



Published in final edited form as:

Stat Methods Med Res. 2020 July ; 29(7): 1891–1912. doi:10.1177/0962280219877520.

Sample size considerations for comparing dynamic treatment regimens in a sequential multiple-assignment randomized trial with a continuous longitudinal outcome

Nicholas J. Seewald¹, Kelley M. Kidwell², Inbal Nahum-Shani³, Tianshuang Wu⁴, James R. McKay⁵, Daniel Almirall^{1,3}

¹Department of Statistics, University of Michigan, Ann Arbor, Michigan, USA

²Department of Biostatistics, University of Michigan, Ann Arbor, Michigan, USA

³Institute for Social Research, University of Michigan, Ann Arbor, MI, USA

⁴AbbVie, North Chicago, Illinois, USA

⁵Department of Psychiatry, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania, USA

Abstract

Clinicians and researchers alike are increasingly interested in how best to personalize interventions. A dynamic treatment regimen (DTR) is a sequence of pre-specified decision rules which can be used to guide the delivery of a sequence of treatments or interventions that are tailored to the changing needs of the individual. The sequential multiple-assignment randomized trial (SMART) is a research tool which allows for the construction of effective DTRs. We derive easy-to-use formulae for computing the total sample size for three common two-stage SMART designs in which the primary aim is to compare mean end-of-study outcomes for two embedded DTRs which recommend different first-stage treatments. The formulae are derived in the context of a regression model which leverages information from a longitudinal outcome collected over the entire study. We show that the sample size formula for a SMART can be written as the product of the sample size formula for a standard two-arm randomized trial, a deflation factor that accounts for the increased statistical efficiency resulting from a longitudinal analysis, and an inflation factor that accounts for the design of a SMART. The SMART design inflation factor is typically a function of the anticipated probability of response to first-stage treatment. We review modeling and estimation for DTR effect analyses using a longitudinal outcome from a SMART, as well as the estimation of standard errors. We also present estimators for the covariance matrix for a variety of common working correlation structures. Methods are motivated using the ENGAGE study, a SMART aimed at developing a DTR for increasing motivation to attend treatments among alcohol- and cocaine-dependent patients.

Declaration of conflicting interests

The authors have no conflicting interests to declare.

Supplemental material

Supplementary material and the R code used to generate the simulation results in this paper is available at https://osf.io/q7zv8/?view_only=f6b35cea8d4a42a7bc369ed3a0c443c3

1 Introduction

Dynamic treatment regimens (DTRs) are sequences of pre-specified decision rules leading to courses of treatment which adapt to a patient's changing needs.¹ DTRs operationalize clinical decision-making by recommending particular treatments or intervention components to certain subsets of patients at specific times.² Consider the following example DTR which was designed to increase engagement with an intensive outpatient rehabilitation program (IOP) for patients with alcohol and/or cocaine dependence: "Within a week of the participant becoming non-engaged in the IOP, provide a phone-based session focusing on helping the patient re-engage in the IOP. At week 8, look back at the participant's engagement pattern over the past eight weeks. If the participant continued to not engage, provide a second phone-based session, this time focusing on facilitating personal choice (i.e., highlighting various treatment options the patient can choose from in addition to IOP). Otherwise, provide no further contact."³ Notice that the DTR recommends intervention strategies for both engaged and non-engaged participants at week 8. Alternative names for DTRs include adaptive treatment strategies^{4,5} and adaptive interventions,^{6,7} among others.

Scientists often have questions about how best to sequence and individualize interventions in the context of a DTR. Sequential, multiple-assignment, randomized trials (SMARTs) are one type of randomized trial design that can be used to answer questions at multiple stages of the development of high-quality DTRs.^{8,9,10} The characteristic feature of a SMART is that some or all participants are randomized more than once, often based on previously-observed covariates. Each randomization corresponds to a critical question regarding the development of a high-quality DTR, typically related to the type, timing, or intensity of treatment. SMARTs have been employed in a variety of fields, including oncology,^{11,12,13} surgery,^{14,15} substance abuse,¹⁶ and autism¹⁷.

Most SMARTs contain an embedded "tailoring variable", a pre-defined covariate observed during treatment which determines whether or how a participant will be randomized in the next stage of the SMART. For example, participants who "respond" to treatment may be re-randomized between different treatment options than participants who do not respond. SMARTs with embedded tailoring variables also contain embedded DTRs; that is, by design, participants in the SMART receive sequences of treatments which are consistent with the recommendations made by one or more DTRs. Note that SMARTs need not contain an embedded tailoring variable; however, we restrict our focus in this manuscript to those that do. We discuss this in more detail in section 2.

The comparison of two embedded DTRs which recommend different first-stage treatments is a common primary aim for a SMART.⁷ There exist data analytic methods for addressing this aim when the outcome is continuous,⁷ survival,¹⁸ binary,¹⁹ cluster-level²⁰ and longitudinal.^{21,22} A key step in designing a SMART, as with any randomized trial, is determining the sample size needed to be able detect a desired effect with given power. However, there is no existing method for determining sample size for such a comparison when the outcome is continuous and longitudinal.

Our primary contribution is tractable sample size formulae for SMARTs with a continuous longitudinal outcome in which the primary aim is an end-of-study comparison of two DTRs which recommend different first-stage treatments. Additionally, we present estimators for parameters in the working covariance matrix used in the analysis methods developed by Lu et al.²¹

In section 2, we provide a brief overview of three common SMART designs and introduce a motivating example. Section 3 reviews the estimation procedure introduced by Lu et al., and extends it by developing estimators for various working covariance structures.²¹ In section 4, we develop and present sample size formulae for SMARTs in which the primary aim is a comparison of two embedded DTRs which recommend different first-stage treatments using a continuous longitudinal outcome. The sample size formulae are evaluated via simulation in section 5.

2 Dynamic Treatment Regimens and Sequential Multiple-Assignment Randomized Trials

A DTR is a sequence of functions (“decision rules”), each of which takes as inputs a person’s history up to the time of the current decision (including baseline covariates, adherence, responses to previous treatments, etc.) and outputs a recommendation for the next treatment.¹⁰ Consider the example DTR in section 1. The recommended first-stage treatment is a phone-based session with a focus on re-engagement with the IOP. At week 8, each participant’s history of engagement is assessed, and an appropriate second-stage treatment is recommended. For participants who have shown a pattern of continued non-engagement, the recommended second-stage treatment is a second phone-based session focusing on personal choice. For all other participants, the DTR recommends no further contact. The tailoring variable is an indicator as to whether or not the participant demonstrated a pattern of continued non-engagement prior to week 8.

We consider two-stage SMARTs in which the primary outcome is continuous and repeatedly measured in participants over the course of the study. Our examples refer to trials in which at least one observation of the outcome is made in each stage, though that is not required for the estimation method presented in section 3. For simplicity, we refer to the tailoring variable as response status to first-stage treatment, and, in the second stage, we describe participants as “responders” or “non-responders”. We denote a DTR embedded in a SMART with a triple of the form (a_1, a_{2R}, a_{2NR}) , where a_1 is an indicator for the recommended first-stage treatment, a_{2R} an indicator for the second-stage treatment recommended for responders, and a_{2NR} the second-stage treatment recommended for non-responders. Throughout, (a_1, a_{2R}, a_{2NR}) is non-random and is used to index the DTRs embedded in a SMART.

We introduce three common two-stage SMART designs in figure 1 which vary in the subsets of participants who are re-randomized after the first stage. Each of these designs contains an embedded tailoring variable, and thus, for the purposes of this manuscript, contains embedded DTRs.

In design I, all participants are re-randomized. There are eight DTRs embedded in this design: for example, the DTR which starts by recommending A, then recommends C for responders and F for non-responders. Using the notation in figure 1, this DTR would be written $(1, 1, -1)$. SMARTs of this form have been run in the fields of drug dependence,^{23,24} smoking cessation,²⁵ and childhood depression,²⁶ among others.

SMARTs using design II restrict the second randomization to only non-responders; that is, only participants who have a certain value of the tailoring variable (here, “non-response”) are re-randomized. This is perhaps the most common SMART design, and it has been utilized in the study of ADHD,²⁷ adolescent marijuana use,²⁸ alcohol and cocaine dependence³, and more. There are four embedded DTRs in this design. Because responders are not re-randomized, a_{2R} is set to zero for all embedded DTRs.

In design III, re-randomization is restricted to only non-responders who receive a particular first-stage treatment. SMARTs of this type have been used to investigate cognition in children with autism spectrum disorder^{17,29} and implementation of a re-engagement program for patients with mental illness.³⁰ There are three DTRs embedded in this design. Note that, as in design II, responders are not re-randomized, so a_{2R} is set to zero for all embedded DTRs. Furthermore, a_{2NR} is set to zero when $a_1 = -1$, as non-responders to treatment B are not re-randomized.

For more information on various SMART designs and case studies for each type, see Lei, et al.³¹

To illustrate our ideas, we use ENGAGE, a SMART designed to study the effects of offering cocaine- and/or alcohol-dependent patients who did not engage in an IOP phone-based sessions either geared toward re-engaging them in an IOP or offering a choice of treatment options.³ The study recruited 500 cocaine- and/or alcohol-dependent adults who were enrolled in an IOP and failed to attend two or more sessions in the first two weeks. ENGAGE is modeled on design II. In the context of figure 1, treatment A was two phone-based motivational interviews focused on re-engaging the participant with the IOP (“MI-IOP”); treatment B was two phone-based motivational interviews geared towards helping the participant choose and engage with an intervention of their choice (“MI-PC”). Participants who exhibited a pattern of continued non-engagement after eight weeks were considered non-responders, and re-randomized to receive either MI-PC (treatments D and G) or no further contact (treatments E and H). Responders were provided no further contact (treatments C and F). Following the coding in figure 1, the example DTR from section 1 is labeled $(1, 0, 1)$.

An important continuous outcome in ENGAGE is “treatment readiness”. This is a measure of a patient’s willingness and ability to commit to active participation in a substance abuse treatment program. The score ranges from 8–40 and is coded so that higher scores indicate greater treatment readiness. Measurements are taken at baseline, and 4, 8, 12, and 24 weeks after program entry.

3 Estimation

We extend the work of Lu and colleagues by offering more detailed guidance on the estimation of model parameters used in computing quantities of interest on which to compare two embedded DTRs.²¹ We first review the method below.

3.1 Marginal Mean Model

Consider a SMART design with embedded DTRs labeled by (a_1, a_{2R}, a_{2NR}) . Suppose we have a longitudinal outcome $Y_i = (Y_{t_1, i}, \dots, Y_{t_T, i})$, $i = 1, \dots, n$, observed such that $Y_{t,i}$ is measured for each of n participants at each of T timepoints $\{t_j: j = 1, \dots, T; t_1 < \dots < t_T\}$. We do not require that these timepoints be equally-spaced, though they must be common to all participants in the study. Define $t^* \in \{t_j\}$ to be the time of the measurement taken immediately before the assessment of response status and second randomization. In ENGAGE, for example, $T = 5, \{t_j\} = \{0, 4, 8, 12, 24\}$, and $t^* = t_3 = 8$. Let X_i be a vector of mean-centered baseline covariates, such as age at baseline, sex, etc., for the i th individual.

We are interested in $E\left[Y_t^{(a_1, a_{2R}, a_{2NR})} | X\right]$, the marginal mean outcome at time t under DTR (a_1, a_{2R}, a_{2NR}) conditional on X . This is the mean outcome at time t had all individuals with characteristics X been offered DTR (a_1, a_{2R}, a_{2NR}) . Recall that a DTR recommends treatments for both responders and non-responders; therefore, $E\left[Y_t^{(a_1, a_{2R}, a_{2NR})} | X\right]$ is marginal over response status. Note that $Y_{t,i}^{(a_1, a_{2R}, a_{2NR})}$ is a potential outcome: the value of the outcome $Y_{t,i}$ that would be observed had participant i been treated according to the DTR (a_1, a_{2R}, a_{2NR}) .

We impose a modeling assumption on $E\left[Y_t^{(a_1, a_{2R}, a_{2NR})} | X\right]$ namely, that $E\left[Y_t^{(a_1, a_{2R}, a_{2NR})} | X\right] = \mu_t^{(a_1, a_{2R}, a_{2NR})}(X; \theta)$, where $\mu_t^{(a_1, a_{2R}, a_{2NR})}(X; \theta)$ is a marginal structural mean model with unknown parameters $\theta = (\eta^\top, \gamma^\top)^\top$. We use η to represent a column vector of parameters indexing baseline covariates, and γ is a column vector of coefficients on terms involving treatment effects; we discuss in more detail below. As noted by Lu and colleagues, the sequential nature of treatment delivery in SMARTs may suggest constraints on the form of $\mu_t^{(a_1, a_{2R}, a_{2NR})}(X; \theta)$ which depend, in part, on the design of the SMART.²¹ For instance, in ENGAGE, at time $t = 0$, no treatments have been assigned, so all DTRs share a common mean. At times $t = 4$ and $t = 8$, the four embedded DTRs differ only by recommended first-stage treatment; thus there are two means of $Y_t^{(a_1, a_{2R}, a_{2NR})}$ at each timepoint. Finally, for times $t > t^* = 8$, each DTR has a different mean $Y_t^{(a_1, a_{2R}, a_{2NR})}$.

An example marginal structural mean model for ENGAGE (and, more generally, design II) is

$$\mu_t^{(a_1, a_{2R}, a_{2NR})}(X_1; \theta) = \eta_1 X_1 + \gamma_0 + \mathbb{1}\{t \leq t^*\}(\gamma_1 t + \gamma_2 a_1 t) + \mathbb{1}\{t > t^*\}(\gamma_1 t^* + \gamma_2 t^* a_1 + \gamma_3(t - t^*) + \gamma_4(t - t^*) a_1 + \gamma_5(t - t^*) a_{2NR} + \gamma_6(t - t^*) a_1 a_{2NR}), \tag{1}$$

where $\mathbb{1}\{E\}$ is the indicator function for the event E .

Using contrast coding, i.e., $\{a_1, a_{2NR}\} \in \{-1, 1\}^2$, we can write

$$2\gamma_2 = E \left[\frac{Y_{t_j}^{(1,0, \cdot)} - Y_{t_k}^{(1,0, \cdot)}}{t_j - t_k} - \frac{Y_{t_j}^{(-1,0, \cdot)} - Y_{t_k}^{(-1,0, \cdot)}}{t_j - t_k} \mid \mathbf{X} \right], t_j, t_k \leq t^*. \tag{2}$$

This represents the difference in slopes of expected treatment readiness in the first stage of the SMART between DTRs starting with different first-stage treatments (second-stage treatment is arbitrary, as $t < t^*$). Also, we can interpret η_1 as the difference in expected outcome $Y_t^{(a_1, a_{2R}, a_{2NR})}$ associated with a one-unit difference in baseline covariate X_1 , marginal over all embedded DTRs.

We present example models for designs I and III in the online supplement. For more on modeling considerations for longitudinal outcomes in SMARTs, see Lu et al.²¹

3.2 Observed Data

Suppose we have data arising from a SMART with n participants. Let $A_{1,j} \in \{-1, 1\}$ be a random variable which indicates first-stage treatment randomly assigned to participant j ($j = 1, \dots, n$), and let $R_j \in \{0, 1\}$ indicate whether the j th participant responded to $A_{1,j}$, in which case $R_j = 1$, or not, so $R_j = 0$. Define $A_{2,j} \in \{-1, 1\}$ to be the randomly-assigned second-stage treatment. Throughout, we use uppercase A to denote random treatment assignments; lowercase a 's are non-random indices used to denote embedded DTRs.

In design II, since only non-responders are re-randomized, we set $A_{2,j} = 0$ for responders; similarly for design III. We observe a continuous outcome $Y_{t,j}$ for each participant at each of T timepoints. In general, the data collected on the j th individual over the course of the study are of the form

$$(X_j, Y_0, A_{1,j}, Y_{[0 < t \leq t^*], j}, R_j, A_{2,j}, Y_{[t > t^*], j}),$$

where $Y_{[u < t \leq v], j}$ is a vector consisting of all values of the outcome observed for the j th participant between times u and v .

3.3 Estimating Equations

Our goal is to estimate and make inferences on θ , the length- p column vector of mean parameters in the marginal structural mean model of interest. For notational convenience, let \mathcal{D} be the set of DTRs embedded in the SMART under study; for instance, in design II,

$$\mathcal{D} = \{(a_1, a_{2R}, a_{2NR}) : a_1 \in \{-1, 1\}, a_{2R} = 0, a_{2NR} \in \{-1, 1\}\}.$$

This creates, which can be corrected using inverse-probability weighting.^{7,32,2}

Let $W^{(d)}(A_{1,i}, R_i, A_{2,i})$ be a weight associated with participant i and DTR $d \in \mathcal{D}$ defined as

$$W^{(d)}(A_{1,i}, R_i, A_{2,i}) = \frac{I^{(d)}(A_{1,i}, R_i, A_{2,i})}{P(A_{1,i} = a_1)P(A_{2,i} = a_2 | A_{1,i} = a_1, R_i)}, \quad (3)$$

where $I^{(d)}(A_{1,i}, R_i, A_{2,i})$ is an indicator of whether participant i is consistent with DTR d . The form of $I^{(d)}(A_{1,i}, R_i, A_{2,i})$ depends on the particular SMART design under study; for each of the designs in figure 1, these expressions are shown in table 1.

We use $W^{(d)}(A_{1,i}, R_i, A_{2,i})$ to account for the facts that, in some SMARTs (e.g., designs II and III) there is known imbalance in the proportion of responders and non-responders consistent with each DTR, and that some (or all) participants are consistent with more than one embedded DTR.

In design II, for example, only non-responders to first-stage treatment are re-randomized; if all randomizations are with probability 0.5, $W^{(1,0,1)}(1, 1, 0) = (.5 \times 1)^{-1} = 2$ and $W^{(1,0,1)}(1, 0, 1) = (.5 \times .5)^{-1} = 4$. Note that in design I, all participants are re-randomized; hence, all participants receive a weight of 4. The analyst may freely substitute $W^{(d)}(A_{1,i}, R_i, A_{2,i}) = I^{(d)}(A_{1,i}, R_i, A_{2,i})$ in this case.

Define $D^{(d)}(X_i) \in \mathbb{R}^{T \times p}$ to be the Jacobian of $\mu^{(d)}(X_i; \theta)$ with respect to θ , i.e., $D^{(d)}(X_i) = \mu^{(d)}(X_i; \theta) / \theta^T$. Let $V^{(d)}(X_i; \tau) \in \mathbb{R}^{T \times T}$ be a working covariance matrix for $Y^{(d)}$, conditional on baseline covariates X , under DTR $d \in \mathcal{D}$. Here, $\tau = (\sigma^T, \rho^T)^T$ is a vector of parameters indexing variance (σ) and correlation (ρ) components of the working covariance structure. We discuss $V^{(d)}(X_i; \tau)$ in detail in section 3.4. We estimate θ by solving the estimating equations

$$0 = \frac{1}{n} \sum_{i=1}^n \sum_{d \in \mathcal{D}} \left[W^{(d)}(A_{1,i}, R_i, A_{2,i}) \cdot D^{(d)}(X_i)^T V^{(d)}(X_i; \tau)^{-1} (Y_i - \mu^{(d)}(X_i; \theta)) \right]. \quad (4)$$

We call the solution to equation (4) $\hat{\theta}$.

Under usual regularity conditions for M -estimators (see, e.g., van der Vaart, theorem 5.4.1)³³ and given data from a SMART (see appendix A), $\hat{\theta}$ is consistent for θ . Furthermore, $\sqrt{n}(\hat{\theta} - \theta)$ has an asymptotic multivariate normal distribution:

$$\sqrt{n}(\hat{\theta} - \theta) \Rightarrow \mathcal{N}(\theta, B^{-1} M B^{-1}),$$

where

$$\mathbf{B}: = \mathbb{E} \left[\sum_{d \in \mathcal{D}} W^{(d)}(A_{1,i}, R_i, A_{2,i}) \mathbf{D}^{(d)}(\mathbf{X}_i)^\top \mathbf{V}^{(d)}(\mathbf{X}_i; \boldsymbol{\tau})^{-1} \mathbf{D}^{(d)}(\mathbf{X}_i) \right] \in \mathbb{R}^{p \times p} \quad (5)$$

and

$$\mathbf{M}: = \mathbb{E} \left[\left(\sum_{d \in \mathcal{D}} W^{(d)}(A_{1,i}, R_i, A_{2,i}) \mathbf{D}^{(d)}(\mathbf{X}_i)^\top \mathbf{V}^{(d)}(\mathbf{X}_i; \boldsymbol{\tau})^{-1} (\mathbf{Y}_i - \boldsymbol{\mu}^{(d)}(\mathbf{X}_i; \boldsymbol{\theta})) \right)^{\otimes 2} \right] \in \mathbb{R}^{p \times p}, \quad (6)$$

with $\mathbf{Z}^{\otimes 2} = \mathbf{Z}\mathbf{Z}^\top$. Proofs of these claims are available in the supplement. Note that $\hat{\boldsymbol{\theta}}$ is consistent for $\boldsymbol{\theta}$ regardless of the chosen structure of $\mathbf{V}^{(d)}(\mathbf{X}; \boldsymbol{\tau})$; however, we conjecture that choices of $\mathbf{V}^{(d)}(\mathbf{X}; \boldsymbol{\tau})$ closer to the true covariance matrix $\text{Var}(\mathbf{Y}^{(d)})$ will yield more efficient estimates.

3.4 Estimation of the Working Covariance Matrix

Decisions regarding the structure of $\mathbf{V}^{(d)}(\mathbf{X}; \boldsymbol{\tau})$ should be made by the scientist according to existing knowledge regarding the within-person covariance structure of $\mathbf{Y}^{(d)}$. In general, for an embedded DTR $d \in \mathcal{D}$, $\mathbf{V}^{(d)}(\mathbf{X}; \boldsymbol{\tau})$ takes the form

$$\mathbf{V}^{(d)}(\mathbf{X}; \boldsymbol{\sigma}, \boldsymbol{\rho}) = \mathbf{S}^{(d)}(\boldsymbol{\sigma})^{1/2} \mathbf{R}^{(d)}(\boldsymbol{\rho}) \mathbf{S}^{(d)}(\boldsymbol{\sigma})^{1/2},$$

where $\mathbf{S}^{(d)}(\boldsymbol{\sigma})^{1/2} \in \mathbb{R}^{T \times T}$ is a diagonal matrix with diagonal entries $\sigma_{t_1}^{(d)}, \dots, \sigma_{t_T}^{(d)}$, where $\sigma_{t_j}^{(d)} = \text{Var}(Y_{t_j}^{(d)})$, and $\mathbf{R}^{(d)}(\boldsymbol{\rho}) \in \mathbb{R}^{T \times T}$ is working correlation matrix for $\mathbf{Y}^{(d)}$. Note that this notation allows for different working covariance structures for each DTR, as well as non-constant variances of the longitudinal outcome across time.

We propose the following procedure to estimate $\mathbf{V}^{(d)}(\mathbf{X}; \boldsymbol{\tau})$. First, estimate $\boldsymbol{\theta}$ by solving equation (4) using the $T \times T$ identity matrix as $\mathbf{V}^{(d)}(\mathbf{X}; \boldsymbol{\tau})$ for all $d \in \mathcal{D}$. Call the solution $\hat{\boldsymbol{\theta}}_{(0)}$. Next, use $\hat{\boldsymbol{\theta}}_{(0)}$ to estimate $\sigma_t^{(d)}$ as follows

$$\left(\hat{\sigma}_t^{(d)} \right)^2 = \frac{\sum_{i=1}^n W^{(d)}(A_{1,i}, R_i, A_{2,i}) \left(Y_{i,t} - \mu_t^{(d)}(\mathbf{X}_i; \hat{\boldsymbol{\theta}}_{(0)}) \right)^2}{\sum_{i=1}^n W^{(d)}(A_{1,i}, R_i, A_{2,i}) - p} \quad (7)$$

where p is the dimension of $\boldsymbol{\theta}$. If the scientist believes that this variance is constant over time for each DTR, the estimator in equation (7) can be averaged over time; one can also average over DTR if one believes the variance is constant across all embedded DTRs. Estimators for $\boldsymbol{\rho}^{(d)}$ vary with choice of correlation structure $\mathbf{R}^{(d)}(\boldsymbol{\rho})$; we present estimators for selected structures in table 2.

Note that the denominator in equation (7) must be positive. For a fixed DTR $d \in \mathcal{D}$, the sum of the weights is, in expectation, the total sample size n (see supplement). Therefore, the denominator is approximately $n - p$, as in the usual mean squared error in multiple

regression, for example. We recommend that analysts choose appropriately parsimonious marginal structural mean models (i.e., $p < n$) to ensure that equation (7) is positive.

To complete the estimation procedure, we again solve equation (4), this time using $\widehat{\mathbf{V}}^{(d)}(\mathbf{X}; \widehat{\boldsymbol{\tau}}) = \mathbf{S}^{(d)}(\widehat{\boldsymbol{\sigma}})^{1/2} \mathbf{R}^{(d)}(\widehat{\boldsymbol{\rho}}) \mathbf{S}^{(d)}(\widehat{\boldsymbol{\sigma}})^{1/2}$ as the working covariance matrix. This process can be further iterated, as suggested by Liang and Zeger;³⁴ we call the final estimate of the model parameters $\widehat{\boldsymbol{\theta}}$.

4 Sample Size Formulae for End-of-Study Comparisons

Often, longitudinal outcomes are collected in trials to improve the efficiency of the primary aim analysis, even when comparing two treatment groups on the mean of some summary measure such as the end-of-study observation.³⁵ Fitting a longitudinal regression model and using the result to estimate the difference in mean summary measure improves efficiency of the comparison by leveraging within-person correlation (see, e.g. Fitzmaurice et al., section 2.5).³⁶ Furthermore, this approach allows investigators to simultaneously address secondary aims using the same regression model.

As with standard randomized clinical trials, a common primary aim of a SMART is the comparison of mean end-of-study outcomes for two embedded DTRs which recommend different first-stage treatments^{12,18,31,37}. We now present sample size formulae for SMARTs with longitudinal outcomes in which this is the primary aim, addressed using the general estimation procedure of section 3; that is, using a regression approach which includes all observed outcome data. We restrict our focus to two-stage SMARTs in which the outcome is observed at three timepoints – baseline, just prior to the second randomization, and at the end of the study – and in which all randomizations occur with probability 0.5. Additionally, we consider a saturated, piecewise-linear mean structure $\boldsymbol{\mu}^{(d)}(\boldsymbol{\theta})$ similar to model (1).

Recall from section 3.1 that $\boldsymbol{\theta}$ is a p -vector of parameters which indexes a marginal structural mean model for the treatment effects in a SMART. Let \mathbf{c} be a length- p contrast vector so that the null hypothesis of interest takes the form

$$H_0: \mathbf{c}^\top \boldsymbol{\theta} = 0,$$

which we will test against an alternative of the form $H_1: \mathbf{c}^\top \boldsymbol{\theta} = \cdot$. To compare mean end-of-study outcomes between two embedded DTRs which recommend different first-stage treatments, the estimand of interest is

$$\mathbf{c}^\top \boldsymbol{\theta} = \mathbb{E} \left[Y_{t_T}^{(1, a_{2R}, a_{2NR})} - Y_{t_T}^{(-1, a'_{2R}, a'_{2NR})} \right], \quad (8)$$

for some choice of a_{2R} , a'_{2R} , a_{2NR} , and a'_{2NR} . For example, to test equality of mean end-of-study outcomes for DTRs (1, 0, 1) and (-1, 0, -1) in design II under model (1) (assuming no \mathbf{X} , $\{t_j\} = \{0, 1, 2\}$, $t^* = 1$), the estimand is the linear combination $\mathbf{c}^\top \boldsymbol{\gamma}$, where $\mathbf{c}^\top = (0, 0, 2, 0, 2, 2, 0)$.

We employ a 1-degree of freedom Wald test. The test statistic is

$$Z = \frac{\sqrt{nc}^\top \hat{\theta}}{\sigma_c},$$

where $\sigma_c = \sqrt{c^\top \mathbf{B}^{-1} \mathbf{M} \mathbf{B}^{-1} c}$. Under the null hypothesis, by asymptotic normality of $\sqrt{n}(\hat{\theta} - \theta)$, the test statistic follows a standard normal distribution.

Define δ to be the standardized effect size as described by Cohen for an end-of-study comparison, i.e.,

$$\delta = \frac{\Delta}{\sigma}, \quad (9)$$

where $\sigma = \text{Var}(Y_t^{(d)})$ (see working assumption A2 below).³⁸

In order to simplify the form of σ_c and obtain tractable sample size formulae, we make the following working assumptions (A1 and A2):

A1 Constrained conditional covariance matrices for DTRs under comparison.

- a.** The variability of $Y_t^{(d)}$ around the DTR mean $\mu_t^{(d)}(\theta)$ among non-responders is no more than the variance of $Y_t^{(d)}$ unconditional on response; i.e.,

$$\mathbb{E}\left[\left(Y_t^{(d)} - \mu_t^{(d)}(\theta)\right)^2 \mid R^{(a1)} = 0\right] \leq \mathbb{E}\left[\left(Y_t^{(d)} - \mu_t^{(d)}(\theta)\right)^2\right],$$

for all $t > t^*$ and DTRs $d \in \mathcal{D}$ under study.

- b.** For times $t_i < t_j < t^*$, response status is uncorrelated with products of residuals; i.e.,

$$\text{Cov}\left(R^{(a1)}, \left(Y_{t_i}^{(d)} - \mu_{t_i}^{(d)}(\theta)\right)\left(Y_{t_j}^{(d)} - \mu_{t_j}^{(d)}(\theta)\right)\right) = 0.$$

for DTRs $d \in \mathcal{D}$ under study.

- c.** The covariance between the end-of-study measurement and the measurements prior to the second stage among responders is less than or equal to the same quantity among non-responders:

$$\text{Cov}\left(Y_t^{(d)}, Y_{t_T}^{(d)} \mid R^{(a1)} = 1\right) \leq \text{Cov}\left(Y_t^{(d)}, Y_{t_T}^{(d)} \mid R^{(a1)} = 0\right)$$

for DTRs $d \in \mathcal{D}$ under study and $t < t^*$.

A2 Exchangeable marginal covariance structure.

The marginal variance of $\mathbf{Y}^{(d)}$ is constant across time and DTR, and has an exchangeable correlation structure with correlation ρ , i.e.,

$$\text{Var}(\mathbf{Y}^{(d)}) = \Sigma = \sigma^2 \begin{bmatrix} 1 & \rho & \cdots & \rho \\ \rho & 1 & \cdots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \cdots & 1 \end{bmatrix}_{T \times T}$$

for all $d \in \mathcal{D}$.

Note that the above are *working* assumptions which we do not believe hold in general. We will see in sections 5 and 6 that sample size formula (10) (given below) is robust to moderate violations of working assumption A1 and that inputs to the formula can be adjusted in a way to accommodate violations of working assumption A2. A working assumption similar to A1(a) is commonly made in developing sample-size formulae for SMARTs with outcomes collected once at the end of the study.^{39,19,20} Working assumptions A1(b) and A1(c) impose further constraints on the covariance of the outcome conditional on response and allow for tractable sample size formulae. We believe working assumption A1(b) is approximately satisfied in most common definitions of response (see the supplement). Working assumption A2 is not strictly necessary, but is used to simplify the sample size formulae and facilitate easier elicitation of parameters. See section 6 for more discussion.

Working assumption A1 arises specifically as a consequence of unequal weights in equation (4) (i.e., when there exists imbalance between responders and non-responders, by design); therefore, the assumption is not necessary in design I, and can be relaxed to apply to only the two DTRs in which non-responders are re-randomized in design III. See appendix B for more details on how this assumption is used. Furthermore, working assumption A2 cannot be satisfied in design I if all eight embedded DTRs have unique means.

Under working assumptions A1 and A2, the minimum-required sample size to detect a standardized effect size δ with power at least $1 - \beta$ and two-sided type-I error α is

$$n \geq \frac{4(z_{1-\alpha/2} + z_{1-\beta})^2}{\delta^2} \cdot (1 - \rho^2) \cdot \text{DE}, \quad (10)$$

where DE is a SMART-specific “design effect” for an end-of-study comparison (see table 3). These design effects are functions of response rates $P(R^{(a_1)} = 1) = r_{a_1}$; if researchers do not have well-informed estimates of these probabilities, they may use a conservative design effect in which $r_{a_1} = 0$ for $a_1 \in \{-1, 1\}$.

Note that the first term in formula (10) is the typical sample size formula for a traditional two-arm randomized trial with a continuous end-of-study outcome and equal randomization probability. The middle term is due to the within-person correlation in the outcome, and is identical to the corresponding correction term for GEE analyses sized to detect a group-by-time interaction when there is no baseline group effect (see, e.g., Fitzmaurice et al., ch. 20³⁶).

5 Simulations

We conducted a variety of simulations to assess the performance of sample size formula (10). We are interested in the empirical power for a comparison of the DTR which recommends only treatments indicated by 1 and the DTR which recommends only treatments indicated by -1 when the study is sized to detect an effect size of δ . In ENGAGE, this might correspond to a comparison of mean end-of-study outcomes under the DTR which recommends MI-IOP in the first stage, no further contact for engagers, and MI-PC in the second stage for continued non-engagers versus the mean end-of-study outcomes under the regimen which recommends MI-PC in the first-stage, then no further contact for both engagers and non-engagers.

We consider four types of scenarios: first, when no assumptions are violated; second, when each of working assumptions A1(a) to A1(c) are violated; finally, when the working correlation structure is misspecified, in violation of working assumption A2. In each scenario, sample sizes are computed based on nominal power $1 - \beta = 0.8$ and two-sided type-I error $\alpha = 0.05$.

We believe sample sizes from formula (10) will be slightly conservative when all assumptions are satisfied, as formula (10) is an interpretable upper bound on a sharper formula given in appendix B and the supplement. For design I, we do not expect power to be affected by violations of working assumption A1, as the assumption arises as a consequence of over- or under-representation of responders and non-responders consistent with a particular DTR (see appendix B). Since there is no such imbalance in design I, working assumption A1 is not applicable. Similarly, in design III, only non-responders to one first-stage treatment are re-randomized, so we expect that empirical power will decrease slightly, but not seriously, when violating working assumption A1. We expect empirical power to suffer most severely when violating this working assumption in design II.

We further conjecture that scenarios in which the true within-person correlation structure of $\mathbf{Y}^{(d)}$ is autoregressive, sample sizes from formula (10) will be very anti-conservative. Under an AR(1) correlation structure, less information about the end-of-study outcome is provided by the baseline measure than would be under an exchangeable correlation structure. Since, by using formula (10), we have assumed more information is available from earlier measurements than is actually the case, we will be underpowered. Similarly, we expect over-estimation of ρ in formula (10) to lead to anti-conservative sample sizes.

5.1 Data Generative Process

For each simulation, the true marginal mean model is as in model (1) for design II; analogous models are used for designs I and III – see the supplement for examples. We do not include baseline covariates \mathbf{X} ; this is a conservative approach, as we believe that adjustment for prognostic covariates typically will increase power (see, eg., Kahan et al.⁴⁰). Estimates of marginal means from ENGAGE were used to inform a reasonable range of “true” means from which to simulate, though the scenarios presented here are not designed to mimic ENGAGE exactly. All simulations take $T = 3$ and values of γ and σ are chosen to achieve $\delta = 0.3$ or $\delta = 0.5$ (“small” and “moderate” effect sizes, respectively).

Data were generated according to a conditional mean model which, when averaged over response, yields the marginal model of interest. Potential outcomes $Y_{t,i}^{(d)}$ were simulated from appropriately-parameterized normal distributions (see appendix C for details); data were “observed” by selecting the potential outcome corresponding to treatment assignment as generated from a Bernoulli(0.5) distribution.

We consider three mechanisms for generating response status. In the first, “ R_{\perp} ”, response is generated from a Bernoulli (r_{a_1}) distribution, where $r_{a_1} = P(R^{(a_1)} = 1)$, independently of all previously-observed data. In the second and third scenarios (“ R_+ ” and “ R_- ”, respectively), response status is still generated from a Bernoulli distribution, but each individual is assigned a probability of response correlated with their observed value of Y_1 . These correlations are either positive or negative, depending on the response model. This is intended to mimic different coding choices for Y , in the sense of responders tending to have higher or lower values of Y_1 than non-responders. For details of how these are generated, see appendix C, which also contains additional details regarding the data generative models used. In the supplement, we present simulation results under additional models for response.

For each simulation scenario, we compute upper and lower bounds on allowable values of $\text{Var}\left(Y_2^{(d)} \mid R^{(a_1)} = 1\right)$, beyond which it is not possible to either achieve the desired marginal variance, or which induces violation of working assumption A1(a). The results shown in the corresponding column of table 4 were generated when responders’ variances were set to 75% of the lower bound beyond which the fixed marginal variance forces

$$E\left[\left(Y_t^{(d)} - \mu_t^{(d)}(\theta)\right)^2 \mid R^{(a_1)} = 0\right] \geq \sigma^2.$$

Violation of working assumption A1(b) was induced by defining response status as

$$R^{(a_1)} = 1 \left\{ Y_1^{(d)} \in \left(-\infty, \kappa_{a_1}^{\text{low}}\right] \cup \left[\kappa_{a_1}^{\text{high}}, \infty\right) \right\}, \quad (11)$$

where κ^{low} and κ^{high} are chosen to be the $r/2$ and $(1 - r/2)$ th quantiles of the $\mathcal{N}(\mu_1^{(d)}, \sigma^2)$ distribution, respectively. This ensures control on response probability while also inducing large positive correlation between $R^{(a_1)}$ and $\left(Y_1^{(d)} - \mu_1^{(d)}\right)^2$.

Violation of working assumption A1(c) was induced by choosing

$\text{Cor}\left(Y_t^{(d)}, Y_2^{(d)} \mid R^{(a_1)} = 1\right) > \text{Cor}\left(Y_t^{(d)}, Y_2^{(d)} \mid R = 0\right)$ while keeping respective variances fixed. In our generative model, it was difficult to exert precise control over these quantities when response was related to prior outcomes; as such, these violations were induced under the R_{\perp} response model.

There exist natural constraints on how much larger than $\text{cov}\left(Y_t^{(d)}, Y_2^{(d)} \mid R = 0\right)$ the responders’ covariance can be while ensuring that (1) all conditional covariance matrices are positive definite and (2) $\text{cov}\left(Y_t^{(d)}, Y_2^{(d)} \mid R = 0\right) \geq 0$ for $t = 0, 1$. These constraints vary with ρ . We

choose $\text{Cor}(Y_1^{(d)}, Y_2^{(d)} | R^{(a_1)} = 1)$ such that $\text{cov}(Y_1^{(d)}, Y_2^{(d)} | R^{(a_1)} = 1)$ is the midpoint between the minimum covariance for which the assumption is violated and the maximum covariance allowed by the aforementioned constraints.

5.2 Simulation Results

Simulation results based on 3,000 simulated data sets are compiled in table 4. We find that sample size formula (10) performs as expected when all assumptions are satisfied. Empirical power is not significantly less than the target power of 0.8, per a one-sided binomial test with level 0.05. The sample size is, as expected, often conservative, particularly when within-person correlation is high.

There may be some concern that, for high within-person correlation, formula (10) is overly conservative; should this concern arise, we recommend use of the sharper formulae presented in the supplement. The difference between the sharper formulae and formula (10) is maximized when $\rho = (1 + \sqrt{5})/2 \approx 0.62$, so we expect to see the largest differences in power between formula (10) and the sharp formula when we set $\rho = 0.6$.

When all working assumptions are satisfied, we see that empirical power for R_+ and R_- scenarios are similar or slightly higher than under the R_{II} model. In general, there do not appear to be practical differences in empirical power between the response models.

As conjectured, violating working assumption A1(a) does not impact empirical power in design I (compare the results to column “ R_+ ”). For design II, empirical power is consistently less than the nominal value when working assumption A1(a) is violated. However, while the empirical power is often statistically significantly less than 0.8, for practical purposes the loss of power is relatively small. For design III, we notice small reductions in power relative to scenarios in which both working assumptions A1 and A2 are satisfied, though the conservative nature of formula (10) appears to protect against more severe loss of power. This suggests that our sample size formula is moderately robust to reasonable violations of A1(a).

For small ρ , we see no meaningful change in empirical power when violating working assumption A1(b). However, as ρ increases, this also leads to increased correlation between response and the other products of first-stage residuals, which increases the severity of the violation. For $\rho = 0.6$, we see noticeable, but not extreme, departures from nominal power. When $\rho = 0.8$, our generative model was not able to violate working assumption A1(b) without also violating working assumption A1(a); as such, we omit those results.

Interestingly, as can be seen in the supplement, defining *non*-response as in equation (11) (i.e., replacing $R^{(a_1)}$ with $1 - R^{(a_1)}$) leads to higher-than-nominal power. When there exists *negative* correlation between response and products of squared first-stage residuals, the form of σ_c^2 derived in appendix B is more conservative, leading to increased power.

Simulation results show that our sample size formula is quite robust to violations of working assumption A1(c) for low-to-moderate within-person correlations; at high correlations, the

empirical power is statistically significantly less than 0.8. However, as with working assumption A1(a), the practical reduction in power is relatively small.

The final column of table 4 suggests that formula (10) is highly sensitive to violations of working assumption A2 in regards to the true correlation structure. In particular, when the true correlation structure is not exchangeable with correlation ρ and is instead AR(1) with correlation ρ , empirical power is substantially lower than the target of 0.8, particularly as ρ increases. This is unsurprising: as our assumed exchangeable ρ increases, the difference between the assumed and actual correlation between the end-of-study measurement and earlier measurements increases, leading to more severe loss of power.

Note that when within-person correlation is high, sample size becomes rather small. Since the method presented here is based on asymptotic normality, we caution the reader that small sample sizes (e.g., $n < 100$) provided by formula (10) may be quite sensitive to violation of the working assumptions.

In figure 2, we examine the effect on empirical power of misspecifying the within-person correlation. Analytically, we see from formula (10) that if the assumed ρ is smaller than the true within-person correlation, the sample size will be conservative. On the other hand, when the assumed ρ in formula (10) is larger than the true correlation, the sample size will be anti-conservative. Figure 2 shows plots of empirical power against the difference between the assumed within-person correlation ρ_{guess} and the true ρ . For small ρ_{guess} , formula (10) appears to be quite robust to misspecification of ρ ; however, as ρ_{guess} increases, the formula becomes highly sensitive to such a violation of working assumption A2. This is supported analytically, since formula (10) is a function of ρ_{guess}^2 .

6 Discussion

We have derived sample size formulae for SMART designs in which the primary aim is a comparison of two embedded DTRs that begin with different first-stage treatments on a continuous, longitudinal outcome. We derived the formulae for three common SMART designs.

The sample size formula is the product of three components: (1) the formula for the minimum sample size for the comparison of two means in a standard two-arm trial (see, e.g., Friedman et al.,⁴¹ page 147), (2) a deflation factor of $1 - \rho^2$ that accounts for the use of a longitudinal outcome, and (3) a SMART-specific “design effect”, an inflation factor that accounts for the SMART design.

The SMART design effect can be interpreted as the cost of conducting the SMART relative to conducting a standard two-arm randomized trial of the two DTRs which comprise the primary aim. The benefit of conducting a SMART (relative to the standard two-arm randomized trial) is the ability to answer additional, secondary questions that are useful for constructing effective DTRs. For example, such questions may focus on one or more of the other pairwise comparisons between DTRs, on whether the first- and second-stage treatments work synergistically to impact outcomes (e.g., a test of the null that $\gamma_6 = 0$ in

model (1)), or may focus on hypothesis-generating analyses that seek to estimate more deeply-tailored DTRs.^{42,43,44}

The formulae are expected to be easy-to-use for both applied statistical workers and clinicians. Indeed, inputs α , β , and n are as in the sample size formula for a standard z -test. Furthermore, estimates of ρ , r_{a1} , and σ are often readily available from the literature or can be estimated using data from prior studies (e.g., prior randomized trials, or external pilot studies).

We make a number of recommendations concerning the use of the formulae; in particular, how best to use the formulae conservatively in the absence of certainty concerning prior estimates of ρ , r_{a1} , and/or the structure of the variance of the repeated measures outcome.

First, in designs II and III, if there is uncertainty concerning the response rate (e.g., response rate estimates are based on data from smaller prior studies), one approach is to err conservatively by assuming a smaller-than-estimated response rate. In both designs, the most conservative approach is to assume a response rate of zero.

Second, as in standard randomized trials in which the primary aim is a pre-post comparison, the required sample size decreases as the hypothesized within-person correlation increases.⁴⁵ Therefore, if the hypothesized ρ is larger than the true ρ , the computed sample size will be anti-conservative, resulting in an under-powered study. Indeed, we see this in the results of the simulation experiment (see figure 2). Here, again, one approach is to err conservatively towards smaller values of ρ .

Finally, working assumption A2 (concerning the variance of the repeated measures outcome) may be seen as overly restrictive in the imposition of an exchangeable correlation structure. For example, studies with a continuous repeated measures outcome may observe an autoregressive correlation structure. However, the exchangeable working assumption can be employed conservatively in the following way: if the hypothesized structure is not exchangeable, one approach is to set ρ in formula (10) to the smallest plausible value (e.g., the within-person correlation between the baseline and end-of-study measurements for an autoregressive structure). Because this approach utilizes a lower bound on the value of the true within-person correlations, it is expected to yield a larger than needed (more conservative) sample size. Similarly, if the true within-person correlation is expected to differ by DTR, one approach is to employ the smallest plausible ρ . As with the third recommendation, these recommendations are not unique to SMARTs; indeed, these strategies may also be used to size standard two-arm randomized trials with repeated measures outcome.

In the case where $\text{Var}(Y_t^{(d)})$ varies with time and/or DTR, we conjecture that power will suffer if a pooled estimate of σ^2 is used when the variance decreases with time. To see this, consider that the standardized effect size δ defined in equation (9) has as a denominator the pooled standard deviation of $Y_2^{(d)}$ across the groups under comparison. Should the estimate of pooled standard deviation be larger than the true value, the variance of $c^T \hat{\theta}$ will increase; since the estimate will be less efficient than hypothesized, power will be lower than

expected. Conversely, we also conjecture that when $\text{Var}(Y_t^{(d)})$ increases with t , the sample size will be conservative using similar reasoning.

The main contribution of this manuscript is the development of sample size formulae for SMARTs in which the primary aim is an end-of-study comparison of two embedded DTRs which recommend different first-stage treatments (so-called “separate-path” DTRs).⁴⁶ It is possible, though, that some trialists may have interest in sizing a SMART for an end-of-study comparison of “shared-path” DTRs; that is, two DTRs which recommend the same first-stage treatment. We believe that, for the comparison of shared-path DTRs, investigators are better set to use a standard sample size calculation to compare the second-stage treatments which differ between the DTRs, then upweighting the result by the proportion of participants expected to be in these groups.

There are a number of interesting ways to build on this manuscript in future methodological work. First, some scientists may be interested in a primary aim comparison that involves other features of the marginal mean trajectory, such as the area under the curve (AUC). Future work could develop formulae for these other primary aim comparisons. An important challenge here is in whether and how to define the standardized effect size δ . Second, an interesting extension of this work is to better understand the cost-benefit trade-off between adding additional sample size versus adding additional measurement occasions to the SMART design. The formulae presented here employ the rather simplistic working assumption that there are $T = 3$ measurement occasions (at baseline, the end of the first stage, and the end of the second stage). Based on limited simulation experiments, sample sizes based on our formulae are expected to perform conservatively when $T > 3$. Future work could develop rules of thumb for how best to allocate additional sample size versus additional measurement occasions given budget constraints (e.g., a fixed total study cost and fixed costs for an additional participant and additional measurement occasion). Third, as the field moves toward simulation-based approaches for sample size calculation, there is a clear need for the development of software that would allow applied statistical workers and clinicians to make fewer (or more flexible) assumptions concerning many of the features of the SMART, or to be more flexible with respect to the design of the SMART. An important challenge here is to make the software general enough to be used across a number of different types of SMART designs (e.g., three stages of randomization), yet not so flexible that it is difficult to use. The benefits of this is the ability to examine the power for various different scientific questions given a single data generative model and for many other types of SMARTs.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgements

This research was supported in part through computational resources and services provided by Advanced Research Computing at the University of Michigan, Ann Arbor. We would also like to thank the two anonymous reviewers for their helpful and insightful comments which led to a much-improved manuscript.

Funding

The authors disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: This work was supported by the Eunice Kennedy Shriver National Institute of Child Health and Human Development [grant number R01HD073975]; the National Institute of Biomedical Imaging and Bioengineering [grant number U54EB020404]; the National Institute of Mental Health [grant number R03MH097954]; the National Institute on Alcohol Abuse and Alcoholism [grant numbers P01AA016821, RC1AA019092]; and the National Institute on Drug Abuse [grant numbers R01DA039901, P50DA039838]. The content of this paper is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

A: Identifiability Assumptions

We make the following assumptions in order to show that equation (4) has mean zero.

I1 Positivity.

The probabilities $P(A_1 = 1)$ and $P(A_2 = 1 \mid A_1, R)$ are non-zero.

I2 Consistency with potential outcomes.⁴⁷

A participant's observed responder status is "consistent" with the participant's corresponding potential responder status under the assigned first-stage treatment; i.e.,

$$R_i = \mathbb{1}\{A_{1,i} = 1\}R_i^{(1)} + \mathbb{1}\{A_{1,i} = -1\}R_i^{(-1)}.$$

Furthermore, a participant's observed repeated measures outcomes are consistent with the participant's corresponding potential repeated measures outcomes under the assigned treatment sequence; see table A1.

Table A1:

Design-specific consistency assumptions. $d \in \mathcal{D}$ indexes embedded DTRs (a_1, a_{2R}, a_{2NR}) .

Design	$Y_{2,i}$
I	$\sum_{d \in \mathcal{D}} \frac{1}{2} \mathbb{1}\{A_{1,i} = a_1\} (R_i \mathbb{1}\{A_{2,i} = a_{2R}\} + (1 - R_i) \mathbb{1}\{A_{2,i} = a_{2NR}\}) Y_{2,i}^{(d)}$
II	$\sum_{d \in \mathcal{D}} \mathbb{1}\{A_{1,i} = a_1\} \left(\frac{1}{2} R_i + (1 - R_i) \mathbb{1}\{A_{2,i} = a_2\} \right) Y_{2,i}^{(d)}$
III	$\sum_{d \in \mathcal{D}} \mathbb{1}\{A_{1,i} = a_1\} \left(\mathbb{1}\{a_1 = -1\} + \mathbb{1}\{a_1 = 1\} \left(\frac{1}{2} R_i + (1 - R_i) \mathbb{1}\{A_{2,i} = a_2\} \right) \right) Y_{2,i}^{(d)}$

The factor of 1/2 for responders in designs II and III accounts for the fact that these participants are consistent with two DTRs. For example in design II, if $R_j = 1$ for some j , $Y_j^{(a_1, 0, 1)} = Y_j^{(a_1, 0, -1)} = Y_j^{(a_1, 0, \cdot)}$.

I3 Sequential randomization.

At each stage in the SMART, observed treatments A_1 and A_2 are assigned independently of future potential outcomes, given the participant's history up to that point. That is,

$$\begin{cases} \{Y_i^{(d)}, R_i^{(a1)}\} \perp\!\!\!\perp A_{1,i} \\ \{Y_i^{(d)}\} \perp\!\!\!\perp A_{2,i} \mid A_{1,i}, R_i \end{cases}$$

Identifiability assumptions I1 and I3 are satisfied by design in a SMART (see, e.g., Lavori and Dawson⁴⁸); identifiability assumption I2 connects the potential outcomes and observed data, and is typically accepted in the analysis of randomized trials.

B: Derivation of Sample Size Formulae

We derive the sample size formulae for comparing two DTRs which recommend different first-stage treatments that are embedded in a SMART in which a continuous longitudinal outcome is collected throughout the study. These formulae are based on the regression analyses described in section 3 and a Wald test.

We consider a SMART in which the outcome is collected at three timepoints: at baseline ($t_1 = 0$), immediately before assessing response/non-response ($t_2 = 1$), and at the end of the study ($t_3 = 2$). We ignore the presence of baseline covariates \mathbf{X} and assume $\boldsymbol{\mu}^{(d)}(\boldsymbol{\theta})$ is piecewise-linear in $\boldsymbol{\theta}$ (see, for example, model (1)). Recall that $\boldsymbol{\theta} \in \mathbb{R}^p$.

Recall from section 4 that we wish to test the null hypothesis $H_0: \mathbf{c}^\top \boldsymbol{\theta} = 0$, where $\mathbf{c} \in \mathbb{R}^p$ is a contrast vector specifying a linear combination of $\boldsymbol{\theta}$. In particular, we are interested in contrasts \mathbf{c} which yield an end-of-study comparison between two embedded DTRs which recommend different first-stage treatments. For example, in design II, the end-of-study comparison of DTRs (1, 0, 1) vs. (-1, 0, -1) is given by $\mathbf{c} = (0, 0, 2, 0, 2, 2, 0)^\top$. Since, here, \mathbf{c} is a vector, this yields a 1-degree of freedom Wald test for which we can use a Z statistic:

$$Z = \frac{\sqrt{nc}^\top \hat{\boldsymbol{\theta}}}{\sigma_c},$$

where $\sigma_c = \sqrt{\mathbf{c}^\top \mathbf{B}^{-1} \mathbf{M} \mathbf{B}^{-1} \mathbf{c}}$. Under H_0 , by asymptotic normality of $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})$, the test statistic follows an asymptotic standard normal distribution. Suppose we wish to test H_0 against the alternative hypothesis $\mathbf{c}^\top \boldsymbol{\theta} = \Delta$. The minimum-required sample size is

$$n \geq (z_{1-\alpha/2} + z_{1-\beta})^2 \frac{\sigma_c^2}{\Delta^2}, \quad (\text{B1})$$

where z_p is the p th quantile of the standard normal distribution. Formula (B1) is a standard result in the clinical trials literature;^{41,49} however, because of the dependence on σ_c , the formula is not useful as written. The goal of the remainder of this appendix is to derive a closed-form upper bound on σ_c so as to obtain a sample size formula in terms of marginal quantities which can be more easily elicited from clinicians, or estimated from the literature. In particular, we want this upper bound to be a multiple of σ^2 , the assumed common marginal variance across time and DTR, so that the final formula will involve Cohen's effect size $\delta = \Delta / \sigma$.

Recall the definitions of $\mathbf{B} \in \mathbb{R}^{p \times p}$ and $\mathbf{M} \in \mathbb{R}^{p \times p}$ in equations (5) and (6), respectively. These quantities depend on $\mathbf{D}^{(d)} \in \mathbb{R}^{T \times p}$, the jacobian of $\boldsymbol{\mu}^{(d)}(\boldsymbol{\theta})$, and $\mathbf{V}^{(d)}(\boldsymbol{\tau}) \in \mathbb{R}^{T \times T}$, the working covariance matrix for $\mathbf{Y}^{(d)}$. By assumed linearity of $\boldsymbol{\mu}^{(d)}(\boldsymbol{\theta})$, $\mathbf{D}^{(d)}$ is a fixed, constant matrix for all d . Furthermore, we assume that the working covariance matrix $\mathbf{V}^{(d)}(\boldsymbol{\tau})$ is correctly specified and satisfies working assumption A2 so that $\mathbf{V}^{(d)}(\boldsymbol{\tau}) = \boldsymbol{\Sigma}$ for all $d \in \mathcal{D}$. Note that $\boldsymbol{\Sigma}$ is non-random.

The estimand in equation (8) is a function of potential outcomes; as written in equations (5) and (6), \mathbf{B} and \mathbf{M} are functions of observed data. We begin by expressing \mathbf{B} in terms of potential outcomes. Under the positivity, consistency, and sequential ignorability conditions (identifiability assumptions I1 to I3), we can apply lemma 4.1 of Murphy et al.¹⁰ so that

$$\begin{aligned} \mathbf{B} &= \sum_{d \in \mathcal{D}} \mathbb{E}_{A_1, R, A_2} \left[\mathbf{W}^{(d)}(A_1, R, A_2) \mathbf{D}^{(d)} (\mathbf{V}^{(d)}(\boldsymbol{\tau}))^{-1} (\mathbf{D}^{(d)})^\top \right] \\ &= \sum_{d \in \mathcal{D}} \mathbf{D}^{(d)} \boldsymbol{\Sigma}^{-1} (\mathbf{D}^{(d)})^\top. \end{aligned} \quad (\text{B2})$$

We now turn our attention to \mathbf{M} . Expanding the outer product inside the expectation, we have

$$\begin{aligned} \mathbf{M} &= \mathbb{E}_{A_1, R, A_2, \mathbf{Y}} \left[\left(\sum_{d \in \mathcal{D}} \mathbf{W}^{(d)}(A_1, R, A_2) \mathbf{D}^{(d)} (\mathbf{V}^{(d)}(\boldsymbol{\tau}))^{-1} (\mathbf{Y} - \boldsymbol{\mu}^{(d)}(\boldsymbol{\theta})) \right)^{\otimes 2} \right] \\ &= \sum_{d \in \mathcal{D}} \mathbb{E}_{A_1, R, A_2, \mathbf{Y}} \left[\left(\mathbf{W}^{(d)}(A_1, R, A_2) \right)^2 (\mathbf{D}^{(d)} \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \boldsymbol{\mu}^{(d)}(\boldsymbol{\theta})))^{\otimes 2} \right] \\ &\quad + \sum_{d \neq d'} \mathbb{E}_{A_1, R, A_2, \mathbf{Y}} \left[\mathbf{W}^{(d)}(A_1, R, A_2) \mathbf{W}^{(d')} (A_1, R, A_2) \mathbf{D}^{(d)} \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \boldsymbol{\mu}^{(d)}(\boldsymbol{\theta})) \right. \\ &\quad \left. (\mathbf{Y} - \boldsymbol{\mu}^{(d')}(\boldsymbol{\theta}))^\top \boldsymbol{\Sigma}^{-1} (\mathbf{D}^{(d')})^\top \right]. \end{aligned} \quad (\text{B3})$$

Consider a single summand of the first term in equation (B3). We can write this as

$$\mathbf{D}^{(d)} \boldsymbol{\Sigma}^{-1} \mathbb{E}_{A_1, R, A_2, \mathbf{Y}} \left[\mathbf{W}^{(d)}(A_1, R, A_2)^2 (\mathbf{Y} - \boldsymbol{\mu}^{(d)}(\boldsymbol{\theta}))^{\otimes 2} \right] \boldsymbol{\Sigma}^{-1} (\mathbf{D}^{(d)})^\top, \quad (\text{B4})$$

where, as before, $\mathbf{Z}^{\otimes 2} = \mathbf{Z}\mathbf{Z}^\top$. The inner expectation is a $T \times T$ matrix, the (i, j) th element of which is

$$\mathbb{E}_{A_1, R, A_2, \mathbf{Y}} \left[\mathbf{W}^{(d)}(A_1, R, A_2)^2 (Y_{ti} - \mu_{ti}^{(d)}(\boldsymbol{\theta})) (Y_{tj} - \mu_{tj}^{(d)}(\boldsymbol{\theta})) \right]. \quad (\text{B5})$$

Notice that the work above is design-independent: \mathbf{B} and \mathbf{M} have the same form as equations (B2) and (B3), respectively, for all designs. Below, we proceed only for design II, but derivations for designs I and III are analogous, substituting appropriate definitions of $\mathbf{W}^{(d)}(A_1, R, A_2)$. Recall that, for design II, when all randomization probabilities are 0.5, $\mathbf{W}^{(d)}(A_1, R, A_2) = 2 \mathbb{1} \{A_1 = a_1^{(d)}\} (R + 2(1 - R) \mathbb{1} \{A_2 = a_2^{(d)}\})$.

It can be shown that we can achieve an upper bound on σ_c which involves σ and ρ by imposing working assumption A1 on the inner expectation in equation (B4) such that all diagonal terms are at least $2 \cdot DE \cdot \sigma^2$ and all off-diagonal terms are at most $2 \cdot DE \cdot \rho\sigma^2$, where DE is the design effect as in table 3.

Consider, for example, $t = 1$. By repeated use of iterated expectation and application of identifiability assumptions I2 and I3, equation (B5) becomes

$$\begin{aligned}
 & E_{Y_{t_0}, A_1, Y_{t_1}, R, A_2, Y_{t_2}} \left[W^{(d)}(A_1, R, A_2)^2 (Y_{t_1} - \mu_{t_1}^{(d)}(\theta))^2 \right] \\
 &= E_{Y_{t_0}, A_1, Y_{t_1}, R, A_2} \left[4 \mathbb{1} \{A_1 = a_1^{(d)}\} (R + 4(1 - R) \mathbb{1} \{A_2 = a_2^{(d)}\}) (Y_{t_1} - \mu_{t_1}^{(d)}(\theta))^2 \right] \\
 &= E_{Y_{t_0}^{(d)}, A_1, Y_{t_1}, R^{(a_1)}, A_2^{(d)}} \left[4 \mathbb{1} \{A_1 = a_1^{(d)}\} (R^{(a_1)} + 4(1 - R^{(a_1)}) \mathbb{1} \{A_2 = a_2^{(d)}\}) \right. \\
 &\quad \left. (Y_{t_1}^{(d)} - \mu_{t_1}^{(d)}(\theta))^2 \right] \tag{B6} \\
 &= E_{S_2(\bar{A}_1)} \left[4 \mathbb{1} \{A_1 = a_1^{(d)}\} (R^{(a_1)} + 4(1 - R^{(a_1)}) E_{A_2 | S_2(\bar{A}_1)} \left[\mathbb{1} \{A_2 = a_2^{(d)}\} \right] \right. \\
 &\quad \left. (Y_{t_1}^{(d)} - \mu_{t_1}^{(d)}(\theta))^2 \right] \\
 &= E_{Y_{t_0}^{(d)}, A_1, Y_{t_1}^{(d)}, R^{(a_1)}} \left[4 \mathbb{1} \{A_1 = a_1^{(d)}\} (2 - R^{(a_1)}) (Y_{t_1}^{(d)} - \mu_{t_1}^{(d)}(\theta))^2 \right].
 \end{aligned}$$

$$= 4 E_{Y_1^{(d)}} \left[(Y_1 - \mu_1^{(d)})^2 \right] - 2 E_{Y_1^{(d)}, R^{(a_1)}} \left[(Y_1 - \mu_1^{(d)})^2 R^{(a_1)} \right] \tag{B7}$$

$$= 4\sigma^2 - 2\text{Cov} \left((Y_1 - \mu_1^{(d)})^2, R^{(a_1)} \right) - 2E \left[R^{(a_1)} \right] E \left[(Y_1 - \mu_1^{(d)})^2 \right] \tag{B8}$$

$$= 2(2 - r_{a_1})\sigma^2. \tag{B9}$$

Equation (B7) follows from equation (B6) by identifiability assumption I3 and smoothing over $Y_{t_0}^{(d)}$, equation (B8) arises from the definition of covariance, and equation (B9) is a consequence of working assumption A1(b).

Similar derivations can be performed for the remaining combinations (i, j) . Under working assumptions A1 and A2, equation (B5) is exactly equal to $2(2 - r_{a_1})\Sigma_{i, j}$ for $i, j \in \{1, 2\}$. For the last diagonal element ($i = j = 3$), equation (B5) is at least $2(2 - r_{a_1})\Sigma_{3, 3}$; the remaining off-diagonal quantities are bounded above by $2(2 - r_{a_1})\Sigma_{i, j}$. This allows us to bound $c^T \mathbf{B}^{-1} \mathbf{M} \mathbf{B}^{-1} c$ above by

$$\begin{aligned}
\mathbf{c}^\top \mathbf{B}^{-1} \mathbf{M} \mathbf{B}^{-1} \mathbf{c} &\leq 2 \cdot \frac{1}{2} ((2 - r_1) + (2 - r_{-1})) \mathbf{c}^\top \mathbf{B}^{-1} \left(\sum_{d \in \mathcal{D}} \mathbf{D}^{(d)} \boldsymbol{\Sigma}^{-1} \right)^{\otimes 2} \mathbf{B}^{-1} \mathbf{c} \\
&= \frac{4\sigma^2(1 - \rho) \left(\rho^2 + 4\rho - \frac{1}{2}(r_1 + r_{-1})(2\rho + 1) + 2 \right)}{1 + \rho}.
\end{aligned} \tag{B10}$$

Plugging equation (B10) into formula (B1) leads to the aforementioned “sharp” sample size formula for design II. Some algebra shows that

$$\sigma_c^2 \leq 4\sigma^2 \cdot (1 - \rho^2) \cdot \frac{1}{2} ((2 - r_1) + (2 - r_{-1})), \tag{B11}$$

which allows for an easy-to-understand sample size formula. Plugging this result into formula (B1), we arrive at formula (10).

C: Details Concerning the Data-Generative Process for Simulations

To construct table 4, we employ two data-generative models. Here, we describe the first, which we believe to be more realistic and which is used to simulate under all scenarios in table 4 except for those in which working assumption A1(c) is violated. A description of the second model, used to violate working assumption A1(c), is available in the supplement.

In general, generating realistic longitudinal data from a SMART is difficult when precise control must be exerted over the marginal covariance structure of the outcomes. As such, the generative model described here is rather complex. We attempt to distill the details in this appendix and provide further details about response status and variance generation in the supplement.

For each scenario described in table 4, we compute the sample size for the trial using formula (10) and the appropriate design effect from table 3. We then, for each “participant” i , generate potential outcomes under each embedded DTR as follows:

$$\begin{aligned}
Y_{0,i}^{(d)} &= \gamma_0 + \epsilon_{0,i} \\
Y_{1,i}^{(d)} \Big| Y_{0,i}^{(d)} &= \gamma_0 + \rho Y_{0,i}^{(d)} + \gamma_1 + \gamma_2 a_1 + \epsilon_{1,i}^{(a_1)} \\
R_i^{(a_1)} \Big| Y_{0,i}^{(d)}, Y_{1,i}^{(d)} &= g_{a_1} \left(Y_{1,i}^{(a_1)} \right) \\
Y_{2,i}^{(d)} \Big| Y_{0,i}^{(d)}, Y_{1,i}^{(d)}, R_i^{(a_1)} &= (1 - c_0(\rho) - c_1(\rho)) \gamma_0 + c_0(\rho) Y_{0,i}^{(d)} + c_1(\rho) Y_{1,i}^{(d)} \\
&\quad + (1 - c_1(\rho)) (\gamma_1 + \gamma_2 a_1^{(d)}) + \gamma_3 + \gamma_4 a_1^{(d)} + \xi^{(d)} \left(R_i^{(a_1)} \right) \\
&\quad + \left(R_i^{(a_1)} - r_{a_1,i} \right) (\lambda_1 + \lambda_2 a_1^{(d)}) + \epsilon_{2,i}^{(d)} \left(R_i^{(a_1)} \right),
\end{aligned} \tag{C1}$$

where $\epsilon_0 \sim \mathcal{N}(0, \sigma^2)$, $\epsilon_1^{(d)} \sim \mathcal{N}(0, (1 - \rho^2)\sigma^2)$, and $\epsilon_2^{(d)} \left(R_i^{(a_1)} \right) \sim \mathcal{N}\left(0, v^{(d)} \left(R_i^{(a_1)} \right)\right)$. $\sigma^2 = \text{Var}(Y_t^{(d)})$ is the assumed constant marginal variance of the outcome as in working assumption A2. Here, $a_j^{(d)}$ is the j th-stage treatment recommended by DTR d .

Table C1:

Design-specific conditional mean components $\xi^{(d)}(R^{(a_1)})$ for model (C1). $r_{a_1, i}$ is the probability of response to first-stage treatment for participant i . Dashes indicate that the corresponding parameter is not used in the model for that design; variances in the lower part of the table span across multiple rows when (non-)responders to the corresponding first-stage treatment are not re-randomized.

Design	$\xi^{(d)}(R^{(a_1)})$
I	$\frac{R^{(a_1)}}{r_{a_1, i}}(\gamma_5 + \gamma_7 a_1^{(d)})a_{2R}^{(d)} + \frac{1 - R^{(a_1)}}{1 - r_{a_1, i}}(\gamma_6 + \gamma_8 a_1^{(d)})a_{2NR}^{(d)}$
II	$\frac{1 - R^{(a_1)}}{1 - r_{a_1, i}}(\gamma_5 + \gamma_6 a_1^{(d)})a_{2NR}^{(d)}$
III	$\frac{1 - R^{(a_1)}}{1 - r_{a_1, i}} \mathbb{1} \{ a_1^{(d)} = 1 \} \cdot \gamma_5 a_{2NR}^{(d)}$

The error terms ϵ_t add the “additional” variance to $Y_j^{(d)}$ which is needed to achieve marginal variance σ^2 . Outcomes $Y_{t,i}^{(d)}$ are generated as functions of an individual’s outcomes at previous timepoints, which induces variance in, say, $Y_2^{(d)}$ when marginalizing over $Y_1^{(d)}$ and $Y_0^{(d)}$. Hence, $v^{(d)}(R^{(a_1)})$ is the additional variance added to the response-conditional end-of-study outcome beyond that which is induced by defining $Y_{2,i}^{(d)} | R_i^{(a_1)}$ as a function of $Y_{0,i}^{(d)}$ and $Y_{1,i}^{(d)}$. This is required to ensure that the marginal variance of $Y_2^{(d)}$ is σ^2 . All errors ϵ are generated independently of one another.

The parameters γ_j are interpreted exactly as in, say, model (1) (see section 3.3 and equation (2)), and index the generative marginal structural mean model. The functions $c_t(\rho)$ control within-person correlation between $Y_{2,i}^{(d)}$ and previously-observed outcomes $Y_{t,i}^{(d)}$, $t < t^*$. For an exchangeable correlation structure, we use $c_0(\rho) = c_1(\rho) = \rho/(1 + \rho)$; for AR(1), $c_0(\rho) = 0$ and $c_1(\rho) = \rho$.

Note that the second-stage outcome $Y_{2,i}^{(d)}$ is generated conditionally on response status, since participants in a SMART can only be a responder or a non-responder to first-stage treatment. Since second-stage treatments in a SMART are often restricted based on response, these treatment effects can typically only be estimated using either responders or non-responders. $\xi^{(d)}(R^{(a_1)})$ is a design-specific function of response which involves marginal parameters γ_j for second-stage treatment effects and their interactions with first-stage treatment effects. Our choices of $\xi^{(d)}(R^{(a_1)})$ are given in table C1. For example, in design II, only non-responders are re-randomized, and so the effect of a_{2NR} should only be simulated among non-responders, and upweighted appropriately to reflect this (see table C1).

The final component of the generative mean structure for $Y_{2,i}^{(d)} | R_i^{(a1)}$ is

$(R_i^{(a1)} - r_{a1,i})(\lambda_1 + \lambda_2 a_1^{(d)})$, where $\tau_{a1,i}$ is the probability of response to first-stage treatment for participant i . The parameters λ_1 and λ_2 control how responders and non-responders differ from the marginal mean at time 2, and cancels to zero when averaged over response status.

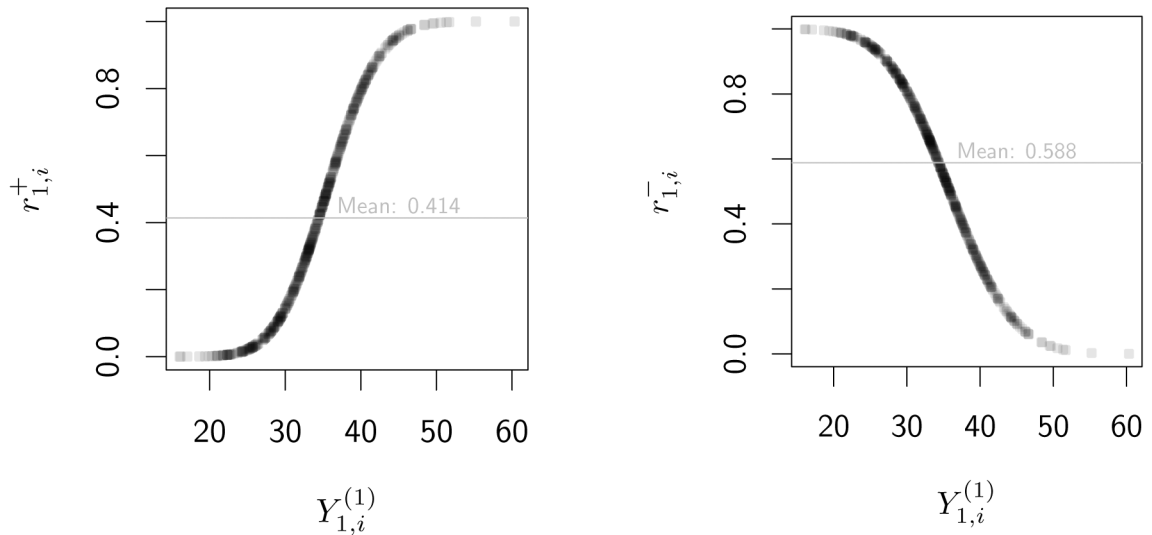
In order to design a realistic generative model, we define response as a function of $Y_{1,i}^{(a1)}$. In table 4, we consider three possible response models: “ R_{II} ”, in which response status is independent of $Y_{1,i}^{(a1)}$, “ R_+ ”, in which $P(R_i^{(a1)} = 1)$ increases with $Y_{1,i}^{(a1)}$, and “ R_- ”, in which $P(R_i^{(a1)} = 1)$ is decreasing in $Y_{1,i}^{(a1)}$. In figure C1, we plot probabilities of response versus $Y_{1,i}^{(d)}$ under both R_+ and R_- response models. We consider both R_+ and R_- to ensure that power is not affected by the choice of coding for Y (e.g., “higher-is-better” vs. “lower-is-better”). More details are provided in the supplement.

Table C2:

Example choices of parameters in data generative model to achieve $\delta = 0.3$ and $\rho = 0.3$ when $r_1 = r_{-1} = 0.4$ under R_+ and working assumptions A1 and A2 are satisfied for each of designs I to III.

Parameter	Design I	Design II	Design III
γ_0	35.0	33.5	35.0
γ_1	-4.0	-0.8	-0.5
γ_2	2.7	0.9	1.0
γ_3	-1.6	-0.8	0.2
γ_4	-1.5	0.4	-0.2
γ_5	0.4	-0.4	0.8
γ_6	-0.4	0.1	-
γ_7	0.4	-	-
γ_8	0.4	-	-
λ_1	0.3	0.1	0.8
λ_2	0.4	-0.5	0.0
σ^2	64	36	64
$v^{(1,1,a2NR)}(1)$	57.62		
$v^{(1,-1,a2NR)}(1)$	57.62	21.12	56.56
$v^{(1,a2R,1)}(0)$	53.07	18.48	55.27
$v^{(1,a2R,-1)}(0)$	53.07	21.64	43.37
$v^{(-1,1,a2NR)}(1)$	57.89		
$v^{(-1,-1,a2NR)}(1)$	57.89	19.39	53.68
$v^{(-1,a2R,1)}(0)$	49.80	14.20	53.97

Parameter	Design I	Design II	Design III
$v^{(-1, a_{2R}, -1)}_{(0)}$	55.78	20.80	

**Figure C1:**

Individual probabilities of response versus potential outcomes at the end of the first stage of a SMART under response models R_+ and R_- . Based on a simulated SMART with $n = 500$ individuals; $a_1 = 1$ was chosen arbitrarily and without loss of generality. Empirical average response rate is plotted as a horizontal line in gray. Darker regions of the curve contain more observations.

(I) Response probabilities in the R_+ model versus (II) Response probabilities in the R_- model versus $Y_{1,i}^{(1)}$ with average response rate $r_1 = 0.4$. Higher values of $Y_{1,i}^{(1)}$ are associated with higher response probabilities of $Y_{1,i}^{(1)}$ are associated with lower response probabilities.

The above models are used to generate potential outcomes $Y_i^{(d)}$ and potential response status $R_i^{(a_1)}$ for each “participant” under each DTR $d = (a_1, a_{2R}, a_{2NR})$. This is done to ensure the generative model satisfies identifiability assumption I2; identifiability assumptions I1 and I3 are satisfied by design. From these potential data, we “observe” data as follows:

1. Choose $Y_{0,i} = Y_{0,i}^{(d)}$. Note that since no treatment has been assigned at time 0, there is only one possible value of Y_0 , so the potential and observed outcomes coincide.
2. Generate $A_{1,i} \sim \text{Bernoulli}(0.5)$.
3. Choose $Y_{1,i} = Y_{1,i}^{(A_{1,i}, \cdot, \cdot)}$, the potential outcome at time 1 corresponding to the DTR(s) which recommend first-stage treatment $A_{1,i}$.
4. Choose $R_i = R_i^{(A_{1,i})}$, the potential response status of the participant under first-stage treatment $A_{1,i}$.

5. Generate $A_{2,i} | R_i \sim \text{Bernoulli}(p_i)$, where p_i is either 0.5 or 0 depending on the value of R_i and the SMART design under consideration. (For example, in design II, $p_i = 0$ for all i such that $R_i = 1$ since responders are not re-randomized.)
6. Choose $Y_{2,i} | R^{(A_1,i)} = Y_{2,i}^{(A_1,i, R^{(A_1,i)} A_{2,i})} (1 - R^{(A_1,i)}) A_{2,i}$, the potential outcome at time 2 had the participant had response status R_i and been treated according to a DTR which recommends $A_{1,i}$ in the first stage and $A_{2,i}$ in the second.

Table C3:

Target and estimated marginal variance matrices $V^{(d)}(\boldsymbol{\tau})$ from the data generative model described in appendix C. The “unstructured estimate” is produced by estimating the variance at each timepoint and for each DTR, and correlation for each DTR using the unstructured estimate in table 2, then averaging over DTRs. The “exchangeable estimate” is computed by assuming variance is constant over time and DTR, and using the exchangeable estimate of ρ from table 2, averaged over DTRs. The exchangeable estimate is used in simulations assuming working assumption A2 is satisfied.

Design	Target Structure	Unstructured Estimate	Exchangeable Estimate
I	$\begin{pmatrix} 64 & 19.2 & 19.2 \\ 19.2 & 64 & 19.2 \\ 19.2 & 19.2 & 64 \end{pmatrix}$	$\begin{pmatrix} 63.9 & 19.3 & 19.1 \\ 19.3 & 63.8 & 18.7 \\ 19.1 & 18.7 & 62.5 \end{pmatrix}$	$\begin{pmatrix} 63.4 & 18.9 & 18.9 \\ 18.9 & 63.4 & 18.9 \\ 18.9 & 18.9 & 63.4 \end{pmatrix}$
II	$\begin{pmatrix} 36 & 10.8 & 10.8 \\ 10.8 & 36 & 10.8 \\ 10.8 & 10.8 & 36 \end{pmatrix}$	$\begin{pmatrix} 35.9 & 10.9 & 11.0 \\ 10.9 & 35.9 & 11.1 \\ 11.0 & 11.1 & 35.8 \end{pmatrix}$	$\begin{pmatrix} 35.9 & 10.9 & 10.9 \\ 10.9 & 35.9 & 10.9 \\ 10.9 & 10.9 & 35.9 \end{pmatrix}$
III	$\begin{pmatrix} 64 & 19.2 & 19.2 \\ 19.2 & 64 & 19.2 \\ 19.2 & 19.2 & 64 \end{pmatrix}$	$\begin{pmatrix} 63.9 & 19.4 & 19.9 \\ 19.4 & 63.7 & 21.3 \\ 19.9 & 21.3 & 63.6 \end{pmatrix}$	$\begin{pmatrix} 63.8 & 20.0 & 20.0 \\ 20.0 & 63.8 & 20.0 \\ 20.0 & 20.0 & 63.8 \end{pmatrix}$

In table C2 we provide values of the parameters chosen for simulations to achieve a standardized effect size $\delta = 0.3$ for the end-of-study comparison of interest. The values of $v^{(d)}(R^{(a_1)})$ are specific to the scenario in which all working assumptions are satisfied, $\rho = 0.3$, average response probabilities are $r_1 = r_{-1} = 0.4$, and when response is computed under the R_+ model. In table C3, we show that the target marginal variance structures are achieved using this data generative model.

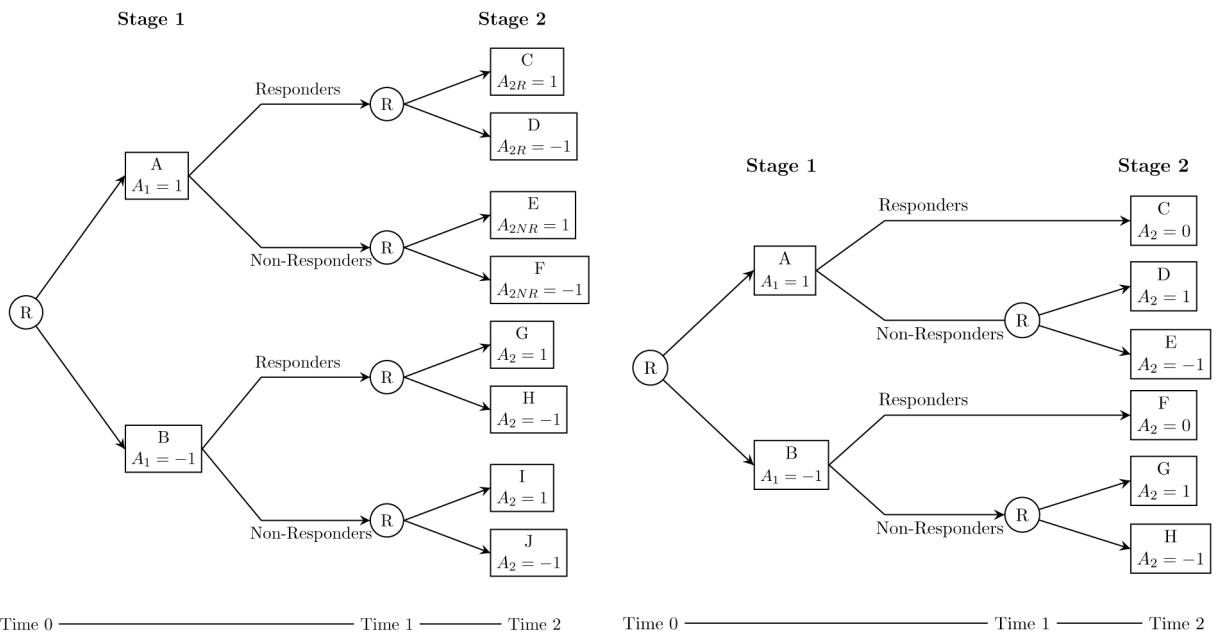
References

- [1]. Kosorok Michael R and Moodie Erica E. M., editors. Adaptive Treatment Strategies in Practice. Society for Industrial and Applied Mathematics, Philadelphia, 2016.
- [2]. Chakraborty Bibhas and Moodie Erica E. M.. Statistical methods for dynamic treatment regimes. Springer, New York, 2013 ISBN 1461474280. doi: 10.1007/978-1-4614-7428-9.
- [3]. McKay James R, Drapkin Michelle L, Van Horn Deborah H A, Lynch Kevin G, Oslin David W, DePhilippis Dominick, Ivey Megan, and Cacciola John S. Effect of patient choice in an adaptive sequential randomization trial of treatment for alcohol and cocaine dependence. J Consult Clin Psychol, 83(6):1021–32, 2015 ISSN 1939–2117. doi: 10.1037/a0039534. URL <http://www.ncbi.nlm.nih.gov/pubmed/26214544>. [PubMed: 26214544]

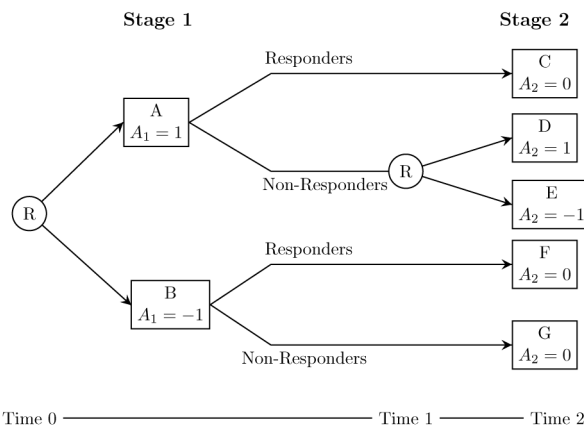
- [4]. Wallace Michael P. and Moodie Erica E. M.. Personalizing medicine: A review of adaptive treatment strategies. *Pharmacoepidem Dr S*, 23(6):580–585, 2014 ISSN 10991557. doi: 10.1002/pds.3606.
- [5]. Ogbagaber Semhar B., Karp Jordan, and Wahed Abdus S.. Design of sequentially randomized trials for testing adaptive treatment strategies. *Stat Med*, 35(6):840–858, 2016 ISSN 10970258. doi: 10.1002/sim.6747. [PubMed: 26412033]
- [6]. Almirall Daniel, Nahum-Shani Inbal, Sherwood Nancy E., and Murphy Susan A.. Introduction to SMART designs for the development of adaptive interventions: with application to weight loss research. *Transl Behav Med*, 4(3):260–274, 2014 ISSN 16139860. doi: 10.1007/s13142-014-0265-0. [PubMed: 25264466]
- [7]. Inbal Nahum-Shani Min Qian, Almirall Daniel, Pelham William E, Gnagy Beth, Fabiano Gregory A, Waxmonsky James G, Yu Jihnee, and Murphy Susan A. Experimental design and primary data analysis methods for comparing adaptive interventions. *Psychol Methods*, 17(4): 457–477, 2012 ISSN 1939–1463. doi: 10.1037/a0029372. [PubMed: 23025433]
- [8]. Lavori Philip W and Dawson Ree. A design for testing clinical strategies: biased individually tailored within-subject randomization. *J R Stat Soc A Stat*, 163(1):29–38, 2000.
- [9]. Lavori Philip W and Dawson Ree. Dynamic treatment regimes: practical design considerations. *Clin Trials*, 1(1):9–20, 2 2004 ISSN 1740–7745 (Print). doi: 10.1191/1740774504cn002oa. [PubMed: 16281458]
- [10]. Murphy Susan A.. An experimental design for the development of adaptive treatment strategies. *Stat Med*, 24(10):1455–1481, 2005 ISSN 02776715. doi: 10.1002/sim.2022. URL 10.1002/sim.2022. [PubMed: 15586395]
- [11]. Auyeung SF, Long Q, Royster EB, Murthy S, McNutt MD, Lawson D, Miller A, Manatunga A, and Musselman DL. Sequential multiple-assignment randomized trial design of neurobehavioral treatment for patients with metastatic malignant melanoma undergoing high-dose interferon-alpha therapy. *Clin Trials*, 6(5):480–490, 10 2009 ISSN 1740–7753; 1740–7745. doi: 10.1177/1740774509344633[doi]. [PubMed: 19786415]
- [12]. Kidwell Kelley M. SMART designs in cancer research: Past, present, and future. *Clin Trials*, 11(4):445–456, 2014. doi: 10.1177/1740774514525691. URL <http://ctj.sagepub.com>. [PubMed: 24733671]
- [13]. Thall Peter F. SMART design, conduct, and analysis in oncology In Kosorok Michael R. and Moodie Erica E. M., editors, *Adapt. Treat. Strateg. Pract*, chapter 4, pages 41–54. Society for Industrial and Applied Mathematics, Philadelphia, 1 edition, 2016.
- [14]. Diegidio Paul, Hermiz Steven, Hibbard Jonathan, Kosorok Michael, and Hultman Charles Scott. Hypertrophic Burn Scar Research: From Quantitative Assessment to Designing Clinical Sequential Multiple Assignment Randomized Trials. *Clin Plast Surg*, 2017 ISSN 00941298. doi: 10.1016/j.cps.2017.05.024. URL 10.1016/j.cps.2017.05.024.
- [15]. Hibbard Jonathan C, Friedstat Jonathan S, Thomas Sonia M, Edkins Renee E, Hultman C Scott, and Kosorok Michael R. LIBERTI: A SMART study in plastic surgery. *Clin Trials*, page 174077451876243, 2018 ISSN 1740–7745. doi: 10.1177/1740774518762435. URL <http://journals.sagepub.com/doi/10.1177/1740774518762435>.
- [16]. Murphy Susan A, Lynch Kevin G, Oslin David, McKay James R, and TenHave Tom. Developing adaptive treatment strategies in substance abuse research. *Drug Alcohol Depen*, 88: S24–S30, 2007.
- [17]. Kasari Connie, Kaiser Ann, Goods Kelly, Nietfeld Jennifer, Mathy Pamela, Landa Rebecca, Murphy Susan, and Almirall Daniel. Communication interventions for minimally verbal children with autism: A sequential multiple assignment randomized trial. *J Am Acad Child Psy*, 53(6):635–646, 2014 ISSN 15275418. doi: 10.1016/j.jaac.2014.01.019. URL 10.1016/j.jaac.2014.01.019.
- [18]. Li Zhiguo and Murphy Susan A.. Sample size formulae for two-stage randomized trials with survival outcomes. *Biometrika*, 98(3):503–518, 2011 ISSN 00063444. doi: 10.1093/biomet/asr019. [PubMed: 22363091]
- [19]. Kidwell Kelley M., Seewald Nicholas J., Tran Qui, Kasari Connie, and Almirall Daniel. Design and Analysis Considerations for Comparing Dynamic Treatment Regimens with Binary

- Outcomes from Sequential Multiple Assignment Randomized Trials. *J Appl Stat*, 2017 ISSN 0266–4763. doi: 10.1080/02664763.2017.1386773. URL 10.1080/02664763.2017.1386773.
- [20]. Necamp Timothy, Kilbourne Amy, and Almirall Daniel. Comparing cluster-level dynamic treatment regimens using sequential, multiple assignment, randomized trials: Regression estimation and sample size considerations. *Stat Methods Med Res*, pages 1–88, 2017 ISSN 0962–2802. doi: 10.1177/0962280217708654. URL <http://journals.sagepub.com/doi/pdf/10.1177/0962280217708654>.
- [21]. Lu Xi, Nahum-Shani Inbal, Kasari Connie, Lynch Kevin G., Oslin David W., Pelham William E., Fabiano Gregory, and Almirall Daniel. Comparing dynamic treatment regimes using repeated-measures outcomes: Modeling considerations in SMART studies. *Stat Med*, 35(10): 1595–1615, 2016 ISSN 10970258. doi: 10.1002/sim.6819. URL 10.1002/sim.6819. [PubMed: 26638988]
- [22]. Li Zhiguo. Comparison of adaptive treatment strategies based on longitudinal outcomes in sequential multiple assignment randomized trials. *Stat Med*, 36(3):403–415, 2017 ISSN 10970258. doi: 10.1002/sim.7136. URL <http://doi.wiley.com/10.1002/sim.7136>. [PubMed: 27646957]
- [23]. Oslin David W.. Managing Alcoholism in People Who Do Not Respond to Naltrexone (EXTEND), 2005 URL <https://clinicaltrials.gov/ct2/show/NCT00115037>.
- [24]. Fitzsimons Heather, Tuten Michelle, O’Grady Kevin, Chisolm Margaret S, and Jones Hendree E. A smart design: Response to reinforcement-based treatment intensity among pregnant, drug-dependent women. *Drug Alcohol Depen*, 156:e69, 11 2015 ISSN 0376–8716. doi: 10.1016/j.drugalcdep.2015.07.1106. URL 10.1016/j.drugalcdep.2015.07.1106.
- [25]. Fu Steven S, Rothman Alexander J, Vock David M, Lindgren Bruce, Almirall Daniel, Begnaud Abbie, Melzer Anne, Schertz Kelsey, Glaeser Susan, Hammett Patrick, and Joseph Anne M. Program for lung cancer screening and tobacco cessation: Study protocol of a sequential, multiple assignment, randomized trial. *Contemp Clin Trials*, 60(July):86–95, 2017 ISSN 15592030. doi: 10.1016/j.cct.2017.07.002. URL 10.1016/j.cct.2017.07.002. [PubMed: 28687349]
- [26]. Eckshtain Dikla. Using SMART Experimental Design to Personalize Treatment for Child Depression, 2013 URL <https://clinicaltrials.gov/ct2/show/NCT01880814>.
- [27]. Pelham William E., Fabiano Gregory A., Waxmonsky James G., Greiner Andrew R., Gnagy Elizabeth M., Pelham William E., Coxe Stefany, Verley Jessica, Bhatia Ira, Hart Katie, Karch Kathryn, Konijnendijk Evelien, Tresco Katy, Nahum-Shani Inbal, and Murphy Susan A.. Treatment Sequencing for Childhood ADHD: A Multiple-Randomization Study of Adaptive Medication and Behavioral Interventions. *J Clin Child Adolesc*, 45(4):396–415, 2016 ISSN 15374416. doi: 10.1080/15374416.2015.1105138. URL
- [28]. Budney Alan J.. Behavioral Treatment of Adolescent Substance Use (SMART), 2014 URL <https://clinicaltrials.gov/ct2/show/NCT02063984>.
- [29]. Almirall Daniel, DiStefano Charlotte, Chang Ya-Chih, Shire Stephanie, Kaiser Ann, Lu Xi, Nahum-Shani Inbal, Landa Rebecca, Mathy Pamela, and Kasari Connie. Longitudinal effects of adaptive interventions with a speech-generating device in minimally verbal children With ASD. *J Clin Child Adolesc*, 4416(August 2017):1–15, 2016 ISSN 1537–4416. doi: 10.1080/15374416.2016.1138407. URL <http://www.tandfonline.com/doi/full/10.1080/15374416.2016.1138407>{% }5Cn.
- [30]. Kilbourne AM, Abraham KM, Goodrich DE, Bowersox NW, Almirall D, Lai Z, and Nord KM. Cluster randomized adaptive implementation trial comparing a standard versus enhanced implementation intervention to improve uptake of an effective re-engagement program for patients with serious mental illness. *Implement Sci*, 8(1):136, 2013 ISSN 1748–5908. doi: 10.1186/1748-5908-8-136. URL <http://www.implementationscience.com/content/8/1/136>. [PubMed: 24252648]
- [31]. Lei Huitan, Nahum-shani Inbal, Lynch Kevin G, Oslin David, and Murphy Susan A.. A “SMART” design for building individualized treatment sequences. *Annu Rev Clin Psycho*, 8: 21–48, 2012 ISSN 1548–5951. doi: 10.1146/annurev-clinpsy-032511-143152. URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3887122>{% }7B{& }7Dtool=pmcentrez{% }7B{& }7Drendertype=abstract.

- [32]. Cole Stephen R. and Hernán Miguel A.. Constructing inverse probability weights for marginal structural models. *Am J Epidemiol*, 168(6):656–664, 2008 ISSN 00029262. doi: 10.1093/aje/kwn164. [PubMed: 18682488]
- [33]. van der Vaart AW. *Asymptotic Statistics Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, 1998. doi: 10.1017/CBO9780511802256.
- [34]. Liang Kung-Yee and Zeger Scott L.. Longitudinal data analysis using generalized linear models. *Biometrika*, 73(1):13–22, 4 1986 ISSN 0006–3444. doi: 10.1093/biomet/73.1.13. URL <http://biomet.oxfordjournals.org/cgi/content/long/73/1/13>.
- [35]. Ahn Chul, Heo Moonseong, and Zhang Song. *Sample Size Calculations for Clustered and Longitudinal Outcomes in Clinical Research Chapman & Hall/CRC Biostatistics Series*. CRC Press, Boca Raton, 2015 ISBN 978-1-4667-5626-3.
- [36]. Fitzmaurice Garrett M., Laird Nan M., and Ware James H.. *Applied Longitudinal Analysis*. Wiley, Hoboken, 2 edition, 2011.
- [37]. Bowden Charles L.. Sequential Multiple Assignment Treatment for Bipolar Disorder. <https://clinicaltrials.gov/ct2/show/NCT01588457>, 3 2017.
- [38]. Cohen J. *Statistical power analysis for the behavioral sciences*, volume 2nd. Routledge, New York, 1988 ISBN 0805802835. doi: 10.1234/12345678.
- [39]. Oetting AI, Levy JA, Weiss RD, and Murphy SA. Statistical Methodology for a SMART Design in the Development of Adaptive Treatment Strategies In Shrout P, Keyes K, and Ornstein K, editors, *Causality Psychopathol. Find. Determ. Disord. their Cures*, pages 179–205. American Psychiatric Publishing, Inc., Arlington, VA, 2011.
- [40]. Kahan Brennan C., Jairath Vipul, Doré Caroline J., and Morris Tim P.. The risks and rewards of covariate adjustment in randomized trials: An assessment of 12 outcomes from 8 studies. *Trials*, 15(139):1–7, 2014 ISSN 17456215. doi: 10.1186/1745-6215-15-139. [PubMed: 24382030]
- [41]. Friedman Lawrence M., Furberg Curt D., and Demets David L.. *Fundamentals of clinical trials*. Springer, New York, 4 edition, 2010 ISBN 9781441915856. doi: 10.1007/978-1-4419-1586-3.
- [42]. Watkins Christopher J. C. H.. Learning from Delayed Rewards. PhD thesis, King’s College, 1989 URL http://www.cs.rhul.ac.uk/~}chrisw/new/{_}thesis.pdf.
- [43]. Nahum-Shani Inbal, Qian Min, Almirall Daniel, Pelham William E., Gnagy Beth, Fabiano Gregory a., Waxmonsky James G., Yu Jihnhee, and Murphy Susan a.. Q-Learning: A Data Analysis Method for Constructing Adaptive Interventions. *Psychol Methods*, 17(4):478–494, 2012 ISSN 1082-989X. doi: 10.1037/a0029373. [PubMed: 23025434]
- [44]. Zhang Yichi, Laber Eric B., Tsiatis Anastasios, and Davidian Marie. Using decision lists to construct interpretable and parsimonious treatment regimes. *Biometrics*, 71(4):895–904, 2015 ISSN 15410420. doi: 10.1111/biom.12354. [PubMed: 26193819]
- [45]. Zhang Song, Cao Jing, and Ahn Chul. A GEE approach to determine sample size for preand post-intervention experiments with dropout. *Comput Stat Data An.*, 69:114–121, 2014 ISSN 01679473. doi: 10.1016/j.csda.2013.07.037. URL 10.1016/j.csda.2013.07.037.
- [46]. Kidwell Kelley M. and Wahed Abdus S.. Weighted log-rank statistic to compare shared-path adaptive treatment strategies. *Biostatistics*, 14(2):299–312, 4 2013 ISSN 1465-4644. doi: 10/gfppsw. [PubMed: 23178734]
- [47]. Robins James M.. Causal Inference from Complex Longitudinal Data In Berkane M, editor, *Latent Var. Model. Appl. to Causality*, volume 120 of *Lecture Notes in Statistics*, pages 69–117. Springer, New York, 1997 ISBN 978-1-4612-1842-5. doi: 10.1007/978-1-4612-1842-54. URL http://link.springer.com/10.1007/978-1-4612-1842-5{_}4.
- [48]. Lavori PW and Dawson R. Introduction to dynamic treatment strategies and sequential multiple assignment randomization. *Clin Trials*, 11(4):393–399, 2014 ISSN 1740–7745. doi: 10.1177/1740774514527651. URL <http://ctj.sagepub.com/cgi/doi/10.1177/1740774514527651>. [PubMed: 24784487]
- [49]. Lachin John M.. Introduction to sample size determination and power analysis for clinical trials. *Control Clin Trials*, 2(2):93–113, 1981 ISSN 01972456. doi: 10.1016/0197-2456(81)90001-5. [PubMed: 7273794]



(I) All participants are re-randomized, regardless of (II) The second randomization is restricted to only response status. non-responders.



(III) The second randomization is restricted to only non-responders to treatment A.

Figure 1: Three commonly-used two-stage SMART designs. Each design varies in choice of which subsets of participants are re-randomized. Circled R indicates randomization, capital letters indicate (potentially non-unique) treatments, and a_{\cdot} provides a coding system used to index embedded DTRs.

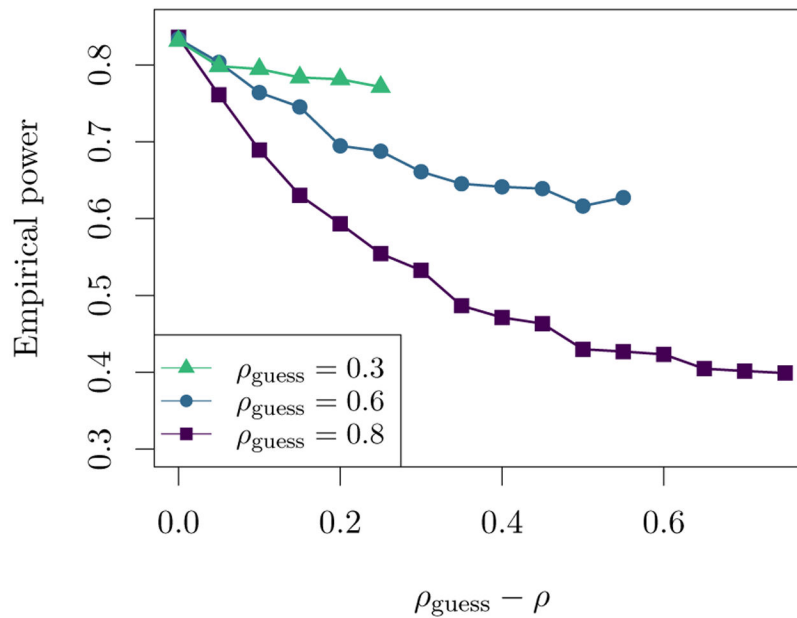


Figure 2: Empirical power versus the difference between the true within-person correlation ρ and hypothesized correlation ρ_{guess} used to compute sample size. Results are shown for design II with a hypothesized response rate of 0.4, and sample size was chosen to detect standardized effect size $\delta = 0.3$ for the comparison of DTRs (1, 0, 1) and (-1, 0, -1). Each point is based on 3000 simulations with target power 0.8 and significance level 0.05. Results are similar for designs I and III and different values of δ and r .

Table 1:Design-specific indicators for consistency with a given DTR $d \in \mathcal{D}$.

Design	$\mathbf{I}^{(d)}(A_{1,i}, R_i, A_{2,i})$
I	$\mathbb{1}\{A_{1,i} = a_1\}(\mathbb{1}\{A_{2,i} = a_{2R}\}R_i + \mathbb{1}\{A_{2,i} = a_{2NR}\}(1 - R_i))$
II	$\mathbb{1}\{A_{1,i} = a_1\}(R_i + \mathbb{1}\{A_{2,i} = a_{2NR}\}(1 - R_i))$
III	$\mathbb{1}\{A_{1,i} = a_1\}(\mathbb{1}\{a_1 = -1\} + \mathbb{1}\{a_1 = 1\})(R_i + \mathbb{1}\{A_{2,i} = a_{2NR}\}(1 - R_i))$

Table 2:

Correlation estimators for selected working correlation structures, assuming constant within-person variance over time. $d \in \mathcal{D}$ is an embedded DTR, $W_i^{(d)}$ is shorthand for $W^{(d)}(A_{1,i}, R_i, A_{2,i})$ and $\hat{e}_{i,t}^{(d)}(\hat{\theta})$ is the estimated residual $Y_{i,t} - \mu_i^{(d)}(\mathbf{X}_i; \hat{\theta})$.

Cor. structure	$\text{Cor}(Y_{t_j}^{(d)}, Y_{t_k}^{(d)})$	Estimator
AR(1)	$\begin{cases} 1 & t_j = t_k \\ (\rho^{(d)})^{ j-k } & t_j \neq t_k \end{cases}$	$\hat{\rho}^{(d)} = \frac{\sum_{i=1}^n W_i^{(d)} \sum_{m=1}^{T-1} \hat{e}_{i,t_m}^{(d)}(\hat{\theta}) \hat{e}_{i,t_{m+1}}^{(d)}(\hat{\theta})}{(\hat{\sigma}^{(d)})^2 \cdot n \cdot (T-1)}$
Exchangeable	$\begin{cases} 1 & t_j = t_k \\ \rho^{(d)} & t_j \neq t_k \end{cases}$	$\hat{\rho}^{(d)} = \frac{\sum_{i=1}^n W_i^{(d)} \sum_{l < m} \hat{e}_{i,t_l}^{(d)}(\hat{\theta}) \hat{e}_{i,t_m}^{(d)}(\hat{\theta})}{(\hat{\sigma}^{(d)})^2 \cdot n \cdot T(T-1)/2}$
Unstructured	$\begin{cases} 1 & t_j = t_k \\ \rho_{t_j, t_k}^{(d)} & t_j \neq t_k \end{cases}$	$\hat{\rho}_{t_j, t_k}^{(d)} = \frac{\sum_{i=1}^n W_i^{(d)} \hat{e}_{i,t_j}^{(d)}(\hat{\theta}) \hat{e}_{i,t_k}^{(d)}(\hat{\theta})}{(\hat{\sigma}^{(d)})^2 \cdot n}$

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 3:

Design effects for sample size formula (10). $r_{a_1} = P(R^{(a_1)} = 1)$ is the response rate to first-stage treatment a_1 .

Design	Design effect	Conservative design effect
I	2	2
II	$\frac{1}{2}(2 - r_1) + \frac{1}{2}(2 - r_{-1})$	2
III	$\frac{1}{2}(3 - r_1)$	$\frac{3}{2}$

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 4:

Sample sizes and empirical power results for an end-of-study comparison of the DTR recommending only treatments indexed by 1 and that which recommends only treatments indicated by -1 . δ is the true standardized effect size as defined in equation (9), r is the common probability of response to first-stage treatment, and ρ is the true exchangeable within-person correlation. n is computed using formula (10) with $\alpha=0.05$ and $\beta=0.2$. R_{\perp} refers to a generative model in which response status is independent of all prior outcomes; R_+ and R_- refer to generative models in which response is positively or negatively correlated with Y_1 , respectively. All violation scenarios assume the R_+ generative model, except working assumption A1(c). Results are the proportion of 3000 Monte Carlo simulations in which we reject $H_0 : c^T \theta = 0$ at the 5% level.

Design	δ	r	ρ	n	Empirical power						
					A1 and A2 satisfied			Violation of A1			Violation of A2
					R_{\perp}	R_+	R_-	A1(a)	A1(b)	A1(c)	True AR(1)
I	0.3	0.4	0.0	698	0.798	0.807	0.803	0.798	0.796	‡	‡
			0.3	635	0.819	0.817	0.800	0.820	0.804	0.815	0.780*
			0.6	447	0.815	0.862	0.773*	0.865	0.817	0.827	0.728*
			0.8	252	0.835	0.925	0.733*	‡	‡	0.840	0.721*
		0.6	0.0	698	0.796	0.799	0.806	0.800	0.791	‡	‡
		0.3	635	0.808	0.813	0.792	0.824	0.805	0.807	0.775*	
		0.6	447	0.833	0.856	0.798	0.859	0.831	0.838	0.727*	
		0.8	252	0.827	0.901	0.758*	‡	‡	0.835	‡	
	0.5	0.4	0.0	252	0.799	0.801	0.798	0.798	0.801	‡	‡
			0.3	229	0.813	0.815	0.797	0.814	0.811	0.814	0.771*
			0.6	161	0.824	0.872	0.789	0.868	0.833	0.843	0.742*
			0.8	91	0.843	0.931	0.734**§	0.926	‡	0.839§	0.725*
		0.6	0.0	252	0.796	0.797	0.810	0.792	0.802	‡	‡
		0.3	229	0.817	0.815	0.808	0.811	0.823	0.823	0.771*	
		0.6	161	0.838	0.859	0.790	0.861	0.832	0.837	0.749*	
		0.8	91	0.835§	0.896	0.765**§	0.896	‡	0.859	‡	
II	0.3	0.4	0.0	559	0.801	0.801	0.808	0.778*	0.803	‡	‡
			0.3	508	0.804	0.813	0.831	0.800	0.797	0.798	0.795
			0.6	358	0.817	0.819	0.834	0.807	0.759*	0.788	0.811
			0.8	201	0.836	0.814	0.836	0.809	‡	0.792	0.806
		0.6	0.0	489	0.804	0.796	0.793	0.736*	0.810	‡	‡
		0.3	445	0.797	0.804	0.818	0.758*	0.795	0.780*	0.804	
		0.6	313	0.824	0.831	0.844	0.793	0.752*	0.770*	0.824	
		0.8	176	0.845	‡	‡	0.754*	‡	0.776*	0.842	
	0.5	0.4	0.0	201	0.801	0.800	0.802	0.768*	0.794	‡	‡
			0.3	183	0.813	0.800	0.819	0.790	0.813	0.796	0.803

Design	δ	r	ρ	n	Empirical power							
					A1 and A2 satisfied			Violation of A1			Violation of A2	
					R_{II}	R_{+}	R_{-}	A1(a)	A1(b)	A1(c)	True AR(1)	
			0.6	129	0.814	0.828	0.833	0.810	0.763 [*]	0.799	0.815	
			0.8	73	0.839	0.841	0.852	0.829	†	0.795	0.804	
		0.6	0.0	176	0.807	0.799	0.796	0.733 [*]	0.808	‡	‡	
			0.3	160	0.816	0.815	0.821	0.767 [*]	0.808	0.802	0.812	
			0.6	113	0.829	0.830	0.837	0.792	0.765 [*]	0.770 [*]	0.817	
			0.8	64	0.845 [§]	†	†	0.783 ^{*§}	†	0.789 [§]	†	
III	0.3	0.4	0.0	454	0.806	0.813	0.806	0.782 [*]	0.794	‡	‡	
			0.3	413	0.815	0.809	0.814	0.789	0.800	0.800	0.775 [*]	
				0.6	291	0.821	0.811	0.818	0.794	0.783 [*]	0.787 [*]	0.687 [*]
				0.8	164	0.824	0.812	0.839	0.812	†	0.802	0.637 [*]
		0.6	0.0	419	0.813	0.814	0.817	0.781 [*]	0.769 [*]	‡	‡	
				0.3	381	0.823	0.812	0.808	0.776 [*]	0.791	0.795	0.771 [*]
				0.6	268	0.823	0.817	0.844	0.807	0.750 [*]	0.754 [*]	0.709 [*]
				0.8	151	0.820	†	†	0.803	†	0.784 [*]	†
		0.5	0.4	0.0	164	0.808	0.804	0.795	0.776 [*]	0.802	‡	‡
				0.3	149	0.822	0.815	0.827	0.811	0.791	0.805	0.789
				0.6	105	0.811	0.810	0.812	0.810	0.798	0.785 [*]	0.698 [*]
				0.8	59	0.838	†	0.823	0.845	†	0.817 [§]	0.684 [*]
		0.6	0.0	151	0.798	0.809	0.803	0.778 [*]	0.772 [*]	‡	‡	
				0.3	138	0.812	0.809	0.814	0.800	0.782 [*]	0.799	0.778 [*]
				0.6	97	0.803 [§]	0.812	0.826 [§]	0.826 [§]	0.762 [*]	0.774 ^{*§}	0.705 ^{*§}
				0.8	55	0.826 [§]	†	†	0.837 [§]	†	0.797 [§]	†

* Statistically significantly less than 0.8 at the 5% level.

† Our data generative model could not accommodate this scenario (see appendix C).

‡ Violation of this working assumption is not applicable when $\rho = 0$.

§ Fewer than 3000 simulations generated data in which all treatment sequences were observed.