

SOFTWARE

Open Access



A web application for the design of multi-arm clinical trials

Michael J. Grayling^{1*}  and James MS. Wason^{1,2}

Abstract

Background: Multi-arm designs provide an effective means of evaluating several treatments within the same clinical trial. Given the large number of treatments now available for testing in many disease areas, it has been argued that their utilisation should increase. However, for any given clinical trial there are numerous possible multi-arm designs that could be used, and choosing between them can be a difficult task. This task is complicated further by a lack of available easy-to-use software for designing multi-arm trials.

Results: To aid the wider implementation of multi-arm clinical trial designs, we have developed a web application for sample size calculation when using a variety of popular multiple comparison corrections. Furthermore, the application supports sample size calculation to control several varieties of power, as well as the determination of optimised arm-wise allocation ratios. It is built using the Shiny package in the R programming language, is free to access on any device with an internet browser, and requires no programming knowledge to use. It incorporates a variety of features to make it easier to use, including help boxes and warning messages. Using design parameters motivated by a recently completed phase II oncology trial, we demonstrate that the application can effectively determine and evaluate complex multi-arm trial designs.

Conclusions: The application provides the core information required by statisticians and clinicians to review the operating characteristics of a chosen multi-arm clinical trial design. The range of designs supported by the application is broader than other currently available software solutions. Its primary limitation, particularly from a regulatory agency point of view, is its lack of validation. However, we present an approach to efficiently confirming its results via simulation.

Keywords: False discovery rate, Familywise error-rate, Multiple comparisons, Optimal design, Power, Sample size

Background

Drug development is becoming an increasingly expensive process, with the estimated average cost per approved new compound now standing at over \$1 bn [1]. In no small part this is due to the high failure rate of clinical trials, in particular in phases II and III. This is particularly true in the field of oncology, where the likelihood of approval from phase I is only 5.1% [2]. Consequently, the clinical research community is constantly seeking new methods that may improve the efficiency of the drug development process.

One possible method, which has received substantial attention in recent years, is the idea to make use

of multi-arm designs that compare several experimental treatments to a shared control group. Several desirable, inter-related, features of such designs have now been described. For example, the number of patients on the control treatment is typically reduced compared to conducting separate two-arm trials, and simultaneously patients are more likely to be randomized to an experimental treatment, which may help with recruitment [3, 4]. Furthermore, the overall required sample size, for the same level of power, will typically be smaller than that which would be required if multiple two-arm trials were conducted [5]. Finally, multi-arm designs offer a fair head-to-head comparison of experimental treatments in the same study [3, 4], and the cost of assessing a treatment in a multi-arm trial is often around half of that for a separate two-arm trial [3].

*Correspondence: michael.grayling@newcastle.ac.uk

¹Population Health Sciences Institute, NE2 4AX Newcastle, UK

Full list of author information is available at the end of the article



Based upon these advantages, and their experiences of utilising such designs in several oncology trials, Parmar et al. [3] make a compelling case for the need for more multi-arm designs to be used in clinical research. We are not aware of any systematic evidence on whether this has now permeated through to practice, but a simple search of PubMed Central suggests it may be the case: 859 articles have included the phrases “multi-arm” and “clinical trial” since 2015, as opposed to just 273 in all years prior to this. Considering this result in combination with the findings of Baron et al. [6], who determined 17.9% of trials published in 2009 were multi-arm, as well as the recent publication of a key guidance document on reporting results from multi-arm trials [7], it is clear that there is now much interest within the trials community in such designs.

However, whilst there are numerous advantages of multi-arm trials, it is important to recognise that determining a suitable design for a multi-arm clinical trial can be a substantially more complex process than for a two-arm trial. In particular, a decision must be made on how to account for the multiple comparisons that will be made. Indeed, whether the final analysis should adjust for multiplicity has been a topic of much debate within the literature. In brief, presented arguments primarily revolve around the fact that failing to account for multiplicity can substantially increase the probability of committing a type-I error. Yet, if a series of two-arm trials were conducted, no adjustment would be made to the significance level used in each trial. For brevity, we will not repeat all further arguments on this issue here, and instead refer the reader to several key discussions on multiplicity [5, 8–18].

For the purposes of what follows in this article, the more important consideration is that when a multiple comparison correction (MCC) is to be used, one of a wide selection must actually be chosen (see, e.g., [19–21] for an overview). MCCs vary widely in their complexity, with Bonferroni’s correction often recommended because of its simplicity [7]. However, other MCCs often perform better in terms of the operating characteristics they impart, as Bonferroni’s correction is known to be conservative [10, 18, 20, 22]. A recent review found that amongst those multi-arm trials that did adjust for multiplicity, 50% used one of the comparatively simple Bonferroni or Dunnett corrections [5]. Thus, there arguably remains the potential for increased efficiency gains to be made in multi-arm trials, if more advanced MCCs can be employed.

Furthermore, regardless of whether a MCC is utilised, there are other complications that must also be addressed in multi-arm trial design, including how to power the trial, and what the allocation ratio to each experimental arm relative to the control arm will be. Indeed, power is not a simple quantity in a multi-arm trial, whilst the literature on how to choose the allocation ratios in an optimal

manner is extensive (see, e.g., [23] for an overview), and deciding whether to specify allocation ratios absolutely, or whether they can be optimised to improve trial efficiency may not be an easy decision.

These considerations imply that user-friendly software for designing multi-arm clinical trials would be a valuable tool in the trials community. It is unfortunate therefore that, as we discuss further later, little software is available to assist with such studies. For this reason, we have developed a web application for multi-arm clinical trial design. We hope that the availability of this application will assist with the utilization of more advanced multi-arm designs in future clinical trials.

Implementation

The web application is written using the Shiny package [24] in the R programming language [25]. It is available as a function in (for off-line local use), and is built using other functions from, the R package multiarm [26]. A vignette is provided for multiarm that gives great detail on its formal statistical specifications. A less technical summary is provided here.

Design setting

It is assumed that outcomes X_{ik} will be accrued from patients $i \in \{1, \dots, n_k\}$ on treatment arms $k \in \{0, \dots, K\}$, with arm $k = 0$ corresponding to a shared control arm, and arms $k \in \{1, \dots, K\}$ to several experimental arms. Later, we provide more information on the precise types of outcome that are currently supported by the web application. The hypotheses of interest are assumed to be $H_k : \tau_k \leq 0$ for $k \in \{1, \dots, K\}$. Here, τ_k corresponds to a treatment effect for experimental arm $k \in \{1, \dots, K\}$ relative to the control arm. Thus, we assume one-sided tests for superiority. Note that in the app, reference is also made to the global null hypothesis, H_G , which we define to be the scenario with $\tau_1 = \dots = \tau_K = 0$.

To test hypothesis H_k , we assume that a Wald test statistic, z_k , is computed

$$z_k = \frac{\hat{\tau}_k}{\sqrt{\text{Var}(\hat{\tau}_k)}} = \hat{\tau}_k I_k^{1/2}, \quad k \in \{1, \dots, K\}.$$

In what follows, we use the notation $\mathbf{z}_k = (z_1, \dots, z_k)^\top \in \mathbb{R}^k$. With this, note that our app supports design in particular scenarios where \mathbf{Z}_k , the random pre-trial value of \mathbf{z}_k , has (at least asymptotically) a k -dimensional multivariate normal (MVN) distribution, with

$$\begin{aligned} \mathbb{E}(Z_l) &= \tau_l I_l^{1/2}, \quad l = 1, \dots, k, \\ \text{Cov}(Z_l, Z_l) &= 1, \quad l \in \{1, \dots, k\}, \\ \text{Cov}(Z_{l_1}, Z_{l_2}) &= I_{l_1}^{1/2} I_{l_2}^{1/2} \text{Cov}(\tau_{l_1}, \tau_{l_2}), \quad l_1 \neq l_2, \quad l_1, l_2 \in \{1, \dots, k\}. \end{aligned}$$

As is discussed further later, this includes normally distributed outcome variable scenarios and, for large sample

sizes, other parametric distributions such as Bernoulli outcome data.

Ultimately, to test the hypotheses, \mathbf{z}_K is converted to a vector of p -values, $\mathbf{p} = (p_1, \dots, p_K)^\top \in [0, 1]^K$, via $p_k = 1 - \Phi_1(z_k, 0, 1)$, for $k \in \{1, \dots, K\}$. Here, $\Phi_n\{(a_1, \dots, a_n)^\top, \boldsymbol{\lambda}, \Sigma\}$ is the cumulative distribution function of an n -dimensional MVN distribution, with mean $\boldsymbol{\lambda}$ and covariance matrix Σ . Precisely

$$\Phi_n\{(a_1, \dots, a_n)^\top, \boldsymbol{\lambda}, \Sigma\} = \int_{-\infty}^{a_1} \dots \int_{-\infty}^{a_n} \phi_n\{\mathbf{x}, \boldsymbol{\lambda}, \Sigma\} dx_1 \dots dx_n,$$

where $\phi_n\{\mathbf{x}, \boldsymbol{\lambda}, \Sigma\}$ is the probability density function of an n -dimensional MVN distribution with mean $\boldsymbol{\lambda}$ and covariance matrix Σ , evaluated at vector $\mathbf{x} = (x_1, \dots, x_n)^\top$.

Then, which null hypotheses are rejected is determined by comparing the p_k to a set of significance thresholds specified based on a chosen MCC, in combination with a nominated significance level $\alpha \in (0, 1)$. Before we describe the currently supported MCCs however, we will first describe the operating characteristics that are currently evaluated by the app.

Operating characteristics

Our app returns a wide selection of statistical operating characteristics that may be of interest when choosing a multi-arm trial design. Specifically, it can compute the following quantities for any nominated multi-arm design and true set of treatment effects

- The conjunctive power (P_{con}): The probability that all of the null hypotheses are rejected, irrespective of whether they are true or false.
- The disjunctive power (P_{dis}): The probability that at least one of the null hypotheses is rejected, irrespective of whether they are true or false.
- The marginal power for arm $k \in \{1, \dots, K\}$ (P_k): The probability that H_k is rejected, irrespective of whether it is true or false.
- The per-hypothesis error-rate (*PHER*): The expected value of the number of type-I errors divided by the number of hypotheses.
- The a -generalised type-I familywise error-rate ($FWER_{Ia}$): The probability that at least $a \in \{1, \dots, K\}$ type-I errors are made. Note that $FWER_{I1}$ is the conventional familywise error-rate ($FWER$); the probability of making at least one type-I error.
- The a -generalised type-II familywise error-rate ($FWER_{IIa}$): The probability that at least $a \in \{1, \dots, K\}$ type-II errors are made.
- The false discovery rate (*FDR*): The expected proportion of type-I errors amongst the rejected hypotheses.

- The false non-discovery rate (*FNDR*): The expected proportion of type-II errors amongst the hypotheses that are not rejected.
- The positive false discovery rate (*pFDR*): The rate that rejections are type-I errors.
- The sensitivity (*Sensitivity*): The expected proportion of the number of correct rejections of the hypotheses to the number of false null hypotheses.
- The specificity (*Specificity*): The expected proportion of the number of correctly not rejected hypotheses to the number of true null hypotheses.

Multiple comparison corrections

Per-hypothesis error-rate control

The most simple method for selecting the significance thresholds against which to compare the p_k , is to compare each to the chosen significance level α . That is, to reject H_k for $k \in \{1, \dots, K\}$ if $p_k \leq \alpha$. This controls the *PHER* to α .

A potential problem with this, however, can be that the statistical operating characteristics of the resulting design may not be desirable (e.g., in terms of $FWER_{I1}$). As discussed earlier, it is for this reason that we may wish to make use of a MCC. Currently, the web application supports the use of a variety of such MCCs, which aim to control either (a) the conventional familywise error-rate, $FWER_{I1}$ (with these techniques sub-divided into single-step, step-down, and step-up corrections) or (b) the *FDR*.

Single-step familywise error-rate control

These MCCs test each of the H_k against a common significance level, $\gamma \in (0, 1)$ say, rejecting H_k if $p_k \leq \gamma$. The currently supported single-step corrections are

- Bonferroni's correction: This sets $\gamma = \alpha/K$ [27].
- Sidak's correction: This sets $\gamma = 1 - (1 - \alpha)^{1/K}$ [28].
- Dunnett's correction: This sets $\gamma = 1 - \Phi_1\{z_D, 0, 1\}$, where z_D is the solution of the following equation

$$\alpha = 1 - \Phi_K\{(z_D, \dots, z_D)^\top, \mathbf{0}_K, \text{Cov}(\mathbf{Z}_K)\},$$

with $\mathbf{0}_n = (0, \dots, 0)^\top \in \mathbb{R}^n$ an n -dimensional vector of zeroes [29].

Note that each of the above specify a γ such that the maximum probability of incorrectly rejecting at least one of the null hypotheses H_k , $k \in \{1, \dots, K\}$, over all possible values of $\boldsymbol{\tau} \in \mathbb{R}^K$ is at most α . This is referred to as strong control of $FWER_{I1}$.

Step-down familywise error-rate control

Step-down MCCs work by ranking the p -values from smallest to largest. We will refer to these ranked p -values by $p_{(1)}, \dots, p_{(K)}$, with associated hypotheses $H_{(1)}, \dots, H_{(K)}$. The $p_{(k)}$ are compared to a vector of significance levels $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_K) \in (0, 1)^K$. Precisely, the

maximal index k such that $p_{(k)} > \gamma_k$ is identified, and then $H_{(1)}, \dots, H_{(k-1)}$ are rejected and $H_{(k)}, \dots, H_{(K)}$ are not rejected. If $k = 1$ then we do not reject any of the null hypotheses, and if no such k exists then we reject all of the null hypotheses. The currently supported step-down corrections are

- Holm-Bonferroni correction: This sets $\gamma_k = \alpha / (K + 1 - k)$ [30].
- Holm-Sidak correction: This sets $\gamma_k = 1 - (1 - \alpha)^{K+1-k}$.
- Step-down Dunnett correction: This can only currently be used when the $\text{Cov}(Z_{k_1}, Z_{k_2})$ are equal for all $k_1 \neq k_2, k_1, k_2 \in \{1, \dots, K\}$. In this case, it sets $\gamma_k = 1 - \Phi_1\{z_{Dk}, 0, 1\}$, where z_{Dk} is the solution to

$$\alpha = 1 - \Phi_{K+1-k}\{(z_{Dk}, \dots, z_{Dk})^\top, \mathbf{0}_{K+1-k}, \text{Cov}(\mathbf{Z}_{K+1-k})\}.$$

Note that the above methods provide strong control of $FWER_{I1}$.

Step-up familywise error-rate control

Step-up MCCs also work by ranking the p -values from smallest to largest, and similarly utilise a vector of significance levels γ . However, here, the largest k such that $p_{(k)} \leq \gamma_k$ is identified. Then, the hypotheses $H_{(1)}, \dots, H_{(k)}$ are rejected, and $H_{(k+1)}, \dots, H_{(K)}$ are not rejected. Currently, one such correction is supported: Hochberg’s correction [31], which sets $\gamma_k = \alpha / (K + 1 - k)$. This method also provides strong control of $FWER_{I1}$.

False discovery rate control

It may be of interest to instead control the FDR , which can offer a compromise between strict $FWER_{I1}$ control and $PHER$ control, especially when we expect a large proportion of the experimental treatments to be effective. Currently, two methods that will control the FDR to at most α over all possible $\tau \in \mathbb{R}^K$ are supported. They function in the same way as the step-up corrections discussed above, with

- Benjamini-Hochberg correction: This sets $\gamma_k = k\alpha / K$ [32].
- Benjamini-Yekutieli correction: This sets [33]:

$$\gamma_k = \frac{k\alpha}{K \left(1 + \frac{1}{2} + \dots + \frac{1}{K}\right)}.$$

Sample size determination

The sample size required by a design to control several types of power to a specified level $1 - \beta$, under certain specific scenarios, can be computed. Precisely, following for example [34], values for ‘interesting’ and ‘uninteresting’ treatment effects, $\delta_1 \in \mathbb{R}^+$ and $\delta_0 \in (-\infty, \delta_1)$ respectively, are specified and the following definitions are made

- The global alternative hypothesis, H_A , is given by $\tau_1 = \dots = \tau_K = \delta_1$.
- The least favourable configuration for experimental arm $k \in \{1, \dots, K\}$, LFC_k , is given by $\tau_k = \delta_1, \tau_1 = \dots = \tau_{k-1} = \tau_{k+1} = \dots = \tau_K = \delta_0$.

Then, the following types of power can be controlled to level $1 - \beta$ by design’s determined using the app

- The conjunctive power under H_A .
- The disjunctive power under H_A .
- The minimum marginal power under the respective LFC_k .

Allocation ratios

One of the primary goals of the app is to aid the choice of values for n_0, \dots, n_K . The app specifically supports the determination of values for these parameters by searching for a suitable n_0 via a one-dimensional root solving algorithm, and then sets $n_k = r_k n_0, r_k \in (0, \infty)$, for $k \in \{1, \dots, K\}$. Here, r_k is the allocation ratio for experimental arm k relative to the control arm.

For this reason, the app also allows the allocation ratios to be specified in a variety of ways: they can be defined explicitly, or alternatively can be determined in an optimal manner. For this optimality problem, many possible optimality criteria have been defined, each with their own merits. Therefore, we refer the reader to Atkinson (2007) [23] for further details of optimal allocation in multi-arm designs. Instead, we simply note that in the web application, the allocation ratios can currently be determined for three such criteria

- A -optimality: Minimizes the trace of the inverse of the information matrix of the design. This results in the minimization of the average variance of the treatment effect estimates.
- D -optimality: Maximizes the determinant of the information matrix of the design. This results in the minimization of the volume of the confidence ellipsoid for the treatment effect estimates.
- E -optimality: Maximizes the minimum eigenvalue of the information matrix. This results in the minimization of the maximum variance of the treatment effect estimates.

The optimal allocation ratios are identified in the app using available closed-form solutions were possible (see [35] for a summary of these), otherwise non-linear programming is employed.

Other design specifications

Finally, the web application also supports the following options

- Plot production: Plots can be produced of (a) all of the operating characteristics quantities listed earlier when $\tau_1 = \dots = \tau_K = \theta$, as well as (b) the P_k when $\tau_k = \theta$ and $\tau_l = \theta - (\delta_1 - \delta_0)$ for $l \neq k$. If these are selected for rendering, the quality of the plots, in terms of the number of values of θ used for line-graph production, can also be controlled.
- Require $n_k \in \mathbb{N}$ for $k \in \{0, \dots, K\}$: By default, the sample size determined for each arm will only be required to be a positive number. In practice, such values need to be integers. This can thus be enforced if desired, with the integer n_k specified by rounding up their determined continuous values.

Supported outcome variables

Normally distributed outcome variables

Currently, the app supports multi-arm trial design for scenarios in which the outcome variables are assumed to be either normally or Bernoulli distributed.

Precisely, for the normal case, it assumes that $X_{ik} \sim N(\mu_k, \sigma_k^2)$, and that σ_k^2 is known for $k \in \{0, \dots, K\}$. Then, for each $k \in \{1, \dots, K\}$

$$\begin{aligned} \tau_k &= \mu_k - \mu_0, \\ \hat{\tau}_k &= \frac{1}{n_k} \sum_{l=1}^{n_k} x_{ik} - \frac{1}{n_0} \sum_{l=1}^{n_0} x_{i0}, \\ I_k &= \frac{1}{\frac{\sigma_0^2}{n_0} + \frac{\sigma_k^2}{n_k}}, \end{aligned}$$

where x_{ik} is the realised value of X_{ik} .

Note that in this case, \mathbf{Z}_K has a MVN distribution, and thus the operating characteristics can be computed exactly and efficiently using MVN integration [36]. Furthermore, the distribution of \mathbf{Z}_K does not depend upon the values of the μ_k , $k \in \{0, \dots, K\}$. Consequently, these parameters play no part in the inputs or outputs of the app.

Bernoulli distributed outcome variables

In this case, $X_{ik} \sim \text{Bern}(\pi_k)$ for response rates π_k , and for each $k \in \{1, \dots, K\}$

$$\begin{aligned} \tau_k &= \pi_k - \pi_0, \\ \hat{\tau}_k &= \frac{1}{n_k} \sum_{l=1}^{n_k} x_{ik} - \frac{1}{n_0} \sum_{l=1}^{n_0} x_{i0}, \\ I_k &= \frac{1}{\frac{\pi_0(1-\pi_0)}{n_0} + \frac{\pi_k(1-\pi_k)}{n_k}}. \end{aligned}$$

Thus, a problem for design determination becomes that the I_k are dependent on the unknown response rates. In practice, this is handled at the analysis stage of a trial by

setting

$$I_k = \frac{1}{\frac{\hat{\pi}_0(1-\hat{\pi}_0)}{n_0} + \frac{\hat{\pi}_k(1-\hat{\pi}_k)}{n_k}},$$

for $\hat{\pi}_k = \sum_{i=1}^{n_k} x_{ik}/n_k$, $k \in \{0, \dots, K\}$. This is the assumption made where required in the app. With this, \mathbf{Z}_K is only asymptotically MVN. Thus, in general it would be important to validate operating characteristics evaluated using MVN integration via simulation.

In addition, note that the above problem also means that the operating characteristics under H_G , H_A , and the LFC_k are not unique without further restriction. Thus, to achieve uniqueness, the app requires a value be specified for π_0 for use in the definition of these scenarios. Moreover, for this reason, the inputs and outputs of functions supporting Bernoulli outcomes make no reference to the τ_k , and work instead directly in terms of the π_k . Finally, note that this problem also means that to determine A -, D -, or E -optimised allocation ratios, a specific set of values for the π_k must be assumed.

In this case, we should also ensure that $\delta_1 \in (0, 1)$ and $\delta_0 \in (-\pi_0, \delta_1)$, for the assumed value of π_0 , since $\pi_k \in [0, 1]$ for $k \in \{1, \dots, K\}$.

Results Support

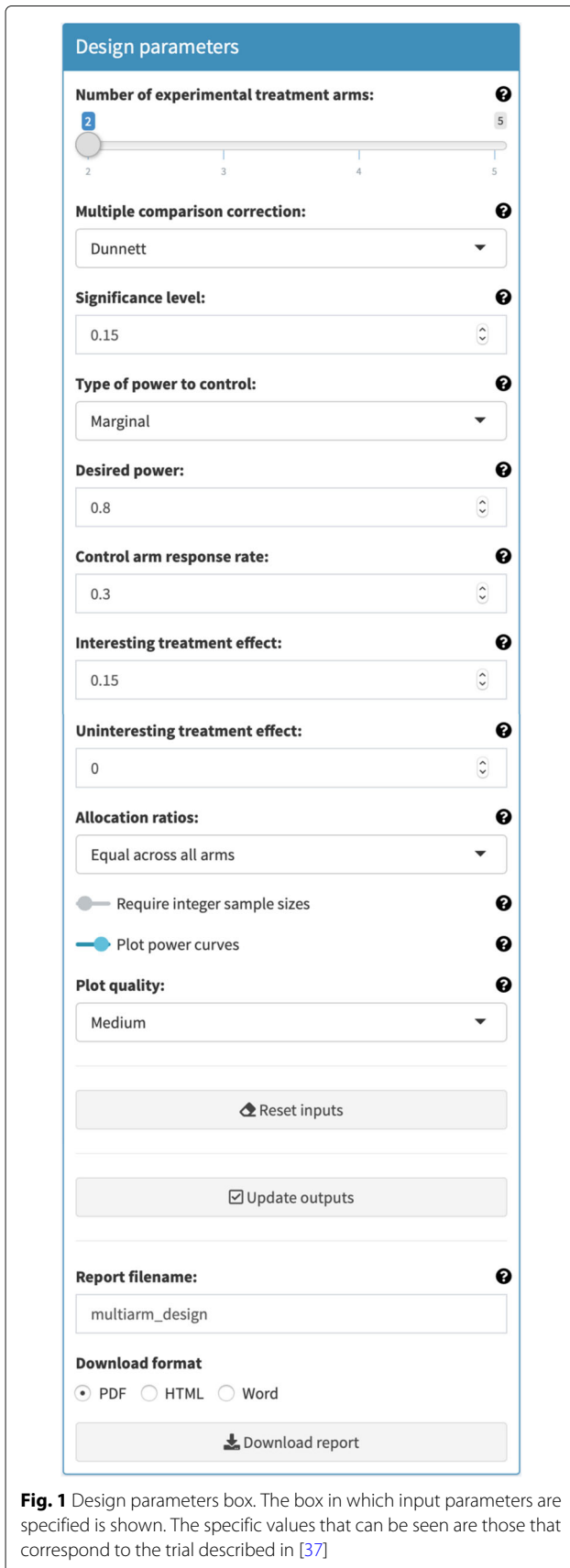
The web application is freely available from <https://mjgrayling.shinyapps.io/multiarm/>. The R code for the application can also be downloaded from <https://github.com/mjg211/multiarm>. Furthermore, as noted earlier, the app is built in to the package multiarm [26], as the function `gui()`, for ease-of-use without internet access. The application has a simple interface, and has the capability to

- Determine the sample required in each arm in a specified multi-arm clinical trial design scenario;
- Summarise and plot the operating characteristics of the identified design;
- Produce a report describing the chosen design scenario, the identified design, and a summary of its operating characteristics.

Inputs

The outputs (i.e., the identified design and its operating characteristics) are determined based upon the following set of user specified inputs (Fig. 1)

1. The number of experimental treatment arms, K .
2. The chosen multiple comparison correction (e.g., Dunnett’s correction).
3. The significance level, α .
4. The type of power to control (e.g., the conjunctive power under H_A).



5. The desired power, $1 - \beta$.
6. For Bernoulli distributed data, the control arm response rate π_0 .
7. The interesting treatment effect, δ_1 .
8. The uninteresting treatment effect, δ_0 .
9. For normally distributed data, the standard deviations, $\sigma_0, \dots, \sigma_K$. These are allocated by first selecting the type of standard deviations (e.g., that they are assumed to be equal across all arms), and then the actual values for the parameters.
10. The allocation ratios (e.g., A-optimal).
11. For Bernoulli distributed data, when searching for optimal allocation ratios, the response rates to assume in the search.
12. Whether the sample size in each arm should be required to be an integer;
13. Whether plots should be produced, and if so the plot quality.

Note that a *Reset inputs* button is provided to simplify returning the inputs to their default values. Once the inputs have been specified as desired, the outputs can be generated by clicking the *Update outputs* button.

Example

Here, we demonstrate specification of the input parameters (Fig. 1), and then subsequent output generation (Figs. 2, 3, and 4), for parameters motivated by a three-arm phase II randomized controlled trial of treatments for myelodysplastic syndrome patients, described in [37]. This trial compared, via a binary primary outcome, two experimental treatments with conventional azacitidine treatment. The trial was designed with $\alpha = 0.15$, $\beta = 0.2$, $\delta_1 = 0.15$, and $\pi_0 = 0.3$. For simplicity, we assume that the familiar Dunnett correction will be used, that $\delta_0 = 0$, and that allocation will be equal across the arms ($r_1 = \dots = r_K = 1$). Finally, we assume it is the minimum marginal power that should be controlled.

Each input widget in Fig. 1 can be seen to have been allocated accordingly based on the description above, whilst we have additionally elected to produce plots (of medium quality), and to not require the arm-wise sample sizes to be integers. Note that in Fig. 1 we can see that the input widgets are supported by help boxes that can be opened by clicking on the small question marks beside them.

Figure 2 then depicts the output to the *Design summary* box once the user clicks on *Update outputs*. Specifically, a summary of the chosen inputs and the identified design is rendered. Furthermore, in Fig. 3 we can see the tables that provide the various statistical quantities under H_G , H_A , the LFC_k , as well as the various treatment effect scenarios that are considered for plot production.

Finally, in Fig. 4 the plots discussed earlier are shown. Observe that horizontal and vertical lines are added at the

Design summary

Design setting

The trial will be designed to compare K experimental treatments to a shared control arm. Response X_{ik} , from patient $i = 1, \dots, n_k$ in arm $k = 0, \dots, K$, will be assumed to be distributed as $X_{ik} \sim \text{Bern}(\pi_k)$. Then, the hypotheses to be tested will be:

$$H_k : \tau_k = \pi_k - \pi_0 \leq 0, k = 1, \dots, K.$$

The *global null hypothesis*, H_G , will be:

$$\pi_0 = \dots = \pi_K.$$

The *global alternative hypothesis*, H_A , will be:

$$\pi_1 = \dots = \pi_K = \pi_0 + \delta_1.$$

The *least favourable configuration* for experimental arm k , LFC_k , will be:

$$\pi_k = \pi_0 + \delta_1, \pi_1 = \dots = \pi_{k-1} = \pi_{k+1} = \dots = \pi_K = \pi_0 + \delta_0.$$

Here, δ_1 and δ_0 are *interesting* and *uninteresting* treatment effects respectively.

Inputs

The following choices were made:

- $K = 2$ experimental treatments will be included in the trial.
- A significance level of $\alpha = 0.15$ will be used, in combination with **Dunnett's correction**.
- The response rate in the control arm will be assumed to be: $\pi_0 = 0.3$.
- The **marginal power for each null hypothesis** will be controlled to level $1 - \beta = 0.8$ under **each of their respective least favourable configurations**.
- The interesting and uninteresting treatment effects will be: $\delta_1 = 0.15$ and $\delta_0 = 0$ respectively.
- The target allocation to each of the experimental arms will be: **the same as the control arm**.
- The sample size in each arm **will not** be required to be an integer.
- Plots **will** be produced.

Outputs

- The total required sample size is: $N = 293,963$.
- The required sample size in each arm is: $(n_0, \dots, n_K) = (97,988, 97,988, 97,988)$.
- Therefore, the realised allocation ratios to the experimental arms are: $(r_1, \dots, r_K) = (1, 1)$.
- The maximum familywise error-rate is: **0.15**.
- The **minimum marginal power** is: **0.8**.
- The following critical threshold should be used with the chosen multiple comparison correction: **0.087**.

Fig. 2 Design summary box. The box in which a summary of the input parameters and of the identified design is rendered is shown. The specific output that can be seen corresponds to the inputs from Fig. 1

values α , $1 - \beta$, δ_1 , and δ_0 respectively. Note that these plots are outputted in a manner to allow the user to zoom in on a particular sub-component if desired.

In all, Figs. 2, 3, and 4 provide a set of outputs with a variety of features that should be anticipated given the chosen input parameters. Firstly, the specification that the allocation to all arms should be equal means that $n_0 = \dots = n_K$. In addition, $FWER_{I1}$ is equal to 0.15 under H_G , and the minimum marginal power is 0.8, as was desired. Moreover, the specification that $r_1 = \dots = r_K$ means that P_{con} and P_{dis} are equal for each of the LFC_k , and $P_1 = P_2$.

Finally, as noted above, and as can be seen in Fig. 1, a *Generate report* button is provided that can produce a copy of the outputs in either PDF (.pdf), HTML (.html), or Word (.docx) format. The user can also nominate a name for this file in the *Report filename* input widget. This allows a record of designs to be stored, presented, and compared to other designs if required. A copy of the report, in PDF form, for the inputs shown in Fig. 1, is given as Additional file 1.

Comparison to other software solutions

In this section we discuss solutions that are available for designing multi-arm trials in a range of popular trial design packages, using this to describe the advantages and disadvantages of our web application.

Firstly, we note that we are unaware of any other code for R that directly facilitates the design of a multi-arm trial: in particular the CRAN Task View for Clinical Trial Design, Monitoring, and Analysis does not list any potential solution [38]. Nonetheless, a multi-arm trial designed to achieve a particular level of marginal power, that controls either the *PHER* or the *FWER* via a single-step MCC, could be identified using one of the many functions available for designing two-arm trials (see, e.g., `power.prop.test()` from the stats package). However, one would not then be able to readily explore the resultant design's operating characteristics. Similar statements hold for Stata [39] and SAS [40], with the `power` command and the PROC POWER procedure respectively enabling the determination and evaluation of two-arm

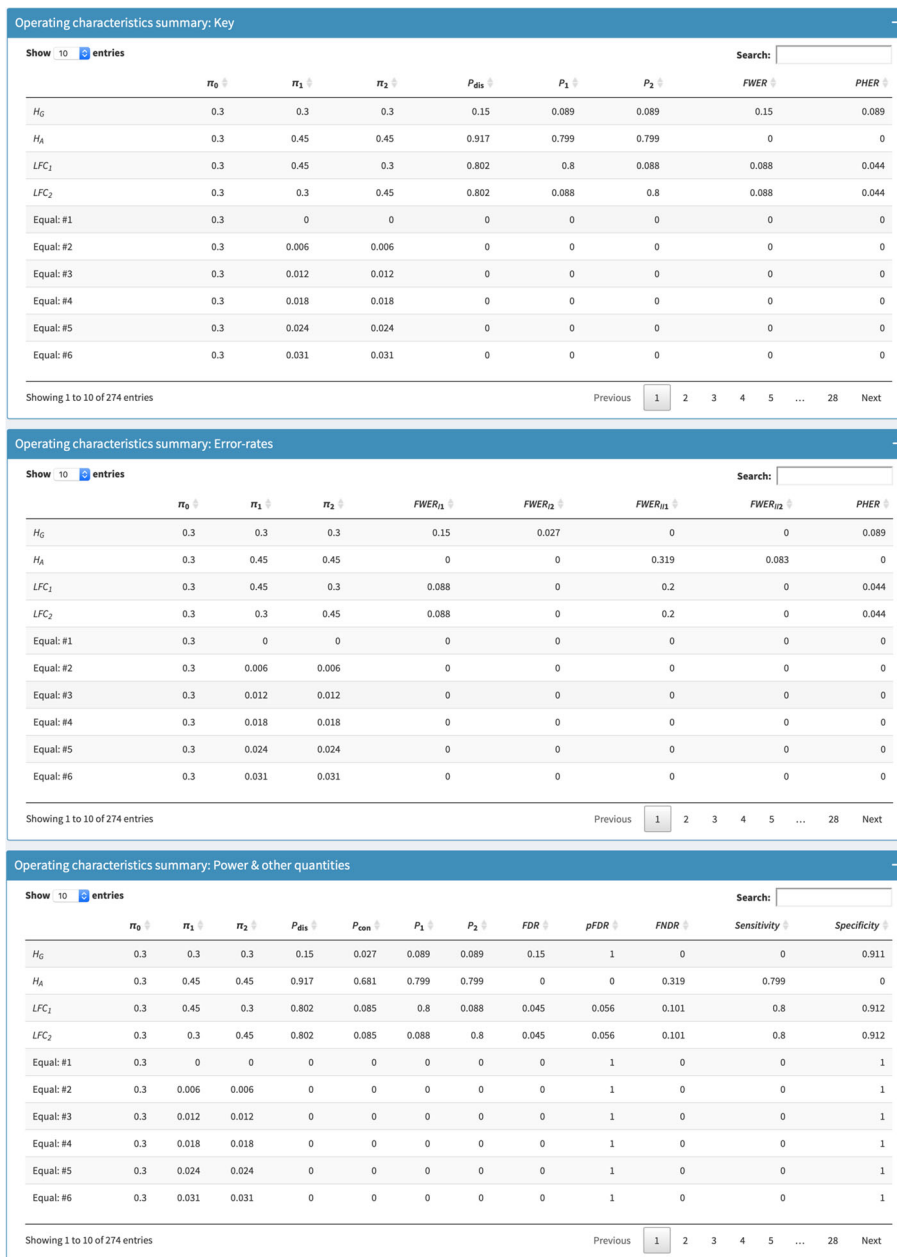


Fig. 3 Operating characteristics summary. The boxes in which a summary of the identified designs operating characteristics is produced is shown. The specific output that can be seen corresponds to the inputs from Fig. 1

trial designs, but neither directly supports multi-arm trial design. Moreover, nQuery [41], to the best of our knowledge does not appear to currently support the design of multi-arm trials.

Direct solutions for certain types of multi-arm trial are available in several other proprietary software packages: namely East [42], FACTS [43], and PASS [44]. Unfortunately, the cost of these packages may be prohibitive to many working within academia. Indeed,

this was our primary motivation for developing the presented web application, and we are only able to comment precisely here on the available functionality in PASS, as we do not have access to either East or FACTS.

Firstly, we note that from the information provided online, the MULTIARM module for East facilitates the determination of a range of multi-arm trial designs. So to does it support their comparison in terms of numerous

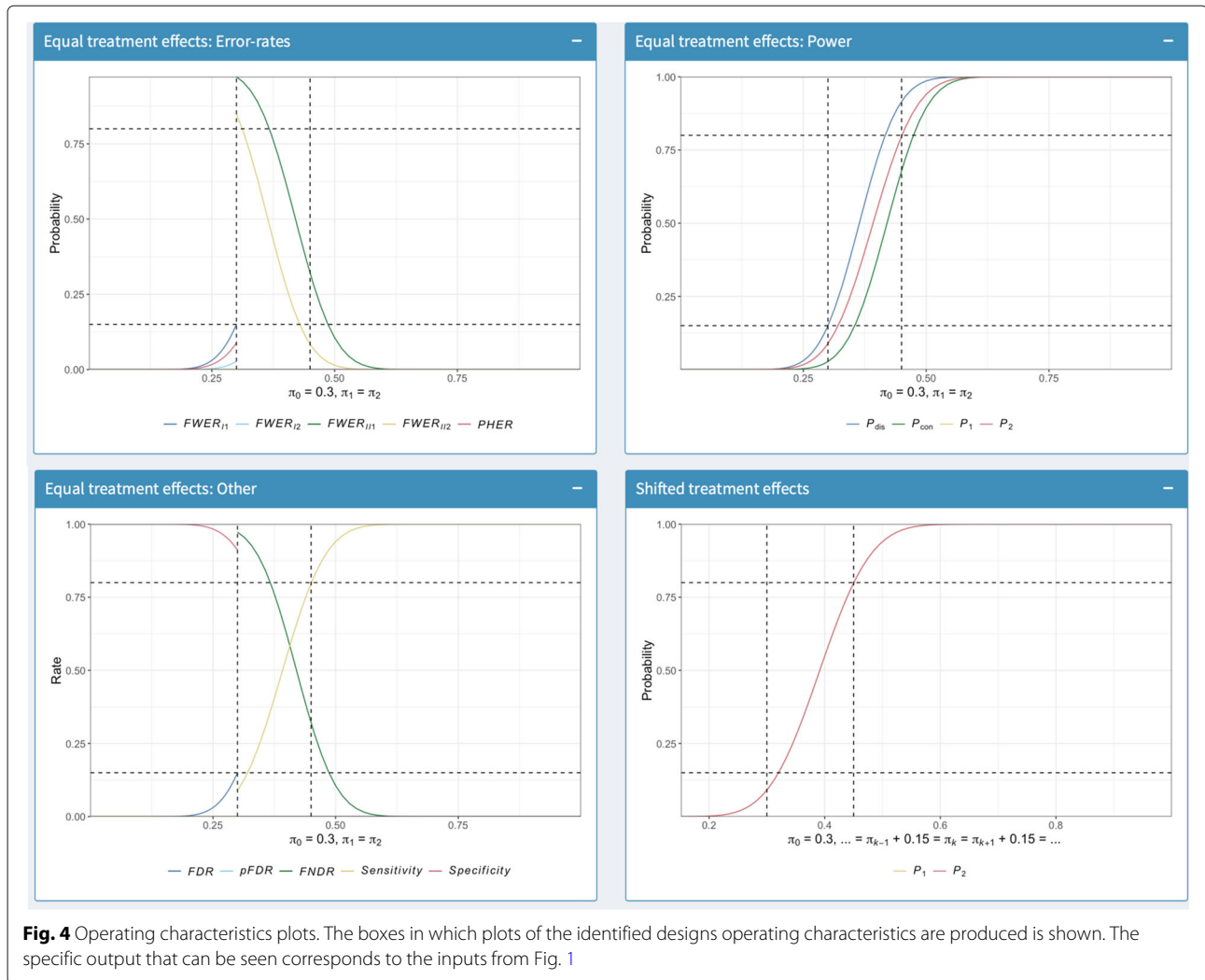


Fig. 4 Operating characteristics plots. The boxes in which plots of the identified designs operating characteristics are produced is shown. The specific output that can be seen corresponds to the inputs from Fig. 1

operating characteristics, including the *FWER* and several varieties of power. It will also produce a selection of insightful plots, handles both continuous and binary outcome variables, and eleven MCCs. Less information is available online about the precise support available in FACTS, but it is stated that its ‘Core’ functionality can handle scenarios with multiple treatment arms. In PASS, support is provided to design a multi-arm trial with Bernoulli outcomes via formula provided in Chow *et al.* (2008) [45]. Specifically, the Bonferroni correction is used to control the *FWER* to a specified level, and the sample size required to achieve a particular level of the minimum marginal power can be computed, under several allocation ratio scenarios. Furthermore, a report is ultimately generated on the calculations performed. PASS also supports similar calculations, using either Dunnett’s or the Kruskal-Wallis MCC, for a vast array of outcome types via simulation (including both Bernoulli and normally distributed outcomes). These calculations explicitly address

the sample size required to control the conjunctive or disjunctive power, and allow for flexible assumptions about the allocation ratios.

Thus, a variety of multi-arm trial designs can be determined using solutions other than our web application. However the cost of these packages may render them unsuitable, particularly in academic departments. This reveals arguably the greatest advantage of our web application: that it is provided under a license that makes it completely free to utilise and modify as a user sees fit. In addition, like the discussed proprietary solutions, our web application allows for calculations via a GUI that contains several features to make it easier to use, without compromising on the type of multi-arm designs that can be determined. In fact, we would argue that our application supports a broader range of multi-arm design scenarios than any other currently available solution.

We feel that there are only two principal limitations of our application. Firstly, MVN integration is utilised by

the application in all instances to determine the statistical operating characteristics of potential multi-arm designs. This makes the execution time for returning outputs with many possible input parameters fast. However, there is an unavoidable complexity in certain multi-arm designs, which may make execution time long. This is particularly true of scenarios with $K \geq 5$. It can also be true of designs that utilise the more complex step-wise MCCs. It is for this reason that the web application places an upper cap in the inputs of $K = 5$, and also returns a warning in scenarios for which a lengthy execution time would be anticipated. Nonetheless, users may have to wait several minutes in certain situations to identify their desired design. In contrast, proprietary solutions may exploit more efficient solutions to reduce execution time, with FACTS in particular noting its use of efficient low-level languages.

More significantly, it is crucial that all software for clinical trial design be validated. Each of the discussed proprietary solutions will almost certainly have gone through more rigorous testing than we are able to achieve. Specifically, it is challenging to validate our results because of the limited freely available software solutions for multi-arm trials. We have compared the output of our application to that of PASS for a variety of supported input parameters, but output for many possible inputs remains difficult to corroborate because of a lack of equivalent available functionality. For this reason, we have carefully followed recommended good-programming practices and perform all statistical calculations within the application by calling functions from the R package `multiarm`, in which the code has been modularised [26].

Furthermore, in this package we have created a function that simulates multi-arm clinical trials that use a given design. This allows us to perform an additional check on our analytical computations. As an example, we demonstrate how to identify the example design discussed above, but under the assumption of normally distributed data with $\sigma_1 = \dots = \sigma_K = 1$:

```
> set.seed(1)
> design <- multiarm::des_ma(K = 2,
+ alpha = 0.15,
+ beta = 0.2,
+ delta1 = 0.15,
+ delta0 = 0,
+ sigma = c(1, 1, 1),
+ ratio = c(1, 1),
+ correction = "dunnett",
+ power = "marginal",
+ integer = T)
```

Then, 100,000 replicate simulations of trials that utilise this design, under H_G , H_A , and the LFC_k , can be calculated with:

```
> simulated <- multiarm::sim_ma(design)
```

Finally, the maximum absolute difference in the operating characteristics of this design, as determined analytically and via simulation can be evaluated as:

```
> max(abs(simulated$sim - design$opchar))
[1] 0.002166331
```

Thus, the maximal difference is within what would be anticipated allowing for simulation error.

In Additional file 2, we demonstrate how we repeated the above for 1000 randomly generated combinations of possible input parameters, thus covering an extremely wide range of supported design scenarios. As above, the analytical operating characteristics returned by the web application in the *Operating characteristics summary* boxes were compared to those based on trial simulation, using 100,000 replicate simulations in each instance. Across all considered scenarios, the maximum absolute difference between the analytical and simulated operating characteristics was just 5×10^{-3} , which is again within what would be anticipated due to simulation error. Consequently, it does appear that our application is functioning as it should. However, it remains that the principal argument for not utilising our application would be to attain a stronger guarantee on the results.

Conclusions

A possible barrier to previous calls for increased use of multi-arm clinical trial designs is a lack of available easy-to-access user-friendly software that facilitates associated sample size calculations. For this reason, we have created an online web application that supports multi-arm trial design determination for a wide selection of possible input parameters. Its use requires no knowledge of statistical programming languages and is facilitated via a simple user interface. Furthermore, we have made the application available on the internet, so that it is readily accessible, and have also made it freely available for download for remote use without an internet connection. Like similar applications that have been released recently for phase I clinical trial design [46, 47], we hope that the availability of this application will assist with the design of future multi-arm studies. As we have discussed, however, users should bear in mind the primary limitation of our application: that it is not validated. Therefore, alternative proprietary solutions may be needed if certain guarantees on outputs are required.

Finally, we note several possible avenues for future development of the web application. Firstly, numerous papers have now provided designs for adaptive multi-arm trials (e.g., [48, 49]), and software for their determination in certain settings [50, 51]. Given the evidential increased interest in such designs [52], allowing for their determination would be a valuable extension to our application. In addition, our web application currently focuses on design

for normally and Bernoulli distributed outcomes. But, time-to-event outcomes are also commonly used in oncology. Permitting such calculations therefore likewise offers a valuable avenue for subsequent versions of the app.

Availability and requirements

Project name: Multi-arm trial web application.

Project home page: <https://mjgrayling.shinyapps.io/multiarm/>.

Operating system(s): Platform independent.

Programming language: R.

Other requirements: Version 3.5.2 or later.

License: MIT.

Any restrictions to use by non-academics: None.

Supplementary information

Supplementary information accompanies this paper at <https://doi.org/10.1186/s12885-020-6525-0>.

Additional file 1: PDF report. A copy of the PDF report generated by clicking the Generate report button in the web application, for the input parameters shown in Fig. 1.

Additional file 2: Analytical vs. simulated operating characteristics comparison. R code to replicate our comparison of the analytical operating characteristics returned by the web application against those based on simulation.

Abbreviations

MCC: Multiple comparison correction; MVN: Multivariate normal

Acknowledgements

Not applicable.

Authors' contributions

MJG and JMSW contributed to conception of the web application. MJG wrote the code for the web application. MJG and JMSW contributed to drafting and revising the manuscript. MJG and JMSW gave final approval of the manuscript submitted for publication.

Funding

This work was supported by the Medical Research Council [grant number MC_UU_00002/6 to JMSW]. The funding body did not have any role in the design of this study, collection, analysis, and interpretation of data, nor in the writing of the manuscript.

Availability of data and materials

Access to the application online is available at <https://mjgrayling.shinyapps.io/multiarm/>. The R code for the application can be downloaded from <https://github.com/mjg211/multiarm>.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

¹Population Health Sciences Institute, NE2 4AX Newcastle, UK. ²MRC Biostatistics Unit, University of Cambridge, CB2 0SR Cambridge, UK.

Received: 20 June 2019 Accepted: 8 January 2020

Published online: 31 January 2020

References

1. J.A. DiMasi, H.G. Grabowski, R.W. Hansen, Innovation in the pharmaceutical industry: new estimates of R&D costs. *J Health Econ.* **47**, 20–33 (2016)
2. Biotechnology Innovation Organization (BIO), Biomedtracker, AMPLION, Clinical development success rates 2006–2015 (2016)
3. M.K.B. Parmar, J. Carpenter, M.R. Sydes, More multiarm randomised trials of superiority are needed. *Lancet.* **384**(9940), 283–4 (2014)
4. T. Jaki, J.M.S. Wason, Multi-arm multi-stage trials can improve the efficiency of finding effective treatments for stroke: a case study. *BMC Cardiovasc Disord.* **18**(1), 215 (2018)
5. J.M.S. Wason, L. Stecher, A.P. Mander, Correcting for multiple-testing in multi-arm trials: is it necessary and is it done? *Trials.* **15**, 364 (2014)
6. G. Baron, E. Perrodeau, I. Boutron, P. Ravaud, Reporting of analyses from randomized controlled trials with multiple arms: a systematic review. *BMC Med.* **11**, 84 (2013)
7. E. Juszczak, D.G. Altman, S. Hopewell, K. Schulz, Reporting of multi-arm parallel-group randomized trials: extension of the CONSORT 2010 statement. *JAMA.* **321**(16), 1610–20 (2019)
8. K.J. Rothman, No adjustments are needed for multiple comparisons. *Epidemiology.* **1**(1), 43–6 (1990)
9. R.J. Cook, V.T. Farewell, Multiplicity considerations in the design and analysis of clinical trials. *J R Stat Soc Ser A.* **159**(1), 93–110 (1996)
10. M.A. Proschan, M.A. Waclawiw, Practical guidelines for multiplicity adjustment in clinical trials. *Control Clin Trials.* **21**(6), 527–39 (2000)
11. R. Bender, S. Lange, Adjusting for multiple testing - when and how? *J Clin Epidemiol.* **54**(4), 343–349 (2001)
12. R.J. Feise, Do multiple outcome measures require p-value adjustment? *BMC Med Res Methodol.* **2**, 8 (2002)
13. M.D. Hughes, Multiplicity in clinical trials. *Encycl Biostat.* **5**, 3446–51 (2005)
14. B. Freidlin, E.L. Korn, R. Gray, A. Martin, Multi-arm clinical trials of new agents: some design considerations. *Clin Cancer Res.* **14**, 4368–4371 (2008)
15. G. Li, M. Taljaard, E.R. Van den Heuvel, M.A.H. Levine, D.J. Cook, G.A. Wells, P.J. Devereaux, L. Thabane, An introduction to multiplicity issues in clinical trials: the what, why, when and how. *Int J Epidemiol.* **46**(2), 746–55 (2016)
16. E.M. Agency, Guideline on Multiplicity Issues in Clinical Trials (2017). https://www.ema.europa.eu/en/documents/scientific-guideline/draft-guideline-multiplicity-issues-clinical-trials_en.pdf. Accessed 17 Jan 2020
17. U. F. D. Administration, Multiple Endpoints in Clinical Trials Guidance for Industry (2017). <https://www.fda.gov/regulatory-information/search-fda-guidance-documents/multiple-endpoints-clinical-trials-guidance-industry>. Accessed 17 Jan 2020
18. D.R. Howard, J.M. Brown, S. Todd, W.M. Gregory, Recommendations on multiple testing adjustment in multi-arm trials with a shared control group. *Stat Methods Med Res.* **27**(5), 1513–30 (2018)
19. Y. Hochberg, A.C. Tamhane, *Multiple Comparison Procedures*. (Wiley, New York, 1987)
20. J.C. Hsu, *Multiple Comparisons*. (Chapman & Hall, London, 1996)
21. F. Bretz, T. Hothorn, P. Westfall, *Multiple Comparisons using R*. (CRC Press, Boca Raton, 2010)
22. A.J. Sankoh, R.B.S. D'Agostino, M.F. Huque, Efficacy endpoint selection and multiplicity adjustment methods in clinical trials with inherent multiple endpoint issues. *Stat Med.* **22**(20), 3133–50 (2003)
23. A. Atkinson, A. Donev, R. Tobias, *Optimum Experimental Designs, with SAS*. (Oxford University Press, Oxford, 2007)
24. W. Chang, J. Cheng, J.J. Allaire, Y. Xie, J. McPherson, shiny: Web Application Framework for R (2019). <https://CRAN.R-project.org/package=shiny>. Accessed 17 Jan 2020
25. R Core Team, *R: a Language and Environment for Statistical Computing*. (R Foundation for Statistical Computing, Vienna, 2018). <https://www.R-project.org/>. Accessed 17 Jan 2020
26. M.J. Grayling, multiarm: Design and analysis of fixed-sample multi-arm clinical trials (2019). <http://www.github.com/mjg211/multiarm/>. Accessed 17 Jan 2020
27. C.E. Bonferroni, Teoria statistica delle classi e calcolo delle probabilità. *Pubblicazioni del R Istituto Superiore di Scienze Economiche e Commerciali di Firenze.* **8**, 3–62 (1936)
28. Z. Šidák, Rectangular confidence regions for the means of multivariate normal distributions. *J Am Stat Assoc.* **62**(318), 626–33 (1967)
29. C.W. Dunnett, A multiple comparison procedure for comparing several treatments with a control. *J Am Stat Assoc.* **50**(272), 1096–121 (1955)
30. S. Holm, A simple sequentially rejective multiple test procedure. *Scand J Stat.* **6**(2), 65–70 (1979)

31. Y. Hochberg, A sharper bonferroni procedure for multiple tests of significance. *Biometrika*. **75**(4), 800–2 (1988)
32. Y. Benjamini, Y. Hochberg, Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Stat Soc Ser B*. **57**(1), 289–300 (1995)
33. Y. Benjamini, D. Yekutieli, The control of the false discovery rate in multiple testing under dependency. *Annals Stat*. **29**(4), 1165–88 (1995)
34. J. Wason, D. Magirr, M. Law, T. Jaki, Some recommendations for multi-arm multi-stage trials. *Stat Methods Med Res*. **25**(2), 716–27 (2016)
35. O. Sverdlov, W.F. Rosenberger, On recent advances in optimal allocation designs in clinical trials. *J Stat Theory Pract*. **7**(4), 753–73 (2013)
36. A. Genz, F. Bretz, T. Miwa, M. X, L. F, S. F, H. T, mvtnorm: Multivariate normal and t distributions. R package version 1.0-10 (2019). <http://CRAN.R-project.org/package=mvtnorm>. Accessed 17 Jan 2020
37. L. Jacob, U. M, S. Boulet, I. Begaj, S. Chevret, Evaluation of a multi-arm multi-stage Bayesian design for phase II drug selection trials - an example in hemato-oncology. *BMC Med Res Methodol*. **16**, 67 (2016)
38. C.RAN. Task View: Clinical Trial Design, Monitoring, and Analysis. <https://cran.r-project.org/web/views/ClinicalTrials.html>. Accessed: 16 Oct 2019
39. Stata. <https://www.stata.com/>. Accessed: 16 Oct 2019
40. SAS. https://www.sas.com/en_gb/home.html. Accessed: 16 Oct 2019
41. nQuery. <https://www.statsols.com/nquery>. Accessed: 16 Oct 2019
42. East. <https://www.cytel.com/software/east>. Accessed: 04 May 2019
43. FACTS. <https://www.berryconsultants.com/software/>. Accessed: 16 Oct 2019
44. PASS. <https://www.ncss.com/software/pass/>. Accessed: 16 Oct 2019
45. S. Chow, H. Wang, J. Shao, *Sample Size Calculations in Clinical Research*. (Chapman & Hall, Boca Raton, 2008)
46. G.M. Wheeler, M.J. Sweeting, A.P. Mander, AplusB: A Web Application for Investigating A + B Designs for Phase I Cancer Clinical Trials. *PLoS ONE*. **11**(7), 0159026 (2016)
47. N.A. Wages, G.R. Petroni, A web tool for designing and conducting phase I trials using the continual reassessment method. *BMC Cancer*. **18**, 133 (2018)
48. D. Magirr, T. Jaki, J. Whitehead, A generalized Dunnett test for multi-arm multi-stage clinical studies with treatment selection. *Biometrika*. **99**(2), 494–501 (2012)
49. J. Wason, N. Stallard, J. Bowden, C. Jennison, A multi-stage drop-the-losers design for multi-arm clinical trials. *Stat Methods Med Res*. **26**(1), 508–24 (2017)
50. F.M.S. Barthel, P. Royston, M.K.B. Parmar, A menu-driven facility for sample-size calculation in novel multiarm, multistage randomized controlled trials with a time-to-event outcome. *Stata J*. **9**(4), 505–23 (2009)
51. T. Jaki, P. Pallmann, D. Magirr, The R package MAMS for designing multi-arm multi-stage clinical trials. *J Stat Softw*. **88**(4), 1–25 (2019)
52. M. Dimairo, E. Coates, P. Pallmann, S. Todd, S.A. Julious, T. Jaki, J. Wason, A.P. Mander, C.J. Weir, F. Koenig, M.K. Walton, K. Biggs, J. Nicholl, T. Hamasaki, M.A. Proschan, J.A. Scott, Y. Ando, D. Hind, D.G. Altman, Development process of a consensus-driven CONSORT extension for randomised trials using an adaptive design. *BMC Med*. **16**, 210 (2018)

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

