# A menu-driven facility for complex sample size calculation in randomized controlled trials with a survival or a binary outcome

Patrick Royston
MRC Clinical Trials Unit, London, UK
patrick.royston@ctu.mrc.ac.uk

Abdel Babiker
MRC Clinical Trials Unit, London, UK
abdel.babiker@ctu.mrc.ac.uk

**Abstract.** We present a menu-driven Stata program for the calculation of sample size or power for complex clinical trials with a survival time or a binary outcome. The features supported include up to six treatment arms, an arbitrary time-to-event distribution, fixed or time-varying hazard ratios, unequal patient allocation, loss to follow-up, staggered patient entry, and crossover of patients from their allocated treatment to an alternative treatment. The computations of sample size and power are based on the logrank test and are done according to the asymptotic distribution of the logrank test statistic, adjusted appropriately for the design features.

**Keywords:** st0013, randomized controlled trials, survival analysis, logrank test, experimental design

## 1 Introduction

Stata includes just one program (`sampsi`) for calculating sample size in randomized controlled studies. It deals only with comparisons between two groups in terms of binary or Normally distributed outcome variables. Many such trials, however, are designed around a survival-time outcome measure such as the time to death, to disease progression, or to healing of a lesion. The trials may compare more than two groups and are subject to loss to follow-up, withdrawal from allocated treatment, and staggered entry. The purpose of the present article and software is to provide a very flexible tool for determining sample size in such studies. Because inevitably there are many potential "options" (in the Stata sense), a conventional ado-file may be dauntingly complex. For this reason, we have provided a menu-driven front end for sample-size calculation. The `Study size` menu is initiated by entering `ssmenu on` in the Stata command window.

For survival-time outcomes, the basic assumption behind the calculations is that the groups will be compared by using the logrank test (Peto and Peto 1972); for example, using the `sts test` command for `st` data in Stata. The logrank test is based on the comparison between observed and (conditionally) expected numbers of events. Power is greatest under the assumption of proportional hazards between the groups. In our implementation, therefore, the null and alternative hypotheses are set up in terms of hazard ratios (HRs). The following study design features have been included:

- Up to six treatment groups.

- Arbitrary baseline time-to-event distribution.

- Time-varying hazard ratios (i.e., nonproportional hazards).

- Arbitrary allocation ratios across groups.

- Loss to follow-up.

- Staggered patient entry.

- Crossover from allocated treatment to alternative treatment.

- Survival analysis by unweighted, Tarone–Ware or Harrington–Fleming versions of the logrank test (Tarone and Ware 1977; Harrington and Fleming 1982). The weights used in the last two versions of the test are the square root of the total number at risk (Tarone–Ware) and the estimated overall survivor function raised to power $I$ (default is $I = 1$) (Harrington–Fleming) at each failure time.

- In addition, two flavors of $\chi^2$ test (unconditional and conditional on the total number of events) are available to compare the proportions of failures at the end of the study. Complex trials with a binary outcome, involving possible loss to follow-up, staggered entry, and treatment crossover, may thereby be designed with the aid of the software.

- Trend test, with specified doses if required.

Sample-size calculation for the comparison between two survival distributions using the logrank test was described in the simplest case of two groups and proportional hazards by Freedman (1982) and Schoenfeld (1982). Extension to more than two groups was given by Ahnn and Anderson (1995). A method incorporating loss to follow-up, staggered entry, and treatment crossover was proposed for two-group designs by Lakatos (1988) and extended to more than two groups by Ahnn and Anderson (1998). The latter methodologies are based on a nonstationary Markov model, which allows for an arbitrary number of switches between treatments. In our approach, only one treatment switch is permitted. We believe that this is quite sufficient in practice and may be supported by certain medical arguments about the effects of treatment. In addition, it allows direct calculation of the expected failure-time distribution adjusted for loss to follow-up, staggered entry, and treatment crossover.

## 2   Design of menu and dialogs

All the features are available from the `Study size` menu and its associated dialogs. When the computations are complete, Stata displays, in the Review window, the command line that generated the results. The computations are performed by an ado-file called `calcssi`, which is provided as part of this article. By recalling the command

from the Review window, editing, and re-executing it, the menu system may also be used as a tutor for the command-driven approach using `calcssi`.

We provide an item within the `Study size` menu for computing sample size for simpler studies with a binary outcome. It extends the facilities available with Stata's `sampsi` command and relies on an ado-file called `calcssbi`, also provided with this article.

When `ssmenu` has been executed using `ssmenu on`, a new item `Study size` appears on the system menu-bar. The menu is turned off by entering `ssmenu off`. `Study size` contains the following four items:

| | |
|---|---|
| `Survival - Basic setup` | Sets up basic design parameters (number of groups etc.) |
| `Survival - Advanced options` | Accesses the more complicated design options |
| `Survival - Compute` | Sets up more design parameters and runs the calculations |
| `Binary outcomes` | For trials with a binary outcome and a simple design |

The first three items deal with the design of potentially complex trials with a survival-time outcome. The fourth is for simpler trials with a binary outcome and is independent of the first three. We will describe these items in turn.

## 2.1 Survival - Basic setup

Having clicked on the `Survival - Basic setup` menu item, you enter the number of groups, the number of periods, and the baseline failure (or survival) distribution. The number of groups must be between two and six. The number of periods means the duration of the trial, in arbitrary time units, from entry of the first patient to the time the analysis is presumed to be carried out. Each period is of length one time unit, and the time units are chosen such that rates of events (failure, withdrawal, loss to follow-up, and entry into the trial) are approximately constant within each period. In the simplest case, you need only enter 1 (the default) as the number of periods. The baseline failure distribution is entered as cumulative probabilities of failure either at the end of each period (if there is more than one period) or only at the end of the last period (however many periods there are). By pressing the appropriate radio button on the next line, you can tell the program that you are entering survival probabilities instead of failure probabilities. These are, of course, just one minus the failure probabilities.

The format for entering the failure distribution is either a single value such as 0.5, representing the cumulative probability at the end of the last period, or values preceded by `p<`*period_number*`>=`, separated by spaces, representing the cumulative probability at the end of each period. You must enter the values in one or the other of these two ways; they may not be mixed. For example, suppose there were four periods, and the cumulative failure probabilities were 0.1, 0.2, 0.3, and 0.35. You would enter these as

`p1=0.1 p2=0.2 p3=0.3 p4=0.35`. If you give values for some of the periods only (e.g., `p1=0.1 p3=0.3`), the program will interpolate and extrapolate failure probabilities for the other periods, assuming piecewise exponential distributions with a constant hazard in each period. If just one period is specified, a constant hazard is assumed so that the failure distribution is exponential.

The default test is a global test of difference between the groups based on a $\chi^2$ distribution with $k-1$ degrees of freedom, where $k$ is the number of treatment groups. A trend test on 1 degree of freedom is also available. You can check the box for `Trend` if you require a trend test. If you have dose levels for the groups, these should be entered in the `Dose` box. If you do not specify dose levels, a linear trend test will be performed, equivalent to specifying doses $1, 2, \ldots, k$.

The default method of analysis in the `Method of sample size calculation` list-box is the unweighted logrank test. The two weighted versions of the test (Tarone–Ware and Harrington–Fleming) are usually used when nonproportional hazards are anticipated. They place different weights on different portions of the survival distribution.

Complex trials with a binary outcome are supported by another method of analysis (binomial, conditional, and unconditional on total events), also available on the `Method of sample size calculation` list-box. The groups are compared in terms of the expected proportions of failures at the end of the study. For consistency with the survival analysis paradigm, these overall event probabilities are specified indirectly in terms of the baseline failure or survival distribution and the hazard ratios for the groups, rather than directly (and more familiarly) as the event probabilities for the groups.

## 2.2   Survival - Advanced options

The `Survival - Advanced options` menu item provides the more complicated design features mentioned in the *Introduction*. For many trial designs, you will not need this item at all. Cumulative probabilities of loss to follow-up may be entered for any group in the trial in the same format as for the failure probabilities above; that is, a single value for the end of the trial *or* values preceded by `p<`*period_number*`>=` for any subset of periods. Likewise, a cumulative probability of withdrawal from allocated treatment is required if you wish to make use of this feature. (We will give an example of this later.) You will then need to specify either the group(s) to which the patients transferred (the so-called `target group on crossover`) or the hazard ratio of failure following withdrawal, compared with the baseline hazard. While the latter is the more flexible option, the former is more likely to be used in practice. For example, control-arm patients who show signs of disease progression or other treatment failure some way into the trial may be transferred to the experimental regimen received by Group 2, and therefore take on the hazard ratio expected in Group 2. This can dilute the expected treatment effect considerably and so increase the sample size.

## 2.3 Survival - Compute

In addition to `Survival - Basic setup`, the `Survival - Compute` menu requires some design parameters to be specified. Most important of these are the hazard ratios in relation to the baseline hazard function implied by (and internally computed from) the baseline failure-time distribution. HRs are entered separately for each group. Typically, group 1 will be the control arm and will be assumed to have an HR of 1. The experimental arm(s) will have HR(s) less than 1, representing an anticipated improvement in survival. If it is known that the hazard ratio is time-varying, time-specific values (HR function) may be given. These are entered in order for each period. If too few HR values are entered for any group, the values during subsequent periods are assumed to equal the last value entered for that group. HR functions must be entered for at least two groups. The default for the remaining groups is the geometric mean of the specified HR functions.

Allocation ratios (weights) are by default taken as equal across the groups. Unequal allocation ratios are entered in an obvious way; e.g., `2 1 1` would assume twice as many patients in group 1 as in groups 2 and 3.

In reality, patients enter trials over time and are not all available at the beginning of the study. You indicate the recruitment duration by filling the `Duration` box with the number of periods it takes to recruit. By default this is 0, meaning (unrealistically) that all patients are recruited at the beginning. The default is to assume steady recruitment (uniform distribution), and this may be changed to a negative exponential distribution with a rate that you enter in the `Shape of distribution` box. `Period weights` refers to the proportion of patients recruited during each period and is by default 1, meaning equal numbers of patients in each period.

The usual two-sided Type I error probability is entered in the `Alpha` box, and the power in the `Power` box. Alternatively, you can specify the sample size (i.e., the total number of patients in all groups) and get the power by pressing the `Specify sample size` radio button.

The `Detailed output` check box gives further information when the program is run by pressing the `Compute` button.

### Technical note: Computation

The computations of sample size and power are based on the logrank test and are done according to the asymptotic distribution of the logrank test statistic, $Q$. $Q$ is defined as $U'V^{-1}U$, where $U$ is the vector of the total observed minus expected number of events in each of the $k$ groups in the design except for the first, and $V$ is the covariance matrix of $U$. The distribution of $Q$ differs under the null and alternative hypotheses. The null hypothesis, $H_0$, is that the survival distributions are identical in the $k$ groups. The alternative hypothesis, $H_1$, is that at least one distribution is different from the others. $Q$ is distributed asymptotically as $\chi^2$ on $k-1$ degrees of freedom under $H_0$ and as noncentral $\chi^2$ on $k-1$ degrees of freedom under local $H_1$. Loosely speaking, "local $H_1$" implies that the hazard ratios between treatment groups are not far from 1. The

noncentrality parameter equals $E\left(U'\right)V^{-1}E\left(U\right)$, the expectation being taken under $H_1$. Under distant $H_1$, the distribution of $Q$ may be approximated by a scaled, noncentral $\chi^2$ on $k-1$ degrees of freedom. The scaling factor and the noncentrality parameter may be determined by the method of moments from the asymptotic mean and variance of $Q$ under distant $H_1$. The test statistic for trend when dose levels are specified is also based on $U$, and is asymptotically normally distributed under $H_0$ and $H_1$ with known mean and variance. Computation is required of the survival functions in each of the $k$ groups allowing for loss to follow-up, treatment crossover, etc. Although lengthy and tedious, the calculations are algebraically fairly straightforward. A report on the methodology is being prepared and will be published in due course. In the meantime, mathematical details are available from the authors on request. The `calcssi` Stata software has been in use in real applications for several years. In the two-group case, `calcssi` gives results that may be up to about 5% higher than those from the method of Lakatos (1988), as implemented in the SAS program SIZE by Shih (1995).

## 2.4   Binary outcomes

The `Binary outcomes` menu provides for the comparison of up to six groups with proportions of events as the outcome. It allows for comparison of more than two groups, for differing allocation ratios, and for a dose/response relationship via a trend test, as with the survival menus. However, no allowance is made for loss to follow-up, crossover, or staggered entry into the study. (These possibilities are provided within the survival analysis framework by the two `binomial` options on the `Method of sample size calculation` list-box in the `Survival - Basic setup` dialog, followed by the features of the `Survival - Advanced options` and `Survival - Compute` dialogs.) The default test does not condition on the total number of observed events. A conditional test which uses the Peto approximation to the odds ratio is available as an option. This approximation is very adequate with small to moderate treatment effect (odds ratio between 0.5 and 2.0), but underestimates the sample size with larger differences between the groups. Unlike with `sampsi`, no continuity correction is included in the computations; therefore, the sample sizes will always be somewhat smaller than those obtained using `sampsi`. The statistical jury is out on whether it is preferable to use a continuity correction for such calculations or not. We chose not to do so.

### Technical note: Choice of test to compare binary outcomes

We will not attempt a review of the rather extensive literature for sample size calculation in the two-group binary outcomes study, but will restrict ourselves to a few comments. A recent review is given by Sahai and Khurshid (1996). Many software program authors choose as their default option (or indeed perhaps their only option) the unconditional test with nonlocal alternatives and *without* continuity correction. A Normal approximation to the binomial distribution is assumed. The addition of a continuity correction may considerably increase the study size and has been criticized as unnecessarily conservative (i.e., potentially wasteful of resources, particularly in small

studies where cost per patient is in some sense high). Our default tends slightly in the direction of conservatism, as may be seen in the examples of the previous subsection, but will give always smaller samples than `sampsi` does.

It may be helpful to bear in mind the principle of being willing to analyze the data according to the same test that was used in the sample-size calculation. For some people, this will rule out use of the conditional test (Peto approximation), which assumes that the total number of events is fixed.

The local alternative option uses the same variance for the observed proportion minus the expected proportion under the null hypothesis and the alternative hypothesis. There is usually little difference between the results with local and nonlocal alternatives. For technical reasons, we prefer the approach with local alternatives because the mathematics extends naturally from two groups to many groups. For practical reasons (rather than for any methodological limitations), our software allows up to six groups. Such extension to many groups is not straightforward with nonlocal alternatives, although a reasonable approximation is used in the software.

In summary, unless there are strong reasons and a clear rationale to do otherwise, we believe that our default choice of test for sample size and power calculations in the two-group and the multi-group situation is satisfactory.

# 3  Examples

## 3.1  A basic survival study design

We will give a hypothetical example with design parameters that are typical of a trial in some types of cancer. We will start with the basic situation, and then elaborate the design to illustrate some of the more complicated options.

There are two arms, control and experimental. From existing data, we know that the two-year survival rate of control-arm patients is approximately 20%. We are hoping that the experimental treatment will improve survival with a hazard ratio of 0.7, representing a two-year survival probability of about 32%.

Figure 1 shows the completed `Basic setup` dialog for the trial. The only value that had to be entered was 0.8 for the baseline cumulative probability of failure at the end of the study (two years).
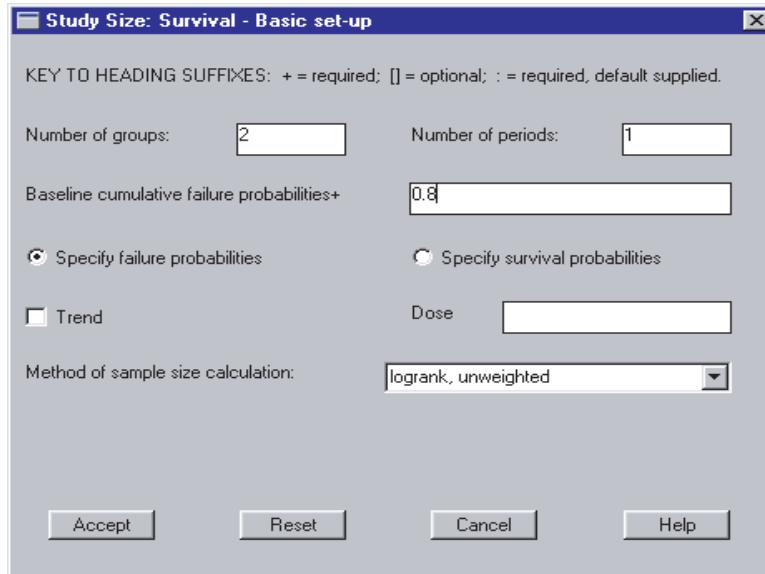
Figure 1: A completed Basic setup screen.

The second menu that must be used is `Survival - Compute`. We fill in the hazard ratios for groups 1 and 2, the recruitment duration as 1, and the power as 0.9 (see Figure 2).
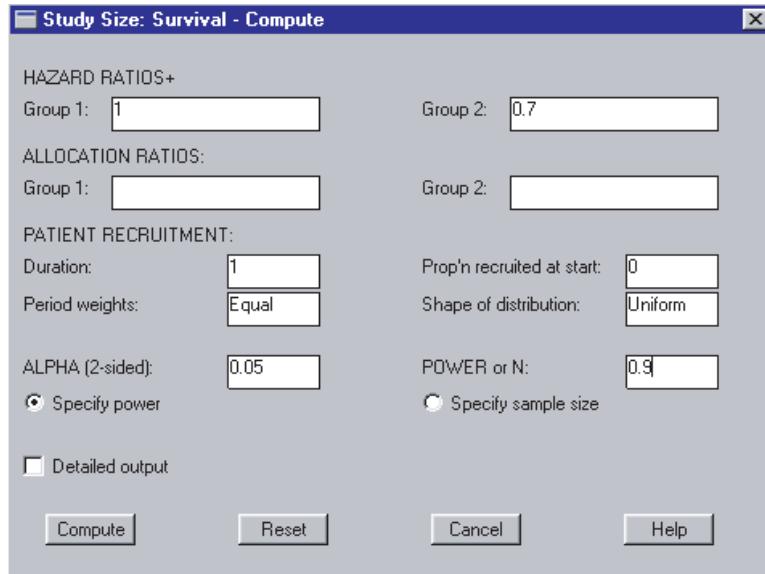


Figure 2: A Survival Computation screen.

On pressing the `Compute` button, the following results are obtained:

```
        Sample size: 2-group comparison
        Unweighted logrank test
        ─────────────────────────────────────────────────
        Allocation ratio:                    Equal group sizes
        Two-sided alpha =                    0.050
        Power =                              0.900

        Total sample size =                       736
        Expected total number of events =         333
```

The total number of patients required is 736 (368 per group). The total number of events is 333.

To recapitulate: we recruit patients at a uniform rate for 2 years (= 1 period), and then we analyze the data at that point. In reality, we would be more likely to follow the patients up for some time before analysis, both for practical reasons of trial management and to accumulate more events. Suppose we were to recruit for 2 years, as before, but then to follow-up for 1 year and to analyze the data at that point, 3 years after initiation of the study. Now, the natural time period for the sample-size calculation is 1 year rather than two. We would specify this design by requiring 3 periods with a recruitment duration of 2 periods (i.e., 2 years). Also, we would have to specify the cumulative failure rate to be 0.8 at 2 years by entering `p2=0.8`, rather than leaving the value as 0.8 which would be interpreted incorrectly as the value at the end of the study (now 3 years). The instantaneous event rate (hazard) in period 3 would be taken to be the same as in the latest specified period, here period 2. For examples of how the hazard is computed from the specified failure probabilities, see the description of the `edf0()` option in the help file for `calcssi`.

On pressing `Compute`, the number of events is found to be essentially unchanged (331 versus 333 before), but the required number of patients is reduced from 736 to 461, a 37% saving. If the duration were increased to 4 years by following up for 2 years, then the sample size would be further reduced to 389. We are trading length of follow-up for number of patients. To obtain a given number of events (dictated by the baseline failure probabilities and the hazard ratio), we can have fewer patients followed up over a longer time or more patients over a shorter time.

## 3.2   A more complicated design

We will now assume that the study is carried out over 4 years (four one-year periods) with recruitment over 2 years, and that the hazard ratio for the experimental treatment arm compared with control varies over time. We will take the HRs to be 0.5, 0.65, 0.8, and 0.9 in periods 1 to 4, respectively. The interpretation is that the new treatment reduces the mortality rate quite markedly in the period immediately following its first administration, but that its efficacy decreases fairly rapidly over time, so that by 4 years it is little better than control.

Before running this example, note that using the last-mentioned setup with recruitment over 2 years, follow-up for 2 years, and a constant HR of 0.7, but otherwise the same parameters, gives $n = 389$ with 332 events.

In the `Survival - Basic setup` dialog, we enter the number of periods to be 4 and the cumulative failure probability at period 2 to be `p2=0.8`. In the `Survival - Compute` dialog, we enter the recruitment duration to be 2, the hazard ratio for group 1 to be `1`, and the hazard ratios for group 2 to be `0.5 0.65 0.8 0.9`. On pressing the `Compute` button, the resulting sample size is 178 with 150 events, about half the value with a constant HR. The reason why the number of events has gone down is that most of the deaths occur in the first period where the hazard ratio is 0.5, quite a lot lower than the constant HR of 0.7 we assumed before. This is a more extreme HR, and hence the required number of events is reduced.

We may further complicate the situation by assuming that a proportion of patients crossover from the control to the experimental therapy arm. Suppose that all patients receive the protocol treatment in the first one-year period. However, 25% of them relapse in the second period and are given the new treatment, and the same happens to some of the survivors in the third and fourth period, perhaps as attempted 'rescue therapy', for example. The cumulative proportions of patients given the new treatment by the end of these periods are now 0.35 and 0.45, respectively. We specify the design by entering the cumulative distribution of the time to crossover as `p1=0 p2=0.25 p3=0.35 p4=0.45`, and the target group for crossover from group 1 to be 2, with no crossover from group 2. These options are entered in the `Survival - Advanced options` dialog, which is shown in Figure 3.
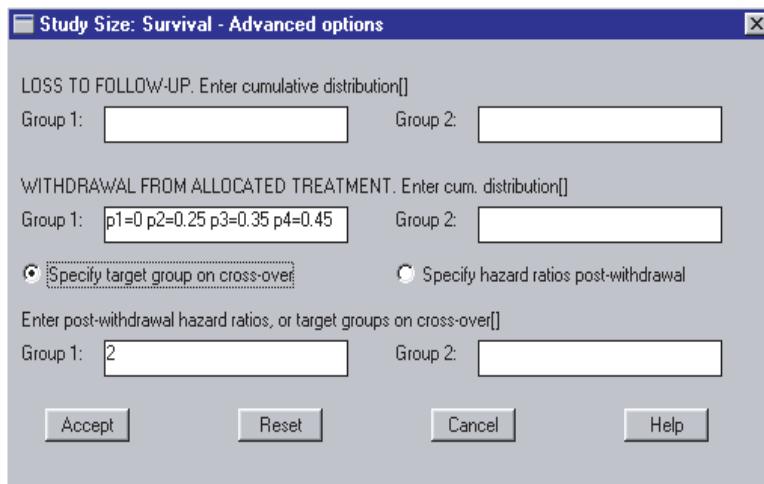


Figure 3: An Advanced option screen.

No data are entered for Group 2 since no crossover is expected. This results in $n = 217$ with 181 events, an increase from 178 with 150 events because the treatment difference between the groups is diluted by the crossover.

Note that the increase in sample size is modest in this example because of the very high event rate. Many events occur early in the trial when little or no crossover has occurred. If the two-year cumulative event rate was `p2=0.4` rather than `p2=0.8`, we would obtain $n = 442$ (204 events) without crossover but $n = 607$ (442 events) with the above crossover specification, a substantial difference.

## 3.3   A study with a binary outcome

We now show results for a simple study with a binary outcome. Suppose there are two groups of equal size, and the hypothesized proportions of an event are 0.2 and 0.3. Figure 4 shows the necessary setup using the `Binary outcomes` dialog.
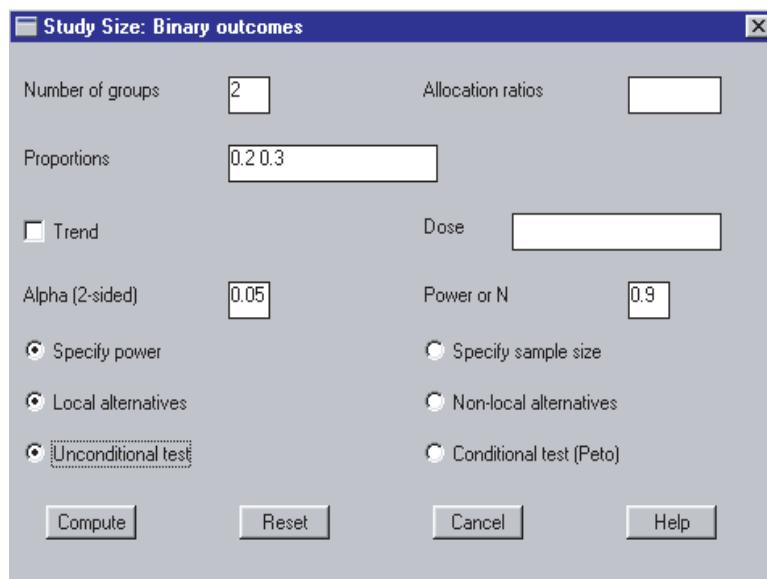


Figure 4: A completed Basic setup screen.

The result of pressing `Compute` with this setup is shown below:

```
          Sample size: 2-groups comparison
       Unconditional comparison of 2 binomial proportions
       ─────────────────────────────────────────────────

Anticipated event probabilities:           0.200, 0.300
Allocation ratios:                          Equal group sizes
Two-sided alpha =                           0.050
Power =                                     0.900

Total sample size =                             789
Expected total number of events =               198
```

A total of 789 subjects are needed. Specifying nonlocal alternatives (by clicking the appropriate radio button) gives a small reduction to 782, whereas specifying the conditional test (Peto) gives a more substantial reduction to 771.

For comparison, `sampsi` gives the following:

```
. sampsi 0.2 0.3, power(0.9)
Estimated sample size for two-sample comparison of proportions
Test Ho: p1 = p2, where p1 is the proportion in population 1
                    and p2 is the proportion in population 2
Assumptions:
          alpha =   0.0500  (two-sided)
          power =   0.9000
             p1 =   0.2000
             p2 =   0.3000
          n2/n1 =   1.00
Estimated required sample sizes:
             n1 =       412
             n2 =       412
```

`sampsi` requires a total of 824 subjects, some 4% more than our default option of local alternatives and the unconditional test. This is due to `sampsi`'s use of an unconditional test with additional application of a continuity correction to the $\chi^2$ statistic. If we apply a continuity correction to our requirement for 789 subjects using the formula given under [R] **sampsi** in the *Stata Reference Manual*, we obtain 414 subjects per group, very similar to `sampsi`.

# 4    Conclusion

It is important that the design of a randomized controlled clinical trial be realistic, allowing for such factors as a tendency for patients' treatment to be switched if the initial treatment is not successful, for possible loss to follow-up, and for the decline of a treatment effect over time, expressible as nonproportional hazards. Basic considerations such as staggered patient entry are also very important and may influence the overall duration of a trial to a major extent. The software provided here should allow researchers to accommodate all of these basic and more advanced design features in a straightforward manner.

# 5    References

Ahnn, S. and S. J. Anderson. 1995. Sample size determination for comparing more than two survival distributions. *Statistics in Medicine* 14: 2273–2282.

——. 1998. Sample size determination in complex clinical trials comparing more than two groups for survival endpoints. *Statistics in Medicine* 17: 2525–2534.

Freedman, L. S. 1982. Tables of the number of patients required in clinical trials using the log-rank test. *Statistics in Medicine* 1: 121–129.

Harrington, D. P. and T. R. Fleming. 1982. A class of rank test procedures for censored survival data. *Biometrika* 69: 553–566.

Lakatos, E. 1988. Sample sizes based on the log-rank statistic in complex clinical trials. *Biometrics* 44: 229–241.

Peto, R. and J. Peto. 1972. Asymptotically efficient rank invariant test procedures (with discussion). *Journal of the Royal Statistical Society, Series B* 135: 185–206.

Sahai, H. and A. Khurshid. 1996. Formulae and tables for the determination of sample sizes and power in clinical trials for testing differences in proportions for the two-sample design: a review. *Statistics in Medicine* 15: 1–21.

Schoenfeld, D. 1982. The asymptotic properties of nonparametric tests for comparing survival distributions. *Biometrika* 68: 316–319.

Shih, J. H. 1995. Sample size calculation for complex clinical trials with survival endpoints. *Controlled Clinical Trials* 16: 395–407.

Tarone, R. E. and J. H. Ware. 1977. On distribution-free tests for equality of survival distributions. *Biometrika* 64: 156–160.

**About the Authors**

Patrick Royston is a medical statistician of 25 years of experience, with a strong interest in biostatistical methodology and in statistical computing and algorithms. At present he works in clinical trials and related research issues in cancer. Currently he is focusing on problems of model building and validation with survival data, including prognostic factors studies, on parametric modeling of survival data and on novel trial designs.

Abdel Babiker is head of the HIV division of the MRC Clinical Trials Unit in London, UK.