

## ABSTRACT

ROSE, ERIC JAMES. Methods for Dynamic Treatment Regimes. (Under the direction of Dr. Eric Laber and Dr. Marie Davidian.)

Precision medicine is a data-driven approach for leveraging patient heterogeneity to improve health care by tailoring treatments for patients based on individual patient characteristics. In practice, clinicians have to make several decisions during the course of treating a disease to adapt to the evolving health status of the patient. A dynamic treatment regime operationalizes clinical decision making as a sequence of decision rules that map up-to-date patient information to a recommended intervention.

Sequential Multiple Assignment Randomized Trials (SMARTs) are considered the gold standard for estimation and evaluation of treatment regimes. In Chapter 2, we derive sample size procedures for a SMART that ensure: (i) sufficient power for comparing the optimal treatment regime with standard of care; and (ii) the estimated optimal regime is within a given tolerance of the true optimal regime with high-probability. We establish asymptotic validity of the proposed procedures and demonstrate their finite sample performance in a series of simulation experiments.

In the context of time-to-event data where the number of decisions depends on each patient's health trajectory, naive application of existing methods for estimating an optimal treatment regime can lead to bias. In Chapter 3, we propose a variant of Q-learning, a regression-based approximate dynamic programming method, to estimate an optimal treatment regime with time-to-event data subject to censoring.

In Chapter 4, we discuss a SMART that was used to study parent messaging strategies for reducing chronic absenteeism in elementary school children. We will compare the adaptive regimes embedded in the trial in addition to estimating an optimal treatment regime in which we tailor interventions based on family and school characteristics.

© Copyright 2019 by Eric James Rose

All Rights Reserved

Methods for Dynamic Treatment Regimes

by  
Eric James Rose

A dissertation submitted to the Graduate Faculty of  
North Carolina State University  
in partial fulfillment of the  
requirements for the Degree of  
Doctor of Philosophy

Statistics

Raleigh, North Carolina

2019

APPROVED BY:

---

Dr. Dennis Boos

---

Dr. Shu Yang

---

Dr. Eric Laber  
Co-chair of Advisory Committee

---

Dr. Marie Davidian  
Co-chair of Advisory Committee

## **DEDICATION**

To my parents, Tom and Alice, and brother, Jason.

## **BIOGRAPHY**

Eric James Rose grew up in Ridgefield, Connecticut where he graduated from Ridgefield High School in 2010. He attended the University of Connecticut where he received a Bachelor of Arts in mathematics and statistics in 2014. He then continued his education at North Carolina State University where he earned a Master of Statistics in 2017.

## **ACKNOWLEDGEMENTS**

I would like to thank my advisors Dr. Eric Laber and Dr. Marie Davidian for all of the help they gave me throughout my graduate studies. Beginning with the courses I took of theirs and continuing with the guidance and encouragement they provided with my research, I am extremely grateful to have had the opportunity to work with them. The faculty and staff in the Department of Statistics at North Carolina State University do an incredible job at creating a great environment for their students. I would also like to thank James Gilman, Nick Kapur, and the rest of the members of Laber Labs who have provided not only helpful feedback on my work, but many friendships.

During my time in graduate school, I have been lucky to make many friendships inside and outside of the department that I credit much of my success to. In particular, the support Susheela and Nikhil Singh and Emily Gower have provided has meant so much to me over the past several years.

Finally, I would like to thank my parents and brother for the love and support they have provided throughout my life. I would have never been able to accomplish this without them.

# TABLE OF CONTENTS

<b>LIST OF TABLES</b> .....	<b>vi</b>
<b>LIST OF FIGURES</b> .....	<b>ix</b>
<b>Chapter 1 Introduction</b> .....	<b>1</b>
1.1 Q-learning .....	2
1.2 Sequential Multiple Assignment Randomized Trials .....	4
<b>Chapter 2 Sample Size Calculations for SMARTs</b> .....	<b>6</b>
2.1 Setup and notation .....	9
2.2 Sample size procedures .....	11
2.2.1 Normality-based sample size procedure .....	11
2.2.2 Projection-based sample size procedure .....	16
2.3 Simulation experiments .....	24
2.4 Discussion .....	30
<b>Chapter 3 Q-learning for Survival Analysis</b> .....	<b>33</b>
3.1 Setup and notation .....	35
3.2 Method .....	38
3.3 Simulation experiments .....	43
3.4 Analysis of STAR*D trial .....	47
3.5 Discussion .....	51
<b>Chapter 4 Parent Messaging Strategies for Student Absenteeism</b> .....	<b>53</b>
4.1 SMART design .....	54
4.2 Setup and notation .....	57
4.3 Analysis .....	58
4.3.1 Synthetic data .....	58
4.3.2 Missing data .....	60
4.3.3 Research Question 1 .....	62
4.3.4 Research Question 2 .....	64
4.3.5 Research Question 3 .....	68
4.3.6 Research Question 4 .....	70
<b>BIBLIOGRAPHY</b> .....	<b>75</b>
<b>APPENDICES</b> .....	<b>84</b>
Appendix A Sample Size Calculations for SMARTs .....	85
A.1 Proofs of technical results .....	85
A.2 Simulation results .....	91
Appendix B Q-learning for Survival Analysis .....	93
B.1 Proofs of technical results .....	93
B.2 Simulation study parameter values .....	95

## LIST OF TABLES

Table 2.1	Estimated power (POW) under a correctly specified generative model using the normality-based sample size procedure at a nominal level of 90. To form a baseline for comparison, $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care. . . . .	27
Table 2.2	Estimated concentration (OPT) under a correctly specified generative model using the normality-based sample size procedure at a nominal level of 90. To form a baseline for comparison, $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care. . . . .	28
Table 2.3	Estimated power (POW) under a model which violated the normality assumptions using the normality-based sample size procedure at a nominal level of 90. To form a baseline for comparison, $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care. . . . .	29
Table 2.4	Estimated concentration (OPT) under a model which violated the normality assumptions using the normality-based sample size procedure at a nominal level of 90. To form a baseline for comparison, $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care. . . . .	30
Table 2.5	Estimated power (POW) under a model for which the normality assumptions hold using the projection-based sample size procedure at a nominal level of 90. To form a baseline for comparison, $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care. . . . .	31
Table 2.6	Estimated concentration (OPT) under a model for which the normality assumptions hold using the projection-based sample size procedure at a nominal level of 90. To form a baseline for comparison, $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care. . . . .	31
Table 2.7	Estimated power (POW) under a model which violated the normality assumptions using the projection-based sample size procedure at a nominal level of 90. To form a baseline for comparison, $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care. . . . .	32
Table 2.8	Estimated concentration (OPT) under a model which violated the normality assumptions using the projection-based sample size procedure at a nominal level of 90. To form a baseline for comparison, $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care. . . . .	32



Table 3.1	Estimated value of the estimated optimal regime and the true value of the estimated optimal regime for the case where the time-to-event and censoring follow a Weibull distribution. . . . .	45
Table 3.2	Estimated value of the estimated optimal regime and the true value of the estimated optimal regime for the case where the time-to-event and censoring have a piecewise hazard model. . . . .	47
Table 3.3	Regression Coefficients for Third Stage Q-function for STAR*D trial	51
Table 3.4	Regression Coefficients for Second Stage Q-function for STAR*D trial	52
Table 3.5	Regression Coefficients for First Stage Q-function for STAR*D trial .	52
Table 4.1	Estimated effect of first-stage intervention on the number of days absent in the fall semester . . . . .	64
Table 4.2	Estimated effect of first-stage intervention on chronic absenteeism in the fall semester . . . . .	64
Table 4.3	Pairwise differences between each adaptive intervention and BAU on the number of days absent during both semesters . . . . .	67
Table 4.4	Pairwise differences between each of the adaptive interventions on the number of days absent during both semesters . . . . .	68
Table 4.5	Pairwise differences between each adaptive intervention and BAU on chronic absenteeism during both semesters . . . . .	68
Table 4.6	Pairwise differences between each of the adaptive interventions on chronic absenteeism during both semesters . . . . .	69
Table 4.7	Differences between second-stage interventions on the number of days absent during both semesters . . . . .	70
Table 4.8	Differences between second-stage interventions on chronic absenteeism during both semesters . . . . .	71
Table 4.9	Parameter estimates for the second stage Q-function . . . . .	74
Table 4.10	Parameter estimates for the first stage Q-function . . . . .	74
Table 4.11	IPW Estimates of expected outcomes of different regimes . . . . .	74
Table A.1	Estimated power (POW) and concentration (OPT) under a correctly specified generative model using the normality-based sample size procedure at a nominal level of 90. To form a baseline for comparison, $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care. . . . .	91
Table A.2	Estimated power (POW) and concentration (OPT) under a model which violated the normality assumptions using the normality-based sample size procedure at a nominal level of 90. To form a baseline for comparison, $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care. . . . .	92
Table A.3	Estimated power (POW) and concentration (OPT) under a model for which the normality assumptions hold using the projection-based sample size procedure at a nominal level of 90. To form a baseline for comparison, $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care. . . . .	92

Table A.4    Estimated power (POW) and concentration (OPT) under a model which violated the normality assumptions using the projection-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care. . . . . 92

## LIST OF FIGURES

Figure 1.1	Example two-stage SMART design . . . . .	5
Figure 4.1	SMART diagram for the parent messaging study . . . . .	56

# CHAPTER

## 1

# INTRODUCTION

Precision medicine focuses on how to tailor treatments for patients based on individual patient characteristics. There is a significant amount of heterogeneity in patients due to differences in attributes such as genetics, lifestyle, and medical history. These differences affect how a patient will respond to the treatments they receive and should therefore be considered in the decision making process for assigning treatments. In precision medicine, we seek to formalize this process with data-driven methods for informing how to best treat patients based on these characteristics.

Conventional methods for studying treatment effects involve conducting a randomized clinical trial in which a set of candidate treatments are randomly assigned to patients. These patients are then followed to determine the expected outcome for each individual treatment option if it was assigned to the entire population of interest. This does not match the clinical

decision making process in practice. Typically clinicians will make several decisions over the course of treating a disease for a single patient. This could be because there are multiple key decision points corresponding to the development of the disease or during regular follow-ups the clinician may determine the patient is not responding to the initial treatment or experiencing an adverse reaction and decide to adjust the treatment plan. This process of synthesizing patient information and assigning treatments sequentially over time is formalized in what is called a dynamic treatment regime.

A dynamic treatment regime operationalizes clinical decision making as a sequence of decision rules that map up-to-date patient information to a recommend intervention. The goal in precision medicine is to estimate an optimal treatment regime where optimal is defined as the regime that maximizes some mean outcome of interest if all patients in the population of interest were to receive treatments via these decision rules.

There have been many different methods proposed for estimating optimal dynamic treatment regimes which can generally be classified into two different categories, direct or indirect methods. Direct estimation methods work by specifying a class of regimes and then maximizing the estimated expected outcome over all regimes within that class. Examples of direct methods include marginal structural mean models (Robins et al., 2008) and outcome weighted learning (Zhao et al., 2012). Indirect methods function by positing models for the outcome of interest and the optimal regime is then determined by assigning the treatment that would maximize the predicted expected outcome. Examples of this include A-learning (Robins, 2004a; Murphy, 2003a) and Q-learning (Murphy, 2005b; Qian et al., 2013). Here we will focus on the use of Q-learning for estimating an optimal dynamic treatment regime.

## **1.1 Q-learning**

Q-learning is an approximate dynamic programming algorithm used for the estimation of an optimal dynamic treatment regime. It was developed in the context of reinforcement

learning for learning an optimal policy for a Markov decision process (Watkins, 1989) and has been used for sequential decision making problems in many different areas of application. It consists of performing a series of regressions each one with a subsequent maximization.

Suppose we are interested in estimating an optimal regime for a situation where we have  $K$  different key decision points. We will assume we observe data of the form  $\{(\mathbf{X}_{1,i}, A_{1,i}, \mathbf{X}_{2,i}, A_{2,i}, \dots, \mathbf{X}_{K,i}, A_{K,i}, Y_i)\}_{i=1}^n$  which is comprised of  $n$  *i.i.d* trajectories of the form  $(\mathbf{X}_1, A_1, \mathbf{X}_2, A_2, \dots, \mathbf{X}_K, A_K, Y)$  such that  $\mathbf{X}_1 \in \mathbb{R}^{p_1}$  is baseline covariate information,  $A_k \in \{-1, 1\}$  is the treatment assigned at stage  $k$ ,  $\mathbf{X}_k \in \mathbb{R}^{p_k}$  is additional information recorded during the course of treatment  $k - 1$ , and  $Y \in \mathbb{R}$  is the outcome of interest coded such that higher values are better. Define  $\mathbf{H}_k = (\mathbf{X}_1, A_1, \dots, \mathbf{X}_k)$  to be the history of a patient at decision point  $k$ . Q-learning works by defining the Q-functions

$$Q_K(\mathbf{h}_K, a_K) = \mathbb{E}(Y | \mathbf{H}_K = \mathbf{h}_K, A_K = a_K)$$

$$Q_k(\mathbf{h}_k, a_k) = \mathbb{E} \left\{ \max_{a_{k+1}} Q_{k+1}(\mathbf{H}_{k+1}, a_{k+1}) | \mathbf{H}_k = \mathbf{h}_k, A_k = a_k \right\} \text{ for } k = 1, \dots, K-1.$$

It can then be shown via dynamic programming that  $\pi_k^{\text{opt}}(\mathbf{h}_k) = \arg \max_{a_k} Q_k(\mathbf{h}_k, a_k)$  (Bellman, 1957). Models  $Q_k(\mathbf{h}_k, a_k; \beta_k)$  are then posited for each Q-function and fit in a backwards iterative manner. First estimate  $\hat{\beta}_K$  by regressing  $Y$  on  $\mathbf{H}_K$  and  $A_K$ . For  $k = K-1, \dots, 1$  the Q-function is estimated by regressing  $\max_{a_{k+1}} Q_{k+1}(\mathbf{H}_{k+1}, a_{k+1}; \hat{\beta}_{k+1})$  on  $\mathbf{H}_k$  and  $A_k$ . This then gives a plug-in estimator for the optimal regime as  $\hat{\pi}_k^{\text{opt}}(\mathbf{h}_k) = \arg \max_{a_k} Q_k(\mathbf{h}_k, a_k; \hat{\beta}_k)$ .

This can be implemented with a large variety of different regression methods for the Q-functions. The choice of regression model determines the form of the decision rule though so choosing simple models such as linear models will lead to simple, interpretable regimes. More flexible, nonparametric regression models could provide a better fit, but may lead to unintelligible regimes.

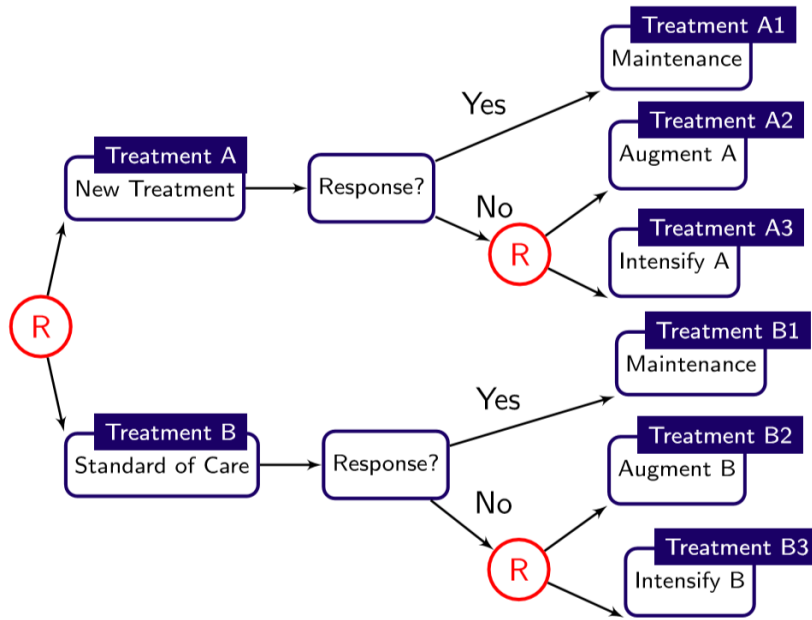
## 1.2 Sequential Multiple Assignment Randomized Trials

To be able to estimate an optimal dynamic treatment regime we need data of the correct form. We cannot take data from different clinical trials and observational studies that each focus on one individual decision point and piece them together. This is because there could be delayed effects of the early treatments that are not seen until later. Therefore we need data on the same set of patients through all the decision points.

Data from observational studies can lead to biased results because of confounding that occurs between the treatment assigned and the patient outcome. For example, sicker patients may be more likely to receive one of the treatment options and therefore worse outcomes for these patients are observed. In the case of multiple decision points we also have time-dependent confounding which further complicates inference. This has led to the preference of using experimental data for the use of studying dynamic treatment regimes.

The Sequential Multiple Assignment Randomized Trial (SMART) is the gold standard for gathering data for the study of sequential decision making for precision medicine. In a SMART, patients are randomized to the different feasible treatment options at each of the potential decision points. This provides high quality data for the estimation of an optimal regime while also mimicking the adaptive nature of assigning treatments that is observed in practice.

One example of a potential SMART with two decision points is shown in figure 1.1. In this example, in the first stage all the patients are randomized to receive either the new treatment which is denoted as treatment A or standard of care, denoted by treatment B. At the end of the first stage we check whether a patient is responding to the initial treatment. For example, a patient may be considered a responder if their symptom severity decreases to below some given threshold. All of the patients that are considered responders maintain their initial treatment during the second stage. If a patient were to receive treatment A during the first stage and is deemed a non-responder at the end of the first-stage, then



**Figure 1.1** Example two-stage SMART design

they are re-randomized to receive either augment A or intensify A during the second stage. Similarly non-responsive patients receiving treatment B are re-randomized in the second stage to augment B or intensify B. This is just one possible design for a SMART among many different potential options. How to determine the best design for a SMART is dictated by the scientific questions of interest that researchers plan on examining using the data collected.

Chapter 2 covers how to conduct sample size calculations for a two-stage SMART where we are trying to size the trial for the estimation of an optimal dynamic treatment regime. Chapter 3 develops an adaptation of Q-learning for when the outcome of interest is a time-to-event. Chapter 4 discusses a SMART that was used to test a parent messaging system to reduce absenteeism in elementary school kids.



## CHAPTER

# 2

# SAMPLE SIZE CALCULATIONS FOR SMARTS

A treatment regime is a sequence of functions, one per stage of clinical intervention, that map up-to-date patient information to a recommended treatment. An optimal treatment regime maximizes the mean of some cumulative clinical outcome when applied to select treatments for individuals in a population of interest (Murphy, 2003b; Robins, 2004b). Thus, an optimal treatment regime leads to better overall healthcare by adapting treatment to the evolving health status of each patient; consequently, optimal treatment regimes have become a primary means of operationalizing precision medicine. Optimal treatment regimes have been estimated across a wide range of application domains including breastfeeding (Moodie et al., 2012), bipolar disorder (Wu et al., 2015; Zhang et al., 2017), cancer (Thall

et al., 2000; Zhao et al., 2011; Wang et al., 2012; Zhang et al., 2012, 2015; Murray et al., 2017), cystic fibrosis (Zhou et al., 2017), diabetes (Ertefaie, 2014; Luckett et al., 2016), depression (Zhao et al., 2012), HIV (Moodie et al., 2007; van der Laan et al., 2005; Cain et al., 2010), smoking cessation (Chakraborty et al., 2009), substance abuse (Nahum-Shani et al., 2017) among others.

Sequential Randomized Multiple Assignment Randomized Trials (SMARTs Lavori and Dawson, 2000, 2004; Murphy, 2005a; Kidwell, 2014) are the gold standard for estimating and evaluating treatment regimes (Murphy et al., 2007; Lei et al., 2012; Chakraborty and Moodie, 2013; Kosorok and Moodie, 2015) and are increasingly common design in clinical and intervention science (PSU Methodology Center, 2017b,a). However, sample size calculations for SMARTs are typically based on power calculations for simple comparisons, e.g., comparison of the mean outcome across two pre-specified treatment sequences (Murphy, 2005a; Lei et al., 2012) and estimation of an optimal treatment regime from data collected in a SMART are almost always (we are not aware of an exception) conducted as part of exploratory, hypothesis-generating analyses. This approach is aligned with the *estimate-and-validate* paradigm wherein: (i) an optimal treatment regime is estimated using data collected in a SMART; and (ii) the performance of the estimated optimal regime is validated in a follow-up trial where the estimated regime is compared head-to-head with standard of care (Murphy, 2005a). This approach is appealing in that it avoids a number of nontrivial technical issues associated with estimating and evaluating a treatment regime using the same data (Robins, 2004b; Moodie et al., 2010; Chakraborty et al., 2010; Song et al., 2015a; Chakraborty et al., 2014; Laber et al., 2014); furthermore, sample size formulae for the comparison of fixed treatment sequences or other commonly used criteria to size SMARTs are straightforward in that they resemble those commonly used in non-sequential randomized trials. Another reason that sample size calculations for SMARTs are often based on simple comparisons is the seemingly widely held belief that sizing a trial to guarantee frequentist operating characteristics for an estimated optimal regime, e.g., providing a performance guarantee

for the estimated regime or powering a comparison of the performance of the optimal regime with standard of care, would either: (i) require a prohibitively large sample size; and/or (ii) rely on unrealistic assumptions about the underlying data-generating model. We provide evidence that for a for a large class of generative models neither of these beliefs appear to be well-founded.

We derive sample size procedures for a SMART that ensure sufficient power in comparing the mean outcome under the optimal regime with standard of care and that the estimated optimal regime will be within a given tolerance of the optimal policy with a given probability. The proposed sample sized procedures we develop here are at two possible extremes in terms of modeling assumptions. Our first procedure imposes significant parametric structure on the data-generating model and consequently we are able to derive sample size formulae that resemble the comparison of two means and require elicitation or estimation of a single scalar parameter. Our second procedure imposes structure only on moments and tail behavior of some components of the data generating model and then uses the bootstrap with oversampling to estimate a sufficient sample size; as this procedure imposes less structure the resulting estimator is more variable and consequently the estimated sample size tends to be larger. One reason for considering these two extremes is that they provide a basis for intermediate procedures that impose as much structure, as appropriate for a given application. We leave such intermediate approaches to future work.

We perceive the proposed work as making the following contributions: (i) it provides the first rigorous yet practical sample size procedures for estimating and evaluating optimal treatment regimes using SMARTs; (ii) it generates new knowledge about how much additional data would be needed to estimate a high-quality regime and consequently a provides a sense of how ‘underpowered’ existing SMARTs are for estimating optimal treatment regimes; and (iii) it provides theoretical guarantees for bootstrap oversampling for sample size calculations that are of independent interest. Furthermore, the proposed criteria used to derive our samples size procedures are closely related to those used in Laber et al. (2015)

to size a single stage two-arm trial to estimate an optimal regime, however, the current procedure applies to multistage trials and provides considerably stronger performance guarantees for the estimated optimal regime.

In Section 2.1, we provide the setup and notation. In Section 2.2, we derive our sample size procedures and state their theoretical properties. In Section 2.3, we evaluate the finite sample performance of the proposed sample size procedures in a series of simulation experiments. A discussion of the proposed methodology and open problems is provided in Section 2.4.

## 2.1 Setup and notation

We consider choosing the sample size,  $n$ , for a two-stage SMART that will produce data,  $\mathcal{D}_n = \{(\mathbf{X}_{1,i}, A_{1,i}, \mathbf{X}_{2,i}, A_{2,i}, Y_i)\}_{i=1}^n$ , which comprises *i.i.d.* trajectories of the form  $(\mathbf{X}_1, A_1, \mathbf{X}_2, A_2, Y)$  where:  $\mathbf{X}_1 \in \mathbb{R}^{p_1}$  denotes baseline subject information;  $A_1 \in \{-1, 1\}$  denotes the first assigned treatment;  $\mathbf{X}_2 \in \mathbb{R}^{p_2}$  denotes subject information collected during the course of the first treatment;  $A_2 \in \{-1, 1\}$  denotes the second assigned treatment; and  $Y \in \mathbb{R}$  denotes the outcome, coded so that higher is better. Define  $\mathbf{H}_1 = \mathbf{X}_1$  and  $\mathbf{H}_2 = (\mathbf{X}_1^\top, A_1, \mathbf{X}_2^\top)^\top$  so that  $\mathbf{H}_t$  denotes the history at time  $t = 1, 2$ . For simplicity, we assume that the trial will employ simple one-to-one randomization so that  $P(A_t = a_t | \mathbf{H}_t) = 1/2$  with probability one for  $a_t \in \{-1, 1\}$ ,  $t = 1, 2$ ; extensions to more complex randomization schemes including those with feasible sets of treatments is straightforward (Schulte et al., 2014). A treatment regime in this context is a pair of functions  $\boldsymbol{\pi} = (\pi_1, \pi_2)$  where  $\pi_t : \text{dom } \mathbf{H}_t \rightarrow \text{dom } A_t$  so that a decision maker following  $\boldsymbol{\pi}$  would recommend treatment  $\pi_t(\mathbf{h}_t)$  to a patient presenting with  $\mathbf{H}_t = \mathbf{h}_t$  at time  $t = 1, 2$ . We define a treatment regime as optimal if it leads to maximal mean outcome if applied to the population from which  $\mathcal{D}_n$  is drawn (other definitions of optimality are possible, see Kosorok and Moodie, 2015; Linn et al., 2016). To formally define an optimal treatment regime, we use potential outcomes (Rubin, 1978; Splawa-Neyman

et al., 1990).

Let  $\mathbf{H}_2^*(a_1)$  denote the potential second history under initial treatment  $a_1$  and  $Y^*(a_1, a_2)$  the potential outcome under treatment sequence  $(a_1, a_2)$ . The potential outcome under a regime  $\pi$  is

$$Y^*(\pi) = \sum_{(a_1, a_2)} Y^*(a_1, a_2) \mathbf{1}_{\pi_1(\mathbf{H}_1)=a_1} \mathbf{1}_{\pi_2\{\mathbf{H}_2^*(a_1)\}=a_2},$$

where  $\mathbf{1}_u$  is an indicator that  $u$  is true. For any regime,  $\pi$ , define  $V(\pi) = \mathbb{E}Y^*(\pi)$ ; an optimal regime,  $\pi^{\text{opt}}$ , satisfies  $V(\pi^{\text{opt}}) \geq V(\pi)$  for all  $\pi$ . Our sample size procedures depend on an estimator of  $\pi^{\text{opt}}$ , in order to construct such an estimator, we make the following assumptions: (C1) sequential ignorability,  $\{\mathbf{H}_2^*(a_1), Y^*(a_1, a_2) : (a_1, a_2) \in \{-1, 1\}^2\} \perp A_t | \mathbf{H}_t$  for  $t = 1, 2$ ; (C2) positivity,  $P(A_t = a_t | \mathbf{H}_t) > 0$  with probability one for each  $a_t \in \{-1, 1\}$  for  $t = 1, 2$ ; and (C3) consistency,  $Y = Y^*(A_1, A_2)$  and  $\mathbf{H}_2 = \mathbf{H}_2^*(A_1)$ . These assumptions are standard in the context of estimating optimal treatment regimes (Robins, 2004b; Schulte et al., 2014; Chakraborty and Moodie, 2013) with (C1) and (C2) holding by design in a SMART.

Under these assumptions, the optimal regime can be characterized in terms of the data-generating model as follows. Define  $Q_2(\mathbf{h}_2, a_2) = \mathbb{E}(Y | \mathbf{H}_2 = \mathbf{h}_2, A_2 = a_2)$  and  $Q_1(\mathbf{h}_1, a_1) = \mathbb{E}\{\max_{a_2} Q_2(\mathbf{H}_2, a_2) | \mathbf{H}_1 = \mathbf{h}_1, A_1 = a_1\}$ , then  $\pi_t^{\text{opt}}(\mathbf{h}_t) = \arg \max_{a_t} Q_t(\mathbf{h}_t, a_t)$  for  $t = 1, 2$  (see Murphy, 2005b; Schulte et al., 2014). Furthermore, it can be seen that  $V(\pi^{\text{opt}}) = \mathbb{E} \max_{a_1} Q_1(\mathbf{H}_1, a_1)$ . Our sample size procedures are based on constructing estimators of  $Q_t(\mathbf{h}_1, a_t)$  for  $t = 1, 2$  and subsequently deriving plug-in estimators of  $\pi$ ; these procedures vary in the structure we impose on these functions. Before describing specific estimators, we state properties of these estimators that we would like to ensure hold with high-probability provided the sample size is sufficiently large.

Let  $\hat{\pi}_n$  denote an estimator of  $\pi^{\text{opt}}$  and let  $B_0 > 0, \gamma, \alpha, \eta, \epsilon, \zeta \in (0, 1)$  be constants. Our goal is to choose  $n$  so that:

(POW) there exists an  $\alpha$ -level test of  $H_0 : V(\pi^{\text{opt}}) \leq B_0$  based on  $\hat{\pi}_n$  that has power at least  $(1 - \gamma) \times 100 + o(1)$  provided  $V(\pi^{\text{opt}}) \geq B_0 + \eta$ ;

$$(OPT) \quad P \left[ \mathbb{E} \{ Y^*(\hat{\pi}_n) | \mathcal{D}_n \} \geq V(\pi^{\text{opt}}) - \epsilon \right] \geq 1 - \zeta + o(1).$$

Condition (POW) ensures sufficient power to test the effectiveness of the optimal treatment regime relative to some baseline expected outcome,  $B_0$ , e.g., the expected outcome under some standard of care. Condition (OPT) ensures that the expected performance of the estimated optimal regime will be near-optimal with high-probability. These conditions are analogous to those used to size a one-stage clinical trial for estimation of an optimal regime except that (OPT) controls the performance of the estimated optimal regime whereas Laber et al. (2015) control the *estimated performance* of the estimated optimal regime. Furthermore, like the one-stage setting, our sample size procedures depend on approximating the sampling distribution of an estimator of  $\mathbb{E} Y^*(\pi^{\text{opt}})$ ; however, as we will later illustrate, constructing a high-quality approximation is markedly more complex in the multistage setting (see also Dawid, 1994; Chakraborty et al., 2009; Moodie et al., 2010; Hirano and Porter, 2012; Chakraborty et al., 2014; Laber et al., 2014; Luedtke and Van Der Laan, 2016).

## 2.2 Sample size procedures

We derive two sample size procedures. The first procedure imposes more parametric structure on the joint distribution of  $(\mathbf{X}_1, A_1, \mathbf{X}_2, A_2, Y)$  than is typical in Q-learning and thereby avoids (or rather assumes away) some of the complexities associated with non-regularity and exceptional laws (Robins, 2004b; Chakraborty et al., 2009; Moodie et al., 2010; Chakraborty et al., 2013, 2014; Laber et al., 2014; Song et al., 2015b). The second proposed procedure does not impose as much parametric structure but at the expense of a more complex and potentially conservative sample size estimator.

### 2.2.1 Normality-based sample size procedure

We make the following assumptions about the generative model:

- (AN1)  $Q_2(\mathbf{h}_2, a_2) = \mathbf{h}_{2,0}^\top \boldsymbol{\beta}_{2,0}^* + a_2 \mathbf{h}_{2,1}^\top \boldsymbol{\beta}_{2,1}^*$ , where  $\mathbf{h}_{2,0} \in \mathbb{R}^{p_{2,0}}$ ,  $\mathbf{h}_{2,1} \in \mathbb{R}^{p_{2,1}}$  are summaries of  $\mathbf{h}_2$  and  $\boldsymbol{\beta}_{2,0}^* \in \mathbb{R}^{p_{2,0}}$ ,  $\boldsymbol{\beta}_{2,1}^* \in \mathbb{R}^{p_{2,1}}$  are unknown parameters;
- (AN2)  $\mathbb{E}(\mathbf{H}_{2,0}^\top \boldsymbol{\beta}_{2,0}^* | \mathbf{H}_1 = \mathbf{h}_1, A_1 = a_1) = \mathbf{h}_{1,0}^\top \boldsymbol{\xi}_{1,0}^* + a_1 \mathbf{h}_{1,1}^\top \boldsymbol{\xi}_{1,1}^*$ , where  $\mathbf{h}_{1,0} \in \mathbb{R}^{p_{1,0}}$ ,  $\mathbf{h}_{1,1} \in \mathbb{R}^{p_{1,1}}$  are summaries of  $\mathbf{h}_1$  and  $\boldsymbol{\xi}_{1,0}^* \in \mathbb{R}^{p_{1,0}}$ ,  $\boldsymbol{\xi}_{1,1}^* \in \mathbb{R}^{p_{1,1}}$  are unknown parameters;
- (AN3)  $\mathbf{H}_{2,1}^\top \boldsymbol{\beta}_{2,1}^* = \mathbf{H}_{1,2}^\top \boldsymbol{\varpi}_{1,2}^* + A_1 \mathbf{H}_{1,3}^\top \boldsymbol{\varpi}_{1,3}^* + \tau^* Z$ , where  $\mathbf{H}_{1,2} \in \mathbb{R}^{p_{1,2}}$ ,  $\mathbf{H}_{1,3} \in \mathbb{R}^{p_{1,3}}$  are summaries of  $\mathbf{H}_1$ ,  $Z$  is a standard normal random variable which is independent of  $\mathbf{H}_1, A_1$ , and  $\tau^* > 0$ ,  $\boldsymbol{\varpi}_{1,2}^* \in \mathbb{R}^{p_{1,2}}$ ,  $\boldsymbol{\varpi}_{1,3}^* \in \mathbb{R}^{p_{1,3}}$ , are unknown parameters;
- (AN4)  $(\mathbf{H}_{1,0}^\top \boldsymbol{\xi}_{1,0}^*, \mathbf{H}_{1,1}^\top \boldsymbol{\xi}_{1,1}^*, \mathbf{H}_{1,2}^\top \boldsymbol{\varpi}_{1,2}^*, \mathbf{H}_{1,3}^\top \boldsymbol{\varpi}_{1,3}^*)^\top \sim \text{Normal}(\boldsymbol{\omega}^*, \boldsymbol{\Omega}^*)$ , where  $\boldsymbol{\omega}^* \in \mathbb{R}^4$  and  $\boldsymbol{\Omega}^* \in \mathbb{R}^{4 \times 4}$  are unknown parameters.

Assumptions (AN1)-(AN3) are similar to those used in interactive Q-learning (IQ-learning Laber et al., 2014) except that in IQ-learning, (AN3) is replaced with a more general location-scale model than the normal linear model used here. The summaries of the history  $\mathbf{h}_t$  for  $t = 1, 2$  can include basis expansions or other non-linear terms as needed. These assumptions were motivated by a desire to create a generative model that is conceptually consistent with the analysis model used in linear Q-learning which remains the most commonly used method for estimating an optimal treatment regime from SMARTs. Assumption (AN4) is not required by IQ-learning as it conditions on  $\mathbf{H}_1$ . The assumption of joint normality could be relaxed, for example, by using a copula or semi-parametric model, but at the expense of more complex expressions that are less amenable to sample size calculations. Therefore we shall not consider such generalizations further.

The following results, which are proved in Appendix A, will be used to inform the construction of an estimator of the optimal treatment regime (see Schulte et al., 2014; Laber et al., 2014, for related expressions). Let  $\Phi$  denote the cumulative distribution function of a standard normal random variable.

**Lemma 2.2.1.** Assume (AN1)-(AN3). For any  $\beta_1 = (\xi_{1,0}^\top, \xi_{1,1}^\top, \varpi_{1,2}^\top, \varpi_{1,3}^\top)^\top$  define

$$Q_1(\mathbf{h}_1, a_1; \beta_1, \tau) = \mathbf{h}_{1,0}^\top \beta_{1,0} + a_1 \mathbf{h}_{1,1}^\top \beta_{1,1} + \frac{2\tau}{\sqrt{2\pi}} \exp \left\{ -\frac{(\mathbf{h}_{1,2}^\top \beta_{1,2} + a_1 \mathbf{h}_{1,3}^\top \beta_{1,3})^2}{2\tau^2} \right\} \\ + (\mathbf{h}_{1,2}^\top \beta_{1,2} + a_1 \mathbf{h}_{1,3}^\top \beta_{1,3}) \left[ 1 - 2\Phi \left\{ -\frac{(\mathbf{h}_{1,2}^\top \beta_{1,2} + a_1 \mathbf{h}_{1,3}^\top \beta_{1,3})}{\tau} \right\} \right].$$

Then,  $Q_1(\mathbf{h}_1, a_1) = Q_1(\mathbf{h}_1, a_1; \beta_1^*, \tau^*)$ .

Let  $W(\mathbf{H}_1, \beta_1) = (\mathbf{H}_{1,0}^\top \xi_{1,0}, \mathbf{H}_{1,1}^\top \xi_{1,1}, \mathbf{H}_{1,2}^\top \varpi_{1,2}, \mathbf{H}_{1,3}^\top \varpi_{1,3})^\top$  and define  $g : \mathbb{R}^4 \rightarrow \mathbb{R}$  as

$$g(\mathbf{v}) = \max_{\rho \in \{-1, 1\}} \left( v_1 + \rho v_2 + \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{(v_3 + \rho v_4)^2}{2} \right\} + (v_3 + \rho v_4) [1 - 2\Phi\{-(v_3 + \rho v_4)\}] \right);$$

it follows from Lemma (2.2.1) that  $\max_{a_1} Q_1(\mathbf{H}_1, a_1) = \tau^* g \{W(\mathbf{H}_1, \beta_1^*)/\tau^*\}$ . Let  $\psi(\mathbf{v}; \omega, \Omega)$  denote the density of a multivariate normal distribution with mean  $\omega \in \mathbb{R}^4$  and covariance  $\Omega \in \mathbb{R}^{4 \times 4}$  and write  $\text{vech}(\Sigma)$  to denote the vector-half operator of symmetric matrix  $\Sigma$  (Henderson and Searle, 1979). The following result shows that  $V(\pi^{\text{opt}})$  is a smooth function of  $\omega^*$ ,  $\tau^*$ , and  $\Omega^*$ .

**Corollary 2.2.2.** Assume (AN1)-(AN4) and let  $g : \mathbb{R}^4 \rightarrow \mathbb{R}$  be defined as above. Then

$$V(\pi^{\text{opt}}) = \nu \{ \tau^*, \omega^*, \text{vech}(\Omega^*) \} = \int_{\mathbb{R}^4} \tau^* g(\mathbf{v}/\tau^*) \psi(\mathbf{v}; \omega^*, \Omega^*) d\mathbf{v}.$$

Thus, given estimators  $\hat{\tau}_n$ ,  $\hat{\omega}_n$ , and  $\hat{\Omega}_n$  of  $\tau^*$ ,  $\omega^*$ , and  $\Omega^*$ , one can use the preceding result to construct a plugin estimator of  $V(\pi^{\text{opt}})$ . We next describe how to construct these estimators.

Let  $\mathbb{P}_n$  denote the empirical measure. Define  $Q_2(\mathbf{h}_2, a_2; \beta_2) = \mathbf{h}_{2,0}^\top \beta_{2,0} + a_2 \mathbf{h}_{2,1}^\top \beta_{2,1}$  and subsequently define  $\hat{\beta}_{2,n} = \arg \min_{\beta_2} \mathbb{P}_n \{ Y - Q_2(\mathbf{H}_2, A_2; \beta_2) \}^2$ . Thus,  $\hat{Q}_{2,n}(\mathbf{h}_2, a_2) = Q_2(\mathbf{h}_2, a_2; \hat{\beta}_{2,n})$  and the estimated optimal rule at the second stage is  $\hat{\pi}_{2,n}(\mathbf{h}_2) = \arg \max_{a_2} \hat{Q}_{2,n}(\mathbf{h}_2, a_2)$ . Define

$$\hat{\xi}_{1,0,n}, \hat{\xi}_{1,1,n} = \arg \min_{\xi_{1,0}, \xi_{1,1}} \mathbb{P}_n \left( \mathbf{H}_{2,0}^\top \hat{\beta}_{2,0,n} - \mathbf{H}_{1,0}^\top \xi_{1,0} - A_1 \mathbf{H}_{1,1}^\top \xi_{1,1} \right)^2 \\ \hat{\varpi}_{1,2,n}, \hat{\varpi}_{1,3,n} = \arg \min_{\varpi_{1,2}, \varpi_{1,3}} \mathbb{P}_n \left( \mathbf{H}_{2,1}^\top \hat{\beta}_{2,1,n} - \mathbf{H}_{1,2}^\top \varpi_{1,2} - A_1 \mathbf{H}_{1,3}^\top \varpi_{1,3} \right)^2$$



so that  $\widehat{\beta}_{1,n} = (\widehat{\xi}_{1,0,n}^\top, \widehat{\xi}_{1,1,n}^\top, \widehat{\omega}_{1,2,n}^\top, \widehat{\omega}_{1,3,n}^\top)^\top$  is the least-squares estimator of  $\beta_1^* = (\xi_{1,0}^{*\top}, \xi_{1,1}^{*\top}, \omega_{1,2}^{*\top}, \omega_{1,3}^{*\top})^\top$ . In addition, define  $\widehat{\tau}_n^2 = \mathbb{P}_n \{ \mathbf{H}_{2,1}^\top \widehat{\beta}_{2,1} - \mathbf{H}_{1,2}^\top \widehat{\omega}_{1,2,n} - A_1 \mathbf{H}_{1,3}^\top \widehat{\omega}_{1,3,n} \}^2$ . The plugin estimator of  $Q_1(\mathbf{h}_1, \mathbf{a}_1)$ , based on (2.2.1), is  $\widehat{Q}_{1,n}(\mathbf{h}_1, \mathbf{a}_1) = Q_1(\mathbf{h}_1, \mathbf{a}_1; \widehat{\beta}_{1,n}, \widehat{\tau}_n)$  and the estimated optimal decision rule at the first stage is  $\widehat{\pi}_{1,n}(\mathbf{h}_1) = \arg \max_{a_1} \widehat{Q}_{1,n}(\mathbf{h}_1, a_1)$ . Furthermore, define  $\widehat{\omega}_n = \mathbb{P}_n W(\mathbf{H}_1, \widehat{\beta}_{1,n})$  and  $\widehat{\Omega}_n = \mathbb{P}_n \{ W(\mathbf{H}_1, \widehat{\beta}_{1,n}) - \widehat{\omega}_n \} \{ W(\mathbf{H}_1, \widehat{\beta}_{1,n}) - \widehat{\omega}_n \}^\top$ . The plugin estimator of  $V(\boldsymbol{\pi}^{\text{opt}})$  is

$$\widehat{V}_n = \nu \{ \widehat{\tau}_n, \widehat{\omega}_n, \text{vech}(\widehat{\Omega}_n) \} = \int_{\mathbb{R}^4} \widehat{\tau}_n g(\mathbf{v}/\widehat{\tau}_n) \psi(\mathbf{v}; \widehat{\omega}_n, \widehat{\Omega}_n) d\mathbf{v}.$$

To establish consistency and asymptotic normality of  $\widehat{V}_n$  we assume:

(AN5)  $\{\tau^*, \omega^{*\top}, \text{vech}(\Omega^*)\}^\top \in \Theta \subseteq \mathbb{R}^{15}$ , where  $\Theta$  is compact;

(AN6)  $\sqrt{n} [\{ \widehat{\tau}_n, \widehat{\omega}_n^\top, \text{vech}(\widehat{\Omega}_n)^\top \}^\top - \{ \tau^*, \omega^{*\top}, \text{vech}(\Omega^*)^\top \}^\top] \rightsquigarrow \text{Normal}(0, \Sigma^*)$ , where  $\Sigma^* \in \mathbb{R}^{15 \times 15}$  is positive definite.

Condition (AN6) follows from moment conditions that are common in  $M$ -estimation; we provide sufficient conditions for (AN6) in Appendix A.

**Lemma 2.2.3.** *Assume (C1)-(C3) and (AN1)-(AN6). Then,*

$$\sqrt{n} \{ \widehat{V}_n - V(\boldsymbol{\pi}^{\text{opt}}) \} \rightsquigarrow \text{Normal}(0, \sigma^{*2}),$$

where  $\sigma^{*2} = \nabla \nu \{ \tau^*, \omega^*, \text{vech}(\Omega^*) \}^\top \Sigma \nabla \nu \{ \tau^*, \omega^*, \text{vech}(\Omega^*) \}$ .

Let  $\widehat{\sigma}_n^2$  be a consistent estimator of  $\sigma^{*2}$  and let  $z_{1-\varrho}$  the  $(1-\varrho)$  quantile of a standard normal distribution, then a test that rejects when  $\sqrt{n} \{ \widehat{V}_n - B_0 \} / \widehat{\sigma}_n \geq z_{1-\alpha}$  is an (asymptotic)  $\alpha$ -level test of  $H_0 : V(\boldsymbol{\pi}^{\text{opt}}) \leq B_0$  with power exceeding  $\Phi(z_\alpha + \sqrt{n}\eta/\sigma^*) + o(1)$  when  $V(\boldsymbol{\pi}^{\text{opt}}) \geq B_0 + \eta$ . Thus, choosing  $n = \lceil (\sigma^{*2}/\eta^2) \{ \Phi^{-1}(1-\gamma) + z_{1-\alpha} \}^2 \rceil$ , satisfies (POW) asymptotically. This expression depends on  $\sigma^*$ , which is unknown in general, thus, a value for  $\sigma^*$  must be elicited from domain experts or estimated from historical data.

The preceding sample size has a familiar form which is unsurprising as it is derived from a test statistic which is asymptotically normal. However, what is perhaps more surprising, is that a similar sample size formula can also be used to ensure that condition (OPT) holds under the following regularity conditions. Define  $\Delta Q_j = Q_j(\mathbf{H}_j, 1) - Q_j(\mathbf{H}_j, -1)$  and  $\Delta \widehat{Q}_{j,n} = \widehat{Q}_{j,n}(\mathbf{H}_j, 1) - \widehat{Q}_{j,n}(\mathbf{H}_j, -1)$  for  $j = 1, 2$ . To select  $n$  so that (OPT) also holds we further assume:

(AN7) there exists positive sequences  $\{c_{n,j}\}_{n \geq 1}$  and  $\{\ell_{n,j}\}_{n \geq 1}$  satisfying  $\liminf_{n \rightarrow \infty} c_{n,j} \geq c_{0,j} > 0$  and  $\liminf_{n \rightarrow \infty} \ell_{n,j} \geq \ell_{0,j} > 0$  such that

$$P \{ \sqrt{n} |\Delta \widehat{Q}_j - \Delta Q_j| > t \} \leq \exp(-c_{n,j} t^{\ell_{n,j}}),$$

for all  $n$  and  $j = 1, 2$ ;

(AN8) there exists  $M_j, \kappa_j > 0$  such that  $P(|\Delta Q_j| \leq \epsilon) \leq M_j \epsilon^{\kappa_j}$  for  $j = 1, 2$  as  $\epsilon \rightarrow 0$ .

The preceding assumptions are relatively mild with (AN7) being weaker than requiring a subexponential tail; e.g., (AN7) and (AN8) would be satisfied if the histories and outcomes are normally distributed. The following result characterizes the concentration of the marginal mean outcome under  $\widehat{\pi}_n$  about  $\pi^{\text{opt}}$  which can subsequently be used to choose a sample size  $n$  that satisfies (OPT).

**Lemma 2.2.4.** *Assume (C1)-(C3) and (AN1)-(AN8). Then there exists  $K$  and  $\delta > 0$  such that*

$$|V(\widehat{\pi}_n) - V(\pi^{\text{opt}})| \leq K n^{-\delta} |\widehat{V}_n - V(\pi^{\text{opt}})| + o_p(1/\sqrt{n}).$$

**Corollary 2.2.5.** *Assume (C1)-(C3) and (AN1)-(AN8). Then setting*

$$n = \left\lceil \left\{ \frac{\Phi^{-1}(1-\zeta)\sigma^*}{\epsilon} \right\}^2 \right\rceil,$$

*satisfies (OPT).*

**Remark 2.2.6.** Given pilot or historical data, one can construct a plug-in estimator of  $\sigma^{*2}$ . In the absence of such data, one can use an elicited value for the variance of  $Y$  under standard care as an *ad hoc* surrogate for  $\sigma^{*2}$ . Heuristic justification for this surrogate is as follows. If the variance of the outcome is at least as large under standard care as it is under the optimal regime *and* the parametric estimator  $\widehat{V}_n$  is at least as efficient as the sample mean of  $n$  observations collected under the optimal policy, then

$$\begin{aligned}
\sigma^{*2} &= \lim_{n \rightarrow \infty} \text{Var}[\sqrt{n}\{\widehat{V}_n - V(\boldsymbol{\pi}^{\text{opt}})\}] \\
&\leq \lim_{n \rightarrow \infty} \text{Var}\left[\frac{1}{\sqrt{n}} \sum_{i=1}^n \{Y_i^*(\boldsymbol{\pi}^{\text{opt}}) - V(\boldsymbol{\pi}^{\text{opt}})\}\right] \\
&\leq \lim_{n \rightarrow \infty} \text{Var}\left[\frac{1}{\sqrt{n}} \sum_{i=1}^n \{Y_i^*(\boldsymbol{\pi}^{\text{soc}}) - V(\boldsymbol{\pi}^{\text{soc}})\}\right] \\
&= \text{Var}\{Y^*(\boldsymbol{\pi}^{\text{soc}})\},
\end{aligned}$$

where  $\boldsymbol{\pi}^{\text{soc}}$  denotes standard of care. If one is unwilling to make the above assumptions, an alternative would be to inflate the elicited value for  $\text{Var}\{Y^*(\boldsymbol{\pi}^{\text{soc}})\}$  by a constant factor.

## 2.2.2 Projection-based sample size procedure

Despite a deluge of new estimators of optimal treatment regimes (see Zhang et al., 2017; Zhou et al., 2017; Laber and Staicu, 2017; Tao and Wang, 2017, and references therein),  $Q$ -learning with linear models remains among the most commonly used methods in practice. This popularity can be partly attributed to: (i) the heavy use of linear models in seminal papers on estimation of optimal treatment regimes (Murphy, 2003b; Robins, 2004b; Murphy, 2005b; Qian and Murphy, 2011); (ii) minimal requirements on the joint distribution of the data-generating model; (iii) theoretical tractability (Chakraborty et al., 2009; Moodie et al., 2010; Chakraborty et al., 2014; Laber et al., 2014); and (iv) good empirical performance even under some forms of misspecification (Schulte et al., 2014). Thus, our second sample size procedure is designed for the setting where analysts plan to estimate the optimal

regime using  $Q$ -learning with linear models. We do not assume that the analysis model is correctly specified nor do we impose any parametric structure on the generative model. However, unlike the procedure described in the preceding section, the resultant sample size procedure derived here relies on quantities that would be difficult to elicit from domain experts, we therefore require that one has suitable pilot data available; such data could be historical or collected as an internal pilot.

We assume that for  $t = 1, 2$  one postulates models of the form  $Q_t(\mathbf{h}_t, a_t; \mu_t) = \mathbf{h}_{t,0}^\top \mu_{t,0} + a_t \mathbf{h}_{t,1}^\top \mu_{t,1}$ , where  $\mathbf{h}_{t,0}, \mathbf{h}_{t,1}$  are summaries of  $\mathbf{h}_t$  and  $\mu_t = (\mu_{t,0}^\top, \mu_{t,1}^\top)^\top$  are unknown parameters. Define  $\mu_2^* = \arg \min_{\mu_2} P \{ Y - Q_2(\mathbf{H}_2, A_2; \mu_2) \}^2$  and  $\mu_1^* = \arg \min_{\mu_1} P \{ \max_{a_2} Q_2(\mathbf{H}_2, a_2; \mu_2^*) - Q_1(\mathbf{H}_1, A_1; \mu_1) \}^2$ . If (AN1) holds then, provided the requisite expectations exist,  $Q_2(\mathbf{h}_2, a_2) = Q_2(\mathbf{h}_2, a_2; \mu_2^*)$ ; however, even if (AN1)-(AN4) hold, it need not follow that  $Q_1(\mathbf{h}_1, a_1) = Q_1(\mathbf{h}_1, a_1; \mu_1^*)$  (Laber et al., 2014). Nevertheless, it is still meaningful to discuss the mean outcome under the estimated optimal regime when these models are misspecified. Define the optimal regime under linear  $Q$ -learning with the above class of models as  $\boldsymbol{\pi}^{Q,\text{opt}} = (\pi_1^{Q,\text{opt}}, \pi_2^{Q,\text{opt}})$  so that  $\pi_t^{Q,\text{opt}}(\mathbf{h}_t) = \arg \max_{a_t} Q_t(\mathbf{h}_t, a_t; \mu_t^*)$ .<sup>1</sup> Define  $\hat{\mu}_{2,n} = \arg \min_{\mu_2} \mathbb{P}_n \{ Y - Q_2(\mathbf{H}_2, A_2; \mu_2) \}^2$  and  $\hat{\mu}_{1,n} = \arg \min_{\mu_1} \{ \max_{a_2} Q_2(\mathbf{H}_2, a_2; \hat{\mu}_{2,n}) - Q_2(\mathbf{H}_1, A_1; \mu_1) \}^2$ . The estimated optimal decision rule at time  $t$  is  $\hat{\pi}_{t,n}^Q(\mathbf{h}_t) = \arg \max_{a_t} Q_t(\mathbf{h}_t, a_t; \hat{\mu}_{t,n})$ .

It is well-known that  $V(\boldsymbol{\pi}^{\text{opt}})$  is not a smooth functional of the generative model and consequently standard approaches for inference, e.g., the bootstrap or series approximations, will not hold without modification (Robins, 2004b; Moodie et al., 2010; Chakraborty et al., 2009, 2013, 2014; Laber et al., 2014; Luedtke and Van Der Laan, 2016). To derive a test which satisfies (POW) we invert a variant of a projection confidence interval (Berger and Boos, 1994; Robins, 2004b) for  $V(\boldsymbol{\pi}^{Q,\text{opt}})$ ; the interval we propose holds regardless of misspecification of the  $Q$ -functions and does not require strong parametric assumptions on the underlying generative model. This approach requires a confidence set for  $(\mu_1^*, \mu_2^*)$

<sup>1</sup>The notation  $\boldsymbol{\pi}^{Q,\text{opt}}$  is a bit misleading in that if the  $Q$ -functions are misspecified (which we allow) it need not follow that  $V(\boldsymbol{\pi}^{Q,\text{opt}}) \geq V^Q(\mu_1, \mu_2)$  for all  $(\mu_1, \mu_2) \in \Theta$ ; we do not assume that  $\boldsymbol{\pi}^{Q,\text{opt}}$  satisfies such an inequality. For additional discussion, see Qian and Murphy (2011) and references therein.

which we construct as follows. Let  $\mathbf{C}_2$  be as in (AN7) and let

$$\widehat{\mathfrak{W}}_{2,n} = \{\mathbb{P}_n \mathbf{C}_2 \mathbf{C}_2^\top\}^{-1} \mathbb{P}_n \mathbf{C}_2 \mathbf{C}_2^\top (Y - \mathbf{C}_2^\top \widehat{\mu}_{2,n}) \{\mathbb{P}_n \mathbf{C}_2 \mathbf{C}_2^\top\}^{-1}$$

so that  $\mathfrak{Z}_{2,n,\varepsilon} = \{\mu_2 : n(\mu_2 - \widehat{\mu}_{2,n})^\top \widehat{\mathfrak{W}}_{2,n}^{-1} (\mu_2 - \widehat{\mu}_{2,n}) \leq \chi_{1-\varepsilon, \dim(\mathbf{C}_2)}\}$  is a Wald-type  $(1 - \varepsilon) \times 100$  confidence set for  $\mu_2^*$ , where  $\chi_{q,k}^2$  is the  $q$ th quantile of a chi-square random variable with  $k$  degrees of freedom. For each  $\mu_2$  define

$$\mu_1^*(\mu_2) = \arg \min_{\mu_1} \mathbb{E} \left\{ \max_{a_2} Q_2(\mathbf{H}_2, a_2; \mu_2) - Q_1(\mathbf{H}_1, A_1; \mu_1) \right\}^2,$$

so that  $\mu_1^*(\mu_2)$  is denotes the population-level for the first-stage  $Q$ -function were it known that  $\mu_2^* = \mu_2$ ; thus,  $\mu_1^* = \mu_1^*(\mu_2^*)$ . Define

$$\widehat{\mu}_{1,n}(\mu_2) = \arg \min_{\mu_1} \mathbb{P}_n \left\{ \max_{a_2} Q_2(\mathbf{H}_2, a_2; \mu_2) - Q_1(\mathbf{H}_1, A_1; \mu_1) \right\}^2,$$

to be the least-squares estimator of  $\mu_1^*(\mu_2)$ . Let  $\mathbf{C}_1 = (\mathbf{H}_{1,0}^\top, A_1 \mathbf{H}_{1,1}^\top)^\top$  and define

$$\mathfrak{W}_{1,n}(\mu_2) = \{\mathbb{P}_n \mathbf{C}_1 \mathbf{C}_1^\top\}^{-1} \mathbb{P}_n \mathbf{C}_1 \mathbf{C}_1^\top \left\{ \max_{a_2} Q_2(\mathbf{H}_2, a_2; \mu_2) - \mathbf{C}_1^\top \widehat{\mu}_{1,n}(\mu_2) \right\} \{\mathbb{P}_n \mathbf{C}_1 \mathbf{C}_1^\top\}^{-1}.$$

A  $(1 - \varepsilon) \times 100\%$  Wald-type confidence set for  $\mu_1^*(\mu_2)$  is

$$\mathfrak{Z}_{1,n,\varepsilon}(\mu_2) = \{\mu_1 : n[\mu_1 - \widehat{\mu}_{1,n}(\mu_2)]^\top \widehat{\mathfrak{W}}_{1,n}^{-1}(\mu_2) [\mu_1 - \widehat{\mu}_{1,n}(\mu_2)] \leq \chi_{1-\varepsilon, \dim(\mathbf{C}_1)}\}.$$

Thus, given  $\varepsilon_1, \varepsilon_2 \in (0, 1)$  with  $\vartheta = \varepsilon_1 + \varepsilon_2 \leq 1$ , a  $(1 - \vartheta) \times 100\%$  confidence set for  $(\mu_1^*, \mu_2^*)$  is

$$\Xi_{n,1-\vartheta} = \{(\mu_1, \mu_2) : \mu_2 \in \mathfrak{Z}_{2,n,1-\varepsilon_2} \text{ and } \mu_1 \in \mathfrak{Z}_{1,n,1-\varepsilon_1}(\mu_2)\}.$$

For each  $\mu_1, \mu_2$  define

$$\begin{aligned} \delta(\mu_1, \mu_2) = & 4Y 1_{A_1} \mathbf{H}_{1,1}^\top \mu_{1,1} > 0 1_{A_2} \mathbf{H}_{2,1}^\top \mu_{2,1} > 0 + 2 \left\{ 1_{A_1} \mathbf{H}_{1,1}^\top \mu_{1,1} \leq 0 - (1/2) \right\} \max_{a_1} Q_1(\mathbf{H}_1, a_1; \mu_1) \\ & + 4 1_{A_1} \mathbf{H}_{1,1}^\top \mu_{1,1} > 0 \left\{ 1_{A_2} \mathbf{H}_{2,1}^\top \mu_{2,1} \leq 0 - (1/2) \right\} \max_{a_2} Q_2(\mathbf{H}_2, a_2; \mu_2), \end{aligned}$$

and subsequently define  $\widehat{V}_n^Q(\mu_1, \mu_2) = \mathbb{P}_n \delta(\mu_1, \mu_2)$  and its population-level analog  $V^Q(\mu_1, \mu_2) = \mathbb{E} \delta(\mu_1, \mu_2)$ . Then,  $\widehat{V}_n(\widehat{\mu}_{n,1}, \widehat{\mu}_{n,2})$  is the augmented inverse probability weighted estimator of  $V(\pi^{Q,\text{opt}})$  (Zhang et al., 2013) and  $V^Q(\mu_1^*, \mu_2^*) = V(\pi^{Q,\text{opt}})$  (see also Qian and Murphy, 2011; Zhao et al., 2015). Define

$$\zeta^2(\mu_1, \mu_2) = \mathbb{E} \left\{ \delta(\mu_1, \mu_2) - \mathbb{E} \delta(\mu_1, \mu_2) \right\}^2$$

and

$$\widehat{\zeta}_n^2 = \mathbb{P}_n \left\{ \delta(\mu_1, \mu_2) - \mathbb{P}_n \delta(\mu_1, \mu_2) \right\}^2.$$

For any fixed  $(\mu_1, \mu_2)$ , it follows that

$$\sqrt{n} \left\{ \widehat{V}_n^Q(\mu_1, \mu_2) - V^Q(\mu_1, \mu_2) \right\} \rightsquigarrow \text{Normal} \left\{ 0, \zeta^2(\mu_1, \mu_2) \right\},$$

provided that  $\mathbb{E} \delta^2(\mu_1, \mu_2) < \infty$ . Choose  $\vartheta_1$  and  $\vartheta_2$  such that  $\vartheta_1 + \vartheta_2 = \alpha$ , then the proposed  $\alpha$ -level test for (POW) rejects when

$$\inf_{(\mu_1, \mu_2) \in \Xi_{n,1-\vartheta_1}} \left[ \widehat{V}_n^Q(\mu_1, \mu_2) - \frac{z_{1-\vartheta_2} \widehat{\zeta}_n(\mu_1, \mu_2)}{\sqrt{n}} \right] \geq B_0.$$

Under the null,  $V(\pi^{Q,\text{opt}}) \leq B_0$ , so that the type I error is bounded above by

$$\begin{aligned}
P \left\{ \inf_{(\mu_1, \mu_2) \in \Xi_{n,1-\vartheta_1}} \left[ \widehat{V}_n^Q(\mu_1, \mu_2) - \frac{z_{1-\vartheta_2} \widehat{\zeta}_n(\mu_1, \mu_2)}{\sqrt{n}} \right] \geq V(\pi^{Q,\text{opt}}) \right\} \\
\leq P \left\{ \widehat{V}_n^Q(\mu_1^*, \mu_2^*) - \frac{z_{1-\vartheta_2} \widehat{\zeta}_n(\mu_1^*, \mu_2^*)}{\sqrt{n}} \geq V(\pi^{Q,\text{opt}}) \right\} + \vartheta_1 + o(1) \\
= P \left\{ \sqrt{n} \left[ \frac{\widehat{V}_n^Q(\mu_1^*, \mu_2^*) - V^Q(\mu_1^*, \mu_2^*)}{\widehat{\zeta}_n(\mu_1^*, \mu_2^*)} \right] \geq z_{1-\vartheta_2} \right\} + \vartheta_1 + o(1) \\
\leq \vartheta_1 + \vartheta_2 + o(1),
\end{aligned}$$

where the first inequality follows from  $P \{(\mu_1^*, \mu_2^*) \in \Xi_{n,1-\vartheta_1}\} \geq 1 - \vartheta_1 + o(1)$ . If  $\widehat{\zeta}_n(\mu_1, \mu_2) > 0$  with probability one for all  $(\mu_1, \mu_2) \in \Theta$ , then it can be seen that the power of the proposed test is

$$\begin{aligned}
P \left\{ \inf_{(\mu_1, \mu_2) \in \Xi_{n,1-\vartheta_1}} \left[ \widehat{V}_n^Q(\mu_1, \mu_2) - \frac{z_{1-\vartheta_2} \widehat{\zeta}_n(\mu_1, \mu_2)}{\sqrt{n}} \right] \geq B_0 \right\} \\
= P \left\{ \inf_{(\mu_1, \mu_2) \in \Xi_{n,1-\vartheta_1}} \left[ \frac{\sqrt{n} \{ \widehat{V}_n^Q(\mu_1, \mu_2) - V^Q(\mu_1, \mu_2) \}}{\widehat{\zeta}_n(\mu_1, \mu_2)} \right. \right. \\
\left. \left. + \frac{\sqrt{n} \{ V^Q(\mu_1, \mu_2) - B_0 \}}{\widehat{\zeta}_n(\mu_1, \mu_2)} \right] \geq z_{1-\vartheta_2} \right\} \tag{2.1}
\end{aligned}$$

$$\begin{aligned}
\geq P \left\{ \inf_{(\mu_1, \mu_2) \in \Xi_{n,1-\vartheta_1}} \left( \frac{\sqrt{n} \{ \widehat{V}_n^Q(\mu_1, \mu_2) - V^Q(\mu_1, \mu_2) \}}{\widehat{\zeta}_n(\mu_1, \mu_2)} \right. \right. \\
\left. \left. + \frac{\min[\sqrt{n} \{ V^Q(\mu_1, \mu_2) - B_0 \}, \sqrt{n} \eta]}{\widehat{\zeta}_n(\mu_1, \mu_2)} \right) \geq z_{1-\vartheta_2} \right\}. \tag{2.2}
\end{aligned}$$

The minimum in (2.2) is analogous to plugging-in the smallest possible difference under the alternative  $V(\pi^{Q,\text{opt}}) - B_0 \geq \eta$ ; see Remark 2.2.9 for additional discussion. The sampling distribution of the test statistic under the alternative is complex and difficult to approximate using series approximations; thus, to estimate a sample size that will yield the desired power, we use the bootstrap.

### 2.2.2.1 Bootstrap power calculation

We assume that one has available pilot data  $\mathcal{D}_{n_0} = \{(\mathbf{X}_{1,i}, A_{1,i}, \mathbf{X}_{2,i}, A_{2,i}, Y_i)\}_{i=1}^{n_0}$  comprising  $n_0$  *i.i.d.* trajectories from the same population from which trial participants will be drawn. We estimate the power required in (POW) using the bootstrap with a resample size of  $n \geq n_0$  and solve for the smallest  $n$  such that the estimated power exceeds a given threshold. In our asymptotic analyses, we let both  $n$  and  $n_0$  diverge to infinity; however, as we anticipate the trial sample size to be much larger than that of the pilot, we focus on an asymptotics in which  $n$  goes to infinity “first.” We assume:

$$(PR1) \quad \mathbb{E}Y^2\|\mathbf{C}_2\|^2 < \infty \text{ and } \mathbb{E}\|\mathbf{C}_1\|^2\|\mathbf{C}_2\|^2 < \infty;$$

$$(PR2) \quad \mathbb{E}\mathbf{C}_1\mathbf{C}_1^\top \text{ and } \mathbb{E}\mathbf{C}_2\mathbf{C}_2^\top \text{ are finite and strictly positive definite;}$$

$$(PR3) \quad \inf_{(\mu_1, \mu_2) \in \Theta} \mathbb{E} \{ \delta(\mu_1, \mu_2) - \mathbb{E}\delta(\mu_1, \mu_2) \}^2 > 0 \text{ and } \sup_{(\mu_1, \mu_2) \in \Theta} \mathbb{E} \{ \delta(\mu_1, \mu_2) - \mathbb{E}\delta(\mu_1, \mu_2) \}^2 < \infty;$$

$$(PR4) \quad \text{the classes } \mathcal{F}_1 = \{ \delta(\mu_1, \mu_2) : (\mu_1, \mu_2) \in \Theta \} \text{ and } \mathcal{F}_2 = \{ \delta^2(\mu_1, \mu_2) : (\mu_1, \mu_2) \in \Theta \} \text{ are Donsker;}$$

$$(PR5) \quad \mathbb{E}\delta(\mu_1, \mu_2) \text{ is uniformly continuous in a neighborhood of } (\mu_1^*, \mu_2^*).$$

The foregoing assumptions are standard in linear  $Q$ -learning and mirror those used in linear regression (Laber et al., 2014).

Let  $\mathbb{P}_{n, n_0}^{(b)}$  denote the bootstrap empirical distribution corresponding to a resample size of  $n$ . For any functional  $Z_n = f(P, \mathbb{P}_{n_0})$  we define its bootstrap analog  $Z_{n_0, n}^{(b)} = f \left\{ \mathbb{P}_{n_0}, \mathbb{P}_{n, n_0}^{(b)} \right\}$ . Let  $P_B$  denote probabilities computed with respect the bootstrap distribution conditional on the pilot data. The bootstrap estimator of the sample size required for (POW) is the



positive integer  $n$  which solves

$$P_B \left\{ \inf_{(\mu_1, \mu_2) \in \Xi_{n_0, n, 1-\vartheta_1}^{(b)}} \left( \frac{\sqrt{n} \{ \widehat{V}_{n_0, n}^{Q(b)}(\mu_1, \mu_2) - \widehat{V}_{n_0}^Q(\mu_1, \mu_2) \}}{\widehat{\zeta}_{n_0, n}^{(b)}(\mu_1, \mu_2)} + \frac{\min \left[ \sqrt{n} \{ \widehat{V}_{n_0}^Q(\mu_1, \mu_2) - B_0 \}, \sqrt{n} \eta \right]}{\widehat{\zeta}_{n_0, n}^{(b)}(\mu_1, \mu_2)} \right) \geq z_{1-\vartheta_2} \right\} \geq 1 - \gamma$$

where  $\vartheta_1 + \vartheta_2 = \alpha$ ; the probability on the left hand side of the inequality can be computed to desired precision by Monte Carlo methods. The following results establish consistency of the bootstrap as  $n_0$  and  $n$  diverge.

**Theorem 2.2.7.** *Assume (C1)-(C3) and (PR1)-(PR4). Let  $\vartheta_1 \in (0, 1)$  be fixed. Let  $\alpha, K > 0$  be arbitrary, then*

$$\lim_{n, n_0 \rightarrow \infty} P \left( \sup_{|v| \leq K} \left| P_B \left[ \inf_{(\mu_1, \mu_2) \in \Xi_{n_0, n, 1-\vartheta_1}^{(b)}} \frac{\sqrt{n} \{ \widehat{V}_{n_0, n}^{Q(b)}(\mu_1, \mu_2) - \widehat{V}_{n_0}^Q(\mu_1, \mu_2) \}}{\widehat{\zeta}_{n_0, n}^{(b)}(\mu_1, \mu_2)} \geq v \right] - P \left[ \inf_{(\mu_1, \mu_2) \in \Xi_{n, 1-\vartheta_1}} \frac{\sqrt{n} \{ \widehat{V}_n^Q(\mu_1, \mu_2) - V^Q(\mu_1, \mu_2) \}}{\widehat{\zeta}_n(\mu_1, \mu_2)} \geq v \right] \right| > \alpha \right) = 0.$$

The preceding result does not include  $\sqrt{n} \{ V^Q(\mu_1, \mu_2) - B_0 \} / \widehat{\zeta}_n(\mu_1, \mu_2)$  (or its bootstrap analog) because, under the alternative, provided  $V(\mu_1, \mu_2) > B_0$  for all  $(\mu_1, \mu_2)$  in a sufficiently small neighborhood of  $(\mu_1^*, \mu_2^*)$ , this term (and its bootstrap analog) will diverge to infinity so that the conclusion of the above theorem will hold trivially. The following result characterizes the limiting tail behavior of this term; the factor of  $\sqrt{n}$  on the right-hand-side of each probability assignment reflects the fact that, under the alternative, we expect  $\sqrt{n} \{ V^Q(\pi^{Q, \text{opt}}) - B_0 \}$  to diverge at rate  $\sqrt{n}$ .

**Theorem 2.2.8.** *Assume (C1)-(C3) and (PR1)-(PR4). Let  $\vartheta_1 \in (0, 1)$  and  $\eta \geq 0$  be fixed. In addition, assume that  $n_0 \rightarrow \infty$  as  $n \rightarrow \infty$  and that there exists  $c > 0$  so that  $\inf_{(\mu_1, \mu_2) \in \Omega} \widehat{\zeta}_n(\mu_1, \mu_2) \geq c$*

c. Let  $\varkappa, K > 0$  be arbitrary then

$$\lim_{n, n_0 \rightarrow \infty} P \left( \sup_{|\nu| \leq K} \left| P_B \left[ \inf_{(\mu_1, \mu_2) \in \Xi_{n_0, n, 1-\vartheta_1}^{(b)}} \frac{\min[\sqrt{n} \{ \widehat{V}_{n_0, n}^{Q, (b)}(\mu_1, \mu_2) - B_0 \}, \sqrt{n} \eta]}{\widehat{\zeta}_{n_0, n}^{(b)}(\mu_1, \mu_2)} \geq \sqrt{n} \nu \right] - P \left[ \inf_{(\mu_1, \mu_2) \in \Xi_{n, 1-\vartheta_1}} \frac{\min[\sqrt{n} \{ \widehat{V}_n^Q(\mu_1, \mu_2) - B_0 \}, \sqrt{n} \eta]}{\widehat{\zeta}_n(\mu_1, \mu_2)} \geq \sqrt{n} \nu \right] \right| > \varkappa \right) = 0.$$

To choose  $n$  so that (OPT) holds asymptotically we make use of the following bound. For any (possibly data-dependent) sequence  $(\tilde{\mu}_{1, n}, \tilde{\mu}_{2, n}) \in \Xi_{n, 1-\vartheta_1}$  such that  $\widehat{V}_n^Q(\mu_1^*, \mu_2^*) \leq \widehat{V}_n^Q(\tilde{\mu}_{1, n}, \tilde{\mu}_{2, n}) + o_P(1/\sqrt{n})$  it follows that

$$P \left[ V^Q(\tilde{\mu}_{1, n}, \tilde{\mu}_{2, n}) \geq V^Q(\pi^{Q, \text{opt}}) + \inf_{(\mu_1, \mu_2) \in \Xi_{n, 1-\vartheta_1}} \{ \widehat{V}_n^Q(\mu_1, \mu_2) - V^Q(\mu_1, \mu_2) \} - \sup_{(\mu_1, \mu_2) \in \Xi_{n, 1-\vartheta_1}} \{ \widehat{V}_n^Q(\mu_1, \mu_2) - V^Q(\mu_1, \mu_2) \} \right] \geq 1 - \vartheta_1 + o(1).$$

Thus, if  $\mathfrak{Q}_{n, 1-\vartheta_2, 1-\vartheta_1}$  is the  $(1 - \vartheta_2)$  quantile of

$$\inf_{(\mu_1, \mu_2) \in \Xi_{n, 1-\vartheta_1}} \{ \widehat{V}_n^Q(\mu_1, \mu_2) - V^Q(\mu_1, \mu_2) \} - \sup_{(\mu_1, \mu_2) \in \Xi_{n, 1-\vartheta_1}} \{ \widehat{V}_n^Q(\mu_1, \mu_2) - V^Q(\mu_1, \mu_2) \}$$

then choosing  $\vartheta_1 + \vartheta_2 \leq \zeta$  and  $n$  such that  $\mathfrak{Q}_{n, 1-\vartheta_1, 1-\vartheta_2} / \sqrt{n} \leq \epsilon$  ensures that (OPT) holds asymptotically. Of course,  $\mathfrak{Q}_{n, 1-\vartheta_1, 1-\vartheta_2}$  is unknown so we estimate it using the bootstrap, i.e., we select  $n$  so that  $\mathfrak{Q}_{n_0, n, 1-\vartheta_1, 1-\vartheta_2}^{(b)} / \sqrt{n} \leq \epsilon$ .

**Remark 2.2.9.** To estimate the power at a given sample size one could use the bootstrap analog of (2.1). Indeed, the preceding theoretical results can be easily modified to hold without the min operation. However, the required sample size derived from (2.1) will be based on an *estimated effect size* rather than the minimal effect size of interest,  $\eta$ . A consequence of using the estimated effect size is that as the true effect size increases the estimated required sample size will decrease keeping the power fixed at (approximately)  $(1 - \gamma) \times 100$ . However, in application, it is desirable to have power  $(1 - \gamma) \times 100$  at effect size  $\eta$  but larger

power if the effect size exceeds  $\eta$ . Taking the minimum, as in (2.2), ensures that the power diverges to one as the true effect size grows large.

## 2.3 Simulation experiments

We examine the finite sample performance of the proposed sample size procedures using a series of simulation experiments. Performance is measured in terms of the proposed criteria (POW) and (OPT). For each generative model, we also compute the number of samples required to compare the mean outcomes under standard care to that under the fixed regimes  $\pi^{i,j}$ ,  $i, j \in \{-1, 1\}$ , where  $\pi_1^{i,j}(\mathbf{h}_1) \equiv i$  and  $\pi_2^{i,j}(\mathbf{h}_2) \equiv j$ . This comparison allows us to evaluate how much the sample size must be inflated to estimate and/or evaluate an optimal dynamic treatment regime relative to the comparison of fixed and embedded regimes (Almirall et al., 2012).

We first consider a generative model in which the assumptions (AN1)-(AN8) for the normality-based sample size procedure hold. This generative model is as follows:

$$\begin{aligned}
\mathbf{X}_1 &\sim N_4\{\mathbf{0}, \Omega_{AR1}(0.5)\}, & \mathbf{H}_{1,0}^T &= (1, X_{1,0}), \\
\mathbf{H}_{1,1}^T &= (1, X_{1,1}), & \mathbf{H}_{1,2}^T &= (1, X_{1,2}), \\
\mathbf{H}_{1,3}^T &= (1, X_{1,3}), & A_1, A_2 &\sim_{i.i.d.} \text{Unif}\{-1, 1\}, \\
\phi_1, \phi_2, \nu &\sim_{i.i.d.} N(0, 1), & X_{2,0} &= \mathbf{H}_{1,0}^T \mu_{1,0}^* + A_1 \mathbf{H}_{1,1}^T \mu_{1,1}^* + \phi_1, \\
X_{2,1} &= \mathbf{H}_{1,2}^T \mu_{2,0}^* + A_1 \mathbf{H}_{1,3}^T \mu_{2,1}^* + \phi_2, & \mathbf{H}_{2,0}^T &= (1, X_{1,0}, A_1, X_{2,0}), \\
\mathbf{H}_{2,1}^T &= (1, X_{1,2}, A_1, X_{2,1}), & Y &= \mathbf{H}_{2,0}^T \beta_{2,0}^* + A_2 \mathbf{H}_{2,1}^T \beta_{2,1}^* + \nu,
\end{aligned}$$

where  $\Omega_{AR1}(0.5)$  is an autoregressive covariance matrix such that  $\{\Omega_{AR1}(0.5)\}_{ij} = 0.5^{|i-j|}$ .

Let  $\pi^{\text{fixed,opt}}$  denote the optimal fixed regime such that  $\pi^{\text{fixed,opt}} = \pi^{i^*, j^*}$ , where  $i^*, j^* = \arg\max_{i, j \in \{-1, 1\}} V(\pi^{i,j})$ . We examine the performance of the proposed methods under parameter values which result in the following relationships between  $\pi^{\text{fixed,opt}}$ ,  $\pi^{\text{opt}}$ , and  $B_0$ :

1.  $V(\pi^{\text{opt}}) = V(\pi^{\text{fixed,opt}}) = B_0 + \eta$ ;

$$2. V(\boldsymbol{\pi}^{\text{opt}}) - 0.5\eta = V(\boldsymbol{\pi}^{\text{fixed,opt}}) = B_0 + \eta;$$

$$3. V(\boldsymbol{\pi}^{\text{opt}}) - \eta = V(\boldsymbol{\pi}^{\text{fixed,opt}}) = B_0 + \eta;$$

$$4. V(\boldsymbol{\pi}^{\text{opt}}) - 2\eta = V(\boldsymbol{\pi}^{\text{fixed,opt}}) = B_0 + \eta.$$

Define  $\Delta = \{V(\boldsymbol{\pi}^{\text{opt}}) - V(\boldsymbol{\pi}^{\text{fixed,opt}})\} / \eta$  to be a measure of benefit associated with implementing an optimal dynamic treatment regime relative to the optimal fixed regime. It can be seen that  $\Delta$  ranges from zero to two across the above scenarios. Parameter settings indexing the generative model which yield these values of  $\Delta$  when  $\eta = 1$  are:

$$1. \mu_{1,0} = (-1, 1), \mu_{1,1} = (4, 1), \mu_{2,0} = (-0.4, -1), \mu_{2,1} = (4, 1), \beta_{2,0}^* = (0.5, 0.5, -1, 1) \text{ and } \beta_{2,1}^* = (1, 0.5, 0.5, 1);$$

$$2. \mu_{1,0} = (-1, 1), \mu_{1,1} = (4, 1), \mu_{2,0} = (-0.4, -1), \mu_{2,1} = (-4, 1), \beta_{2,0}^* = (0.5, 0.5, -1, 1) \text{ and } \beta_{2,1}^* = (1, -0.9, 0.5, 1);$$

$$3. \mu_{1,0} = (-1, 1), \mu_{1,1} = (4, 1), \mu_{2,0} = (-0.4, -1), \mu_{2,1} = (-4, 1), \beta_{2,0}^* = (0.5, 0.5, -1, 1) \text{ and } \beta_{2,1}^* = (1, -1.75, 0.5, 1);$$

$$4. \mu_{1,0} = (-1, 1), \mu_{1,1} = (4, 1), \mu_{2,0} = (-0.4, -1), \mu_{2,1} = (-4, 1), \beta_{2,0}^* = (0.5, 0.5, -1, 1) \text{ and } \beta_{2,1}^* = (1, -3.25, 0.5, 1).$$

Let  $\gamma = 0.1$ ,  $\alpha = 0.05$ ,  $\zeta = 0.1$ , and  $\epsilon = 0.3$ . Thus, if (POW) holds, then an  $\alpha$ -level test of  $H_0 : V(\boldsymbol{\pi}^{\text{opt}}) \leq B_0$  will have approximately 90% power, and if (OPT) holds, then  $P[\mathbb{E}\{Y^*(\hat{\boldsymbol{\pi}}_n) | \mathcal{D}_n\} \geq V(\boldsymbol{\pi}^{\text{opt}}) - 0.3] \geq 0.9 + o(1)$ . Recall that the normality-based sample size procedure requires specification of  $\sigma^*$ . We consider three possibilities: (i)  $\sigma^*$  is known, e.g., correctly elicited from domain experts; (ii)  $\sigma^*$  is estimated using a pilot study of  $n_0$  patients; and (iii)  $\sigma^*$  is estimated using the *ad hoc* procedure presented in Remark 2.2.9 using a sample of  $n_0$  subjects treated under standard care wherein we assume that patients are assigned the optimal treatment 80% of the time and suboptimal treatment the remaining 20% of the time. To form a baseline for comparison, we also compute the sample size

required to power a test of the null  $V(\pi^{\text{fixed,opt}}) \leq B_0$  against the alternative  $V(\pi^{\text{fixed,opt}}) > B_0$  where it is assumed that  $\pi^{\text{fixed,opt}}$  is known *a priori* as is  $\text{Var}\{Y^*(\pi^{\text{fixed,opt}})\}$ ; this reflects the common practice of comparing a fixed regime against standard of care or another fixed regime. All results are based on 500 Monte Carlo replications.

Table 2.1 displays the average estimated sample size and its operating characteristics across the four settings of the proposed generative model and three approaches to selecting  $\sigma^*$  when sizing for just condition (POW). Table 2.2 displays the same results when sizing for condition (OPT). A table of results when sizing for both jointly is contained in Appendix A. In the case where  $\Delta = 0$ , the optimal treatment regime provides no benefit over the optimal fixed regime, thus, this setting reflects a worst-case in terms of the conservatism of sizing for (POW) and (OPT) rather than simply sizing to identify the optimal fixed regime; in this case, the proposed sample size procedure attains nominal levels for (POW) and (OPT) at the cost of an inflated sample size. However, as  $\Delta$  increases, so that the benefits of personalizing treatment relative to a fixed regime also increase, it can be seen that the sample size required for (POW) and (OPT) can (perhaps surprisingly) be considerably smaller than required for identifying an optimal embedded regime provided that one has a high-quality estimate for  $\sigma^*$  either through elicitation or a pilot study. Table 2.2 has power 1.0 for all cases considered which is a consequence of using an upper bound on the difference between the value of the estimated regime and the optimal regime. One could potentially explore data-adaptive adjustments, e.g., the double bootstrap, to reduce this excess power.

We also examined the performance of the normality-based sample size when the postulated modeling assumptions are violated. For these simulations, we used the following

**Table 2.1** Estimated power (POW) under a correctly specified generative model using the normality-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care.

$\Delta$	Method for $\sigma^*$	$n_0$	(POW)	(OPT)	$\hat{n}^{\text{fixed}}$	$\mathbb{E}\hat{n}$	Med( $\hat{n}$ )	SD $\hat{n}$
0	known	50	90.0	-	74	130	—	—
0	pilot study	50	87.2	-	74	100.57	99	26.74
0	surrogate	50	98.6	-	74	186.44	185	44.65
0.5	known	50	100	-	111	124	—	—
0.5	pilot study	50	99.8	-	111	132.14	131	29.39
0.5	surrogate	50	100	-	111	164.06	159	46.60
1	known	50	100	-	151	134	—	—
1	pilot study	50	100	-	151	160.43	158	33.11
1	surrogate	50	100	-	151	203.78	199.5	62.16
2	known	50	100	-	251	165	—	—
2	pilot study	50	100	-	251	235.47	233	46.14
2	surrogate	50	100	-	251	280.04	275	82.16

generative model:

$$\begin{aligned}
 X_1 &\sim N(0, 1), & A_1, A_2 &\sim_{i.i.d.} \text{Unif}\{-1, 1\}, \\
 \phi &\sim t_3, & X_2 &= \mu_0^* + \mu_1^* X_1 + \mu_2^* A_1 + \mu_3^* A_1 X_1 + \mu_4^* X_1^2 + \phi, \\
 \mathbf{H}_{2,0} &= (1, X_1, A_1, X_1 A_1, X_2), & \mathbf{H}_{2,1} &= (1, A_1, X_2), \\
 v &\sim N(0, 1), & Y &= \mathbf{H}_{2,0}^T \beta_{2,0}^* + A_2 \mathbf{H}_{2,1}^T \beta_{2,1}^* + v.
 \end{aligned}$$

As previously, we let  $V(\pi^{\text{opt}}) - \Delta\eta = V(\pi^{\text{fixed, opt}}) = B_0 + \eta$  and consider  $\Delta \in \{0, 0.5, 1, 2\}$ . We set  $\mu^* = (1, 0.5, 0.5, 0.1, 1)$  and choose  $\beta^*$  as follows:

1.  $\Delta = 0$ ,  $\beta_{2,0}^* = (1, 0.5, 0.5, 0.5, 1.5)$ ,  $\beta_{2,1}^* = (-1, -1, 0)$ ;
2.  $\Delta = 0.5$ ,  $\beta_{2,0}^* = (1, 0.5, 0.5, 0.5, 1.5)$ ,  $\beta_{2,1}^* = (-1, -1, 0.55)$ ;
3.  $\Delta = 1$ ,  $\beta_{2,0}^* = (1, 0.5, 0.5, 1, 1.5)$ ,  $\beta_{2,1}^* = (-1, -1, 0.65)$ ;
4.  $\Delta = 2$ ,  $\beta_{2,0}^* = (1, 0.5, 0.5, 2.3, 1.5)$ ,  $\beta_{2,1}^* = (-1, -1, 0.71)$ .

**Table 2.2** Estimated concentration (OPT) under a correctly specified generative model using the normality-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care.

$\Delta$	Method for $\sigma^*$	$n_0$	(POW)	(OPT)	$\hat{n}^{\text{fixed}}$	$\mathbb{E}\hat{n}$	Med( $\hat{n}$ )	SD $\hat{n}$
0	known	50	-	100	74	277	—	—
0	pilot study	50	-	100	74	204.51	199.5	60.60
0	surrogate	50	-	100	74	394.63	390	90.04
0.5	known	50	-	100	111	263	—	—
0.5	pilot study	50	-	100	111	269.91	264	62.58
0.5	surrogate	50	-	100	111	353.35	343.5	98.09
1	known	50	-	100	151	285	—	—
1	pilot study	50	-	100	151	335.69	332.5	69.93
1	surrogate	50	-	100	151	431.93	420	130.70
2	known	50	-	100	251	352	—	—
2	pilot study	50	-	100	251	495.52	491	92.61
2	surrogate	50	-	100	251	605.82	578.5	176.10

The average sample size and operating characteristics of the normality-based sample size procedure when sizing for condition (POW) are displayed in Table 2.3 whereas the average sample size and characteristics for condition (OPT) are in Table 2.4. Results for sizing to guarantee both conditions jointly are contained in Appendix A. The proposed method continues to attain nominal levels for  $\Delta \geq 1$ , but is underpowered when using a pilot study to estimate  $\sigma^*$  and there is little or no benefit to the optimal regime over the optimal embedded regime.

We also applied the projection-based sample size procedure to the two classes of generative models described above. Each Monte Carlo replication consists of the following steps. We first generate a pilot study of size  $n_0$ . The bootstrap method described in Section 3.2 is used to calculate the minimum sample size  $\hat{n}(\mathcal{D}_{n_0})$  to achieve power  $(1 - \gamma) \times 100\%$  using 100 bootstrap replications across a grid of potential sample sizes and then using nonlinear least squares to regress the estimated power on the sample sizes. Let  $\gamma = 0.1$ ,  $\vartheta_1 = 0.01$ , and  $\vartheta_2 = 0.04$ . Which corresponds to 90% power for a test with 5% significance level based on

**Table 2.3** Estimated power (POW) under a model which violated the normality assumptions using the normality-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care.

$\Delta$	Method for $\sigma^*$	$n_0$	(POW)	(OPT)	$\hat{n}^{\text{fixed}}$	$\mathbb{E}\hat{n}$	Med( $\hat{n}$ )	SD $\hat{n}$
0	known	50	99.2	-	70	275	—	—
0	pilot study	50	72.4	-	70	77.96	59	74.25
0	surrogate	50	92.8	-	70	158.91	148	65.31
0.5	known	50	100	-	35	228	—	—
0.5	pilot study	50	80.0	-	35	65.14	46	79.77
0.5	surrogate	50	98.2	-	35	151.95	131	78.04
1	known	50	100	-	47	296	—	—
1	pilot study	50	89.0	-	47	77.77	56.5	72.61
1	surrogate	50	99.6	-	47	194.20	166	281.04
2	known	50	100	-	103	407	—	—
2	pilot study	50	99.4	-	103	139.88	118.5	90.88
2	surrogate	50	100	-	103	252.01	227.5	110.40

a confidence interval for  $V(\pi^{\text{opt}})$  that is being constructed using a 99% confidence set for  $(\mu_1^*, \mu_2^*)$  and a 96% interval for  $V(\mu_1, \mu_2)$  for each fixed value of  $(\mu_1, \mu_2)$ . Table 2.5 displays the results under the normal generative model when sizing for condition (POW); in some cases the pilot study shows no benefit to tailoring treatment, i.e.,  $\widehat{V}_{n_0} \leq B_0$ , in which case  $\widehat{n}(\mathcal{D}_0) = +\infty$ . Table 2.6 displays the results when sizing for condition (OPT) under the normal generative model.

Table 2.7 show the results for when the projection-based method when sizing for condition (POW) is applied to the data generating model for which the normality assumptions do not hold. The results of sizing for condition (OPT) for the model which the normality assumptions do not hold is contained in Table 2.8. It can be seen that for  $\Delta \geq 0.50$  the proposed procedure attains nominal power for (POW) for both generative models.



**Table 2.4** Estimated concentration (OPT) under a model which violated the normality assumptions using the normality-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care.

$\Delta$	Method for $\sigma^*$	$n_0$	(POW)	(OPT)	$\hat{n}^{\text{fixed}}$	$\mathbb{E}\hat{n}$	Med( $\hat{n}$ )	SD $\hat{n}$
0	known	50	-	100	70	586	—	—
0	pilot study	50	-	89.4	70	161.96	122	147.53
0	surrogate	50	-	100	70	339.32	317.5	134.60
0.5	known	50	-	100	35	485	—	—
0.5	pilot study	50	-	85.6	35	119.20	87	118.01
0.5	surrogate	50	-	100	35	322.16	285.5	190.63
1	known	50	-	100	47	630	—	—
1	pilot study	50	-	0.88	47	186.41	142.5	159.20
1	surrogate	50	-	0.99	47	390.41	344.5	191.12
2	known	50	-	100	103	866	—	—
2	pilot study	50	-	99.6	103	296.34	241	211.56
2	surrogate	50	-	100	103	537.83	490.5	277.45

## 2.4 Discussion

We proposed two sample size procedures for two-stage SMARTs when the objective is estimation and evaluation of an optimal dynamic treatment regime. These procedures can be used to design SMARTs or conduct power analyses for observational studies. Furthermore, a comparison of the sample size required for construction of a high-quality estimator of an optimal treatment regime with the sample size required for comparison of fixed treatment sequences (or another simple comparison commonly used to size a SMART) can generate new insights into the cost of precision medicine in a given problem domain.

The proposed procedures were developed under two extremes in terms of the structure imposed on the underlying generative model. At one extreme, we assumed correctly specified parametric models for several functionals of the generative model including the optimal regime; and, at the other extreme, we only imposed moment conditions on a possibly misspecified analysis model. There is large class of intermediate models that could be

**Table 2.5** Estimated power (POW) under a model for which the normality assumptions hold using the projection-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care.

$\Delta$	$n_0$	(POW)	(OPT)	$\hat{n}^{\text{fixed}}$	$\mathbb{E}\hat{n}(\mathcal{D}_{n_0})$	$\text{Med}\{\hat{n}(\mathcal{D}_{n_0})\}$	$\text{SD}\hat{n}(\mathcal{D}_{n_0})$	$P\{\hat{n}(\mathcal{D}_0) = \infty\}$
0	50	85.17	-	74	381.94	342.5	112.46	0.42
0.5	50	99.13	-	111	316.01	292.5	83.89	0.21
1	50	99.73	-	151	312.39	299	71.96	0.17
2	50	100	-	251	364.95	361.5	62.78	0.06

**Table 2.6** Estimated concentration (OPT) under a model for which the normality assumptions hold using the projection-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care.

$\Delta$	$n_0$	(POW)	(OPT)	$\hat{n}^{\text{fixed}}$	$\mathbb{E}\hat{n}(\mathcal{D}_{n_0})$	$\text{Med}\{\hat{n}(\mathcal{D}_{n_0})\}$	$\text{SD}\hat{n}(\mathcal{D}_{n_0})$
0	50	-	100	74	114.73	113	9.72
0.5	50	-	100	111	107	106	9.59
1	50	-	99.8	151	112.22	111	10.69
2	50	-	99.4	251	124.3	123	13.22

constructed from these two base approaches. Furthermore, while the proposed approaches focused on regression-based estimators they can be extended to classification-based or direct-search estimators (Orellana et al., 2010; Zhang et al., 2012; Zhao et al., 2012; Zhang et al., 2012, 2013; Zhao et al., 2015; Zhou et al., 2017; Zhao et al., 2015; Laber and Zhao, 2015) which are becoming increasingly popular; we leave the details of this extension to future work.

**Table 2.7** Estimated power (POW) under a model which violated the normality assumptions using the projection-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care.

$\Delta$	$n_0$	(POW)	(OPT)	$\hat{n}^{\text{fixed}}$	$\mathbb{E}\hat{n}(\mathcal{D}_{n_0})$	$\text{Med}\{\hat{n}(\mathcal{D}_{n_0})\}$	$\text{SD}\hat{n}(\mathcal{D}_{n_0})$	$P\{\hat{n}(\mathcal{D}_0) = \infty\}$
0	50	84.28	-	70	482.18	401	250.04	0.34
0.5	50	92.96	-	35	527.52	491	208.04	0.41
1	50	98.55	-	47	556.37	507.5	227.93	0.28
2	50	99.49	-	103	594.68	536	281.35	0.14

**Table 2.8** Estimated concentration (OPT) under a model which violated the normality assumptions using the projection-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care.

$\Delta$	$n_0$	(POW)	(OPT)	$\hat{n}^{\text{fixed}}$	$\mathbb{E}\hat{n}(\mathcal{D}_{n_0})$	$\text{Med}\{\hat{n}(\mathcal{D}_{n_0})\}$	$\text{SD}\hat{n}(\mathcal{D}_{n_0})$
0	50	-	96.6	70	84.13	83	11.66
0.5	50	-	93.6	35	88.93	88	12.08
1	50	-	95.6	47	90.29	88	14.16
2	50	-	97.6	103	90.53	88	22.73

## CHAPTER

### 3

# Q-LEARNING FOR SURVIVAL ANALYSIS

A dynamic treatment regime is a sequential list of decision rules that map up-to-date patient information to a recommended treatment. An optimal dynamic treatment regime is one that maximizes the mean outcome of interest for patients in the population of interest when assigned treatments suggested by following the regime (Murphy, 2003b; Robins, 2004b). One common approach for the estimation of an optimal regime is Q-learning, a regression-based approximate dynamic programming method (Murphy, 2005b; Schulte et al., 2014). This has been applied in many different contexts such as ADHD (Nahum-Shani et al., 2012), depression (Song et al., 2015a), schizophrenia (Shortreed et al., 2010), smoking cessation (Chakraborty et al., 2010), and dosing for chronic pain (Laber et al., 2018).

In the context of chronic disease research such as cancer, HIV/AIDS, and cardiovascular disease, the outcome of interest is commonly a time to an event. Examples of potential

time-to-events of interest include overall survival time, disease-free survival time, and progression-free survival time. When dealing with time-to-event outcomes we frequently have data in which the event is not observed for some patients and is therefore considered to be censored. This could occur due to administrative censoring in which the event has not occurred by the time the study concludes. Some patients could also drop out of the study or there could be a competing risk that prevents the event of interest from being observed causing censoring. In the context of multiple decision regimes, a patient may experience the event of interest or be censored before reaching all of the possible decision points being considered. These issues will cause bias in a naive application of Q-learning to time-to-event data.

There has been some previous work on the estimation of dynamic treatment regimes for time-to-event outcomes. Bai et al. (2017) and Zhao et al. (2015) proposed using a value search estimator from a classification perspective for a one-stage regime. Hager et al. (2018) use value search for an augmented inverse probability weighted estimator from a classification perspective for estimation of an optimal two-stage regime. Xu et al. (2016) use a Bayesian framework for estimating the value of multi-stage dynamic treatment regimes based on disease status. They do not consider the estimation of an optimal regime though. Jiang et al. (2017) propose a value search estimator using kernel smoothing for multiple decision points. They assume that the timing of the decision points are prespecified and the same for all patients.

Goldberg and Kosorok (2012) proposed adjusting Q-learning by using inverse-probability-of-censoring weighting to correct for the bias. They assumed that censoring was independent of patient history though. Zhao et al. (2011) also use a Q-learning framework with modified support vector regression to adjust for bias due to censoring which is again assumed to be independent of patient history. In this paper we will propose a general framework for correcting for bias in the estimation of an optimal dynamic treatment regime using Q-learning via inverse weighting. We will allow censoring to occur at any point after

the initial treatment and depend on the patient history. The timing of the decision points will also be allowed to vary between patients.

In Section 3.1, we will provide setup and notation. In Section 3.2, we present the inverse weighting Q-learning algorithm. In Section 3.3, we will evaluate the performance of our algorithm with a series of simulation experiments. In Section 3.4, we demonstrate the use of our proposed algorithm on the STAR\*D trial which was a study of the effectiveness of depression treatments for patients with a major depressive disorder. In Section 3.5, a discussion of the methodology is provided.

### 3.1 Setup and notation

We will consider the situation where there is a some fixed maximum number,  $K$ , of decision points that any patient could reach. We will assume we have data,  $\mathcal{D}_n = \{(\kappa_i, \mathcal{T}_{1,i}, \mathbf{X}_{1,i}, A_{1,i}, \dots, \mathcal{T}_{\kappa_i,i}, \mathbf{X}_{\kappa_i,i}, A_{\kappa_i,i}, U_i, \Delta_i)\}_{i=1}^n$  which comprises *i.i.d* trajectories of the form  $(\kappa, \mathcal{T}_1, \mathbf{X}_1, A_1, \dots, \mathcal{T}_\kappa, \mathbf{X}_\kappa, A_\kappa, U, \Delta)$  such that  $\kappa$  denotes the number of observed decision points for the patient,  $\mathcal{T}_k$  denotes the time of decision point  $k$ ,  $\mathbf{X}_1 \in \mathbb{R}^{p_1}$  denotes baseline patient information,  $\mathbf{X}_k \in \mathbb{R}^{p_k}$  denotes covariates recorded during the course of the  $(k-1)^{\text{st}}$  treatment for  $k > 1$ ,  $A_k \in \mathcal{A}_k$  denotes the treatment assigned at decision point  $k$ ,  $U$  denotes the time to event or censoring, and  $\Delta$  is an indicator for whether the event was observed. Let  $T$  be the potential time-to-event and  $C$  be the potential censoring time of a patient. Then  $U = \min(T, C)$  and  $\Delta = I(T \leq C)$ .

Define  $\mathbf{H}_k = \{(\mathcal{T}_j, \mathbf{X}_j, A_j)I(\kappa \geq j) \text{ for } j = 1, \dots, k-1, (\mathcal{T}_k, \mathbf{X}_k)I(\kappa \geq k), \kappa I(\kappa < k), (U, \Delta)I(\kappa < k)\}$  for  $k = 1, \dots, K$  to be the history of a patient available to a decision maker at decision point  $k$ . Define  $\mathbf{H}(u)$  to be the patient history up to time  $u$  such that  $\mathbf{H}(u) = \{I(\kappa \geq j, u \geq \mathcal{T}_j), (\mathcal{T}_j, \mathbf{X}_j, A_j)I(\kappa \geq j, u \geq \mathcal{T}_j) \text{ for } j = 1, \dots, K, I(u > U), (U, \Delta)I(u > U)\}$ . A treatment regime is given by a sequence of  $K$  functions  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$  where  $\pi_k : \text{dom } \mathbf{H}_k \rightarrow \text{dom } A_k$  so that a decision maker assigning treatments by following  $\boldsymbol{\pi}$  would assign treatment  $\pi_k(\mathbf{h}_k)$

to a patient with history  $\mathbf{H}_k = \mathbf{h}_k$  at time  $k = 1, \dots, K$ . We define a treatment regime as optimal if it maximizes the mean outcome when applied to the population of interest. We will use potential outcomes to formally define the value of a regime.

Let  $\bar{\mathbf{A}}_k = (A_1, \dots, A_k)$  denote a sequence of the first  $k$  treatments. Define  $\mathbf{H}_k^*(\bar{\mathbf{a}}_{k-1})$  to be the potential history at the  $k^{\text{th}}$  decision point if assigned treatments  $\bar{\mathbf{A}}_{k-1} = \bar{\mathbf{a}}_{k-1}$  previously. Let  $\mathcal{T}_k^*(\bar{\mathbf{a}}_{k-1})$  be the potential time of the  $k^{\text{th}}$  decision point. We will let the censoring be considered an external process and will thus not have a potential outcome. Define  $\mathcal{X}^*(\bar{\mathbf{a}}_K)$  to be the potential number of decision points a patient would reach if they could not be censored and  $\mathcal{X}_k^*(\bar{\mathbf{a}}_k)$  will denote an indicator for whether a patient reaches the  $k^{\text{th}}$  decision point without experiencing the event. Let  $T^*(\bar{\mathbf{a}}_{\mathcal{X}^*(\bar{\mathbf{a}}_K)})$  be the potential time-to-event. Let  $T^*(\pi)$  be the potential time-to-event for a patient if they were to receive treatments dictated by regime  $\pi$  at each decision point.

Let  $S_k^*(\bar{\mathbf{a}}_{k-1})$  denote the potential time of the  $k^{\text{th}}$  decision point or the potential event time if it occurs after the  $(k-1)^{\text{st}}$  and before the  $k^{\text{th}}$  decision point. If the potential event occurs before the  $(k-1)^{\text{st}}$  decision point then let  $S_k^*(\bar{\mathbf{a}}_{k-1}) = 0$ . Let  $\Upsilon_k^*(\bar{\mathbf{a}}_{k-1})$  denote whether the potential event occurs between the  $(k-1)^{\text{st}}$  and  $k^{\text{th}}$  decision point such that

$$\begin{aligned} S_k^*(\bar{\mathbf{a}}_{k-1}) &= I\{\mathcal{X}_{k-1}^*(\bar{\mathbf{a}}_{k-2}) = 1, \mathcal{X}_k^*(\bar{\mathbf{a}}_{k-1}) = 0\} T^*(\bar{\mathbf{a}}_{k-1}) + I\{\mathcal{X}_k^*(\bar{\mathbf{a}}_{k-1}) = 1\} \mathcal{T}_k^*(\bar{\mathbf{a}}_{k-1}) \\ \Upsilon_k^*(\bar{\mathbf{a}}_{k-1}) &= 1 \quad \text{if } \mathcal{X}_{k-1}^*(\bar{\mathbf{a}}_{k-2}) = 1 \text{ and } \mathcal{X}_k^*(\bar{\mathbf{a}}_{k-1}) = 0 \\ &= 0 \quad \text{if } \mathcal{X}_{k-1}^*(\bar{\mathbf{a}}_{k-2}) = 0 \\ &= -1 \quad \text{if } \mathcal{X}_k^*(\bar{\mathbf{a}}_{k-1}) = 1. \end{aligned}$$

Define the set of all possible potential outcomes as  $W^* = \{\mathcal{X}^*(\bar{\mathbf{a}}), \mathcal{T}_2^*(a_1), X_2^*(a_1), \mathcal{T}_3^*(\bar{\mathbf{a}}_2), X_3^*(\bar{\mathbf{a}}_2), \dots, \mathcal{T}_{\mathcal{X}^*(\bar{\mathbf{a}})}^*(\bar{\mathbf{a}}_{\mathcal{X}^*(\bar{\mathbf{a}})-1}), X_{\mathcal{X}^*(\bar{\mathbf{a}})}^*(\bar{\mathbf{a}}_{\mathcal{X}^*(\bar{\mathbf{a}})-1}), T^*(\bar{\mathbf{a}}_{\mathcal{X}^*(\bar{\mathbf{a}})})\}$ , for all  $\bar{\mathbf{a}} \in \mathcal{A}$ . To be able to construct an estimator for the value of a regime from the observed data we will make the following causal assumptions (C1 - C3). (C1) Consistency states that the observed data is equivalent to the potential outcomes for the treatments that were actually assigned.

Therefore

$$\begin{aligned}\kappa &= \mathcal{X}^*(\bar{\mathbf{A}}) \text{ and } U = T^*(\bar{\mathbf{A}}_\kappa) \text{ when } \Delta = 1 \\ \kappa &= \max_j \{j : \mathcal{X}_j^*(\bar{\mathbf{A}}_{j-1}) = 1 \text{ and } U > \mathcal{T}_j^*(\bar{\mathbf{A}}_{j-1})\} \text{ when } \Delta = 0 \\ X_k &= X_k^*(\bar{\mathbf{A}}_{k-1}), \mathcal{T}_k = \mathcal{T}_k^*(\bar{\mathbf{A}}_{k-1}) \text{ for } k = 1, \dots, \kappa.\end{aligned}$$

(C2) Sequential ignorability states that the treatment assigned at each decision point is independent of the set of all potential outcomes conditional on the patient history. This can be written as  $W^* \perp A_k | H_k, \kappa \geq k$  for  $k = 1, \dots, K$ . (C3) Positivity states that there is a non-zero probability of each feasible treatment being assigned for a given patient history and a non-zero probability of not being censored before any decision point. Therefore

$$P(A_k = a_k | H_k = h_k, \kappa \geq k) > 0$$

for all feasible  $a_k \in \mathcal{A}_k$  and for all  $\mathbf{h}_k$  that satisfies  $P(\mathbf{H}_k = \mathbf{h}_k, \kappa \geq k) > 0$ . The non-zero probability of not being censored can be expressed formally as

$$P\{U \geq u | S_{k+1}^*(\bar{\mathbf{a}}_k) = u, H_k = h_k, A_k = a_k, \kappa \geq k\} > 0$$

for feasible  $a_1 \in \mathcal{A}_1$  and history  $\mathbf{h}_1$  such that  $P\{S_{k+1}^*(\bar{\mathbf{a}}_k) = u, \mathbf{H}_k = \mathbf{h}_k, A_k = a_k, \kappa \geq k\} > 0$ .

Define

$$\lambda_C\{u | \mathbf{H}(u)\} = \lim_{du \rightarrow 0} du^{-1} P\{u \leq U < u + du, \Delta = 0 | U \geq u, \mathbf{H}(u)\}$$

to be the cause-specific hazard function for being censored at time  $u$ . We will assume that the censoring is noninformative in that it is independent of the set of potential outcomes conditional on the history. Therefore  $W^* \perp C | \mathbf{H}_k, A_k, \kappa \geq k$ , for  $k = 1, \dots, K$ . This implies that  $\lambda_C\{u | \mathbf{H}(u), W^*\} = \lambda_C\{u | \mathbf{H}(u)\}$ .

Define  $V(\boldsymbol{\pi}) = \mathbb{E}[f\{T^*(\boldsymbol{\pi})\}]$  to be the value of a fixed regime  $\boldsymbol{\pi}$  where  $f(\cdot)$  is a known,



monotone function. The choice of  $f(\cdot)$  depends on the scientific question of interest. Some common quantities of interest for time-to-events are the survival time where  $f(t) = t$ , the probability of survival until time  $r$  where  $f(t) = I(t \geq r)$ , and the restricted lifetime where  $f(t) = \min(t, L)$ . Our goal is to estimate a regime  $\pi^{\text{opt}}$  such that  $V(\pi^{\text{opt}}) \geq V(\pi)$  for all  $\pi$ .

## 3.2 Method

Q-learning is a regression-based approximate dynamic programming method for estimating an optimal treatment regime. Suppose we have data on patients receiving  $K$  treatments without censoring and a final outcome  $Y \in \mathbb{R}$  coded such that higher values are better. In general Q-learning works by defining the Q-functions

$$Q_K(\mathbf{h}_K, a_K) = \mathbb{E}(Y | \mathbf{H}_K = \mathbf{h}_K, A_K = a_K)$$

$$Q_k(\mathbf{h}_k, a_k) = \mathbb{E} \left\{ \max_{a_{k+1}} Q_{k+1}(\mathbf{H}_{k+1}, a_{k+1}) | \mathbf{H}_k = \mathbf{h}_k, A_k = a_k \right\} \text{ for } k = 1, \dots, K-1.$$

It can then be shown via dynamic programming that  $\pi_k^{\text{opt}}(\mathbf{h}_k) = \arg \max_{a_k} Q_k(\mathbf{h}_k, a_k)$  (Bellman, 1957). Models  $Q_k(\mathbf{h}_k, a_k; \beta_k)$  are then posited for each Q-function and fit in a backwards iterative manner. First estimate  $\hat{\beta}_K$  by regressing  $Y$  on  $\mathbf{H}_K$  and  $A_K$ . For  $k = K-1, \dots, 1$  the Q-function is estimated by regressing  $\max_{a_{k+1}} Q_{k+1}(\mathbf{H}_{k+1}, a_{k+1}; \hat{\beta}_{k+1})$  on  $\mathbf{H}_k$  and  $A_k$ . This then gives a plug-in estimator for the optimal regimes as  $\hat{\pi}_k^{\text{opt}}(\mathbf{h}_k) = \arg \max_{a_k} Q_k(\mathbf{h}_k, a_k; \hat{\beta}_k)$ . Our proposed method for time-to-event outcomes will be based around the same framework with the addition of inverse weighting to correct for bias due to censoring.

For each individual  $i$ , let

$$\begin{aligned}
S_{ki}^\dagger &= I(\kappa_i = k-1)U_i + I(\kappa_i \geq k)\mathcal{T}_{k,i} \\
\Upsilon_{ki}^\dagger &= 1 \quad \text{if } \kappa_i = k-1, \Delta_i = 1 \\
&= 0 \quad \text{if } \kappa_i = k-1, \Delta_i = 0 \\
&= -1 \quad \text{if } \kappa_i \geq k.
\end{aligned}$$

Then  $S_{ki}^\dagger$  is either the time of the  $k^{\text{th}}$  decision point or the time to event/censoring if it occurs before for a patient who has reached the  $(k-1)^{\text{st}}$  decision point and  $\Upsilon_{ki}^\dagger$  is an indicator for whether a patient observes the event, is censored, or reaches the  $k^{\text{th}}$  decision point. Define  $\lambda_C(u|\mathbf{h}_k, a_k)$  to be the hazard function for being censored at time  $u$  given  $\mathbf{H}_k = \mathbf{h}_k$  and  $A_k = a_k$ . We can then define the probability of not being censored until after time  $u$  after reaching decision point  $k$  as

$$\mathcal{K}_k(u|\mathbf{h}_k, a_k) = \exp\left\{-\int_{\tau_k}^u \lambda_C(w|\mathbf{h}_k, a_k)dw\right\}, \quad \tau_k \leq u \leq S_{k+1}^\dagger$$

where  $\mathcal{T}_k = \tau_k$  is contained in  $\mathbf{H}_k = \mathbf{h}_k$  for  $\kappa \geq k$ . Define

$$Q_K(\mathbf{h}_K, a_K) = \mathbb{E}\left\{\frac{\Delta f(U)}{\mathcal{K}_K(U|\mathbf{h}_K, a_K)} \mid \mathbf{H}_K = \mathbf{h}_K, A_K = a_K, \kappa = K\right\}.$$

Then

$$Q_K(\mathbf{h}_K, a_K) = \mathbb{E}[f\{T^*(\bar{a}_K)\} \mid \mathbf{H}_K = \mathbf{h}_K, A_K = a_K, \kappa = K].$$

Define  $\Psi(\mathbf{h}_k) \subseteq \mathcal{A}_k$  to be the feasible treatment options at decision point  $k$  for a patient

with history  $\mathbf{h}_k$ . Therefore  $\pi_K^{\text{opt}}(\mathbf{h}_K) = \arg \max_{a_K \in \Psi(\mathbf{h}_K)} Q_K(\mathbf{h}_K, a_K)$ . Note that

$$\begin{aligned} Q_K(\mathbf{h}_K, a_K) &= \mathbb{E}[f\{T^*(\bar{a}_K)\} | \mathbf{H}_K = \mathbf{h}_K, A_K = a_K, \kappa = K] \\ &= \mathbb{E}[f(\tau_K) + f\{T^*(\bar{a}_K)\} - f(\tau_K) | \mathbf{H}_K = \mathbf{h}_K, A_K = a_K, \kappa = K] \\ &= f(\tau_K) + \mathbb{E}[f\{T^*(\bar{a}_K)\} - f(\tau_K) | \mathbf{H}_K = \mathbf{h}_K, A_K = a_K, \kappa = K]. \end{aligned}$$

Then let

$$Q_K^R(\mathbf{h}_K, a_K) = Q_K(\mathbf{h}_K, a_K) - f(\tau_K)$$

be the expected remaining time to event after reaching decision point  $K$ . Then if we posit a model for  $Q_K^R(\mathbf{h}_K, a_K)$  that is strictly non-negative we will respect that  $f(\cdot)$  is monotone and nondecreasing. Let

$$Q_K^R(\mathbf{h}_K, a_K; \beta_K) = \exp\{g_K(\mathbf{h}_K, a_K; \beta_K)\}.$$

Define

$$\begin{aligned} V_K(\mathbf{h}_K) &= \max_{a_K \in \Psi(\mathbf{h}_K)} Q_K(\mathbf{h}_K, a_K) \\ &= f(\tau_K) + \max_{a_K \in \Psi(\mathbf{h}_K)} Q_K^R(\mathbf{h}_K, a_K). \end{aligned}$$

Therefore  $\pi_K^{\text{opt}}(\mathbf{h}_K) = \arg \max_{a_K \in \Psi(\mathbf{h}_K)} Q_K^R(\mathbf{h}_K, a_K)$  and

$$V_K(\mathbf{h}_K) = \mathbb{E}(f[T^*\{\bar{a}_{K-1}, \pi_K^{\text{opt}}(\mathbf{h}_K)\}] | \mathbf{H}_K = \mathbf{h}_K, \kappa = K).$$

The estimated optimal regime at decision point  $K$  is then given by

$$\hat{\pi}_{Q,K}^{\text{opt}}(\mathbf{h}_K) = \arg \max_{a_K \in \Psi(\mathbf{h}_K)} Q_K^R(\mathbf{h}_K, a_K; \hat{\beta}_K).$$

To be able to estimate  $\beta_K$  though we also need an estimate for the probability of not being censored,  $\mathcal{K}_K(U | \mathbf{h}_K, a_K)$ . Model  $\lambda_C\{u | \mathbf{H}(u)\}$  using a proportional hazards model of

the form

$$\lambda_C\{u|\mathbf{H}(u); \beta_C\} = \lambda_{C0}(u) \exp[g_C\{u, \mathbf{H}(u); \xi_C\}]$$

for an unspecified baseline hazard function  $\lambda_{C0}(\cdot)$  and a function  $g_C\{u, \mathbf{H}(u); \xi_C\}$  that depends on the history and a finite-dimensional parameter  $\xi_C$ . Let  $\beta_C = \{\lambda_{C0}(t), \xi_C^T\}^T$ . For example a posited model for  $g_C\{u, \mathbf{H}(u); \xi_C\}$  could be given by

$$g_C\{u, \mathbf{H}(u); \xi_C\} = \xi_{C,1}^T \mathbf{H}_1 + \xi_{C,2} A_1 + \sum_{k=2}^K I(\kappa \geq k, \mathcal{T}_k \leq u) (\xi_{C,2k-1}^T \mathbf{H}_k + \xi_{C,2k} A_k).$$

We can then fit a model to the data to estimate  $\lambda_C\{u|\mathbf{H}(u); \hat{\beta}_C\}$  and then substitute the fitted model into our expression for  $\mathcal{K}_k(u|\mathbf{h}_k, a_k)$  to get

$$\hat{\mathcal{K}}_k(u|\mathbf{h}_k, a_k) = \exp\left\{-\int_{\tau_k}^u \lambda_C(w|\mathbf{h}_k, a_k; \hat{\beta}_C) dw\right\}, \quad k = 1, \dots, K.$$

We estimate  $\hat{\beta}_K$  by solving the inverse probability of censoring M-estimating equation

$$\sum_{\kappa_i=K} \left[ \frac{\Delta_i}{\hat{\mathcal{K}}_K(U_i|\mathbf{H}_{K_i}, A_{K_i})} \frac{\partial Q_K^R(\mathbf{H}_{K_i}, A_{K_i}; \beta_K)}{\partial \beta_K} \{f(U_i) - f(\mathcal{T}_{K_i}) - Q_K^R(\mathbf{H}_{K_i}, A_{K_i}; \beta_K)\} \right] = 0.$$

Proof that the above is an unbiased estimating equation can be found in Appendix B. Now move to decision  $K - 1$ . Let

$$Q_{K-1}(\mathbf{h}_{K-1}, a_{K-1}) = \mathbb{E} \left\{ \frac{I(\Upsilon_K^\dagger = 1)f(S_K^\dagger) + I(\Upsilon_K^\dagger = -1)V_K(\mathbf{H}_K)}{\mathcal{K}_{K-1}(S_K^\dagger|\mathbf{h}_{K-1}, a_{K-1})} \right. \\ \left. | \mathbf{H}_{K-1} = \mathbf{h}_{K-1}, A_{K-1} = a_{K-1}, \kappa \geq K - 1 \right\}.$$

We then have that

$$Q_{K-1}(\mathbf{h}_{K-1}, a_{K-1}) = \mathbb{E} \left( f \left[ T^* \{ \bar{a}_{K-1}, \pi_K^{\text{opt}}(\mathbf{h}_K) \} \right] \middle| \mathbf{H}_{K-1} = \mathbf{h}_{K-1}, A_{K-1} = a_{K-1}, \kappa \geq K - 1 \right).$$

Form pseudo outcomes for individuals  $i$  such that  $\kappa_i \geq K-1$  and  $\Upsilon_{K-1,i}^\dagger \neq 0$ ,

$$\tilde{V}_{Ki}^R = I(\Upsilon_{Ki}^\dagger = 1)f(S_{Ki}^\dagger) + I(\Upsilon_{Ki}^\dagger = -1)[f(S_{Ki}^\dagger) + \max_{a_K} Q_K^R\{\mathbf{h}_K, a_K; \hat{\beta}_K\}] - f(\mathcal{T}_{K-1,i}).$$

Then estimate  $\hat{\beta}_{K-1}$  by solving the unbiased M-estimating equation

$$\sum_{\kappa_i \geq K-1} \left[ \frac{I(\Upsilon_{K,i}^\dagger \neq 0)}{\hat{\mathcal{K}}_{K-1}(S_{Ki}^\dagger | \mathbf{H}_{K-1,i}, A_{K-1,i})} \frac{\partial Q_{K-1}^R(\mathbf{H}_{K-1,i}, A_{K-1,i}; \beta_{K-1})}{\partial \beta_{K-1}} \right. \\ \left. \{ \tilde{V}_{Ki}^R - Q_{K-1}^R(\mathbf{H}_{K-1,i}, A_{K-1,i}; \beta_{K-1}) \} \right] = 0.$$

The estimated optimal regime at decision point  $K-1$  is then

$$\hat{\pi}_{Q,K-1}^{\text{opt}}(\mathbf{h}_{K-1}) = \arg \max_{a_{K-1} \in \Psi(\mathbf{h}_{K-1})} Q_{K-1}^R(\mathbf{h}_{K-1}, a_{K-1}; \hat{\beta}_{K-1}).$$

We then continue the same procedure for  $k = K-2$  down to  $k = 1$ . In general for  $k = K-1, \dots, 1$ ,

$$Q_k(\mathbf{h}_k, a_k) = \mathbb{E} \left\{ \frac{I(\Upsilon_{k+1,i}^\dagger = 1)f(S_{k+1,i}^\dagger) + I(\Upsilon_{k+1,i}^\dagger = -1)V_{k+1}(\mathbf{H}_{k+1})}{\mathcal{K}_k(S_{k+1,i}^\dagger | \mathbf{h}_k, a_k)} \right. \\ \left. | \mathbf{H}_k = \mathbf{h}_k, A_k = a_k, \kappa \geq k \right\}.$$

We estimate  $\beta_k$  by solving

$$\sum_{\kappa_i \geq k} \left[ \frac{I(\Upsilon_{k+1,i}^\dagger \neq 0)}{\hat{\mathcal{K}}_k(S_{k+1,i}^\dagger | \mathbf{H}_{k,i}, A_{k,i})} \frac{\partial Q_k^R(\mathbf{H}_{k,i}, A_{k,i}; \beta_k)}{\partial \beta_k} \{ \tilde{V}_{k+1,i}^R - Q_k^R(\mathbf{H}_{k,i}, A_{k,i}; \beta_k) \} \right] = 0$$

such that for individuals  $i$  such that  $\kappa_i \geq k$  and  $\Upsilon_{k,i}^\dagger \neq 0$ ,

$$\tilde{V}_{k+1,i}^R = I(\Upsilon_{k+1,i}^\dagger = 1)f(S_{k+1,i}^\dagger) + I(\Upsilon_{k+1,i}^\dagger = -1)[f(S_{k+1,i}^\dagger) + \max_{a_{k+1}} Q_{k+1}^R\{\mathbf{h}_{k+1}, a_{k+1}; \hat{\beta}_{k+1}\}] - f(\mathcal{T}_{k,i}).$$

The estimated optimal regime at decision point  $k$  is then given by

$$\hat{\pi}_{Q,k}^{\text{opt}}(\mathbf{h}_k) = \arg \max_{a_k \in \Psi(\mathbf{h}_k)} Q_k^R(\mathbf{h}_k, a_k; \hat{\beta}_k).$$

Let  $\hat{\pi}_Q^{\text{opt}} = (\hat{\pi}_{Q,1}^{\text{opt}}, \dots, \hat{\pi}_{Q,K}^{\text{opt}})$ . An estimator for the value of the optimal regime is then given by

$$\hat{V}_Q(\pi^{\text{opt}}) = n^{-1} \sum_{i=1}^n \tilde{V}_{1,i}^R.$$

### 3.3 Simulation experiments

We will evaluate the performance of our proposed method using a series of simulation experiments. We will consider two different generative models for the data. The first will be the case in which the survival and censoring time follows a Weibull distribution. The other will be a piecewise constant hazard model for both the survival and censoring time. For each generative model we will compare scenarios with different levels of censoring. We will estimate the optimal treatment regime using three different methods. The first will be the proposed survival Q-learning. We will also compare the results with performing Q-learning by assuming the censoring is independent of the patient history and inverse weighting by using the Kaplan-Meier estimate of the probability of not being censored. The third method will be the naive use of Q-learning for a general outcome using only the patients that were not censored.

We will perform 1000 replications of the simulation study. For each replication an estimated optimal regime will be calculated using each of the three methods and the value of the optimal regime will be estimated. The true value of each estimated regime will then be calculated for each estimated regime using the g-computation algorithm with 1000 Monte-Carlo replications. Performance will be evaluated by examining the mean and standard deviations of the estimated value and true value of the estimated regimes across the 1000 replications. We will let the sample size be equal to 1000 for all of the simulations.

The first generative model we will consider is one in which the survival time, censoring time, and time to the next decision point all have a Weibull distribution. The generative model is as follows:

$$X_1 \sim N(0, 1), \quad A_1 \sim \text{Unif}\{-1, 1\}, \quad \mathcal{T}_1 = 0$$

$$\mathbf{H}_{1,0}^T = \mathbf{H}_{1,1}^T = (1, X_1).$$

Generate

$$C_1 = \text{Weibull}\{\text{shape} = \gamma_{C_1}, \text{scale} = \max(1, \mathbf{H}_{1,0}^T \zeta_{1,0} + A_1 \mathbf{H}_{1,1}^T \zeta_{1,1})\}$$

$$T_1 = \text{Weibull}\{\text{shape} = \gamma_{T_1}, \text{scale} = \max(1, \mathbf{H}_{1,0}^T \xi_{1,0} + A_1 \mathbf{H}_{1,1}^T \xi_{1,1})\}$$

$$\mathcal{T}_2 = \text{Weibull}\{\text{shape} = \gamma_{\mathcal{T}_2}, \text{scale} = \max(1, \mathbf{H}_{1,0}^T \omega_{1,0} + A_1 \mathbf{H}_{1,1}^T \omega_{1,1})\}.$$

If  $C_1 = \min(C_1, T_1, \mathcal{T}_2)$  then the individual is censored at time  $C_1$ . If  $T_1 = \min(C_1, T_1, \mathcal{T}_2)$  the individual has the event observed at time  $T_1$ . If  $\mathcal{T}_2 = \min(C_1, T_1, \mathcal{T}_2)$  then the patient has reached the second decision point. Then

$$X_2 = \mathbf{H}_{1,0}^T \mu_{1,0} + A_1 \mathbf{H}_{1,1}^T \mu_{1,1} + \phi_1, \quad A_2 \sim \text{Unif}\{-1, 1\}$$

$$\mathbf{H}_{2,0}^T = \mathbf{H}_{2,1}^T = (1, X_1, A_1, X_2)$$

and we generate

$$C_2 = \text{Weibull}\{\text{shape} = \gamma_{C_2}, \text{scale} = \max(1, \mathbf{H}_{2,0}^T \zeta_{2,0} + A_2 \mathbf{H}_{2,1}^T \zeta_{2,1})\}$$

$$T_2 = \text{Weibull}\{\text{shape} = \gamma_{T_2}, \text{scale} = \max(1, \mathbf{H}_{2,0}^T \xi_{2,0} + A_2 \mathbf{H}_{2,1}^T \xi_{2,1})\}$$

$$\mathcal{T}_3 = \text{Weibull}\{\text{shape} = \gamma_{\mathcal{T}_3}, \text{scale} = \max(1, \mathbf{H}_{2,0}^T \omega_{2,0} + A_2 \mathbf{H}_{2,1}^T \omega_{2,1})\}.$$

This identical pattern is continued for all 3 stages of treatments. Refer to Appendix B for

the parameter values used. The results are summarized in table 3.1. Note that the true performance of the estimated regimes for all of the methods are very similar to one another. The estimated value of the regime is close to the actual value for our proposed survival Q-learning method while the other competing methods underestimate the true value. As the degree of censoring increases this underestimation becomes even more severe. Our proposed method does lead to an increase in the variance of the estimated value.

**Table 3.1** Estimated value of the estimated optimal regime and the true value of the estimated optimal regime for the case where the time-to-event and censoring follow a Weibull distribution.

Method	Censoring	$\mathbb{E}V(\hat{\pi}^{\text{opt}})$	$\text{SD}V(\hat{\pi}^{\text{opt}})$	$\mathbb{E}\hat{V}(\hat{\pi}^{\text{opt}})$	$\text{SD}\hat{V}(\hat{\pi}^{\text{opt}})$
Survival Q-learn	21.1%	16.426	0.182	16.844	0.862
KM Q-learn	21.1%	16.485	0.168	15.909	0.578
Naive Q-learn	21.1%	16.405	0.163	15.379	0.571
Survival Q-learn	33.7%	16.468	0.192	16.390	1.075
KM Q-learn	33.7%	16.514	0.170	15.231	0.733
Naive Q-learn	33.7%	16.408	0.176	14.446	0.715
Survival Q-learn	46.5%	16.461	0.202	16.032	1.219
KM Q-learn	46.5%	16.501	0.191	14.875	0.854
Naive Q-learn	46.5%	16.367	0.185	13.810	0.849

For the second generative model we will let the survival time, censoring time, and time to the next decision point all have a piecewise constant hazard model. The generative model is then given by:

$$X_1 \sim N(0, 1), \quad A_1 \sim \text{Unif}\{-1, 1\}, \quad \mathcal{T}_1 = 0$$

$$\mathbf{H}_{1,0}^T = \mathbf{H}_{1,1}^T = (1, X_1).$$

Let  $(t_1, t_2, \dots, t_J)$  be time points such that the hazard is constant within each interval  $(t_{j-1}, t_j)$ .



The hazard function for survival time is given by

$$\lambda_{T_1}(u|\mathbf{h}_1, a_1) = \begin{cases} j \max(0.5, \mathbf{h}_{1,0}^T \xi_{1,0} + a_1 \mathbf{h}_{1,1}^T \xi_{1,1})/J & \text{if } t_{j-1} < u \leq t_j \\ \max(0.5, \mathbf{h}_{1,0}^T \xi_{1,0} + a_1 \mathbf{h}_{1,1}^T \xi_{1,1}) & \text{if } u > t_j. \end{cases}$$

Similarly the cause-specific hazard function for the censoring time and time of the next decision point are given by

$$\lambda_{C_1}(u|\mathbf{h}_1, a_1) = \begin{cases} j \max(0.5, \mathbf{h}_{1,0}^T \zeta_{1,0} + a_1 \mathbf{h}_{1,1}^T \zeta_{1,1})/J & \text{if } t_{j-1} < u \leq t_j \\ \max(0.5, \mathbf{h}_{1,0}^T \zeta_{1,0} + a_1 \mathbf{h}_{1,1}^T \zeta_{1,1}) & \text{if } u > t_j. \end{cases}$$

and

$$\lambda_{\mathcal{T}_2}(u|\mathbf{h}_1, a_1) = \begin{cases} j \max(0.5, \mathbf{h}_{1,0}^T \omega_{1,0} + a_1 \mathbf{h}_{1,1}^T \omega_{1,1})/J & \text{if } t_{j-1} < u \leq t_j \\ \max(0.5, \mathbf{h}_{1,0}^T \omega_{1,0} + a_1 \mathbf{h}_{1,1}^T \omega_{1,1}) & \text{if } u > t_j. \end{cases}$$

We generate a time of event, censoring time, and time of the second decision point using each hazard function. As with the previous generative model if  $C_1 = \min(C_1, T_1, \mathcal{T}_2)$  the individual is censored at time  $C_1$ . If  $T_1 = \min(C_1, T_1, \mathcal{T}_2)$  then the individual has the event observed at time  $T_1$  and if  $\mathcal{T}_2 = \min(C_1, T_1, \mathcal{T}_2)$  then the patient has reached the second decision point. If the patient reaches the second stage decision point the intervening information collected and second stage decision are generated as

$$X_2 = \mathbf{H}_{1,0}^T \mu_{1,0} + A_1 \mathbf{H}_{1,1}^T \mu_{1,1} + \phi_1, \quad A_2 \sim \text{Unif}\{-1, 1\}$$

$$\mathbf{H}_{2,0}^T = \mathbf{H}_{2,1}^T = (1, X_1, A_1, X_2).$$

The piecewise hazard functions all have the same form as before and the generative model follows the same pattern through 3 stages of treatments. The parameter values used are contained in Appendix B. The results are summarized in table 3.2. As with the first generative

model, the true performance of the estimated regimes for all of the methods are very similar to one another. For this generative model all of the estimated values are less than the true value of the estimated regime. This is expected as the hazard model is not correctly specified which means our method does not guarantee an unbiased estimate of the true value. In this scenario, our proposed method still is the least biased of all the tested methods. For this generative model there is again a small increase in the variance of the estimated value for our proposed survival Q-learning method.

**Table 3.2** Estimated value of the estimated optimal regime and the true value of the estimated optimal regime for the case where the time-to-event and censoring have a piecewise hazard model.

Method	Censoring	$\mathbb{E}V(\hat{\pi}^{\text{opt}})$	$\text{SD}V(\hat{\pi}^{\text{opt}})$	$\mathbb{E}\hat{V}(\hat{\pi}^{\text{opt}})$	$\text{SD}\hat{V}(\hat{\pi}^{\text{opt}})$
Survival Q-learn	29.4%	12.941	0.126	9.294	0.324
KM Q-learn	29.4%	13.020	0.123	8.555	0.196
Naive Q-learn	29.4%	12.781	0.088	6.988	0.137
Survival Q-learn	40.8%	12.986	0.126	8.849	0.305
KM Q-learn	40.8%	13.125	0.106	8.191	0.190
Naive Q-learn	40.8%	12.745	0.090	6.271	0.129

### 3.4 Analysis of STAR\*D trial

We will also demonstrate usage of our method to the Sequenced Treatment Alternatives to Relieve Depression (STAR\*D) trial (Rush et al., 2003). STAR\*D was a randomized study conducted to compare the effectiveness of different treatment regimes for patients with Major Depressive Disorder (MDD). Effectiveness of treatments was assessed using the *Quick Inventory of Depressive Symptomatology* (QIDS) score which is a measure of the current severity of depression symptoms (Rush et al., 2003). This was used with both clinician and self-reported versions of the QIDS. The scores range from 0 to 27 such that higher values

correspond to worse depression symptoms. A patient with a clinician assessed QIDS score of less than or equal to 5 was considered to be in remission.

The study enrolled 4041 patients, each of which could potentially receive four different stages of treatments each of which could last up to 14 weeks. At any point during the course of a treatment a patient could request to move to the next stage if they could not tolerate the side effects of their current treatment. All of the patients were treated initially with citalopram (CIT) which belongs to the selective serotonin reuptake inhibitor (SSRI) class of antidepressants. If a patient did not enter remission in 14 weeks they progressed to level 2.

In level 2 patients were eligible to receive one of seven different treatment options. We will follow the convention used by Chakraborty and Moodie (2013) to classify the different treatment options into two different categories, SSRI and non-SSRI. In the second stage the potential SSRI options are setraline (SER), CIT+bupropion (BUP), CIT + buspirone (BUS), and CIT + cognitive psychotherapy (CT). The potential non-SSRI options are velafaxine (VEN), BUP, and CT. Patients who were assigned to CIT+CT or CT in level 2 were eligible to move to level 2A in which their level 2 treatment was supplemented with either VEN or BUP. For the analysis we will consider level 2A to be combined with level 2 and any patient who received a SSRI during 2 or 2A will be considered to receive an SSRI in the combined level.

Treatments in level 3 are also classified as either SSRIs or non-SSRIs. The SSRI treatments are an augmentation of any of the level 2 SSRI treatments with lithium (LI) or thyroid hormone (THY). The non-SSRI options are mirtazapine (MIRT), nontriptyline (NTP), or augmenting any of the level 2 non-SSRI treatments with LI or THY. In level 4 the treatment options were tranylcypromine (TCP) or MIRT+VEN. At the beginning of each level, each patient had the option to give a preference to either switch or augment the treatment given in the previous stage. If they gave a preference they were randomized between the options that were consistent with their preference.

For the analysis we will consider level 2 and 2a to be stage 1 for Q-learning since all patients received the same treatment in level 1. Similarly, level 3 and 4 will be stage 2 and 3

respectively. The outcome of interest will be given by the time to remission of depression symptoms which will be denoted by  $T$ . There was a significant amount of censoring due to drop out at each stage of the trial. At each of the stages we will include the QIDS score at the beginning of the stage, the slope of the QIDS score over the previous stage, and the patient preference as covariates in the analysis. This follows the analysis conducted in Pineau et al. (2007) and Chakraborty et al. (2013). Let  $X_{j1}$  be the beginning QIDS score,  $X_{j2}$  be the QIDS slope, and  $X_{j3}$  be the patient preference for stages  $j = 1, 2, 3$ . The patient preference is coded such that  $X_{j3} = 1$  if they prefer to switch from their previous treatment and  $X_{j3} = -1$  otherwise. Let the  $A_j$  denote the treatment for each stage  $j$  such that  $A_j = 1$  if the patient was assigned a SSRI and  $A_j = -1$  if the patient was assigned a non-SSRI.

Only 92 patients reach the third stage and of those only 16 go into remission before the end of the study. Due to the limited sample size we will restrict the covariates in the model to just depend on the features recorded during the course of the second stage treatment as well as the previous treatments assigned. Our model for the Q-function will then be given by

$$Q_3^R(\mathbf{h}_3, \mathbf{a}_3; \beta_3, \psi_3) = \exp\{\beta_{30} + \beta_{31}A_1 + \beta_{32}A_2 + \beta_{33}X_{31} + \beta_{34}X_{32} + \beta_{35}X_{33} + (\psi_{30} + \psi_{31}A_1 + \psi_{32}A_2 + \psi_{33}X_{31} + \psi_{34}X_{32})A_3\}$$

and the cause-specific hazard function for censoring will be modeled using a cox proportional hazards model of the form

$$\lambda_C(u|\mathbf{h}_3, \mathbf{a}_3) = \lambda_{C30}(u) \exp\{\xi_{3,0} + \xi_{3,1}X_{11} + \xi_{3,2}X_{12} + \xi_{3,3}X_{13} + \xi_{3,4}A_1 + \xi_{3,5}X_{21} + \xi_{3,6}X_{22} + \xi_{3,7}X_{23} + \xi_{3,8}A_2 + \xi_{3,9}X_{31} + \xi_{3,10}X_{32} + \xi_{3,11}X_{33} + \xi_{3,12}A_3\}.$$

Since the goal is to minimize the time-to-remission the estimated optimal rule at the third

stage is given

$$\hat{\pi}_3(\mathbf{h}_3) = -\text{sign}(\hat{\psi}_{30} + \hat{\psi}_{31}X_{31} + \hat{\psi}_{32}X_{32} + \hat{\psi}_{33}X_{33}).$$

After fitting both models we then form the pseudo outcomes for the third stage. Since we are trying to minimize the time-to-event, the max operator will be replaced by the min resulting in the pseudo outcome being given by

$$\tilde{V}_{3i}^R = I(\Upsilon_{3i}^\dagger = 1)f(S_{3i}^\dagger) + I(\Upsilon_{3i}^\dagger = -1)[f(S_{3i}^\dagger) + \min_{a_3} Q_3^R\{\mathbf{h}_3, a_3; \hat{\beta}_3, \hat{\psi}_3\}] - f(\mathcal{T}_{2,i}).$$

For the second and first stage Q-functions we will use the following models:

$$Q_2^R(\mathbf{h}_3, a_3; \beta_2, \psi_2) = \exp\{\beta_{20} + \beta_{21}X_{11} + \beta_{22}X_{12} + \beta_{23}X_{13} + \beta_{24}A_1 + \beta_{25}X_{21} + \beta_{26}X_{22} + \beta_{27}X_{23} + (\psi_{20} + \psi_{21}X_{11} + \psi_{22}X_{12} + \psi_{23}X_{13} + \psi_{24}A_1 + \psi_{25}X_{21} + \psi_{26}X_{22} + \psi_{27}X_{23})A_2\}$$

and

$$Q_1^R(\mathbf{h}_3, a_3; \beta_1, \psi_1) = \exp\{\beta_{10} + \beta_{11}X_{11} + \beta_{12}X_{12} + \beta_{13}X_{13} + (\psi_{10} + \psi_{11}X_{11} + \psi_{12}X_{12} + \psi_{13}X_{13})A_1\}.$$

We will model the cause specific hazard function for censoring for  $\mathcal{T}_2 \leq u < \mathcal{T}_3$  and  $u < \mathcal{T}_2$  as

$$\lambda_C(u|\mathbf{h}_2, a_2) = \lambda_{C20}(u) \exp\{\xi_{2,0} + \xi_{2,1}X_{11} + \xi_{2,2}X_{12} + \xi_{2,3}X_{13} + \xi_{2,4}A_1 + \xi_{2,5}X_{21} + \xi_{2,6}X_{22} + \xi_{2,7}X_{23} + \xi_{2,8}A_2\}$$

and

$$\lambda_C(u|\mathbf{h}_1, a_1) = \lambda_{C10}(u) \exp\{\xi_{1,0} + \xi_{1,1}X_{11} + \xi_{1,2}X_{12} + \xi_{1,3}X_{13} + \xi_{1,4}A_1\}.$$

The estimated optimal decision rule at the first two stages are then given by:

$$\hat{\pi}_2(\mathbf{h}_2) = -\text{sign}(\hat{\psi}_{20} + \hat{\psi}_{21}X_{11} + \hat{\psi}_{22}X_{12} + \hat{\psi}_{23}X_{13} + \hat{\psi}_{24}A_1 + \hat{\psi}_{25}X_{21} + \hat{\psi}_{26}X_{22} + \hat{\psi}_{27}X_{23})$$

$$\hat{\pi}_1(\mathbf{h}_1) = -\text{sign}(\hat{\psi}_{10} + \hat{\psi}_{11}X_{11} + \hat{\psi}_{12}X_{12} + \hat{\psi}_{13}X_{13}).$$

The resulting estimates of the regression coefficient are in tables 3.3 - 3.5. The estimated time to remission of depression symptoms of the estimated optimal treatment regime is given by 82.05 days. For the patients in the clinical trial that were not censored, the mean number of days until remission was 136 days. This demonstrates that tailoring treatments for patients can potentially lead to a significant improvement in patient outcomes.

**Table 3.3** Regression Coefficients for Third Stage Q-function for STAR\*D trial

Coefficient	Variable	Estimate
$\beta_{30}$	Intercept	2.56
$\beta_{31}$	Treatment 1	0.28
$\beta_{32}$	Treatment 2	0.24
$\beta_{33}$	QIDS at beginning of stage 3	0.14
$\beta_{34}$	QIDS slope during stage 2	-0.79
$\beta_{35}$	Preference at stage 3	-0.30
$\psi_{30}$	Treatment 3	-0.38
$\psi_{31}$	Treatment 3 $\times$ Treatment 1	0.34
$\psi_{32}$	Treatment 3 $\times$ Treatment 2	-0.05
$\psi_{33}$	Treatment 3 $\times$ QIDS at beginning of stage 3	0.02
$\psi_{34}$	Treatment 3 $\times$ QIDS slope during stage 2	-5.17

### 3.5 Discussion

We proposed an adaptation of Q-learning for the estimation of an optimal dynamic treatment regime when our outcome of interest is a time-to-event. This procedure is a gener-

**Table 3.4** Regression Coefficients for Second Stage Q-function for STAR\*D trial

Coefficient	Variable	Estimate
$\beta_{20}$	Intercept	3.92
$\beta_{21}$	QIDS at beginning of stage 1	0.03
$\beta_{22}$	QIDS slope during level 1	-0.67
$\beta_{23}$	Preference at stage 2	0.04
$\beta_{24}$	Treatment 1	-0.09
$\beta_{25}$	QIDS at beginning of stage 2	0.03
$\beta_{26}$	QIDS slope during stage 1	-0.89
$\beta_{27}$	Preference at stage 2	0.15
$\psi_{20}$	Treatment 2	0.29
$\psi_{21}$	Treatment 2 $\times$ QIDS at beginning of stage 1	-0.00
$\psi_{22}$	Treatment 2 $\times$ QIDS slope during level 1	-1.36
$\psi_{23}$	Treatment 2 $\times$ Preference at stage 2	0.08
$\psi_{24}$	Treatment 2 $\times$ QIDS at beginning of stage 2	-0.00
$\psi_{25}$	Treatment 2 $\times$ QIDS slope during stage 1	-1.04
$\psi_{26}$	Treatment 2 $\times$ Preference at stage 2	0.16

**Table 3.5** Regression Coefficients for First Stage Q-function for STAR\*D trial

Coefficient	Variable	Estimate
$\beta_{10}$	Intercept	3.49
$\beta_{11}$	QIDS at beginning of stage 1	0.07
$\beta_{12}$	QIDS slope during level 1	0.07
$\beta_{13}$	Preference at stage 2	-0.05
$\psi_{10}$	Treatment 1	-0.38
$\psi_{11}$	Treatment 1 $\times$ QIDS at beginning of stage 1	0.02
$\psi_{12}$	Treatment 1 $\times$ QIDS slope during level 1	-0.59
$\psi_{13}$	Treatment 1 $\times$ Preference at stage 2	0.04

alization of the method proposed in Goldberg and Kosorok (2012) in which we no longer assume that the censoring is independent of patient history. Here we only assume we have non informative censoring which is a weaker assumption. This method requires correctly specifying the hazard and analysis models though to get an unbiased estimate of the optimal regime and its value.

## CHAPTER

# 4

## PARENT MESSAGING STRATEGIES FOR STUDENT ABSENTEEISM

Elementary school attendance is very important to children's development. Chronic absenteeism in young students has been found to be highly correlated with significantly lower academic outcomes (Ehrlich et al., 2014; Ginsburg et al., 2014), future academic success (Chang and Romero, 2008), and increased substance abuse (Henry and Thornberry, 2010). These negative effects have also been found to be even more significant for students from lower socioeconomic classes (Ready, 2010). There is also evidence that there are significant spillover effects and students are negatively impacted academically if their classmates are chronically absent (Gottfried, 2019).

Chronic absenteeism is influenced by many different factors. These include barriers to



attendance such as lack of transportation and illness as well as the student or parent having an aversion to school for reasons such as academic struggles or bullying (Henderson et al., 2014). Families also do not always realize the importance of attendance in school for young children (Chang and Romero, 2008). There is a belief that increased parent engagement could help reduce absenteeism (Ginsburg et al., 2014).

## **4.1 SMART design**

This study was conducted by the Institute of Education Sciences along with the American Institutes for Research to test a parent text messaging system to keep parents informed about their child's attendance as well as the importance of attending school for young students. A Sequential Multiple Assignment Randomized Trial (SMART) was used to compare different interventions and examine whether tailoring interventions based on individual family characteristics could lead to significant reductions in chronic absenteeism.

The study was a two-stage SMART with the first stage running during the fall semester and the second stage during the spring. Families were randomized to three different interventions during the first stage. These were Basic Information Messaging - Positive (BIM-pos), Basic Information Messaging - Negative (BIM-neg), and Business As Usual (BAU). Both of the BIM-pos and BIM-neg interventions include a same-day notification to the parent when their child is absent, a same-day message when the student returns to school, and weekly texts on the importance of attendance in school. The difference is for BIM-pos the messages focus on the positive benefits of regular attendance and BIM-neg focuses on the negative consequences of missing school. BAU is the control group and these families receive no additional support.

The first stage begins in October and runs through December. At the end of the fall semester the response status of each student will be assessed. A student will be considered responsive if they were not chronically absent during the fall semester where chronic absen-

teeism is defined as missing more than 10% of the days they were enrolled. In the second stage the students that were considered responsive during the first stage will continue to receive the same intervention. The students that are non-responsive will be randomized to augment their first stage treatment with Student Staff Outreach (SSO) or Goal Commitment Messaging (GCM). Student Staff Outreach involves a staff member reaching out to the parents regularly to help identify with the parent reasons for the chronic absenteeism as well as suggesting methods to help improve attendance. Goal Commitment Messaging works by sending a text message to the parent at the beginning of the week that prompts them to reply and commit to a goal of having perfect attendance for the week. At the end of the week another message is sent that either praises the parent for achieving their goal or encourages them to improve when the goal is not met. These students will continue to receive the same messages from the first-stage intervention they were assigned to in addition to the augmentation interventions during the second stage. All of the students that were assigned to BAU during the first stage will continue with BAU during the second stage no matter what their response status is. The second stage runs from January through the end of the school year in June. The design of the SMART is shown in figure 4.1.

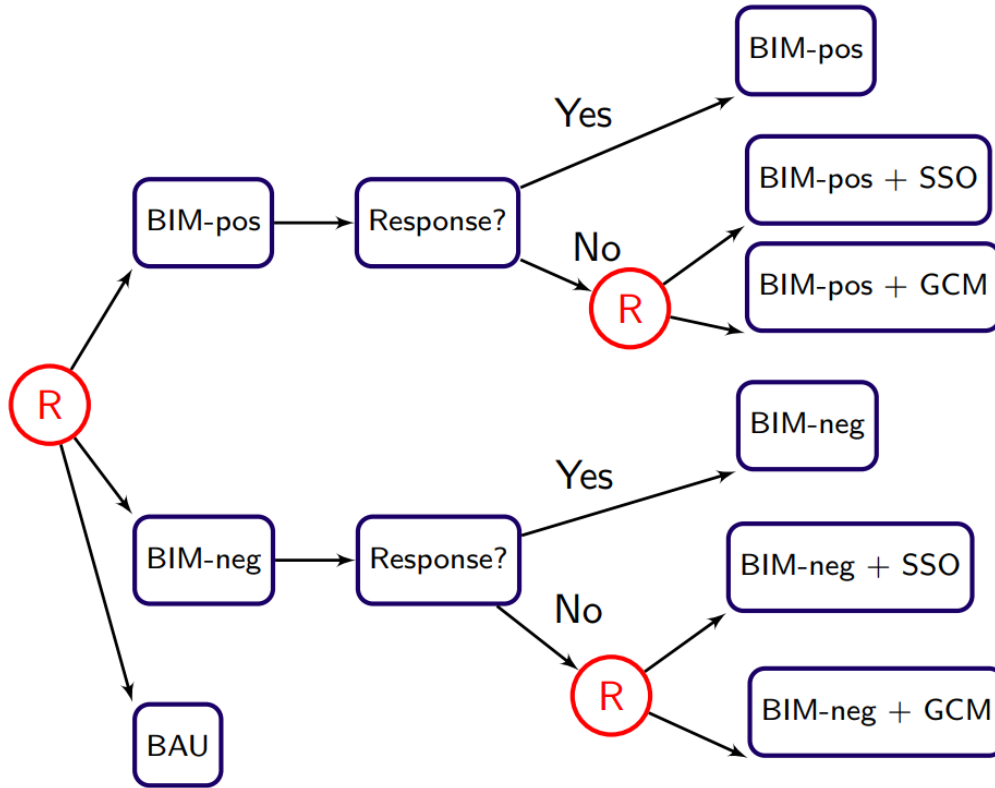
This design leads to five different adaptive interventions embedded within the trial given by:

AI-1: Assign BIM-pos in the fall. If the student is non-responsive in the fall, augment BIM-pos with SSO in the spring and if they are responsive continue with BIM-pos

AI-2: Assign BIM-pos in the fall. If the student is non-responsive in the fall, augment BIM-pos with GCM in the spring and if they are responsive continue with BIM-pos

AI-3: Assign BIM-neg in the fall. If the student is non-responsive in the fall, augment BIM-neg with SSO in the spring and if they are responsive continue with BIM-neg

AI-4: Assign BIM-neg in the fall. If the student is non-responsive in the fall, augment BIM-neg with GCM in the spring and if they are responsive continue with BIM-neg



**Figure 4.1** SMART diagram for the parent messaging study

BAU: Assign BAU to the family for both semesters.

There are several different outcomes of interest that will be used to assess the effectiveness of the different interventions. These outcomes fall into two different categories, attendance outcomes that assess how much the students are attending school and achievement outcomes that assess the academic performance of the students. The attendance outcomes will include an indicator for whether the student was chronically absent and the number of days the student was absent from school. The achievement outcomes are the students' scores from a standardized test in reading and math administered at the end of the school year.

The study was designed to answer four different research questions on the effectiveness of the interventions being tested in the trial. They are as follows:

- RQ1: Are there significant differences in student attendance and achievement scores

between the first stage interventions during the fall semester?

- RQ2: Is there a significant difference in student outcomes between the five different treatment strategies embedded within the trial?
- RQ3: Is there a significant difference in attendance and achievement scores between assigning SSO and GCM for the families that were non-responsive to the first stage intervention?
- RQ4: Is there a benefit to tailoring the interventions based on individual family and school characteristics?

## 4.2 Setup and notation

The data collected from this trial will have the following form,  $\mathcal{D}_n = \{(\mathbf{X}_{ij}^{(1)}, A_{ij}^{(1)}, Y_{ij}^{(1)}, Z_{ij}^{(1)}, R_{ij}, A_{ij}^{(2)}, Y_{ij}^{(2)}, Z_{ij}^{(2)})\}_{i=1, j=1}^{n_j, m}$ , which is given by *i.i.d.* trajectories of the form  $(\mathbf{X}^{(1)}, A^{(1)}, Y^{(1)}, Z^{(1)}, R, A^{(2)}, Y^{(2)}, Z^{(2)})$ . There are  $j = 1, \dots, m$  different schools each of which has  $n_j$  different families participating in the trial. Let  $n = \sum_{j=1}^m n_j$  denote the total number of students in the trial. Let  $\mathbf{X}^{(1)} \in \mathbb{R}^{p_1}$  denote baseline family and school characteristics which in this case will contain the students race, gender, grade, school ID, disability status, English language learner status, sms availability, whether the student is at-risk for chronic absenteeism based on previous years attendance, percentage of days absent during the one month enrollment period before the study starts, and an indicator for whether they were chronically absent during the enrollment period.  $A^{(1)} \in \{-1, 0, 1\}$  denotes the first stage treatment assigned coded such that  $A^{(1)} = 0$  if they were assigned to BAU,  $A^{(1)} = -1$  if assigned to BIM-neg, and  $A^{(1)} = 1$  if assigned to BIM-pos. Let  $R$  denote the responder status of the student at the end of the fall coded such that  $R = 1$  if the student responded and  $R = 0$  if they didn't. For notational convenience let  $R = 1$  for all students assigned to BAU during the first stage.  $A^{(2)} \in \{-1, 1\}$  denotes the second stage treatment assigned coded such that  $A^{(2)} = -1$  if

the student was assigned to GCM and  $A^{(2)} = 1$  if the student was assigned to SSO. Note that  $A^{(2)}$  is nested within  $A^{(1)}$  so  $A^{(2)}$  has different meanings for different values of  $A^{(1)}$ . Let  $Y^{(1)}$  denote the number of days the student was absent in the fall and  $Y^{(2)}$  denote the number of days absent in the spring. Define  $Y = Y^{(1)} + Y^{(2)}$  to be the total number of days the student is absent from school.  $Z^{(1)}$  denotes an indicator for whether the student was chronically absent in the fall and  $Z^{(2)}$  is a indicator for whether they were chronically absent in the spring such that being chronically absent is defined as missing more than 10% of the possible days of school. Define  $Z$  to be an indicator for whether the student was chronically absent for the entire school year.

We will make the same standard causal inference assumptions that were made in chapters 2 and 3. This includes sequential ignorability, positivity, and consistency. Positivity and sequential ignorability hold by design in a SMART. Implicit in the consistency assumption is that there is no interference between students within the trial. That is to say that the outcome of a student is not impacted by the interventions assigned to other students in the trial. This assumption is reasonable in this context though as attendance decisions are determined by the parents for elementary school age children and these decisions are not believed to be impacted by the decisions of other parents in the same school.

## **4.3 Analysis**

### **4.3.1 Synthetic data**

Due to the confidentiality of this study the results of the analysis presented here will use simulated data that has the same form of the actual data. Therefore the conclusions we will make do not reflect the actual results from the study. We will also omit the analysis of the achievement score outcomes, but the methods used are identical to the analysis of the continuous number of days absent outcome. To generate the simulated data we postulate parametric models for the underlying data generating process. We will estimate

the parameters of the model using the actual data and then use the estimated models to generate a new synthetic data set.

For the baseline data we note that the demographic variables such as race, gender, and grade are correlated to one another, but not correlated when conditioned on the school the student attends. Therefore the first step will be to simulate the school for each student which we will assume follows a multinomial distribution with the probability of being from a school given by the sample proportion of students that attended that school. The next step is to simulate the demographic variables which will be given by race, gender, and an indicator for whether the student is in kindergarten. Gender and the kindergarten indicator will be modeled as a binomial random variables with probabilities determined by the within school sample proportions after conditioning on the school the student attends. Race will follow a multinomial distribution again conditional on the school with the probabilities determined by the within school sample proportions.

We will then model the binary student with disability status (SWD) as

$$\text{logit} \{P(\text{SWD}_{ij} = 1)\} = \beta_0 + \beta_1^T X_{ij}^D + \gamma_j$$

such that  $X_{ij}^D$  is a vector of the demographic baseline variables and  $\gamma_j$  is a fixed effect for school  $j$ . Identical models will be used for the students English language learner status and whether they own a SMS capable phone. We will then use the predicted probabilities from the estimated model to simulate data for each of the variables. The at risk at baseline indicator (RSK) will be modeled as

$$\text{logit} \{P(\text{RSK}_{ij} = 1)\} = \beta_0 + \beta_1^T X_{ij}^D + \beta_2 \text{SWD}_{ij} + \beta_3 \text{ELL}_{ij} + \beta_4 \text{SMS}_{ij} + \gamma_j$$

where SWD is the student with disability indicator, ELL is the English language learner status, and SMS is an indicator for SMS phone capabilities. We again will simulate data using the estimated model that's been fit to the original data. The last baseline variable is

the percentage of days the student was absent during September (PCT) before the study begins at the beginning of October. This will be modeled as

$$\text{PCT}_{ij} = \beta_0 + \beta_1^T X_{ij}^D + \beta_2 \text{SWD}_{ij} + \beta_3 \text{ELL}_{ij} + \beta_4 \text{SMS}_{ij} + \beta_5 \text{RSK}_{ij} + \gamma_j + \epsilon$$

where  $\epsilon$  is a random error that is assumed to be normally distributed with mean 0 and variance  $\sigma^2$ .

The next step is to model and simulate the number of days absent in the fall for each student. We will split the data into three separate groups based on the assigned first stage intervention and fit a linear model to the baseline covariates with a fixed effect for school. Therefore for each group we have that

$$Y_{ij}^{(1)} = \beta_0 + \xi^T X_{ij}^{(1)} + \gamma_j + \epsilon.$$

The responder status will then be determined by the percentage of days absent in the fall. The non-responsive students will be randomized between Student Staff Outreach (SSO) and Goal Commitment Messaging (GCM). The number of days absent in the spring will also be modeled using a linear model of the same form as the fall with the number of days absent in the fall added as an additional covariate.

### 4.3.2 Missing data

There is substantial amounts of missing data within this trial due to students dropping out of the study midway through. This was predominantly due to students switching schools during the year causing them to leave the study and therefore have missing outcome data. There is also a small amount of missing baseline data. We will assume that the missing values are missing at random.

For this analysis we will adapt the multiple imputation approach proposed in Shortreed et al. (2014) for handling missing data in SMARTs. Since we are interested in multiple

analyses with this data, by using multiple imputation we will be able to construct the imputed data sets once and use them for all analyses which will ensure consistency across the different analyses. This approach does not use a specification of the full conditional distribution as what is normally used for multiple imputation by chained equations (MICE), but instead uses a forward specification.

The time-ordered nested conditional imputation model used imputes the data sequentially by using only the data that is available at that time point for each imputation. We first impute the missing baseline data using only the baseline data. This is done using a conditional specification of each of the baseline variables. The baseline variables are given by student with disability status, baseline chronic absenteeism, English language learner status, kindergarten status, SMS availability, at risk status, and race. The binary variables are modeled using logistic regression models and the categorical variable, race, is modeled using a polytomous logistic regression model. The values are imputed from simulated draws of the posterior predictive distribution for each variable.

We then impute the number of days absent in the fall semester. This is done by first splitting the semester in half with the first half running from October 1st to November 12th. Any student who was not enrolled for the entire period will be considered missing. The model for the percentage of days absent will contain all of the baseline variables as well as interactions between the baseline variables and the first-stage intervention. The percentage of days absent during the first period is then imputed using predictive mean matching. We then impute the days absent in the second half of the fall in the same manner with the number of days absent in the fall and its interaction with the first stage intervention included in the model.

The response status for the students with missing outcome data is then ascertained based on the imputed values and they are then randomized to student staff outreach or goal commitment messaging if they were non-responsive to the first stage intervention. The number of days absent in the spring is then imputed in an identical manner to the fall



with the semester being split at April 1st.

### 4.3.3 Research Question 1

The first research question focuses on the main effect of the first stage interventions on attendance during the fall semester. We would like to test the three pairwise comparisons of BIM-pos, BIM-neg, and BAU.

To estimate the effect of assigning BIM-pos or BIM-neg we will fit the following model:

$$Y_{ij}^{(1)} = \zeta_0 + \zeta_1 \text{BIM}_{ij} + \xi^T \mathbf{X}_{ij}^{(1)} + \gamma_j + \nu_{ij}$$

where  $Y_{ij}^{(1)}$  is the number of days absent in the fall for student  $i$  in school  $j$ ,  $\text{BIM}_{ij}$  is an indicator for whether they were assigned to either BIM-pos or BIM-neg,  $\mathbf{X}_{ij}^{(1)}$  is a vector of baseline covariates that are mean centered,  $\gamma_j$  is a fixed effect for school  $j$ , and  $\nu_{ij}$  is a random error for student  $i$  in school  $j$  that is assumed to be normally distributed with mean 0 and variance  $\sigma^2$ .

When the outcome of interest is an indicator for whether a student is chronically absent the model will be given by:

$$P(Z_{ij}^{(1)} = 1) = \Phi(\zeta_0 + \zeta_1 \text{BIM}_{ij} + \xi^T \mathbf{X}_{ij}^{(1)} + \gamma_j)$$

where  $\Phi(\cdot)$  is the cumulative distribution function for a standard normal distribution. To test if assigning one of either BIM-pos or BIM-neg has a significant effect for both outcomes, we will test whether  $\zeta_1$  is significantly different from zero using a Wald test.

We will also compare assigning BIM-pos, BIM-neg, and BAU. The model for the number of days absent with each of the individual effects will be given by

$$Y_{ij}^{(1)} = \zeta_0 + \zeta_1 \text{POS}_{ij} + \zeta_2 \text{NEG}_{ij} + \xi^T \mathbf{X}_{ij}^{(1)} + \gamma_j + \nu_{ij}$$

where  $Y_{ij}^{(1)}$  is the number of days absent in the fall for student  $i$  in school  $j$ ,  $POS_{ij}$  is an indicator for whether they were assigned to BIM-pos during the fall,  $NEG_{ij}$  is an indicator for whether they were assigned to BIM-neg,  $\mathbf{X}_{ij}^{(1)}$  is a vector of mean centered baseline covariates,  $\gamma_j$  is a fixed effect for school  $j$ , and  $\nu_{ij}$  is a random error for student  $i$  in school  $j$  that is assumed to be normally distributed with mean 0 and variance  $\sigma^2$ .

Similar to before when the outcome of interest is an indicator for whether a student is chronically absent the model will be given by:

$$P(Z_{ij}^{(1)} = 1) = \Phi(\zeta_0 + \zeta_1 POS_{ij} + \zeta_2 NEG_{ij} + \xi^T \mathbf{X}_{ij}^{(1)} + \gamma_j)$$

where  $\Phi(\cdot)$  is the cumulative distribution function for a standard normal distribution.

To test if assigning BIM-pos has a significant effect for both outcomes, we will test whether  $\zeta_1$  is significantly different from zero using a Wald test. Similarly to test whether BIM-neg has a significant effect we will test if  $\zeta_2$  is significantly different than zero. To test if there is a difference in the effect of BIM-pos and BIM-neg we will test if  $\zeta_1 - \zeta_2$  is significantly different than zero.

Table 4.1 displays the estimated mean number of days absent in the fall semester for those that were assigned to BAU, BIM+, BIM-, or one of either BIM+ or BIM- as well as the results of the test for the pairwise differences between the different intervention options. We use a Bonferroni correction to adjust for multiple comparisons.

Table 4.2 displays the predicted probability of a student being chronically absent in the fall semester if they were assigned to each of the different interventions in the fall. We again use a Bonferroni correction to adjust for multiple comparisons.

We can see that there is a significant reduction in the number of days absent and probability of chronic absenteeism by assigning either BIM-pos or BIM-neg. There was not a significant difference between BIM-pos and BIM-neg for outcomes in the fall semester.

**Table 4.1** Estimated effect of first-stage intervention on the number of days absent in the fall semester

Comparison	BIM+ or BIM-	BIM+	BIM-	BAU	Difference	Effect Size	p-value
BIM+ or BIM- vs. BAU	3.160			3.444	0.284	0.119	<0.001
BIM+ vs. BAU		3.155		3.444	0.289	0.122	<0.001
BIM- vs. BAU			3.166	3.444	0.278	0.117	<0.001
BIM+ vs. BIM-		3.155	3.166		0.011	0.005	0.776

**Table 4.2** Estimated effect of first-stage intervention on chronic absenteeism in the fall semester

Comparison	BIM+ or BIM-	BIM+	BIM-	BAU	Difference	Effect Size	p-value
BIM+ or BIM- vs. BAU	0.176			0.216	0.040	0.101	<0.001
BIM+ vs. BAU		0.173		0.216	0.043	0.108	<0.001
BIM- vs. BAU			0.178	0.216	0.038	0.095	<0.001
BIM+ vs. BIM-		0.173	0.178		0.005	0.013	0.435

### 4.3.4 Research Question 2

The second research question examines the difference between the four adaptive regimes embedded in the trial and BAU. We want to test if any of the adaptive regimes are more effective than BAU as well as if there is a difference between any of the adaptive regimes.

Let  $\pi = (\pi_1, \pi_2)$  denote a treatment regime such that  $\pi_t$  is a function that maps the available information on a student to one of the possible interventions available to that student. There are five different treatment regimes embedded in the design of the trial which we will denote by:

1.  $\pi^{\text{BIM+,SSO}}$  which assigns BIM-pos to all families at the initial stage. During the second stage, those that responded to the first stage intervention continue with BIM-pos and those that are not responsive are given SSO in addition to BIM-pos.

2.  $\pi^{\text{BIM}^+, \text{GCM}}$  which assigns BIM-pos to all families at the initial stage. During the second stage, those that responded to the first stage intervention continue with BIM-pos and those that are not responsive are given GCM in addition to BIM-pos.
3.  $\pi^{\text{BIM}^-, \text{SSO}}$  which assigns BIM-neg to all families at the initial stage. During the second stage, those that responded to the first stage intervention continue with BIM-neg and those that are not responsive are given SSO in addition to BIM-neg.
4.  $\pi^{\text{BIM}^-, \text{GCM}}$  which assigns BIM-neg to all families at the initial stage. During the second stage, those that responded to the first stage intervention continue with BIM-neg and those that are not responsive are given GCM in addition to BIM-neg.
5.  $\pi^{\text{BAU}}$  which assigns BAU at all stages.

Consider the students in the trial that received interventions that are consistent with following the treatment regime  $\pi^{\text{BIM}^+, \text{SSO}}$ . This includes all the students who received BIM-pos in the first stage that were responsive at the end of the fall and therefore continued with BIM-pos as well as those that received BIM-pos in the fall, were non-responsive and then re-randomized to receive SSO in the spring. The probability of a non-responsive student receiving SSO is equal to 0.5. This causes those students who were considered non-responsive to be underrepresented in the data. We will correct for this by inverse weighting by the probability of receiving the intervention that they were assigned. Let  $W_{ij}$  denote the weight for student  $i$  in school  $j$ . For the students that were responsive to BIM-pos or BIM-neg and those that were assigned to BAU will have  $W_{ij} = 3$  since the probability of receiving the initial treatment that they did is equal to  $1/3$ . For the students that received BIM-pos or BIM-neg and were non-responsive will have  $W_{ij} = 6$ . This is because the probability of receiving their first stage intervention is  $1/3$  as before and the probability of receiving their second stage intervention is  $1/2$ .

If we wanted to estimate the mean response of a single adaptive intervention we could limit the data set to only those students that were consistent with that adaptive intervention

and fit a model using weighted least squares. If a student was assigned to BIM-pos in the first stage and was responsive at the end of the fall, they are consistent with following the treatment regimes  $\pi^{\text{BIM+,SSO}}$  and  $\pi^{\text{BIM+,GCM}}$ . We therefore need to make further adjustments to properly account for this. We will make a new data set that takes each of the individuals that were responsive to the first stage intervention and creates an identical copy of their observation. For these individuals we will let  $A_{ij}^{(2)} = 1$  for one of the observations and  $A_{ij}^{(2)} = -1$  for the other.

The model for the number of days absent is then given by

$$Y_{ij} = \beta_0 + \beta_1 \text{BAU}_{ij} + \beta_2 A_{ij}^{(1)} + \beta_3 A_{ij}^{(2)} + \beta_4 A_{ij}^{(1)} A_{ij}^{(2)} + \xi^T \mathbf{X}_{ij}^{(1)} + \gamma_j + \nu_{ij}$$

where  $\text{BAU}_{ij}$  is an indicator for whether student  $i$  in school  $j$  was assigned to BAU. The regression coefficients are estimated using weighted least squares by minimizing

$$\sum_{j=1}^m \sum_{i=1}^{n_j} W_{ij} (Y_{ij} - \beta_0 - \beta_1 \text{BAU}_{ij} - \beta_2 A_{ij}^{(1)} - \beta_3 A_{ij}^{(2)} - \beta_4 A_{ij}^{(1)} A_{ij}^{(2)} - \xi^T \mathbf{X}_{ij}^{(1)} - \gamma_j).$$

The expected number of days absent for each of the difference treatment regimes embedded in the trial can then be expressed as a linear combination of the coefficients such that

1.  $\mu^{\text{BIM+,SSO}} = \beta_0 + \beta_2 + \beta_3 + \beta_4$
2.  $\mu^{\text{BIM+,GCM}} = \beta_0 + \beta_2 - \beta_3 - \beta_4$
3.  $\mu^{\text{BIM-,SSO}} = \beta_0 - \beta_2 + \beta_3 - \beta_4$
4.  $\mu^{\text{BIM-,GCM}} = \beta_0 - \beta_2 - \beta_3 + \beta_4$
5.  $\mu^{\text{BAU}} = \beta_0 + \beta_1.$

We can then test all of the pairwise differences between all of the adaptive treatment regimes. For example, to test if  $\mu^{\text{BIM+,SSO}}$  is different than  $\mu^{\text{BAU}}$  is then equivalent to testing

if  $\beta_2 + \beta_3 + \beta_4 - \beta_1$  is different than zero. We will use robust sandwich standard errors to properly take into account the fact that some of the observations were replicated.

When the outcome of interest is whether the student was chronically absent the model will be given by

$$P(Z_{ij} = 1) = \Phi(\beta_0 + \beta_1 \text{BAU}_{ij} + \beta_2 A_{ij}^{(1)} + \beta_3 A_{ij}^{(2)} + \beta_4 A_{ij}^{(1)} A_{ij}^{(2)} + \zeta^T \mathbf{X}_{ij}^{(1)} + \gamma_j)$$

where  $\Phi$  is the cumulative distribution function of a standard normal random variable.

Table 4.3 and 4.4 display the expected number of days absent for each of the different adaptive interventions embedded in the trial. Table 4.3 displays the results of testing all of the pairwise differences with BAU and table 4.4 displays the tests for the remaining differences. We use a Bonferroni correction to adjust for multiple comparisons.

**Table 4.3** Pairwise differences between each adaptive intervention and BAU on the number of days absent during both semesters

Comparison	AI-1	AI-2	AI-3	AI-4	BAU	Diff.	Effect Size	p-value
AI-1 vs. BAU	10.550				11.410	0.860	0.123	<0.001
AI-2 vs. BAU		10.562			11.410	0.847	0.124	<0.001
AI-3 vs. BAU			10.468		11.410	0.941	0.136	<0.001
AI-4 vs. BAU				10.620	11.410	0.790	0.112	<0.001

Table 4.5 and 4.6 display the predicted probability of a student being chronically absent through the entire school year if assigned to each of the different adaptive interventions embedded in the trial. Table 4.5 displays the results of testing all of the pairwise differences with BAU and table 4.6 displays the tests for the remaining differences. We use a Bonferroni correction to adjust for multiple comparisons.

We can see that there is a significant decrease in the number of days absent and the probability of being chronically absent for all of the embedded interventions when compared

**Table 4.4** Pairwise differences between each of the adaptive interventions on the number of days absent during both semesters

Comparison	AI-1	AI-2	AI-3	AI-4	Difference	Effect Size	p-value
AI-1 vs. AI-2	10.550	10.562			0.013	0.002	0.895
AI-1 vs. AI-3	10.550		10.468		0.081	0.011	0.507
AI-1 vs. AI-4	10.550			10.620	0.070	0.010	0.572
AI-2 vs. AI-3		10.562	10.468		0.094	0.013	0.435
AI-2 vs. AI-4		10.562		10.620	0.057	0.008	0.638
AI-3 vs. AI-4			10.468	10.620	0.151	0.026	0.126

**Table 4.5** Pairwise differences between each adaptive intervention and BAU on chronic absenteeism during both semesters

Comparison	AI-1	AI-2	AI-3	AI-4	BAU	Diff.	Effect Size	p-value
AI-1 vs. BAU	0.167				0.251	0.083	0.206	<0.001
AI-2 vs. BAU		0.170			0.251	0.080	0.198	<0.001
AI-3 vs. BAU			0.167		0.251	0.083	0.206	<0.001
AI-4 vs. BAU				0.179	0.251	0.072	0.175	<0.001

to BAU. There is not a significant difference between any of the four adaptive interventions though.

### 4.3.5 Research Question 3

The third research question asks, for the students that were non-responsive to the first stage intervention, is it better to assign Student Staff Outreach (SSO) or Goal Commitment Messaging (GCM)? This analysis is based on the results of research question 2. If  $\pi^{\text{BIM}+, \text{SSO}}$  leads to significantly less missed days than  $\pi^{\text{BIM}+, \text{GCM}}$  and  $\pi^{\text{BIM}-, \text{SSO}}$  is also significantly better than  $\pi^{\text{BIM}-, \text{GCM}}$ , we could conclude that SSO is preferable to GCM in the second stage for students that are non-responsive. Likewise if  $\pi^{\text{BIM}+, \text{GCM}}$  is preferable to  $\pi^{\text{BIM}+, \text{SSO}}$  and  $\pi^{\text{BIM}-, \text{GCM}}$  is better than  $\pi^{\text{BIM}-, \text{SSO}}$  we could conclude that GCM is the best intervention for

**Table 4.6** Pairwise differences between each of the adaptive interventions on chronic absenteeism during both semesters

Comparison	AI-1	AI-2	AI-3	AI-4	Difference	Effect Size	p-value
AI-1 vs. AI-2	0.167	0.170			0.003	0.007	0.649
AI-1 vs. AI-3	0.167		0.167		0.000	0.001	0.978
AI-1 vs. AI-4	0.167			0.179	0.012	0.031	0.142
AI-2 vs. AI-3		0.170	0.167		0.003	0.008	0.694
AI-2 vs. AI-4		0.170		0.179	0.009	0.023	0.263
AI-3 vs. AI-4			0.167	0.179	0.012	0.031	0.062

students who were non-responsive after the first stage.

It may not be the case that one of SSO and GCM is always preferable no matter what first stage intervention was assigned. There is also interest in comparing if there is a difference between all regimes that assign SSO to non-responsive students and the ones that assign GCM when averaging over the possible first stage treatments. The model for the number of days absent for this analysis will be identical to the one used in research question 2 so we again have that

$$Y_{ij} = \beta_0 + \beta_1 \text{BAU}_{ij} + \beta_2 A_{ij}^{(1)} + \beta_3 A_{ij}^{(2)} + \beta_4 A_{ij}^{(1)} A_{ij}^{(2)} + \xi^T X_{ij}^{(1)} + \gamma_j + \nu_{ij}$$

and

$$P(Z_{ij} = 1) = \Phi(\beta_0 + \beta_1 \text{BAU}_{ij} + \beta_2 A_{ij}^{(1)} + \beta_3 A_{ij}^{(2)} + \beta_4 A_{ij}^{(1)} A_{ij}^{(2)} + \xi^T \mathbf{X}_{ij}^{(1)} + \gamma_j)$$

when we are interested in the probability of a student being chronically absent. We will also fit the model using weighted least squares and duplicating the observations for the responsive students as we did for research question 2.

We then have that the average outcome for regimes that assign School Staff Outreach to



non-responsive students is given by

$$\frac{\mu^{BIM+,SSO} + \mu^{BIM-,SSO}}{2} = \beta_0 + \beta_3$$

and the average outcome for regimes that assign Goal Commitment Messaging to non-responsive students is given by

$$\frac{\mu^{BIM+,SSO} + \mu^{BIM-,SSO}}{2} = \beta_0 - \beta_3.$$

Therefore to test if there is a difference between SSO and GCM we can test if  $\beta_3$  is significantly different than zero.

Table 4.7 contains the estimated mean number of days absent for regimes that assign Student Staff Outreach or Goal Commitment Messaging to non-responders during the second stage and the results of the test for if there is a significant difference between them. Table 4.8 displays the predicted probabilities of a student being chronically absent for the two different regimes. We can see that there is not a significant difference for both outcomes of interest.

**Table 4.7** Differences between second-stage interventions on the number of days absent during both semesters

Comparison	SSO	GCM	Difference	Effect Size	p-value
SSO vs. GCM	10.509	10.591	0.082	0.014	0.234

### 4.3.6 Research Question 4

The fourth research question examines how to tailor interventions based on individual family and school characteristics. For this analysis we will use Q-learning to estimate an

**Table 4.8** Differences between second-stage interventions on chronic absenteeism during both semesters

Comparison	SSO	GCM	Difference	Effect Size	p-value
SSO vs. GCM	0.167	0.174	0.007	0.019	0.100

optimal dynamic treatment regime.

Define  $\mathbf{H}_{ij}^{(t)}$  to be the history of the student available to the decision maker at stage  $t$ . Therefore  $\mathbf{H}_{ij}^{(1)} = \mathbf{X}_{ij}^{(1)}$  and  $\mathbf{H}_{ij}^{(2)} = (\mathbf{X}_{ij}^{(1)}, A_{ij}^{(1)}, Y_{ij}^{(1)}, Z_{ij}^{(1)}, R_{ij})$ . As before  $\boldsymbol{\pi} = (\pi_1, \pi_2)$  denotes a treatment regime such that  $\pi_t : \text{dom } \mathbf{H}^{(t)} \rightarrow \text{dom } A^{(t)}$ . Define  $V(\boldsymbol{\pi}) = \mathbb{E} Y^*(\boldsymbol{\pi})$ . Our goal is to estimate a regime  $\boldsymbol{\pi}^{\text{opt}}$  such that  $V(\boldsymbol{\pi}^{\text{opt}}) \leq V(\boldsymbol{\pi})$ . Note that in previous chapters we wanted to find the regime,  $\boldsymbol{\pi}$ , that maximizes  $V(\boldsymbol{\pi})$ , but in this application the outcome,  $Y$ , is the number of days absent for a student which we hope to be as small as possible.

Let  $\psi_t(\mathbf{h}^{(t)})$  be the feasible intervention options for a student with history  $\mathbf{h}^{(t)}$ . Define  $Q_2(\mathbf{h}^{(2)}, a^{(2)}) = \mathbb{E}\{Y | \mathbf{H}^{(2)} = \mathbf{h}^{(2)}, A^{(2)} = a^{(2)}\}$  for all  $(\mathbf{h}^{(2)}, a^{(2)}) \in \text{dom } \mathbf{H}^{(2)} \times \text{dom } A^{(2)}$  such that  $a^{(2)} \in \psi_2(\mathbf{h}^{(2)})$ . Define

$$Q_1(\mathbf{h}^{(1)}, a^{(1)}) = \mathbb{E} \left\{ \min_{a^{(2)} \in \psi_2(\mathbf{H}^{(2)})} Q_2(\mathbf{H}^{(2)}, a^{(2)} | \mathbf{H}^{(1)} = \mathbf{h}^{(1)}, A^{(1)} = a^{(1)}) \right\}$$

for all  $(\mathbf{h}^{(1)}, a^{(1)}) \in \text{dom } \mathbf{H}^{(1)} \times \text{dom } A^{(1)}$  such that  $a^{(1)} \in \psi_1(\mathbf{h}^{(1)})$ . We can show via dynamic programming that  $\boldsymbol{\pi}^{\text{opt}} = (\pi_1^{\text{opt}}, \pi_2^{\text{opt}})$  where  $\pi_t^{\text{opt}} = \arg \min_{a^{(t)}} Q_t(\mathbf{h}^{(t)}, a^{(t)})$  then fulfills  $V(\boldsymbol{\pi}^{\text{opt}}) \leq V(\boldsymbol{\pi})$  for all treatment regimes  $\boldsymbol{\pi}$ .

As discussed in previous chapters, Q-learning is a regression based approximate dynamic programming algorithm that works by constructing estimators for each of the separate Q-functions. This then leads to a plug-in estimate for the optimal treatment regime. For this study, only the non-responsive students are re-randomized in the second stage which leads to some minor adjustments to how the algorithm was presented in Chapters 2 and 3.

We will postulate linear models for the Q-functions. Therefore the model for the second stage Q-function will be given by

$$Q_2(\mathbf{h}^{(2)}, \mathbf{a}^{(2)}; \beta_2) = \beta_{20} + \beta_{21}x^{(1)} + \beta_{22}a^{(1)} + \beta_{23}x^{(2)} + \beta_{24}a^{(2)} + \beta_{25}x^{(1)}a^{(2)} + \beta_{26}a^{(1)}a^{(2)} + \beta_{27}x^{(2)}a^{(2)}.$$

where  $x^{(1)}$  is an indicator for whether the student was considered at risk at baseline and  $x^{(2)}$  is the number of days absent in the fall. Then the estimator  $\hat{\beta}_2$  is given

$$\hat{\beta}_2 = \arg \min_{\beta_2} \sum_{j=1}^m \sum_{i=1}^{n_j} \{Y_{ij}^{(2)} - Q_2(\mathbf{H}_{ij}^{(2)}, A_{ij}^{(2)}; \beta_2)\} (1 - R_{ij}).$$

The estimated optimal rule for non-responsive students during the second stage is then given by

$$\begin{aligned} \hat{\pi}_2(\mathbf{h}^{(2)}) &= \arg \min_{a^{(2)}} Q_2(\mathbf{h}^{(2)}, \mathbf{a}^{(2)}; \hat{\beta}_2) \\ &= -\text{sign}(\hat{\beta}_{24} + \hat{\beta}_{25}x^{(1)} + \hat{\beta}_{26}a^{(1)} + \hat{\beta}_{27}x^{(2)}). \end{aligned}$$

Define the pseudo outcome,  $\tilde{Y}$ , to be the predicted number of days missing for a student that follows the estimated optimal regime during the second stage. Therefore

$$\tilde{Y}_{ij} = R_{ij}Y_{ij} + (1 - R_{ij}) \left\{ Y_{ij}^{(1)} + \min_{a^{(2)}} Q_2(\mathbf{H}_{ij}^{(2)}, \mathbf{a}^{(2)}; \hat{\beta}_2) \right\}.$$

The model for the first stage Q-function will be given by

$$Q_1(\mathbf{h}^{(1)}, \mathbf{a}^{(1)}; \beta_1) = \beta_{10} + \beta_{11}x^{(1)} + \beta_{12}a^{(1)} + \beta_{13}x^{(1)}a^{(1)}$$

and we can estimate  $\beta_1$  by

$$\hat{\beta}_1 = \arg \min_{\beta_1} \sum_{j=1}^m \sum_{i=1}^{n_j} \{ \tilde{Y}_{ij} - Q_1(\mathbf{H}_{ij}^{(1)}, A_{ij}^{(1)}; \beta_1) \}.$$

We can then estimate the optimal first stage regime by

$$\begin{aligned}\hat{\pi}_1(\mathbf{h}_1) &= \arg \min_{a^{(1)}} Q_1(\mathbf{h}^{(1)}, a^{(1)}; \hat{\beta}_1) \\ &= -\text{sign}(\hat{\beta}_{12} + \hat{\beta}_{13}x^{(1)}).\end{aligned}$$

Table 4.9 displays the estimated coefficients and 95% confidence intervals for the second stage Q-function. Table 4.10 shows the estimated coefficients and 95% confidence intervals for the first stage Q-function. For the first stage coefficients the confidence intervals are calculated using the m-out-of-n bootstrap proposed in Chakraborty et al. (2013) to obtain valid intervals despite the non-regularity. Note that the confidence intervals for all of the treatment interaction terms contain zero. Therefore there is not significant evidence that tailoring interventions based on these characteristics would provide a decrease in the number of days absent.

Table 4.11 displays inverse probability weighted estimates of the expected number of days absent for the estimated optimal regime as well as the adaptive interventions embedded within the regime. 95% confidence intervals are also provided that were calculated using the non-parametric bootstrap. Note that the bootstrap may perform poorly in this context as the IPW is not a smooth estimator. These results also indicate that there is not enough evidence to conclude that tailoring interventions would lead to a decrease in absenteeism.

**Table 4.9** Parameter estimates for the second stage Q-function

Variable	Coefficient Estimate	95% Lower CI	95 % Upper CI
Intercept	2.888	1.648	4.128
Treatment 1	-0.202	-0.456	0.053
At Risk	3.953	3.430	4.476
Fall Absenteeism	1.064	0.889	1.238
Treatment 2	-0.272	-1.512	0.968
Trt 2 * Trt 1	0.020	-0.235	0.274
Trt 2 * At Risk	-0.282	-0.805	0.241
Trt 2 * Fall Absenteeism	0.036	-0.138	0.210

**Table 4.10** Parameter estimates for the first stage Q-function

Variable	Coefficient Estimate	95% Lower CI	95 % Upper CI
Intercept	8.586	8.161	8.490
At Risk	7.267	6.635	7.357
Treatment 1	0.058	-0.118	0.249
Trt 1 * At Risk	-0.133	-0.615	0.179

**Table 4.11** IPW Estimates of expected outcomes of different regimes

Regime	Mean Days Absent	95% Lower CI	95% Upper CI
$\hat{\pi}$	10.161	9.850	10.496
$\pi^{BIM+,SSO}$	10.561	10.229	10.862
$\pi^{BIM+,GCM}$	10.504	10.182	10.846
$\pi^{BIM-,SSO}$	10.418	10.107	10.737
$\pi^{BIM-,GCM}$	10.663	10.355	11.021

## BIBLIOGRAPHY

- Almirall, D., S. N. Compton, M. Gunlicks-Stoessel, N. Duan, and S. A. Murphy (2012). Designing a pilot sequential multiple assignment randomized trial for developing an adaptive treatment strategy. *Statistics in medicine* 31(17), 1887–1902.
- Bai, X., A. A. Tsiatis, W. Lu, and R. Song (2017). Optimal treatment regimes for survival endpoints using a locally-efficient doubly-robust estimator from a classification perspective. *Lifetime Data Analysis* 23(4), 585–604.
- Bellman, R. (1957). *Dynamic Programming*. Princeton University Press.
- Berger, R. L. and D. D. Boos (1994). P values maximized over a confidence set for the nuisance parameter. *Journal of the American Statistical Association* 89(427), 1012–1016.
- Cain, L. E., J. M. Robins, E. Lanoy, R. Logan, D. Costagliola, and M. A. Hernán (2010). When to start treatment? a systematic approach to the comparison of dynamic regimes using observational data. *The international journal of biostatistics* 6(2).
- Chakraborty, B., E. B. Laber, and Y. Zhao (2013). Inference for optimal dynamic treatment regimes using an adaptive m-out-of-n bootstrap scheme. *Biometrics* 69(3), 714–723.
- Chakraborty, B., E. B. Laber, and Y.-Q. Zhao (2014). Inference about the expected performance of a data-driven dynamic treatment regime. *Clinical Trials* 11(4), 408–417.
- Chakraborty, B. and E. E. Moodie (2013). *Statistical Methods for Dynamic Treatment Regimes*. Springer.
- Chakraborty, B., S. Murphy, and V. Strecher (2009). Inference for non-regular parameters in optimal dynamic treatment regimes. *Statistical Methods in Medical Research* 19(3).
- Chakraborty, B., V. Strecher, and S. Murphy (2010). Inference for nonregular parameters

- in optimal dynamic treatment regimes. *Statistical Methods in Medical Research* 19(3), 1–27.
- Chang, H. N. and M. Romero (2008). *Present, Engaged, and Accounted For: The critical Importance of Addressing Chronic Absence in the Early Grades*. New York, NY: National Center for Children in Poverty.
- Dawid, A. (1994). Selection paradoxes of bayesian inference. *Lecture Notes-Monograph Series*, 211–220.
- Ehrlich, S. B., J. A. Gwynne, A. S. Pareja, E. M. Allensworth, P. Moore, S. Jagesic, and E. Sorice (2014). *Preschool Attendance in Chicago Public Schools: Relationship with Learning Outcomes and Reasons for Absences*. Chicago: University of Chicago Consortium on Chicago School Research.
- Ertefaie, A. (2014). Constructing dynamic treatment regimes in infinite-horizon settings. *arXiv preprint arXiv:1406.0764*, 1–26.
- Ginsburg, A., P. Jordan, and H. Chang (2014). *Absences Add Up: How School Attendance Influences Student Success*. San Francisco: Attendance Works.
- Goldberg, Y. and M. R. Kosorok (2012). Q-learning with censored data. *Annals of statistics* 40(1), 529.
- Gottfried, M. A. (2019). Chronic absenteeism in the classroom context: Effects on achievement. *Urban Education* 54(1), 3–34.
- Hager, R., A. A. Tsiatis, and M. Davidian (2018). Optimal two-stage dynamic treatment regimes from a classification perspective with censored survival data. *Biometrics* 74, 1180–1192.
- Henderson, H. V. and S. Searle (1979). Vec and vech operators for matrices, with some uses in jacobians and multivariate statistics. *Canadian Journal of Statistics* 7(1), 65–81.

- Henderson, T., C. Hill, and K. Norton (2014). *The Connection Between Missing School and Health: A Review of Chronic Absenteeism and Student Health in Oregon*. Upstream Public Health.
- Henry, K. L. and T. P. Thornberry (2010). Truancy and escalation of substance use during adolescence. *Journal of Studies on Alcohol and Drugs* 71(1), 115–124.
- Hirano, K. and J. R. Porter (2012). Impossibility results for nondifferentiable functionals. *Econometrica* 80(4), 1769–1790.
- Jiang, R., W. Lu, R. Song, and M. Davidian (2017). On estimation of optimal treatment regimes for maximizing t-year survival probability. *Journal of the Royal Statistical Society, Series B* 79, 1165–1185.
- Kidwell, K. M. (2014). Smart designs in cancer research: Past, present, and future. *Clinical Trials* 11(4), 445–456.
- Kosorok, M. R. (2008). *Introduction to empirical processes and semiparametric inference*. Springer.
- Kosorok, M. R. and E. E. Moodie (2015). *Adaptive treatment strategies in practice: planning trials and analyzing data for personalized medicine*. SIAM.
- Laber, E., K. Linn, and L. Stefanski (2014). Interactive model building for q-learning. *Biometrika* 101(4), 831–847.
- Laber, E. and Y. Zhao (2015). Tree-based methods for estimating individualized treatment regimes. *Biometrika* 102(3), 501–514.
- Laber, E. B., D. J. Lizotte, M. Qian, W. E. Pelham, and S. A. Murphy (2014). Dynamic treatment regimes: Technical challenges and applications. *Electronic journal of statistics* 8(1), 1225.
- Laber, E. B. and A.-M. Staicu (2017). Functional feature construction for individualized treatment regimes. *Journal of the American Statistical Association* (just-accepted).



- Laber, E. B., F. Wu, C. Munera, I. Lipkovich, S. Colucci, and S. Ripa (2018). Identifying optimal dosage regimes under safety constraints: An application to long term opioid treatment of chronic pain. *Statistics in Medicine* 37(9), 1407–1418.
- Laber, E. B., Y.-Q. Zhao, T. Regh, M. Davidian, A. Tsiatis, J. B. Stanford, D. Zeng, R. Song, and M. R. Kosorok (2015). Using pilot data to size a two-arm randomized trial to find a nearly optimal personalized treatment strategy. *Statistics in medicine*.
- Lavori, P. and R. Dawson (2000). A design for testing clinical strategies: biased adaptive within-subject randomization. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 163(1), 29–38.
- Lavori, P. W. and R. Dawson (2004). Dynamic treatment regimes: practical design considerations. *Clinical trials* 1(1), 9–20.
- Lei, H., I. Nahum-Shani, K. Lynch, D. Oslin, and S. Murphy (2012). A “smart” design for building individualized treatment sequences. *Annual Review of Clinical Psychology* 8, 21–48.
- Linn, K. A., E. B. Laber, and L. A. Stefanski (2016). Interactive q-learning for quantiles. *Journal of the American Statistical Association* (just-accepted), 1–37.
- Luckett, D. J., E. B. Laber, A. R. Kahkoska, D. M. Maahs, E. Mayer-Davis, and M. R. Kosorok (2016). Estimating dynamic treatment regimes in mobile health using v-learning. *arXiv preprint arXiv:1611.03531*.
- Luedtke, A. R. and M. J. Van Der Laan (2016). Statistical inference for the mean outcome under a possibly non-unique optimal treatment strategy. *The Annals of Statistics* 44(2), 713–742.
- Moodie, E., T. Richardson, and D. Stephens (2007). Demystifying optimal dynamic treatment regimes. *Biometrics* 63(2), 447–455.

- Moodie, E., T. Richardson, and D. Stephens (2010). Estimating optimal dynamic regimes: Correcting bias under the null. *Biometrics* 63(2), 447–455.
- Moodie, E. E., B. Chakraborty, and M. S. Kramer (2012). Q-learning for estimating optimal dynamic treatment rules from observational data. *Canadian Journal of Statistics* 40(4), 629–645.
- Murphy, S. (2005a). An experimental design for the development of adaptive treatment strategies. *Statistics in medicine* 24(10), 1455–1481.
- Murphy, S. A. (2003a). Optimal dynamic treatment regimes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 65(2), 331–355.
- Murphy, S. A. (2003b). Optimal dynamic treatment regimes. *Journal of the Royal Statistical Society, Series B* 65(2), 331–366.
- Murphy, S. A. (2005b, Jul). A generalization error for Q-learning. *Journal of Machine Learning Research* 6, 1073–1097.
- Murphy, S. A., K. G. Lynch, D. Oslin, J. R. McKay, and T. TenHave (2007). Developing adaptive treatment strategies in substance abuse research. *Drug and alcohol dependence* 88, S24–S30.
- Murray, T., Y. Yuan, and P. Thall (2017). A bayesian machine learning approach for optimizing dynamic treatment regimes. *Journal of the American Statistical Association Early view*, 1–37.
- Nahum-Shani, I., A. Ertefaie, X. L. Lu, K. G. Lynch, J. R. McKay, D. W. Oslin, and D. Almirall (2017). A smart data analysis method for constructing adaptive treatment strategies for substance use disorders. *Addiction* 112(5), 901–909.

- Nahum-Shani, I., M. Qian, D. Almirall, W. E. Pelham, B. Gnagy, G. A. Fabiano, J. G. Waxmonsky, J. Yu, and S. A. Murphy (2012). Q-learning: A data analysis method for constructing adaptive interventions. *Psychological methods* 17(4), 478.
- Orellana, L., A. Rotnitzky, and J. Robins (2010). Dynamic regime marginal structural mean models for estimation of optimal dynamic treatment regimes, part i: Main content. *Int. Jrn. of Biostatistics* 6(2), 1–49.
- Pineau, J., M. G. Bellemare, A. J. Rush, A. Ghizaru, and S. A. Murphy (2007). Constructing evidence-based treatment strategies using methods from computer science. *Drug Alcohol Dependence* 88, S52–S60.
- Praestgaard, J. and J. A. Wellner (1993). Exchangeably weighted bootstraps of the general empirical process. *The Annals of Probability* 21(4), 2053–2086.
- PSU Methodology Center, T. (2017a, July). Nih program announcements.
- PSU Methodology Center, T. (2017b, July). Smart studies.
- Qian, M. and S. Murphy (2011). Performance Guarantees for Individualized Treatment Rules. *The Annals of Statistics* 39(2), 1180–1210.
- Qian, M., I. Nahum-Shani, and S. A. Murphy (2013). Dynamic treatment regimes. In *Modern Clinical Trial Analysis*, pp. 127–148. Springer.
- Ready, D. D. (2010). Socioeconomic disadvantage, school attendance, and early cognitive development: The differential effects of school exposure. *Sociology of Education* 83(4), 271–286.
- Robins, J. (2004a). Optimal structural nested models for optimal sequential decisions. In *Proceedings of the Second Seattle Symposium in Biostatistics: Analysis of Correlated Data*.

- Robins, J., L. Orellana, and A. Rotnitzky (2008). Estimation and extrapolation of optimal treatment and testing strategies. *Statistics in Medicine* 27(23), 4678–4721.
- Robins, J. M. (2004b). Optimal structural nested models for optimal sequential decisions. In *Proceedings of the Second Seattle Symposium on Biostatistics*, pp. 189–326. Springer.
- Rubin, D. (1978). Bayesian inference for causal effects: The role of randomization. *The Annals of Statistics* 6(1), 34–58.
- Rush, A., M. Trivedi, and M. Fava (2003). Depression, IV: STAR\* D treatment trial for depression. *American Journal of Psychiatry* 160(2), 237.
- Rush, A., M. Trivedi, H. Ibrahim, T. Carmody, B. Arnow, D. Klein, J. Markowitz, P. Ninan, S. Kornstein, R. Manber, M. Thase, J. Kocsis, and M. Keller (2003). The 16-Item quick inventory of depressive symptomatology (QIDS), clinician rating (QIDS-C), and self-report (QIDS-SR): a psychometric evaluation in patients with chronic major depression. *Biological Psychiatry* 54(5), 573–583.
- Schulte, P., A. Tsiatis, E. Laber, and M. Davidian (2014). Q- and a-learning methods for estimating optimal dynamic treatment regimes. *Statistical Science* 29(4), 640–661.
- Shortreed, S., E. B. Laber, D. J. Lizotte, T. S. Stroup, J. Pineau, and S. A. Murphy (2010). Informing sequential clinical decision-making through reinforcement learning: an empirical study. *Machine Learning* 84(1–2), 109–136.
- Shortreed, S. M., E. Laber, T. S. Stroup, J. Pineau, and S. A. Murphy (2014). A multiple imputation strategy for sequential multiple assignment randomized trials. *Statistics in Medicine* 33(24), 4204–4214.
- Song, R., W. Wang, D. Zeng, and M. Kosorok (2015a). Penalized q-learning for dynamic treatment regimes. *Statistica Sinica* 25(3), 901–920.

- Song, R., W. Wang, D. Zeng, and M. R. Kosorok (2015b). Penalized q-learning for dynamic treatment regimens. *Statistica Sinica* 25(3), 901.
- Splawa-Neyman, J., D. Dabrowska, T. Speed, et al. (1990). On the application of probability theory to agricultural experiments. essay on principles. section 9. *Statistical Science* 5(4), 465–472.
- Tao, Y. and L. Wang (2017). Adaptive contrast weighted learning for multi-stage multi-treatment decision-making. *Biometrics* 73(1), 145–155.
- Thall, P. F., R. E. Millikan, H.-G. Sung, et al. (2000). Evaluating multiple treatment courses in clinical trials. *Statistics in medicine* 19(8), 1011–1028.
- van der Laan, M. J., M. L. Petersen, and M. M. Joffe (2005). History-adjusted marginal structural models and statically-optimal dynamic treatment regimens. *The International Journal of Biostatistics* 1(1).
- Wang, L., A. Rotnitzky, X. Lin, R. E. Millikan, and P. F. Thall (2012). Evaluation of viable dynamic treatment regimes in a sequentially randomized trial of advanced prostate cancer. *Journal of the American Statistical Association* 107(498), 493–508.
- Watkins, C. (1989). *Learning from Delayed Rewards*. Ph. D. thesis, Cambridge University.
- Wu, F., E. B. Laber, I. A. Lipkovich, and E. Severus (2015). Who will benefit from antidepressants in the acute treatment of bipolar depression? a reanalysis of the step-bd study by sachs et al. 2007, using q-learning. *International journal of bipolar disorders* 3(1), 7.
- Xu, Y., P. Müller, A. S. Wahed, and P. F. Thall (2016). Bayesian nonparametric estimation for dynamic treatment regimes with sequential transition times. *Journal of the American Statistical Association* 111(515), 921–950.
- Zhang, B., A. A. Tsiatis, M. Davidian, M. Zhang, and E. Laber (2012). Estimating optimal treatment regimes from a classification perspective. *Stat* 1(1), 103–114.

- Zhang, B., A. A. Tsiatis, E. B. Laber, and M. Davidian (2012). A robust method for estimating optimal treatment regimes. *Biometrics* 68(4), 1010–1018.
- Zhang, B., A. A. Tsiatis, E. B. Laber, and M. Davidian (2013). Robust estimation of optimal dynamic treatment regimes for sequential treatment decisions. *Biometrika* 100(3), 681–694.
- Zhang, Y., E. Laber, A. Tsiatis, and M. Davidian (2017). List-based treatment regimes. *Journal of the American Statistical Association Early view*, 1–25.
- Zhang, Y., E. B. Laber, A. Tsiatis, and M. Davidian (2015). Using decision lists to construct interpretable and parsimonious treatment regimes. *Biometrics* 71(4), 895–904.
- Zhao, Y., D. Zeng, E. B. Laber, and M. R. Kosorok (2015). New statistical learning methods for estimating optimal dynamic treatment regimes. *Journal of the American Statistical Association* 110(510), 583–598.
- Zhao, Y., D. Zeng, E. B. Laber, R. Song, M. Yuan, and M. R. Kosorok (2015). Doubly robust learning for estimating individualized treatment with censored data. *Biometrika* 102(1), 151–168.
- Zhao, Y., D. Zeng, A. J. Rush, and M. R. Kosorok (2012). Estimating individualized treatment rules using outcome weighted learning. *Journal of the American Statistical Association* 107(499), 1106–1118.
- Zhao, Y., D. Zeng, M. A. Socinski, and M. R. Kosorok (2011). Reinforcement learning strategies for clinical trials in nonsmall cell lung cancer. *Biometrics* 67(4), 1422–1433.
- Zhou, X., N. Mayer-Hamblett, U. Khan, and M. R. Kosorok (2017). Residual weighted learning for estimating individualized treatment rules. *Journal of the American Statistical Association* 112(517), 169–187.

## **APPENDICES**

## APPENDIX

### A

# SAMPLE SIZE CALCULATIONS FOR SMARTS

## A.1 Proofs of technical results

*Proof of Lemma 3.1.* It can be shown by direct calculation that if  $W \sim \text{Normal}(\mu, \omega^2)$  then

$$\begin{aligned}\mathbb{E}|W| &= \int \frac{|\mu + \omega w|}{\sqrt{2\pi}} \exp\left(-\frac{w^2}{2}\right) dw \\ &= \mu \left\{1 - \Phi\left(-\frac{\mu}{\omega}\right)\right\} + \frac{\sqrt{2}\omega}{\sqrt{\pi}} \exp\left(-\frac{\mu^2}{2\omega^2}\right).\end{aligned}$$



Write

$$\begin{aligned}
Q_1(h_1, a_1) &= \mathbb{E} \left\{ \max_{a_2} Q_2(\mathbf{H}_2, a_2) \middle| \mathbf{H}_1 = \mathbf{h}_1, A_1 = a_1 \right\} \\
&= \mathbb{E} \left( \mathbf{H}_{2,1}^\top \boldsymbol{\beta}_{2,0}^* + |\mathbf{H}_{2,1}^\top \boldsymbol{\beta}_{2,1}^*| \middle| \mathbf{H}_1 = \mathbf{h}_1, A_1 = a_1 \right) \\
&= \mathbf{h}_{1,0}^\top \boldsymbol{\xi}_{1,0}^* + a_1 \mathbf{h}_{1,1}^\top \boldsymbol{\xi}_{1,1}^* \\
&\quad + \mathbb{E} \left( |\mathbf{H}_{1,2}^\top \boldsymbol{\varpi}_{1,2}^* + A_1 \mathbf{H}_{1,3}^\top \boldsymbol{\varpi}_{1,3}^* + \tau^* Z| \middle| \mathbf{H}_1 = \mathbf{h}_1, A_1 = a_1 \right) \\
&= \mathbf{h}_{1,0}^\top \boldsymbol{\xi}_{1,0}^* + a_1 \mathbf{h}_{1,1}^\top \boldsymbol{\xi}_{1,1}^* \\
&\quad + \int \frac{|\mathbf{h}_{1,2}^\top \boldsymbol{\varpi}_{1,2}^* + a_1 \mathbf{h}_{1,3}^\top \boldsymbol{\varpi}_{1,3}^* + \tau^* z|}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) dz,
\end{aligned}$$

from which the result follows. ■

*Proof of Lemma 3.3.* The function  $\nu$  is continuously differentiable in a neighborhood of  $\{\tau^*, \omega^{*\top}, \text{vech}(\Omega^*)^\top\}^\top$ . Thus, the result follows immediately from a Taylor series expansion, (AN6), and Slutsky's theorem; that is,

$$\begin{aligned}
\sqrt{n} \{ \widehat{V}_n - V(\boldsymbol{\pi}^{\text{opt}}) \} &= \sqrt{n} [ \nu \{ \widehat{\tau}_n, \widehat{\omega}_n, \text{vech}(\widehat{\Omega}_n) \} - \nu \{ \tau^*, \omega^*, \text{vech}(\Omega^*) \} ] \\
&= \nabla \nu \{ \widetilde{\tau}_n, \widetilde{\omega}_n, \text{vech}(\widetilde{\Omega}_n) \}^\top \sqrt{n} \left\{ \begin{array}{c} \widehat{\tau}_n - \tau^* \\ \widehat{\omega}_n - \omega^* \\ \text{vech}(\widehat{\Omega}_n) - \text{vech}(\Omega^*) \end{array} \right\} \\
&\rightsquigarrow \text{Normal}(0, \Sigma^*).
\end{aligned}$$

■

*Proof of Lemma 3.4.* We show that  $\sqrt{n} |V(\widehat{\boldsymbol{\pi}}_n) - V(\boldsymbol{\pi}^{\text{Q,opt}})| = O(n^{-\delta})$ . For any regime  $\boldsymbol{\pi} =$

$(\pi_1, \pi_2)$  write

$$\begin{aligned}
V(\pi) &= \mathbb{E} \left[ \mathbb{E} \left\{ \mathbb{E} \left( Y | \mathbf{H}_2, A_2 \right) \Big|_{A_2=\pi_2(\mathbf{H}_2)} \Big|_{\mathbf{H}_1, A_1} \right\} \Big|_{A_1=\pi_1(\mathbf{H}_1)} \right] \\
&= \mathbb{E} \left( \mathbb{E} [Q_2 \{ \mathbf{H}_2, \pi_2(\mathbf{H}_2) \} | \mathbf{H}_1, A_1] \Big|_{A_1=\pi_1(\mathbf{H}_1)} \right) \\
&= \mathbb{E} \left( \mathbb{E} \left\{ \max_{a_2} Q_2(\mathbf{H}_2, a_2) | \mathbf{H}_1, A_1 \right\} \Big|_{A_1=\pi_1(\mathbf{H}_1)} \right) \\
&\quad - \mathbb{E} \left\{ \mathbb{E} \left( |\Delta Q_2| 1_{\pi_2(\mathbf{H}_2) \neq \pi_2^{Q, \text{opt}}(\mathbf{H}_2)} | \mathbf{H}_1, A_1 \right) \Big|_{A_1=\pi_1(\mathbf{H}_1)} \right\} \\
&= \mathbb{E} [Q_1 \{ \mathbf{H}_1, \pi_1(\mathbf{H}_1) \}] \\
&\quad - \mathbb{E} \left\{ \mathbb{E} \left( |\Delta Q_2| 1_{\pi_2(\mathbf{H}_2) \neq \pi_2^{Q, \text{opt}}(\mathbf{H}_2)} | \mathbf{H}_1, A_1 \right) \Big|_{A_1=\pi_1(\mathbf{H}_1)} \right\} \\
&= \mathbb{E} \left\{ \max_{a_1} Q_1(\mathbf{H}_1, a_1) \right\} - \mathbb{E} \left( |\Delta Q_1| 1_{\pi_1(\mathbf{H}_1) \neq \pi_1^{Q, \text{opt}}(\mathbf{H}_1)} \right) \\
&\quad - \mathbb{E} \left\{ \mathbb{E} \left( |\Delta Q_2| 1_{\pi_2(\mathbf{H}_2) \neq \pi_2^{Q, \text{opt}}(\mathbf{H}_2)} | \mathbf{H}_1, A_1 \right) \Big|_{A_1=\pi_1(\mathbf{H}_1)} \right\} \\
&= V(\pi^{Q, \text{opt}}) - R_1(\pi_1) - R_2(\pi_1, \pi_2).
\end{aligned}$$

Thus,  $|V(\hat{\pi}_n) - V(\pi^{Q, \text{opt}})| \leq R_1(\hat{\pi}_{1,n}) + R_2(\hat{\pi}_{1,n}, \hat{\pi}_{2,n})$ . We next derive rates of convergence for these remainder terms.

We will make use of the fact that  $R_2(\hat{\pi}_{1,n}, \hat{\pi}_{2,n}) \leq 2\mathbb{E}|\Delta Q_2| 1_{\Delta \hat{Q}_{2,n} \Delta Q_2 < 0}$  which follows from writing

$$\begin{aligned}
R_2(\hat{\pi}_{1,n}, \hat{\pi}_{2,n}) &= \mathbb{E} \left\{ \mathbb{E} \left( |\Delta Q_2| 1_{\Delta \hat{Q}_{2,n} \Delta Q_2 < 0} | \mathbf{H}_1, A_1 \right) \Big|_{A_1=\hat{\pi}_{1,n}(\mathbf{H}_1)} \right\} \\
&\leq \mathbb{E} \left\{ \mathbb{E} \left( |\Delta Q_2| 1_{\Delta \hat{Q}_{2,n} \Delta Q_2 < 0} | \mathbf{H}_1, A_1 = 1 \right) \right\} \\
&\quad + \mathbb{E} \left\{ \mathbb{E} \left( |\Delta Q_2| 1_{\Delta \hat{Q}_{2,n} \Delta Q_2 < 0} | \mathbf{H}_1, A_1 = -1 \right) \right\} \\
&= 2\mathbb{E}|\Delta Q_2| 1_{\Delta \hat{Q}_{2,n} \Delta Q_2 < 0},
\end{aligned}$$

where we have used  $P(A_1 = 1|\mathbf{H}_1) = P(A_1 = -1|\mathbf{H}_1) = 1/2$  with probability one. Let  $\epsilon_n$  be any nonnegative sequence diverging to  $\infty$  with  $n$ , then

$$\begin{aligned}
2\mathbb{E}|\Delta Q_2|1_{\Delta\widehat{Q}_{2,n}\Delta Q_2 < 0} &= 2\mathbb{E}|\Delta Q_2|1_{\Delta\widehat{Q}_{2,n}\Delta Q_2 < 0}1_{|\Delta Q_2| \leq \epsilon_n/\sqrt{n}} \\
&\quad + 2\mathbb{E}|\Delta Q_2|1_{\Delta\widehat{Q}_{2,n}\Delta Q_2 < 0}1_{|\Delta Q_2| > \epsilon_n/\sqrt{n}} \\
&\leq 2(\epsilon_n/\sqrt{n})P(|\Delta Q_2| \leq \epsilon_n/\sqrt{n}) \\
&\quad + 2\mathbb{E}|\Delta Q_2|1_{\Delta\widehat{Q}_{2,n}\Delta Q_2 < 0}1_{|\Delta Q_2| > \epsilon_n/\sqrt{n}} \\
&\leq 2M\epsilon_n^{1+\kappa}/n^{1/2+\kappa/2} \\
&\quad + 2\mathbb{E}|\Delta Q_2|1_{\Delta\widehat{Q}_{2,n}\Delta Q_2 < 0}1_{|\Delta Q_2| > \epsilon_n/\sqrt{n}} \\
&= 2M\epsilon_n^{1+\kappa}/n^{1/2+\kappa/2} \\
&\quad + 2\mathbb{E}|\Delta Q_2|1_{(\Delta\widehat{Q}_{2,n}-\Delta Q_2)\text{sgn}(\Delta Q_2)+|\Delta Q_2| < 0}1_{|\Delta Q_2| > \epsilon_n/\sqrt{n}} \\
&\leq 2M\epsilon_n^{1+\kappa}/n^{1/2+\kappa/2} + 2\sqrt{\mathbb{E}|\Delta Q_2|^2}P\{\sqrt{n}|\Delta\widehat{Q}_{2,n}-\Delta Q_2| > \epsilon_n\} \\
&\leq 2M\epsilon_n^{1+\kappa}/n^{1/2+\kappa/2} + 2\sqrt{\mathbb{E}|\Delta Q_2|^2}\exp(-c_{n,2}\epsilon_n^{\ell_{n,2}}),
\end{aligned}$$

where the second inequality follows from (AN8), the third inequality follows from Cauchy-Schwartz, and the fourth from (AN7). Choose  $\epsilon_n = \{c_{n,2}^{-1} \log(n^{1/2+\kappa/2})\}^{1/\ell_{n,2}}$  to obtain that  $\sqrt{n}R_2(\widehat{\pi}_{1,n}, \widehat{\pi}_{n,2}) = O_p(n^{-\kappa/2})$ . An identical argument can be used to derive an analogous bound on  $\sqrt{n}R_1(\widehat{\pi}_{1,n})$ .  $\blacksquare$

*Proof of Theorem 3.7.* Define  $\mathcal{F}_1 = \{\delta(\mu_1, \mu_2) : (\mu_1, \mu_2) \in \Theta\}$  and  $\mathcal{F}_2 = \{\delta^2(\mu_1, \mu_2) : (\mu_1, \mu_2) \in \Theta\}$ ; then both are Donsker by (PR4). By Corollary 2.2 of Praestgaard and Wellner (1993) it follows that  $\sqrt{n}(\mathbb{P}_{n_0,n}^{(b)} - \mathbb{P}_{n_0}) \stackrel{P}{\underset{B}{\rightrightarrows}} \mathbb{G}$  in  $\ell^\infty(\mathcal{F}_1)$  and  $\ell^\infty(\mathcal{F}_2)$ , where  $\mathbb{G}$  is a Brownian bridge with mean zero and covariance  $\text{Cov}\{\mathbb{G}(f_1), \mathbb{G}(f_2)\} = P(f_1 - Pf_1)(f_2 - Pf_2)$ . Similarly, it follows that  $\mathbb{P}_{n_0,n} \stackrel{P}{\underset{B}{\rightrightarrows}} P$  in  $\ell^\infty(\mathcal{F}_1)$  and  $\ell^\infty(\mathcal{F}_2)$ . Applying Slutsky's theorem for empirical processes (e.g.,

Theorem 7.15 in Kosorok, 2008) it follows that

$$\left\{ \begin{array}{c} \sqrt{n}(\mathbb{P}_{n_0,n} - \mathbb{P}_{n_0}) \\ \mathbb{P}_{n_0,n} \\ \mathbb{P}_{n_0,n} \end{array} \right\} \overset{P}{\rightsquigarrow} \overset{B}{\left( \begin{array}{c} \mathbb{G} \\ P_1 \\ P_2 \end{array} \right)}$$

in  $\ell^\infty(\mathcal{F}_1 \cup \mathcal{F}_2)$ , where  $P_1 = \{Pf : f \in \mathcal{F}_1\}$  and  $P_2 = \{Pf : f \in \mathcal{F}_2\}$ . Define  $\mathcal{F}_3 = \{f^2 - (Pf)^2 : f \in \mathcal{F}_1\}$ . Application of the bootstrap continuous mapping theorem ensures that

$$\left\{ \begin{array}{c} \sqrt{n}(\mathbb{P}_{n_0,n} - \mathbb{P}_{n_0}) \\ \mathbb{P}_{n_0,n} \end{array} \right\} \overset{P}{\rightsquigarrow} \overset{B}{\left( \begin{array}{c} \mathbb{G} \\ P_3 \end{array} \right)}$$

in  $\ell^\infty(\mathcal{F}_1 \cup \mathcal{F}_3)$ , where  $P_3 = \{Pf : f \in \mathcal{F}_3\}$ . The desired result follows from the bootstrap continuous mapping theorem.  $\blacksquare$

*Proof of Theorem 3.8.* Let  $\Delta$  denote symmetric set difference, then we have  $P_B \widehat{\Theta}_{n_0,n}^{(b)} \Delta \widehat{\theta}_n \rightarrow 0$  with probability one. Write  $Z_{n_0,n}^{(b)} = O_{P_B} \{f(n_0, n)\}$  to mean that  $\lim_{M \rightarrow \infty} P_B(|Z_{n_0,n}^{(b)}|/f(n_0, n)| \leq M)$  converges to one in probability. It can be seen that

$$\begin{aligned} \sqrt{n} \{ \widehat{V}_{n_0,n}^{Q(b)}(\mu_1, \mu_2) - B_0 \} &= \sqrt{n} \{ \widehat{V}_{n_0,n}^{Q(b)}(\mu_1, \mu_2) - \widehat{V}_{n_0}^Q(\mu_1, \mu_2) \} + \sqrt{n} \{ \widehat{V}_{n_0}^Q(\mu_1, \mu_2) - V^Q(\mu_1, \mu_2) \} \\ &\quad + \sqrt{n} \{ V^Q(\mu_1, \mu_2) - B_0 \} \\ &= O_{P_B}(1) + O_{P_B}(\sqrt{n}/\sqrt{n_0}) + \sqrt{n} \{ V^Q(\mu_1, \mu_2) - B_0 \}. \end{aligned}$$

Thus, provided that  $V^Q(\mu_1^*, \mu_2^*) > B_0 + \eta$ , it follows that

$$\begin{aligned} &\inf_{(\mu_1, \mu_2) \in \Xi_{n_0,n,1-\theta_1}^{(b)}} \frac{\min \left[ \sqrt{n} \{ \widehat{V}_{n_0,n}^{Q(b)}(\mu_1, \mu_2) - B_0 \}, \sqrt{n} \eta \right]}{\widehat{\zeta}_{n_0,n}^{(b)}(\mu_2, \mu_2)} \\ &= \inf_{(\mu_1, \mu_2) \in \Xi_{n_0,n,1-\theta_1}^{(b)}} \frac{\min \left[ \sqrt{n} \{ V^Q(\mu_1, \mu_2) - B_0 \}, \sqrt{n} \eta \right]}{\widehat{\zeta}_{n_0,n}^{(b)}(\mu_2, \mu_2)} + O_{P_B} \{ \max(1, \sqrt{n}/\sqrt{n_0}) \}, \end{aligned}$$

where, by appeal to (PR5), the leading term on the right hand side of the above equality

diverges to  $+\infty$  at rate  $\sqrt{n}$  and thereby dominates the second term. Similarly, it can be seen that

$$\begin{aligned} & \inf_{(\mu_1, \mu_2) \in \Xi_{n,1-\theta_1}} \frac{\min[\sqrt{n} \{ \widehat{V}_n^Q(\mu_1, \mu_2) - B_0 \}, \sqrt{n}\eta]}{\widehat{\zeta}_n(\mu_2, \mu_2)} \\ &= \inf_{(\mu_1, \mu_2) \in \Xi_{n,1-\theta_1}} \frac{\min[\sqrt{n} \{ V^Q(\mu_1, \mu_2) - B_0 \}, \sqrt{n}\eta]}{\widehat{\zeta}_n(\mu_2, \mu_2)} + O_p(1), \end{aligned}$$

where the leading term diverges to  $+\infty$  at rate  $\sqrt{n}$  and thereby dominates the second term.

The desired result follows from Slutsky's theorem. ■

## A.2 Simulation results

**Table A.1** Estimated power (POW) and concentration (OPT) under a correctly specified generative model using the normality-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care.

$\Delta$	Method for $\sigma^*$	$n_0$	(POW)	(OPT)	$\hat{n}^{\text{fixed}}$	$\mathbb{E}\hat{n}$	Med( $\hat{n}$ )	SD $\hat{n}$
0	known	50	99.8	100	74	277	—	—
0	pilot study	50	97.6	100	74	204.51	199.5	60.60
0	surrogate	50	100	100	74	394.63	390	90.04
0.5	known	50	100	100	111	263	—	—
0.5	pilot study	50	100	100	111	269.91	264	62.58
0.5	surrogate	50	100	100	111	353.35	343.5	98.09
1	known	50	100	100	151	285	—	—
1	pilot study	50	100	100	151	335.69	332.5	69.93
1	surrogate	50	100	100	151	431.93	420	130.70
2	known	50	100	100	251	352	—	—
2	pilot study	50	100	100	251	495.52	491	92.61
2	surrogate	50	100	100	251	605.82	578.5	176.10

**Table A.2** Estimated power (POW) and concentration (OPT) under a model which violated the normality assumptions using the normality-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care.

$\Delta$	Method for $\sigma^*$	$n_0$	(POW)	(OPT)	$\hat{n}^{\text{fixed}}$	$\mathbb{E}\hat{n}$	Med( $\hat{n}$ )	SD $\hat{n}$
0	known	50	100	100	70	586	—	—
0	pilot study	50	79.0	86.4	70	153.17	115.5	145.73
0	surrogate	50	99.4	100	70	339.32	317.5	134.60
0.5	known	50	100	100	35	485	—	—
0.5	pilot study	50	90.2	85.0	35	132.75	90.5	145.46
0.5	surrogate	50	99.8	100	35	322.16	285.5	190.63
1	known	50	100	100	47	630	—	—
1	pilot study	50	95.6	88	47	186.41	142.5	159.20
1	surrogate	50	100	99.0	47	390.41	344.5	191.12
2	known	50	100	100	103	866	—	—
2	pilot study	50	100	99.6	103	296.34	241	211.56
2	surrogate	50	100	100	103	537.83	490.5	277.45

**Table A.3** Estimated power (POW) and concentration (OPT) under a model for which the normality assumptions hold using the projection-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care.

$\Delta$	$n_0$	(POW)	(OPT)	$\hat{n}^{\text{fixed}}$	$\mathbb{E}\hat{n}(\mathcal{D}_{n_0})$	Med $\{\hat{n}(\mathcal{D}_{n_0})\}$	SD $\hat{n}(\mathcal{D}_{n_0})$	$P\{\hat{n}(\mathcal{D}_0) = \infty\}$
0	50	84.58	100	74	387.72	346	107.78	0.48
0.5	50	97.54	100	111	304.47	285	82.12	0.24
1	50	100	100	151	314.61	299.5	71.15	0.08
2	50	100	100	251	366.88	364.5	64.13	0.03

**Table A.4** Estimated power (POW) and concentration (OPT) under a model which violated the normality assumptions using the projection-based sample size procedure at a nominal level of 90. To form a baseline for comparison,  $\hat{n}^{\text{fixed}}$ , shows the required sample size to compare the optimal embedded regime with standard of care.

$\Delta$	$n_0$	(POW)	(OPT)	$\hat{n}^{\text{fixed}}$	$\mathbb{E}\hat{n}(\mathcal{D}_{n_0})$	Med $\{\hat{n}(\mathcal{D}_{n_0})\}$	SD $\hat{n}(\mathcal{D}_{n_0})$	$P\{\hat{n}(\mathcal{D}_0) = \infty\}$
0	50	84.12	100	70	569.31	491.5	270.15	0.41
0.5	50	93.65	100	35	531.89	503	206.85	0.32
1	50	100	100	47	586.21	532.5	254.04	0.21
2	50	100	100	103	581.8	524	251.39	0.03

## APPENDIX

### B

# Q-LEARNING FOR SURVIVAL ANALYSIS

## B.1 Proofs of technical results

Proof that

$$\sum_{\kappa_i=K} \left[ \frac{\Delta_i}{\hat{\mathcal{K}}_K(U_i | \mathbf{H}_{K_i}, A_{K_i})} \frac{\partial Q_K^R(\mathbf{H}_{K_i}, A_{K_i}; \beta_K)}{\partial \beta_K} \{f(U_i) - f(\mathcal{T}_{K_i}) - Q_K^R(\mathbf{H}_{K_i}, A_{K_i}; \beta_K)\} \right] = 0.$$

is an unbiased estimating equation.

To show this we will show that

$$\mathbb{E} \left[ \frac{I(\kappa = K)\Delta}{\mathcal{K}_K(U | \mathbf{H}_K, A_K)} \frac{\partial Q_K^R(\mathbf{H}_K, A_K; \beta_K)}{\partial \beta_K} \{f(U) - f(\mathcal{T}_K) - Q_K^R(\mathbf{H}_K, A_K; \beta_K)\} \right] = 0.$$



We then have that

$$\begin{aligned}
& \mathbb{E} \left[ \frac{I(\kappa = K)\Delta}{\mathcal{K}_K(U|\mathbf{H}_K, A_K)} \frac{\partial Q_K^R(\mathbf{H}_K, A_K; \beta_K)}{\partial \beta_K} \{f(U) - f(\mathcal{T}_K) - Q_K^R(\mathbf{H}_K, A_K; \beta_K)\} \right] \\
&= \mathbb{E} \left( \mathbb{E} \left[ \frac{I(\kappa = K)\Delta}{\mathcal{K}_K(U|\mathbf{H}_K, A_K)} \frac{\partial Q_K^R(\mathbf{H}_K, A_K; \beta_K)}{\partial \beta_K} \{f(U) - f(\mathcal{T}_K) - Q_K^R(\mathbf{H}_K, A_K; \beta_K)\} \right. \right. \\
&\quad \left. \left. \middle| \mathbf{H}_K, A_K, \kappa = K \right] \right) \\
&= \mathbb{E} \left( I(\kappa = K) \frac{\partial Q_K^R(\mathbf{H}_K, A_K; \beta_K)}{\partial \beta_K} \left[ \mathbb{E} \left\{ \frac{f(U)\Delta}{\mathcal{K}_K(U|\mathbf{H}_K, A_K)} \middle| \mathbf{H}_K, A_K, \kappa = K \right\} \right. \right. \\
&\quad \left. \left. - \mathbb{E} \left\{ \frac{f(\mathcal{T}_K)\Delta + Q_K^R(\mathbf{H}_K, A_K; \beta_K)\Delta}{\mathcal{K}_K(U|\mathbf{H}_K, A_K)} \middle| \mathbf{H}_K, A_K, \kappa = K \right\} \right] \right).
\end{aligned}$$

Note that by definition

$$\mathbb{E} \left\{ \frac{f(U)\Delta}{\mathcal{K}_K(U|\mathbf{H}_K, A_K)} \middle| \mathbf{H}_K, A_K, \kappa = K \right\} = Q_K(\mathbf{H}_K, A_K)$$

and

$$\mathbb{E} \left\{ \frac{\Delta}{\mathcal{K}_K(U|\mathbf{H}_K, A_K)} \middle| \mathbf{H}_K, A_K, \kappa = K \right\} = 1.$$

Therefore the above becomes

$$\begin{aligned}
&= \mathbb{E} \left( I(\kappa = K) \frac{\partial Q_K^R(\mathbf{H}_K, A_K; \beta_K)}{\partial \beta_K} \{Q_K(\mathbf{H}_K, A_K) - f(\mathcal{T}_K) - Q_K^R(\mathbf{H}_K, A_K; \beta_K)\} \right) \\
&= 0.
\end{aligned}$$

## B.2 Simulation study parameter values

For the Weibull generative model the following parameter values were used for the case with a low level of censoring:

$$\begin{aligned}
 \mu_{1,0} &= (0, 0.25) & \mu_{1,1} &= (0.25, -1) \\
 \mu_{2,0} &= (0, 0.25, 0.25, 0.25) & \mu_{2,1} &= (0.25, 0.25, 0.25, -1) \\
 \xi_{1,0} &= (7, 1) & \xi_{1,1} &= (-3, 1) \\
 \xi_{2,0} &= (7, 1, 1, 1) & \xi_{2,1} &= (0, 1, 1, 1) \\
 \xi_{3,0} &= (7, 1, 1, 1, 1, 1) & \xi_{3,1} &= (-3, 1, 1, 1, 2, 2) \\
 \zeta_{1,0} &= (7, 1) & \zeta_{1,1} &= (0, 0) \\
 \zeta_{2,0} &= (8, 1, 1, 1) & \zeta_{2,1} &= (0, 0, 0, 0) \\
 \zeta_{3,0} &= (12, 1, 1, 1, 1, 1) & \zeta_{3,1} &= (0, 0, 0, 0, 0, 0) \\
 \omega_{1,0} &= (3, 1) & \omega_{1,1} &= (0, 0) \\
 \omega_{2,0} &= (3, 1, 1, 1) & \omega_{2,1} &= (0, 0, 0, 0) \\
 \omega_{3,0} &= (3, 1, 1, 1, 1, 1) & \omega_{3,1} &= (0, 0, 0, 0, 0, 0) \\
 \gamma_{C_t} &= 5 & \gamma_{T_t} &= 5 & \gamma_{\mathcal{T}_t} &= 10 \quad \forall t.
 \end{aligned}$$

To increase the degree of censoring, the intercept in  $\zeta_{t,0}$  for all  $t$  is decreased while everything else remains the same. For the medium degree of censoring  $\zeta_{t,0}$  changes to

$$\begin{aligned}
 \zeta_{1,0} &= (5, 1) \\
 \zeta_{2,0} &= (6, 1, 1, 1) \\
 \zeta_{3,0} &= (10, 1, 1, 1, 1, 1)
 \end{aligned}$$

and for the high degree of censoring  $\zeta_{t,0}$  changes to

$$\zeta_{1,0} = (4, 1)$$

$$\zeta_{2,0} = (5, 1, 1, 1)$$

$$\zeta_{3,0} = (9, 1, 1, 1, 1, 1).$$

For the piecewise hazard generative model the parameter values for the case of low level censoring is given by:

$$\mu_{1,0} = (0, 0.25) \quad \mu_{1,1} = (0.25, -1)$$

$$\mu_{2,0} = (0, 0.25, 0.25, 0.25) \quad \mu_{2,1} = (0.25, 0.25, 0.25, -1)$$

$$\xi_{1,0} = (1, 0.5) \quad \xi_{1,1} = (-1, 0.3)$$

$$\xi_{2,0} = (1, 0.5, 0.5, 0.5) \quad \xi_{2,1} = (0, 0.3, 0.3, 0.3)$$

$$\xi_{3,0} = (2, 0.3, 0.3, 0.3, 0.3, 0.3) \quad \xi_{3,1} = (-2, 0.5, 0.5, 0.5, 1, 1)$$

$$\zeta_{1,0} = (0.2, 0.2) \quad \zeta_{1,1} = (0, 0)$$

$$\zeta_{2,0} = (0.2, 0.2, 0.2, 0.2) \quad \zeta_{2,1} = (0, 0, 0, 0)$$

$$\zeta_{3,0} = (0.2, 0.2, 0.2, 0.2, 0.2, 0.2) \quad \zeta_{3,1} = (0, 0, 0, 0, 0, 0)$$

$$\omega_{1,0} = (9, 3) \quad \omega_{1,1} = (0, 0)$$

$$\omega_{2,0} = (9, 3, 3, 3) \quad \omega_{2,1} = (0, 0, 0, 0)$$

$$\omega_{3,0} = (9, 3, 3, 3, 3, 3) \quad \omega_{3,1} = (0, 0, 0, 0, 0, 0).$$

For the high level censoring case the values of  $\zeta_{t,0}$  are changed to

$$\zeta_{1,0} = (1, 0.2)$$

$$\zeta_{2,0} = (1, 0.2, 0.2, 0.2)$$

$$\zeta_{3,0} = (1, 0.2, 0.2, 0.2, 0.2, 0.2).$$