the inferior treatment, and we may lose the benefits of the adaptive design (note that Lachin [21, Figure 1] shows minimal loss in power even for N_A/N around 0.75). In principle, power is just one factor to consider in planning a clinical trial, and the relative ethical cost of a type I and type II error may be different in a trial involving adaptive allocation. One might be willing to sacrifice some power by using an adaptive design to assign more patients to the better treatment. Large sample inference does sacrifice some power if sample proportions are far from 0.5, principally because we are ignoring the important information in N_A as a sufficient statistic. Exact unconditional inference incorporates the information in N_A , and hence should be more powerful. In the context of sequential tests, the increased sample size required to maintain the same power for the RPW rule and equal allocation still allows fewer expected failures under RPW allocation [22].

Example 3: AZT in Maternal–Fetal HIV Transmission

A clinical trial comparing AZT to placebo evaluated whether AZT reduces the risk of maternal–fetal HIV transmission [23]. The trial used standard randomization, administering AZT to 239 pregnant women and placebo to 238. The endpoint was time-to-diagnosis of HIV in the newborn. The results were compelling: 60 newborns were HIV-positive in the placebo group and only 20 in the AZT group.

Yao and Wei [11] performed an extensive simulation study to see hypothetically how an adaptive design with delayed response would have affected power. They found significant reductions in treatment failures with minimal loss of power [11]. A more simplistic analysis showed that, if the response had been immediate and binary, an RPW rule would have skewed the allocation to 25:75 in favor of AZT, with a reduction of seven treatment failures [3]. To maintain similar power, the adaptive design would have had to recruit an additional 18% patients, yet the reduction in treatment failures would still have been substantial.

Exact Inference

We can write a computer networking algorithm to find the exact distribution of statistics like equation (1). Such an unconditional statistic is then a function of the three sufficient statistics (S_A, S_B, N_A) and is likely to be more powerful than a statistic that conditions on N_A . With the advent of parallel processing, we have found that exact computations for $N = 100$ or more are possible. We can also ascertain exact power, which is useful in planning studies. Wei et al. [4] wrote similar code to find exact unconditional confidence intervals for the simple difference and odds ratio, and applied it for $n = 12$.

For multicenter clinical trials, a test of homogeneity of treatment effects (e.g., equal odds ratios) among clinics is desirable. An exact test for $K 2 \times 2$ tables was recently programmed for the RPW rule [24].

Programming exact inference requires extensive computational complexity. As another option, investigators can simulate the distribution of the test statistics much more easily, and with high-speed computers, thousands of replications are possible, making results very accurate. An advantage to simulations is that they can incorporate delayed-response mechanisms.

The Fluoxetine Trial (1994)

In one analysis, Tamura et al. [7] discussed finding the exact joint distribution of four quantities under a randomization model: time-to-surrogate response, response to primary outcome, surrogate response, and HAMD₁₇ change. They then computed equation (1) from this joint distribution. Because the number of patients was too large for them to do this efficiently, they simulated the joint distribution assuming a fixed delay in surrogate response [7].

Sequential Tests

Little in the literature has addressed adaptive designs and sequential testing, and yet it is standard for clinical trials to be monitored for efficacy, and investigators often apply early stopping rules. Although the RPW rule may be most appropriate for short-term trials, investigators incorporate early stopping rules to ensure that no more patients than necessary be assigned to inferior treatments. A recent paper [22] simulates a triangular test with Christmas tree correction [25] using RPW allocation. The authors find that power seems to be unaffected by the allocation rule and that the RPW rule allows fewer treatment failures, on average, than equal allocation. One can find a more theoretical treatment of group sequential monitoring with the RPW rule in Bandyopadhyay and Biswas [26], in which the authors derive and tabulate operating characteristics.

Bootstrap Confidence Intervals

A recent paper discusses simple bootstrap confidence intervals for p_A and p_B under RPW allocation [10]. Obtaining our observed estimates p_A and p_B from the data, we simulate the RPW rule K times, using p_A and p_B as the underlying response probabilities. Then we compute the K estimates, $\hat{p}_{A1}^*, \dots, \hat{p}_{AK}^*$ and

Table 4 Conditions Under Which the RPW Rule Is Reasonable

-
- The therapies have been evaluated previously for toxicity.
 - Response is binary.
 - Delay in response is moderate, allowing adapting to take place.
 - Sample sizes are moderate (at least 50 subjects).
 - Duration of the trial is limited and recruitment can take place during the entire trial.
 - The trial is carefully planned with extensive computations done under different models and initial urn compositions.
 - The experimental therapy is expected to have significant benefits to public health if it proves effective.
-

$\hat{p}_{b1}^*, \dots, \hat{p}_{bK}^*$ from the K replications of the simulation. We order these, and take the $K\alpha/2$ th and $K(1 - \alpha)/2$ th ones as the upper and lower bounds of the interval on p_A and p_B , respectively. Simulations show that these confidence intervals have nearly the perfect coverage of $100(1 - \alpha)\%$. Similarly, one can compute bootstrap confidence intervals on functions of p_A and p_B , such as the odds ratio, by ordering the simulated measures.

One advantage of these bootstrap confidence intervals is that they are based on simulation, so that one can incorporate prestratification, staggered entry, censoring, and delayed response. Exact inference cannot deal with these mechanisms.

The Fluoxetine Trial (1994)

In the fluoxetine trial, a delay in surrogate response was observed and differed with respect to the response (but not with respect to stratum or treatment group). Time-to-surrogate response was approximately normal with mean 43 days and variance 122 days for responders and approximately uniform on the interval (20 days, 75 days) for nonresponders [10]. We incorporated stochastic delay mechanisms into the analysis, using the bootstrap procedure above on the surrogate outcome. For confidence intervals within each stratum, see Rosenberger and Hu [10].

CONCLUSIONS

Our first example, the ECMO trial, was atypical of adaptive designs and should not constitute a reason to neglect adaptive designs in modern clinical trials. A more realistic example of a modern multicenter clinical trial is the fluoxetine trial, example 2. The principal importance of the fluoxetine trial was logistical. The Lilly team has demonstrated quite effectively that investigators can employ suitable information technology for adaptive randomization, and that clinicians and subjects can be willing and enthusiastic to undertake such a trial.

Palmer and Rosenberger discuss ethical conditions under which adaptive designs generally are appropriate [27]. This paper has addressed more logistical and statistical considerations. Table 4 summarizes conditions under which the RPW rule is a reasonable alternative to a standard clinical trial design.

The RPW rule is a viable alternative to equal allocation in clinical trials with binary response, provided that statisticians plan very thoughtfully and intensively at the design phase. Such planning should involve theoretical questions, such as exploring the variability of the design by exact computations or else by simulation if there is delayed response. If extensive data on safety are not available before the trial, investigators should consider starting with equal allocation until safety data are available. This should mitigate the potential for assigning more patients to a treatment that is more efficacious but also more toxic. In any case, starting the urn with more than one ball of each color should lessen the probability of a success run of the type seen in the ECMO trial. More than likely, most RPW trials would be short-term; most longer-term trials would involve time-to-response as an endpoint. Clearly, however, we can minimize variability in the urn process by using sufficiently large samples.

If at all possible, investigators should carry out exact unconditional inference or simulate the exact distribution, or employ a bootstrap method if confidence intervals are desired. If response is delayed, one can incorporate the delayed response-mechanism via simulation.

Professor Rosenberger's research is supported by grant R29-51017-04 from the National Institute of Diabetes and Digestive and Kidney Diseases. Deepak Agarwal and Padmanabhan Seshaiyer programmed the simulations and exact computations; they were also supported by the grant.

REFERENCES

1. Wei LJ, Durham SD. The randomized play-the-winner rule in medical trials. *J Am Stat Assoc* 1978;73:840-843.
2. Flournoy N, Rosenberger WF (eds). *Adaptive Designs*. Hayward, Calif: Institute of Mathematical Statistics; 1995.
3. Rosenberger WF. New directions in adaptive designs. *Stat Sci* 1996;11:137-149.
4. Bartlett RH, Roloff DW, Cornell RG, et al. Extracorporeal circulation in neonatal respiratory failure: A prospective randomized study. *Pediatrics* 1985;76:479-487.
5. Royall RM. Ethics and statistics in randomized clinical trial. *Stat Sci* 1991;6:52-62.
6. Rosenberger WF, Lachin JM. The use of response-adaptive designs in clinical trials. *Control Clin Trials* 1993;14:471-484.
7. Tamura RN, Faries DE, Andersen JS, et al. A case study of an adaptive clinical trial in the treatment of out-patients with depressive disorder. *J Am Stat Assoc* 1994; 89:768-776.
8. Rout CC, Rocke DA, Levin J, et al. A reevaluation of the role of crystalloid preload in the prevention of hypotension associated with spinal anesthesia for elective Cesarean section. *Anesthesiology* 1993;79:262-269.
9. UK Collaborative ECMO Trial Group. UK Collaborative randomized trial of neonatal extracorporeal membrane oxygenation. *Lancet* 1996;348:75-82.
10. Rosenberger WF, Hu F. Bootstrap methods for adaptive designs. *Stat Med* (in press) 1999.
11. Yao Q, Wei LJ. Play the winner for phase II/III clinical trials. *Stat Med* 1996; 15:2413-2423.
12. Hallstrom A, Brooks MM, Peckova M. Logrank, play the winner, power and ethics. *Stat Med* 1996;15:2135-2142.
13. Rosenberger WF, Seshaiyer P. Adaptive survival trials. *J Biopharm Stat* 1997; 7:617-624.

14. Wei LJ, Smythe RT, Lin DY, et al. Statistical inference with data-dependent treatment allocation rules. *J Am Stat Assoc* 1990;85:156–162.
15. Rosenberger WF, Sriram TN. Estimation for an adaptive allocation design. *J Stat Planning Inf* 1997;59:309–319.
16. Coad DS. Sequential tests for an unstable response variable. *Biometrika* 1991; 78:113–121.
17. Coad DS. A comparative study of some data-dependent allocation rules for Bernoulli data. *J Stat Comput Simulation* 1992;40:219–231.
18. Athreya KB, Karlin S. Limit theorems for the split times of branching processes. *J Math Mech* 1967;17:257–277.
19. Melfi V, Page C. Variability in adaptive designs for estimation of success probabilities. In: Flournoy N, Rosenberger WF, Wong WK, eds. *New Developments and Applications of Experimental Designs*. Hayward, Calif: Institute of Mathematical Statistics; 1998.
20. Matthews PC, Rosenberger WF. Variance in randomized play-the-winner clinical trials. *Stat Probab Lett* 1997;35:233–240.
21. Lachin JM. Statistical properties of randomization in clinical trials. *Control Clin Trials* 1988;9:289–311.
22. Coad DS, Rosenberger WF. A comparison of the randomised play-the-winner rule and the triangular test for clinical trials with binary responses. *Stat Med* 1999; 18:761–769.
23. Connor EM, Sperling RS, Gelber R, et al. Reduction of maternal-infant transmission of human immunodeficiency virus type 1 with zidovudine treatment. *New Engl J Med* 1994;331:1173–1180.
24. Chen L, Tamura RN. Analysis of homogeneity of treatment effect in adaptive multicenter clinical trials. *J Biopharm Stat* 1998;8:55–68.
25. Whitehead J. *The Design and Analysis of Sequential Clinical Trials*. Chichester: Wiley; 1997.
26. Bandyopadhyay U, Biswas A. Some sequential tests in clinical trials based on randomized play-the-winner rule. *Calcutta Stat Assoc Bull* 1997;47:67–89.
27. Palmer CR, Rosenberger WF. Ethics and practice: Alternative designs for phase III randomized clinical trials. *Control Clin Trials* 1999;20:172–186.