

The DURATIONS randomised trial design: Estimation targets, analysis methods and operating characteristics

Clinical Trials
2020, Vol. 17(6) 644–653
© The Author(s) 2020



Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/1740774520944377
journals.sagepub.com/home/ctj



Matteo Quartagno , James R Carpenter, A Sarah Walker, Michelle Clements and Mahesh KB Parmar

Abstract

Background: Designing trials to reduce treatment duration is important in several therapeutic areas, including tuberculosis and bacterial infections. We recently proposed a new randomised trial design to overcome some of the limitations of standard two-arm non-inferiority trials. This DURATIONS design involves randomising patients to a number of duration arms and modelling the so-called ‘duration-response curve’. This article investigates the operating characteristics (type-I and type-2 errors) of different statistical methods of drawing inference from the estimated curve.

Methods: Our first estimation target is the shortest duration non-inferior to the control (maximum) duration within a specific risk difference margin. We compare different methods of estimating this quantity, including using model confidence bands, the delta method and bootstrap. We then explore the generalisability of results to estimation targets which focus on absolute event rates, risk ratio and gradient of the curve.

Results: We show through simulations that, in most scenarios and for most of the estimation targets, using the bootstrap to estimate variability around the target duration leads to good results for DURATIONS design-appropriate quantities analogous to power and type-I error. Using model confidence bands is not recommended, while the delta method leads to inflated type-I error in some scenarios, particularly when the optimal duration is very close to one of the randomised durations.

Conclusions: Using the bootstrap to estimate the optimal duration in a DURATIONS design has good operating characteristics in a wide range of scenarios and can be used with confidence by researchers wishing to design a DURATIONS trial to reduce treatment duration. Uncertainty around several different targets can be estimated with this bootstrap approach.

Keywords

Durations, estimand, estimation methods, operating characteristics

Introduction

In several therapeutic areas, it is important to identify the optimal duration of treatment, defined as the shortest duration providing an acceptable efficacy. For example, reducing antibiotic treatment duration has been suggested as a way of combatting antimicrobial resistance,¹ but this has to be done while maintaining high cure rates. Furthermore, shorter treatment durations often increase adherence, reduce side effects and will be more cost-effective, provided they do not lead to an increased risk of relapse.

We recently proposed the DURATIONS randomised trial design² as an improvement over standard non-inferiority trials³ for identifying the shortest

acceptable treatment duration.⁴ Its main attraction is that it does not involve selection of a single shorter duration to test against a control, which is often chosen based on very limited prior evidence; instead, patients are randomised to multiple durations, enabling the relationship between duration and response to be directly

MRC Clinical Trials Unit, University College London Institute for Clinical Trials and Methodology, London, UK

Corresponding author:

Matteo Quartagno, MRC Clinical Trials Unit, University College London Institute for Clinical Trials and Methodology, 90 High Holborn, Second Floor, London WC1V 6LJ, UK.
Email: m.quartagno@ucl.ac.uk

estimated using pre-specified flexible regression models. In a wide range of examples, randomising 500 patients to 5–7 durations enables the underlying duration-response curve to be estimated within 5% average absolute error in 95% of simulations.²

The DURATIONS design moves away from a binary hypothesis testing paradigm so that the trial outcome is not a treatment difference between two fixed durations, but rather the whole estimated duration-response curve. In applications, we wish to use this curve to inform decision-making, and hence it is essential to understand the properties of various ways of drawing inference from the curve. This article, therefore, first defines DURATION-design quantities analogous to power and type-I error and then compares different strategies of inference for these quantities through extensive simulations.

Methods

Suppose there is a treatment that is known to be highly effective compared to placebo, and that is usually prescribed as standard-of-care to patients with a particular disease or condition for a fixed course duration; examples include antibiotics for bacterial infections and direct acting antivirals against Hepatitis C. In addition, suppose the recommended treatment course is D_{MAX} days, although we believe that shorter durations might be similarly effective so that we suspect the shortest effective duration might be as small as D_{MIN} .

Now, suppose we want to design a trial to identify the shortest effective duration and assume that the primary outcome of the trial is cure, a binary variable equal to 1 if the patient recovers from their condition or 0 otherwise. Using a DURATIONS randomised trial design,² we randomise N , say 500, patients to multiple, for example 7, duration arms, between and including D_{MAX} and D_{MIN} . After observing the responses, we fit a pre-specified fractional polynomial logistic model with cure as outcome and duration as the only covariate, with up to two polynomial terms. This gives us an estimated duration-response curve, similar to the black curve in Figure 1.

When provided with this curve, how should clinicians and policy-makers choose what is the ‘optimal’ duration to prescribe? A simple choice is to target the duration leading to at most a fixed loss of efficacy (risk difference δ) compared to the control (maximum) duration, for example, 5% less. Hence, if the estimated control cure rate at D_{MAX} was $\pi_{D_{MAX}}$, our objective would be to find the minimum duration D^* corresponding to a cure rate of at least $\pi_{D^*} = \pi_{D_{MAX}} - \delta$. The rationale for this choice closely corresponds to that for the choice of margin in non-inferiority trials, that is, it answers the question ‘what is the minimum treatment efficacy that we would like to preserve, considering the ancillary

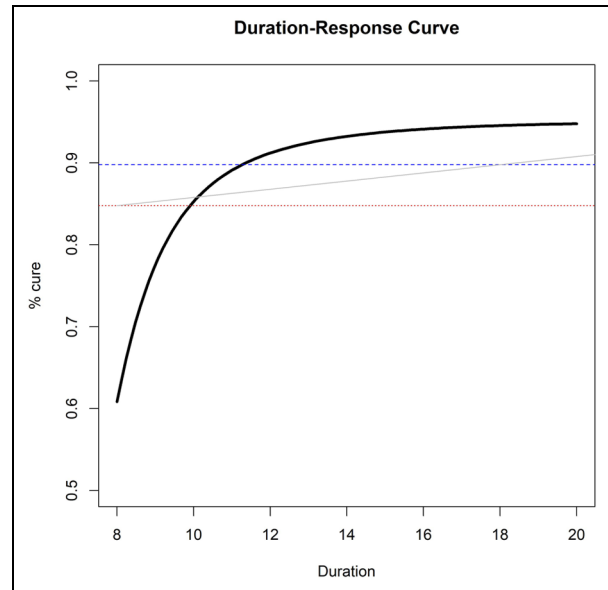


Figure 1. Example of estimated duration-response curve (solid, black), drawn against three possible non-inferiority margins or ‘acceptability frontiers’ (dotted, red = 10% fixed risk difference; dashed, blue = 5% fixed risk difference; solid, grey = duration-specific acceptability frontier).

advantages of the active treatment (i.e. of reducing treatment duration)?’

In this section, we propose different ways of estimating such a D^* from the duration-response curve, and we describe a simulation study to investigate their operating characteristics. Finally, we extend to different possible objectives.

Model confidence bands

The simplest approach is to extrapolate D^* from the duration-response curve, selecting the duration at which $\pi_{D^*} = \pi_{D_{MAX}} - \delta$. However, the fitted curve is just the most likely estimated from the sample data. Thus, there is a non-negligible probability that the true D^* is below this value. Instead, we will typically wish to keep type-1 error below 2.5%, so that we do not recommend a duration that is not long enough more than 2.5% of the times.

The first and most naïve method of controlling type-1 error uses the lower bound of the pointwise confidence bands around the curve, looking for the duration D^* satisfying $\pi_{D^*} - 1.96 \cdot SE_{\pi_{D^*}} = \pi_{D_{MAX}} - \delta$. To avoid recommending non-integer durations, which might make little sense in applications, we round the selected duration D_L^* up to the next whole number (i.e. using the ceiling function); hence, we select the duration D_I that satisfies

$$D_I = \min_{D_i} (D_i > D_L^*)$$

With D_i being the set of integer durations.

Table 1. Properties of different methods to estimate confidence intervals.

Method	Targets difference in efficacy between two points	Addresses model selection uncertainty	Addresses multiplicity issue	Analysis method used in simulations
Confidence bands	No	No	No	Fixed-2 fractional polynomials (gamlss package in R)
Delta method	Yes	No	No	Fixed-2 fractional polynomials (gamlss)
Bootstrap CI	Yes	Yes	No	Fixed-2 fractional polynomials (gamlss)
Bootstrap duration CI	Yes	Yes	Yes	Both standard (mfp package) and fixed-2 (gamlss) fractional polynomials

Delta method CI

To compare two specific points on the duration-response curve, we need to estimate a confidence interval around the difference in outcome between the two points. In our application, we want to compare each shorter duration D_i against the longest duration D_{MAX} ; hence, we want to estimate the confidence intervals around $\pi_{D_{MAX}} - \pi_{D_i}$ for every D_i (including any integer duration randomised to or not randomised to below D_{MAX}). Let us call these confidence intervals $CI_{D_{MAX}-D_i}^{1-\alpha} = (CI_{D_{MAX}-D_i}^l, CI_{D_{MAX}-D_i}^u)$, with α being the confidence level. These can be estimated using the delta-method, which gives a good approximation based on Taylor series expansions (details in Appendix A). Inference can then be drawn by selecting the minimum duration D^* for which $CI_{D_{MAX}-D^*}^u < \delta$.

Bootstrap CI

Alternatively, rather than relying on these approximations, given the observed dataset with N observations, we can sample with replacement N observations M times, generating M bootstrap samples. We then fit the fractional polynomial model on each of the M datasets and estimate in each bootstrap sample $\pi_{D_{MAX}} - \pi_{D_i}$ for every D_i . We then calculate bootstrap confidence intervals around these quantities and use these to choose the optimal D^* .

Bootstrap duration CI

Another way of using bootstrapping is to directly estimate a confidence interval around D^* . M bootstrap samples are selected similarly to the previous section, but instead of estimating $\pi_{D_{MAX}} - \pi_{D_i}$ in each sample, we estimate the corresponding D^* , denoted D_m^* , with m indexing the bootstrap sample. Then a bootstrap confidence interval $(D^{*,l}, D^{*,u})$ can be constructed from the bootstrap mean estimate \widehat{D}^* and its standard error, and the recommended duration would be

$$D_I = \min_i (D_i > D^{*,u})$$

Theoretical comparison of methods

Table 1 provides an overview of the properties of different methods. The attractiveness of the confidence bands method comes from its simplicity, as it is probably the method most researchers would naturally use to estimate D^* from the duration-response curve. However, it has several limitations, the most important being that the pointwise confidence bands for the curve are not the same as the confidence interval for the difference between two specific points on the curve.

While the delta method addresses this problem, it is affected by at least one other issue, at least when using fractional polynomials⁵ as the flexible regression method: specifically, that model selection variability should be taken into account.⁶ The two bootstrap methods are theoretically appealing strategies to solve this problem, as repeating the fractional polynomial selection step for each bootstrap sample is one approach to address model selection variability.

With the delta method and the bootstrap CI method, we estimate multiple confidence intervals around the curve and compare each upper bound against the maximum tolerable risk difference. Hence, we may theoretically run into a multiplicity issue. However, our repeated tests are performed to solve an equation, rather than to formally compare different duration arms, making this less problematic. Bauer et al.⁷ discussed a similar version of this serial gate-keeping strategy, showing that it has strong control over type 1 error.

When the model used to estimate the duration-response curve is correct, we expect the bootstrap duration CI method to estimate a confidence interval for D^* that covers the true value at the nominal level. The assumption about model correctness is important; our choice of fractional polynomials as the preferred analysis method was originally driven by the fact that in many situations, we are not confident what the true underlying model is, and hence flexible models are preferred. While the standard fractional polynomial algorithm was built as a parsimonious modelling technique,⁸ in our application, with a single covariate and a reasonable number of expected events, a modi-

fied algorithm selecting *exactly* two polynomial terms rather than a *maximum of* two is likely to be preferable and makes it easier to satisfy the assumption that the model used is approximately correct. Henceforth, we refer to this method as the fixed-2 fractional polynomial analysis.

Alternative objectives

Thus far, we assumed that the trial objective was to find the minimum duration D^* such that $\pi_{D^*} \geq \pi_{D_{MAX}} - \delta$. This clearly corresponds with standard non-inferiority trials designed with margin δ_{NI} on the risk difference scale: in both, we assume that if the cure rate in the intervention arm is poorer only within a certain limit, then its secondary advantages make it preferable to the control arm. Now we discuss different objectives. Some are again in common with standard non-inferiority trials, while others are specific to our design.

Fixed rate. Instead of considering a fixed difference in cure rate compared to control, one might simply want to identify the duration D^* leading to a specific cure rate π^*

$$D^* = \min_{D_i} (\pi_{D_i} \geq \pi^*)$$

This may be reasonable when there is good prior information on the expected cure rate $\pi_{D_{MAX}}$ at the control duration D_{MAX} , although one can never be sure that the population recruited will reflect historical controls. The bootstrap CI and bootstrap duration CI methods could be used as analysis methods for this estimation target.

Fixed risk ratio. Analogously to a standard non-inferiority trial, the margin δ could also be defined on a relative scale; for example, as the proportion of the cure rate at D_{MAX} that should be preserved to consider the shorter duration preferable

$$\min_{D_i} (\pi_{D_i} \geq \delta \cdot \pi_{D_{MAX}})$$

For example, if we wanted to preserve 90% of the effect of treatment at D_{MAX} , we would choose $\delta = 0.9$.

The acceptability frontier. While non-inferiority trials naturally have a single non-inferiority margin δ_{NI} , in DURATIONS trials it might instead be logical to have different margins for each specific duration D_i . For example, for a duration that is a half of D_{MAX} , we might be happy to tolerate a slightly larger δ than for a duration equal to two-thirds of D_{MAX} , assuming that advantages of shorter durations generally increase as the duration reduces. Hence, the objective of the trial might be to find the duration D^* such that

$$D^* = \min_{D_i} (\pi_{D_i} \geq \pi_{D_{MAX}} - \delta(D_i))$$

where $\delta(D_i)$ is the function that indicates the acceptable loss in cure rate for each duration D_i . We call this the *Acceptability Frontier*, a possible example of which is the grey line in Figure 1.

Maximum acceptable gradient. Instead of targeting a specific cure rate, or difference in cure rate, investigators might be interested in the gradient of the duration-response curve. If we expect the curve to asymptote after a certain duration D^* , then we might define such D^* as the point after which the gradient of the function is always below a certain threshold δ

$$D^* = \min_{D_i} (\nabla \pi(D) \leq \delta, \forall D > D_i)$$

For example, in the scenarios above, we found that $\delta = 2\%$ was a reasonable value for this threshold.

Simulations

We compared the methods presented above in a simulation study designed using the recently proposed ‘Aims-Data generating mechanisms-Estimands-Methods-Performance measure’ framework.⁹

Aims. The main aim is to compare different strategies of drawing inference from the duration-response curve. Specific questions are: are bootstrap methods necessary, or does the simpler delta-method suffice? Is the fixed-2 fractional polynomial analysis preferable? How problematic is the multiplicity issue with the bootstrap CI method?

Data-generating mechanisms. The DURATIONS design aims to be resilient to the true underlying duration-response relationship. Consequently, we generate data under sixteen different scenarios that are listed and plotted in the additional material online, including the eight scenarios in (Quartagno et al, 2018). These reflect a wide range of possible duration-response relationships, including both those generated by fractional polynomials and those generated from sigmoid functions (which are not strictly within the fractional polynomial paradigm). Across scenarios, the optimal duration’s proximity to the nearest integer varies, to explore the effect this may have on results. For comparisons, we re-scale the x-axis so that the minimum duration considered (D_{MIN}) is 8 days, and the maximum D_{MAX} is 20 days, with seven randomised arms evenly spread between them (8, 10, 12, 14, 16, 18 and 20). We generate 1000 data sets from each scenario of $N = 500$ individuals. Using the same process, we additionally generate 200 datasets of $N = 750$ and $N = 1000$ individuals, to explore the sensitivity of results to total sample size.

Estimation targets. We initially assume the estimation target of interest is the duration D^* leading to a cure rate that is $\delta\%$ less than at the longest (and currently used) duration D_{MAX} . This approximately corresponds to the shortest duration non-inferior to standard-of-care within a non-inferiority margin δ on the risk difference scale. In our base case, we consider $\delta = 10\%$, a margin that has been recommended for example by the US Food and Drug Administration (FDA) for non-inferiority trials of antibiotics for adult bacterial community-acquired pneumonia.¹⁰ However, as in some applications, this may be considered too large a margin, we additionally explore the sensitivity of results using $\delta = 5\%$. This risk difference-based estimation target might not be optimal in all applications, and hence later we explore other estimation targets.

Methods. We compare the four different methods of estimating \widehat{D}^* described above. For the bootstrap methods, we use $M = 500$ replications. For the bootstrap CI method, we build confidence intervals using the BCa method,¹¹ which is based on the percentiles of the bootstrap distribution, but with adjustments to account for both bias and skewness. For the bootstrap duration CI method, we use the percentile confidence interval instead, as BCa is not applicable in some scenarios where the bootstrap distribution is truncated at 8 days. In addition, we compare the standard (i.e. selecting up to two polynomial terms) and fixed-2 (i.e. selecting exactly two terms) fractional polynomial algorithms for the bootstrap duration CI method.

Performance measures. The choice of performance measures is not straightforward. Unlike most randomised trial designs, our DURATIONS randomised design does not define a null hypothesis to test. This is because the longest treatment duration we randomise patients to is already known to be (usually highly) effective, and our aim is only to find the most appropriate shorter duration. Under very simple assumptions, that is, that the duration-response curve is continuous and monotonic, we know that either there is an integer duration D_I in the interval $[D_{MIN}; D_{MAX}]$ such that $\pi_{D_I} = \pi_{D_{MAX}} - \delta$, or that even the shortest duration is sufficiently effective, that is, $\pi_{D_I} = \pi_{D_{MIN}} \geq \pi_{D_{MAX}} - \delta$.

Given this, and having defined our estimation target, we define three main performance measures of interest:

1. **Type-1 Error:** in this context, we define type-1 error as the probability that our trial ends up recommending a duration \widehat{D}_I that is insufficiently effective.

$$Type1Error = P(\pi_{\widehat{D}_I} < \pi_{D_{MAX}} - \delta)$$

2. **Optimal Power:** there are several ways to define power compared to standard hypothesis testing. We define *optimal power* as the probability that the trial ends up recommending the actual optimal duration, that is the minimum effective (integer) duration

$$OptimalPower = P(\widehat{D}_I = \min_{D_i} (\pi_{D_i} \geq \pi_{D_{MAX}} - \delta))$$

3. **Acceptable Power:** designing a trial to reach high levels of Optimal Power might be difficult, but one might be interested in simply finding an effective duration that is shorter than the recommended one, even if it is not necessarily the *minimum* effective duration. We define the *acceptable power* as the probability that our trial identifies one such duration

$$AcceptablePower = P(\pi_{\widehat{D}_I} \geq \pi_{D_{MAX}} - \delta)$$

In addition, we measure performance by considering the distribution of all the recommended durations \widehat{D}_I as this allows us to quantify how far from the optimal duration we are both when a type I error occurs or when we recommend an acceptable (but not optimal) duration.

Different estimation targets. We performed additional simulations to explore the operating characteristics of the bootstrap duration CI method with each of the different estimation targets. We analysed the same 1000 datasets generated for the 16 scenarios in our base case simulations, with $N = 500$. As the target cure rate, we chose in each scenario the value 10% less than the true cure rate at the longest duration. For the fixed risk ratio, we used $\delta = 0.9$. The acceptability frontier was a linear function, such that the largest acceptable risk difference at D_{MIN} was 10%, and at 18 days 5% (Figure 1). Finally, the maximum acceptable gradient was taken to be 2%. For this last estimation target, we did not use bootstrap to estimate uncertainty, but we simply considered the single point estimate in the data.

Results

Figure 2 compares type-1 error and acceptable power across the different methods. It is immediately clear that type-1 error is not adequately controlled under certain scenarios. These scenarios are mainly those where the curvature at the optimal duration is positive, that is, those for which steepness of the curve is increasing at the optimal duration. These are arguably not the most likely scenarios in our settings, as we expect to be investigating part of the duration-response curve where the curve is asymptotic; nevertheless, it is preferable to use bootstrap methods that provide better inference by

Table 2. Results from using the Bootstrap CI and Bootstrap duration CI methods, with the fixed-2 fractional polynomial analysis method and base-case design parameters (estimation target = shortest duration non-inferior to 20 days within 10% risk difference). In italics: scenarios for which the upper bound of the Monte-Carlo confidence interval for type-I error was not strictly controlled within 2.5%, with Monte-Carlo confidence interval. See additional material Supplemental Tables A-B for other methods.

	Acceptable power (%)	Optimal power (%)	Type-I error (%) (target value: 2.5%)	True min duration	Estimated minimum duration recomm.	Estimated 2.5th Perc duration recomm.	Estimated median duration recomm.
Bootstrap CI:							
Scenario 1	98.2	9.8	1.8	13.1	12	14	16
Scenario 2	99.2	16.5	0.8	14.5	14	15	16
Scenario 3	94.5	10	1.6	15.9	15	16	18
Scenario 4	100	82.7	0	8.0	8	8	8
Scenario 5	99.6	5.4	0.4	9.7	9	10	12
Scenario 6	99.7	3.2	0.3	10.8	10	11	14
Scenario 7	94.4	30.7	4.3 (3.0,5.6)	16.2	14	16	18
Scenario 8	87	12.9	10.1(8.2,12)	15.0	8	15	17
Scenario 9	100	2.6	0	12.6	13	13	14
Scenario 10	99.9	0.8	0	15.2	16	17	18
Scenario 11	79	8.6	0.1	16.8	16	17	18
Scenario 12	98.9	29.9	1.1	11.2	11	12	13
Scenario 13	95.5	26.6	4.5 (3.2,5.8)	8.1	8	8	10
Scenario 14	56.2	8.6	3.8 (2.6,5.0)	15.0	8	15	17
Scenario 15	99.3	4.8	0.7	12.5	12	13	15
Scenario 16	99.8	4.9	0.2	12.0	11	12	14
Bootstrap duration CI:							
Scenario 1	97.7	9.9	2.3	13.1	11	14	16
Scenario 2	99.7	14.4	0.3	14.5	14	15	16
Scenario 3	99	5.3	1	15.9	14	16	18
Scenario 4	100	86.1	0	8.0	8	8	8
Scenario 5	99.9	5.4	0.1	9.7	9	10	12
Scenario 6	99.7	3.5	0.3	10.8	10	11	14
Scenario 7	95.2	46.8	4.8 (3.5,6.1)	16.2	15	16	17
Scenario 8	93.3	15.7	6.7 (5.2,8.2)	15.0	8	15	17
Scenario 9	100	2.3	0	12.6	13	14	14
Scenario 10	100	0.6	0	15.2	16	17	17
Scenario 11	99.9	9.8	0.1	16.8	16	17	18
Scenario 12	99	40.3	1	11.2	11	12	13
Scenario 13	96.3	29	3.7 (2.5,4.9)	8.1	8	8	10
Scenario 14	95.5	5.7	4.5 (3.2,5.8)	15.0	8	15	18
Scenario 15	99.3	5	0.7	12.5	12	13	15
Scenario 16	99.8	4.6	0.2	12.0	11	12	14

taking into account model selection variability. Differences between the two bootstrap methods are less marked (Table 2, and Supplemental Figure b in additional online material), although type-I error is generally slightly lower using the bootstrap duration CI method.

In terms of the analysis model, the fixed-2 fractional polynomial method, which is used in all analysis methods except for the last reported method in Figure 2 (and Supplemental Table E in additional material), is preferable with regards to type-I error. This is because it selects the best fitting model, while standard fractional polynomials seek a parsimonious model. Similarly less satisfactory results were reported using standard fractional polynomials with the bootstrap CI method as well (results not shown). There are only a few scenarios

under which the two bootstrap methods fail to control the type-1 error within 2.5%, namely scenarios 7, 8, 13 and 14 (Figure 3). However, this figure shows that the difference between the estimated minimum duration and the actual minimum acceptable duration is small (≤ 1 day) in these scenarios and that the 2.5th percentiles of recommended durations (and associated cure rate) are very close to the optimal duration (and corresponding cure rate).

In terms of power, differences are not as pronounced, and all methods achieve very good acceptable power ($>90\%$) under most scenarios. Of note, the simulated sample size ($N = 500$) was determined to estimate the duration-response curve within a certain average absolute error (5%), and not to achieve a specific power under any of our scenarios. As expected,

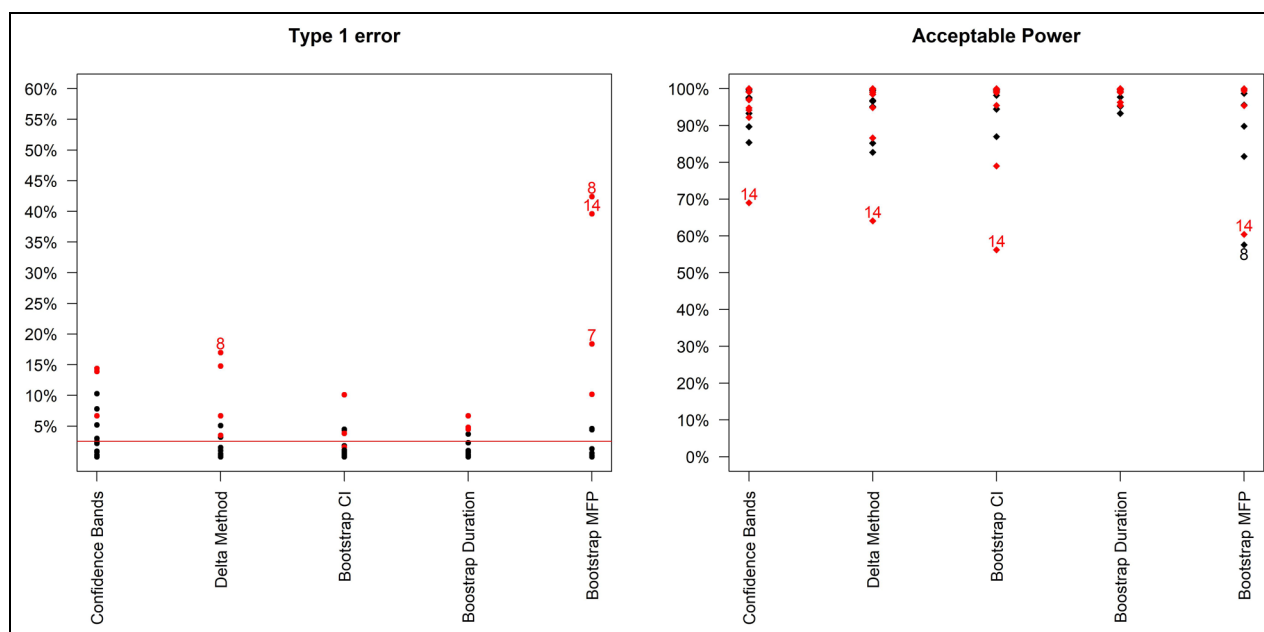


Figure 2. Type-I error and acceptable power (probability of recommending any sufficiently effective shorter duration) of the 5 analysis methods across the 16 simulation scenarios. Bootstrap MFP uses the Bootstrap duration CI method, but using the standard fractional polynomial approach as the analysis method (as in R package mfp). Scenarios leading to type I error $> 15\%$ or Acceptable Power $< 70\%$ are indexed. In addition, scenarios where the curvature is positive at the optimal duration are in red.

optimal power is always substantially lower, particularly in scenarios where the cure rate at the optimal duration is very close to the minimum acceptable cure rate. Again, differences between different methods are not substantial.

The only scenario where optimal power exceeds 80% is Scenario 4 (Table 2, Supplemental Tables A-D in additional material), that is, constant cure rate at every duration. In this scenario, optimal power can be broadly interpreted as the power of a non-inferiority trial of D_{MAX} against D_{MIN} with a $\delta\%$ non-inferiority margin, where the two durations have the same true cure rate.

Finally, conclusions are similar using $\delta = 5\%$ (Supplemental Figure a in additional material), although power is lower in all scenarios. In particular, acceptable power remains quite low for Scenario 14, only crossing the 80% threshold for $N = 1000$. Optimal power remains very low for all scenarios and sample sizes.

Different estimation targets. Results (Supplemental Tables D-G and Figure c in the additional material) show reasonable performance for the first three alternative estimation targets. For the maximum gradient target, type-1 errors are large, and further analysis refinements would be necessary to use this.

Analysis example

Given the simulation results, here we sketch our favoured three-step approach to analysing a

DURATIONS trial. First, an acceptability frontier should be defined, answering the question ‘what would the non-inferiority margin be for a trial comparing the longest duration to each shorter one?’. In this example, we assume a reasonable non-inferiority margin is 10% cure rate difference compared to the control duration, as in the base-case scenario of our simulations (Figure 4). Second, we run our fixed-2 fractional polynomial algorithm to estimate the duration-response curve (black solid line in Figure 4). Third, we use the bootstrap (BCa method)¹¹ to find the confidence interval either around our estimated optimal duration \widehat{D}^* or the differences in cure rate from the control at each duration (left and right panel of Figure 4 respectively). In this example, both methods recommend $\widehat{D}^* = 13$ days as optimal duration. Both methods are valid; the Bootstrap Duration CI method (left panel) led to slightly better results in our simulations, but the Bootstrap CI method (right panel) has the advantage that it can be compared against any non-inferiority margin or acceptability frontier.

Discussion

In this article, we have compared different strategies for drawing inference from a duration-response curve estimated using a DURATIONS randomised trial design. We defined quantities analogous to type-1 error and power in this scenario and found a method based on bootstrap re-sampling – to estimate the duration

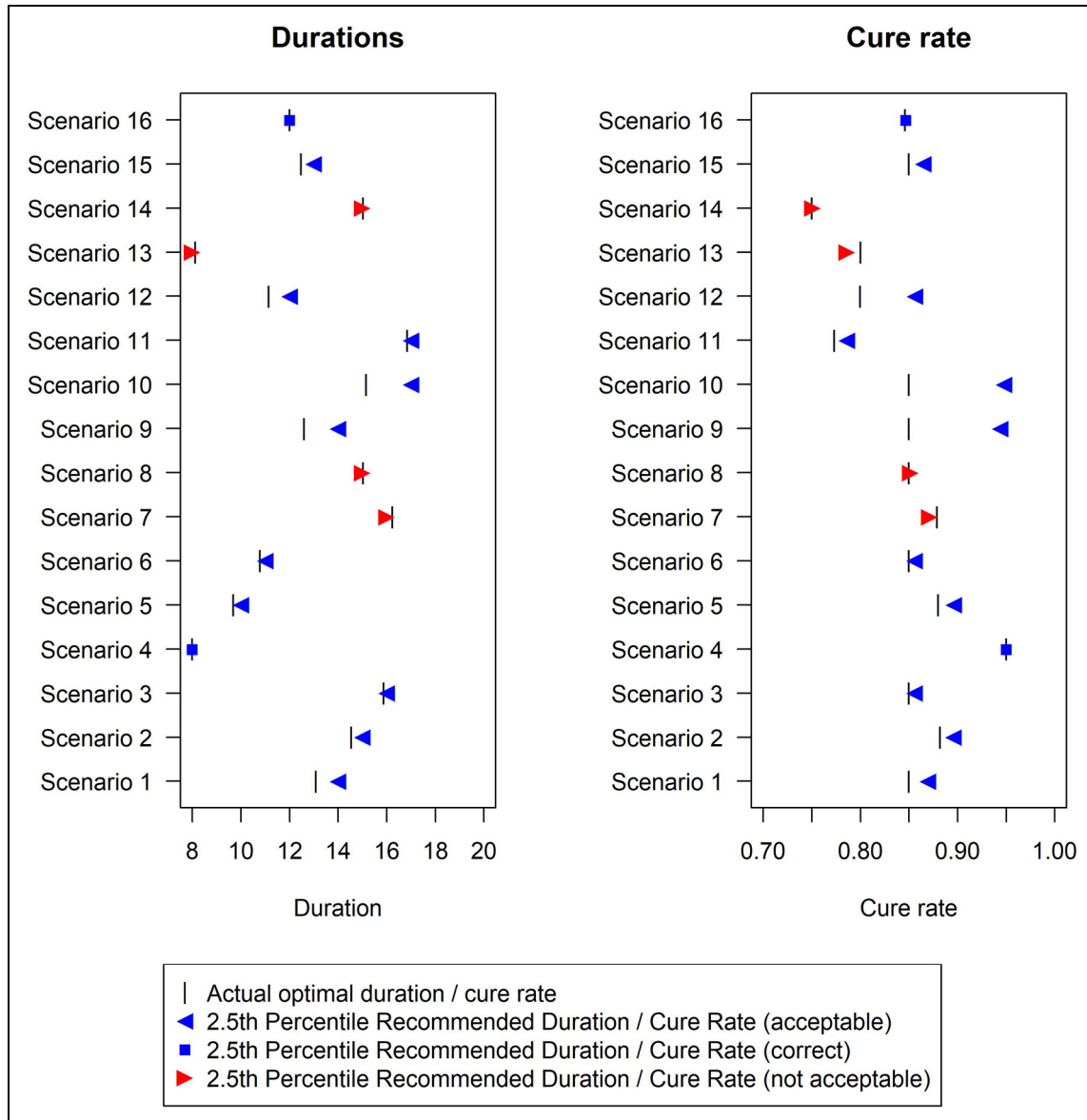


Figure 3. Duration recommended and associated cure rate across the 16 simulation scenarios using the Bootstrap duration CI method, with base-case design parameters (estimation target = shortest duration non-inferior to 20 days within 10% risk difference). The vertical bar indicates the true minimum effective duration.

associated with a specific cure rate difference from control – has good inferential properties when combined with a fixed-2 fractional polynomial analysis method. This is therefore our recommended approach.

One issue with the standard non-inferiority design for identifying optimal treatment duration is the potential for so-called ‘bio-creep’, that is, the erosion of efficacy from sequentially testing for non-inferiority shorter durations, iteratively updating the control duration to one previously shown to be non-inferior.¹² One advantage of the DURATIONS design is that it avoids this problem, as long as all the durations are evaluated in the same trial. Another advantage is its resilience; in

a standard non-inferiority trial, whenever a single design parameter turns out to have been badly misjudged, the whole trial can quickly lose power or interpretability. By contrast, the DURATIONS design has been developed to be flexible enough to maintain good properties against a wide range of duration-response curves.

Extensions

Design and analysis of randomised trials often intersect, and hence what is an analysis decision (how to analyse and draw inference from the observed data) also has

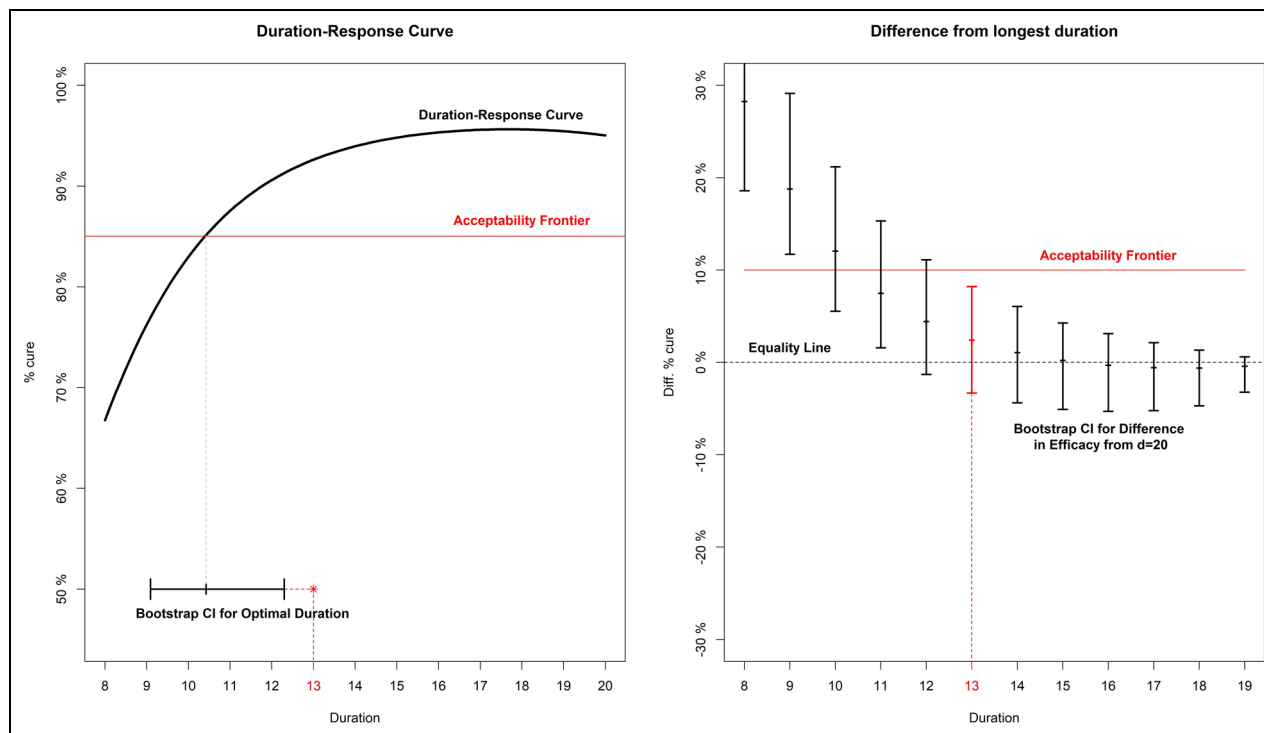


Figure 4. Analysis example for a hypothetical trial. On the left panel, the duration-response curve is estimated and then a bootstrap CI is built around the point where it crosses the acceptability frontier. On the right panel, bootstrap CIs are built around the difference in efficacy (cure rate) between each arm and the longest ($d = 20$).

design implications (how to best design a trial that we aim to analyse in a particular way). Hence, future work will investigate how to design a DURATIONS randomised trial that we aim to analyse with methods presented here, particularly considering the key challenge that the error rates depend on the true duration-response curve, which is unknown at the point of design.

We focused on binary outcomes, but similar methods could be easily used for continuous and survival outcomes, with the only additional complication of having to derive a suitable estimation target. When developing the target based on the acceptability frontier, we assumed that this frontier was subjectively drawn by the analyst, based on assumptions about the trade-offs between shortening duration and loss of treatment efficacy. An alternative could be to derive the acceptability frontier using available data on the secondary advantages of shorter durations. For example, if the goal was to design a DURATIONS trial in Hepatitis C, the acceptability frontier could be built using data on costs.

Although the methods presented here were motivated by trials whose aim was to optimise treatment duration in large Phase III-type trials, the same methods could be used to optimise treatment dose in smaller Phase II-type trials. There is a large literature on

methods for these smaller dose-finding trials, and some of the methods introduced in those settings could be adapted to our design. For example, several papers have investigated dose selection using gatekeeping strategies.^{7,13,14} One important difference is that in our application, there is generally already an accepted maximum duration, D_{MAX} , which the aim is to reduce. Of note, our approach uses a ceiling function to round the optimal duration estimated, a conservative choice to control type I error, possibly at the cost of losing some optimal power. In phase II settings, more liberal choices (e.g. standard rounding) could be made.

Future work could include investigating the effect of including additional covariates in the fractional polynomial model, for example age or sex, if there was evidence that the optimal duration might vary depending on these factors. It could also investigate using an adaptive design, to allow for closure of poorly performing arms (i.e. those at the lowest durations).⁴

For tuberculosis and related settings, where the optimal duration might be investigated for a new drug or new regimen, it is important to investigate the best way to include a formal comparison with an independent control treatment of fixed duration, for example with the standard 6-month tuberculosis treatment course with rifampicin, isoniazid, pyrazinamide and ethambutol.

Conclusion

We recently proposed the DURATIONS randomised trial design as an alternative to the standard two-arm non-inferiority design when the goal is to optimise treatment duration. In this article, we have investigated the operating characteristics of various methods of drawing inference from the duration-response curve and found that a method based on bootstrap to estimate a duration associated with a specific cure rate has good properties and is an appealing choice. Using this analysis method, DURATIONS randomised trials could help identify better treatment durations in an optimal way across many illnesses.⁴

Code

The code used for the simulations will be made available on the GitHub page of the first author.


Declaration of conflicting interests

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

Funding

The author(s) disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: This work was supported by the Medical Research Council [MC_UU_12023/29]. ASW is an NIHR Senior Investigator. The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR, or the Department of Health.

ORCID iD

Matteo Quartagno  <https://orcid.org/0000-0003-4446-0730>

Supplemental material

Supplemental material for this article is available online.

References

1. Llewelyn MJ, Fitzpatrick JM, Darwin E, et al. The antibiotic course has had its day. *BMJ* 2017; 358: j3418.
2. Quartagno M, Walker AS, Carpenter JR, et al. Rethinking non-inferiority: a practical trial design for optimising treatment duration. *Clin Trials* 2018; 15(5): 477–488.
3. Rehal S, Morris TP, Fielding K, et al. Non-inferiority trials: are they inferior? A systematic review of reporting in major medical journals. *BMJ Open* 2016; 6(10): 012594.
4. Pouwels KB, Yin M, Butler CC, et al. Optimising trial designs to identify appropriate antibiotic treatment durations. *BMC Med* 2019; 17: 1–7.
5. Royston P, Ambler G and Sauerbrei W. The use of fractional polynomials to model continuous risk variables in epidemiology. *Int J Epidemiol* 1999; 28(5): 964–974.
6. Faes C, Aerts M, Geys H, et al. Model averaging using fractional polynomials to estimate a safe level of exposure. *Risk Anal* 2007; 27(1): 111–123.
7. Bauer P, Röhm J, Maurer W, et al. Testing strategies in multi-dose experiments including active control. *Stat Med* 1998; 17: 2133–2148.
8. Royston P and Altman DG. Regression using fractional polynomials of continuous covariates: parsimonious parametric modelling. *J R Stat Soc Ser C* 1994; 43: 429–467.
9. Morris TP, White IR and Crowther MJ. Using simulation studies to evaluate statistical methods. *Stat Med* 2019; 39: 2074–2102.
10. Fleming TR and Powers JH. Issues in noninferiority trials: the evidence in community-acquired pneumonia. *Clin Infect Dis* 2008; 47: 108–120.
11. Carpenter J and Bithell J. Bootstrap confidence intervals: when, which, what? A practical guide for medical statisticians. *Stat Med* 2000; 19: 1141–1164.
12. Macfadden DR and Hanage WP. Potential for erosion of efficacy in non-inferiority trials of decreasing duration of antibiotic therapy. *Clin Infect Dis* 2019; 69(7): 1262.
13. Bretz F, Maurer W, Brannath W, et al. A graphical approach to sequentially rejective multiple test procedures. *Stat Med* 2009; 28(4): 586–604.
14. Dmitrienko A, Millen B, Brechenmacher T, et al. Development of gatekeeping strategies in confirmatory clinical trials. *Biom J* 2011; 53(6): 875–893.