

ABSTRACT

DONG, LIN. Semiparametric Methods for Decision Making and Causal Effect Generalization. (Under the direction of Dr. Eric Laber and Dr. Shu Yang).

Semiparametric method is a broad term that encompasses many estimation and modeling approaches. In the simplest terms, all such approaches assume that some features of the data can be constructed parametrically, and the others can be modeled nonparametrically. We study its usage in the context of dynamic treatment regimes, causal effect generalization, and predictive modeling. Chapter 1 gives an overview of the research work. Chapter 2 is dedicated to the use of semiparametric method in dynamic treatment regimes, where we present an estimating equation framework to tackle the missing data issue in estimating optimal treatment regimes. We propose an augmented weighting estimating equation approach that applies to a broad class of estimators including both Q-learning and outcome weighted learning. We establish consistency and demonstrate its advantages in finite samples using simulation experiments and application to a schizophrenia study. In Chapter 3, we develop a new semiparametric framework to evaluate the average treatment effect of a target population integrating the complementary features of the randomized clinical trial and real world evidence studies. We establish related asymptotic results and derive the semiparametric efficiency bound under this setting. We propose an augmented calibration weighting estimator that achieves such bound when the nuisance models are correctly specified. We apply our proposed methods to estimate the effect of adjuvant chemotherapy in early-stage non-small-cell lung cancer where we utilize the data from trial CALGB 9633 and the National Cancer Database. In Chapter 4, we focus on the use of semiparametric methods in predictive modeling. We introduce a new class of neural network architecture to represent some of the widely-used semiparametric regression models, including additive models, single index and additive index models. A hypothesis testing based model ranking is proposed and validated by simulation experiments. These developments help fill the methodology gaps that remain in the emerging fields of precision medicine, data integration and supervised learning.

© Copyright 2019 by Lin Dong

All Rights Reserved

Semiparametric Methods for Decision Making and Causal Effect Generalization

by
Lin Dong

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. Eric Laber
Co-chair of Advisory Committee

Dr. Shu Yang
Co-chair of Advisory Committee

Dr. Rui Song

Dr. Yair Goldberg
External Member

Dr. Jamian (Krishna) Pacifici

DEDICATION

To my family and friends.

BIOGRAPHY

Lin Dong (Dǒng Lín in Chinese pinyin) was born and grew up in Beijing, China. She obtained her B.S. in Statistics and Operations Research from Hong Kong Baptist University in May, 2014. She then joined North Carolina State University to pursue her Ph.D. in Statistics. Lin will defend her dissertation in June 2019.

ACKNOWLEDGEMENTS

I would like to express my deepest appreciation to my advisors. I am deeply indebted to Dr. Eric Laber for his continuous support and guidance. He is a great advisor and leader, who has been providing me and Laber-Labs with valuable resources. I thank him for giving me the chances to tackle various tasks and challenges, which greatly broaden my horizon. I would like to extend my deepest gratitude to Dr. Shu Yang for her kindly guidance and patience that cannot be underestimated. She is a great scholar and is always prompt on answering my questions. The completion of my dissertation would not have been possible without the support of them.

I would like to extend my sincere thanks to my committee members, Dr. Rui Song, Dr. Yair Goldberg and Dr. Jamian Pacific for their advises and time devoted to serve in the committee. I must also thank Dr. Vijay Nair and Dr. Joel Vaughan at Wells Fargo for their constructive advises in the structured neural network project. I would like to thank Dr. Shannon Holloway for helping me polish the writing of my dissertation. I also wish to thank Dr. Sujit Ghosh for his guidance in a biosimilar project.

The department of statistics at North Carolina State University is a great place to conduct collaborative research. Many thanks also to the entire faculty and staffs in the department for offering such a friendly environment.

I owe my long overdue thanks to my professors in HKBU. My gratitude goes to Dr. Lixing Zhu, Dr. Man-Lai Tang and Dr. Xiaonan Wu for encouraging me to pursue a doctoral degree and their recommendations.

I want to express my sincere appreciation to my Laber-Labs lab-mates for the valuable discussions we had over the years. Special thanks to the friends I met during my graduate study and internships. We had a lot of good-times together.

Last but not the least, I thank my parents for their everlasting love and support and I thank Zhentao for his constant companion.

TABLE OF CONTENTS

LIST OF TABLES	vii
LIST OF FIGURES	ix
Chapter 1 Introduction	1
1.1 Semiparametric methods	1
1.2 Overview of research work	2
Chapter 2 Dynamic treatment regimes with missing data	5
2.1 Introduction	5
2.2 Setup and notation	7
2.3 Estimation with complete data	9
2.3.1 Q -learning with complete data	10
2.3.2 Outcome weighted learning with complete data	12
2.4 Estimation with incomplete data	18
2.4.1 Missingness mechanism	18
2.4.2 Inverse probability weighted estimating equations	19
2.4.3 Augmented inverse probability weighted estimating equations	20
2.5 Simulation and Data Application	22
2.5.1 Simulation Studies	22
2.5.2 CATIE Trial Analysis	24
2.6 Discussion	30
Chapter 3 Integrative analysis of randomized clinical trials with real world evidence studies	31
3.1 Introduction	31
3.2 Basic setup	34
3.2.1 Notation: causal effect and two data sources	34
3.2.2 Identification assumptions	35
3.2.3 Existing estimation methods	37
3.3 Calibration weighting estimators	38
3.4 Semiparametric efficient estimator when Y and A are available in RWE	41
3.4.1 Augmented calibration weighting estimator	41
3.4.2 Semiparametric models by the method of sieves	45
3.5 Simulation	48
3.6 Real data application	55
3.7 Concluding remarks	57
Chapter 4 Model building with structured neural networks	58
4.1 Introduction	58
4.2 Set-up and notations	60

4.3	Structured neural networks	61
4.3.1	Architecture of structured neural networks	61
4.3.2	Regularization	65
4.3.3	Computation	66
4.4	Prediction error based model ranking	66
4.4.1	Prediction error	66
4.4.2	Hypothesis test	67
4.5	Simulation studies	70
4.5.1	Comparative study of sNN and semiparametric regression models	70
4.5.2	Simulation studies on model ranking	71
4.6	Concluding remarks and future work	74
BIBLIOGRAPHY		81
APPENDICES		94
Appendix A	Appendix for Chapter 2	95
A.1	IPWCC estimating equations	95
A.2	AIPWCC estimating equations	96
A.3	Simulation settings	97
A.4	Proof of Theorem 1	97
A.5	Proof of Theorem 2	99
Appendix B	Appendix for Chapter 3	102
B.1	Additional treatment calibration	102
B.2	Proofs	104
B.2.1	Proof of Theorem 3	104
B.2.2	Proof of Theorem 4	111
B.2.3	Proof of Theorem 5	114
B.2.4	Proof of Theorem 5 and Theorem 6	116
B.3	Conditions for the sieves estimator	117
B.4	Additional simulation study	118
B.4.1	Comparison of the CW estimators	118
B.4.2	Increased sample sizes	119
Appendix C	Appendix for Chapter 4	125
C.1	sNN specifications	125
C.2	Ridge functions	126

LIST OF TABLES

Table 2.1	Cross-validated value estimates of the optimal regimes estimated using different methods.	29
Table 3.1	Simulation results for continuous outcome bias of point estimates, Monte Carlo variance, relative bias of bootstrap variance estimate, and the coverage rate of 95% Wald confidence interval.	51
Table 3.2	Simulation results for binary outcome: bias of point estimates, Monte Carlo variance, relative bias of bootstrap variance estimate, and the coverage rate of 95% Wald confidence interval.	53
Table 3.3	Number of patients by treatment of the CALGB 9633 trial sample and the NCDB sample.	55
Table 3.4	Covariate and outcome means comparison of the CALGB 9633 trial sample and the NCDB sample.	57
Table 3.5	Point estimate, standard error and 95% Wald confidence interval of the causal risk difference between adjuvant chemotherapy and observation based on the CALGB 9633 trial sample and the NCDB sample.	57
Table 4.1	Summary of structured neural networks and their model structures.	62
Table 4.2	Comparison on MSE (S.E.) of sNN with semiparametric regressions: (a). SIM-Net verses PPR with one ridge function; (b). AM-Net verses GAM; (c). AIM-Net verses PPR with three ridge functions.	78
Table 4.3	Summary of Case 1: point estimate of prediction risk, 95% coverage of the Wald-type confidence intervals for each model; mean of p-values, empirical type I error rate and power for each pairwise comparison w/o Holm adjustment; percentage of the estimated best set contains the true best model w/o Holm adjustment and the average length of the the estimated best set.	79
Table 4.4	Summary of Case 2: point estimate of prediction risk and 95% coverage of the Wald-type confidence interval; mean of p-value, empirical type I error rate and power w/o Holm adjustment; percentage of the estimated best set contains the true best model w/o Holm adjustment and the average length of the the estimated best set.	79
Table 4.5	Summary of Case 3: point estimate of prediction risk and 95% coverage of the Wald-type confidence interval; mean of p-value, empirical type I error rate and power w/o Holm adjustment; percentage of the estimated best set contains the true best model w/o Holm adjustment and the average length of the the estimated best set.	80
Table B.1	Simulation results for continuous outcome with population size $N = 500000$: bias of point estimates, Monte Carlo variance, relative bias of bootstrap variance estimate, and the coverage rate of 95% Wald confidence interval.	121

Table B.2 Simulation results for continuous outcome with population size $N = 500000$: bias of point estimates, Monte Carlo variance, relative bias of bootstrap variance estimate, and the coverage rate of 95% Wald confidence interval. . 123

LIST OF FIGURES

Figure 2.1	Relative value of Q-learning with IPWCC estimator, AIPWCC estimator and MI when the missingness model is correctly specified. The length of the corresponding vertical bar is the Monte Carlo standard deviation of the relative value estimates	25
Figure 2.2	Relative value of outcome weighted learning with IPWCC estimator, AIPWCC estimator and MI when the missingness model is correctly specified. The length of the corresponding vertical bar is the Monte Carlo standard deviation of the relative value estimates.	26
Figure 2.3	Relative value of Q-learning with IPWCC estimator, AIPWCC estimator and MI when the missingness model is misspecified. The length of the corresponding vertical bar is the Monte Carlo standard deviation of the relative value estimates.	27
Figure 2.4	Relative value of outcome weighted learning with IPWCC estimator, AIPWCC estimator and MI when the missingness model is misspecified. The length of the corresponding vertical bar is the Monte Carlo standard deviation of the relative value estimates.	28
Figure 2.5	Missingness pattern of CATIE study after cleaning.	29
Figure 3.1	Demonstration of the sampling and treatment assignment regimes for the RCT and RWE samples within the target population.	36
Figure 3.2	Boxplot of estimators for continuous outcome under four model specification scenarios.	52
Figure 3.3	Boxplot of estimators for binary outcome under four model specification scenarios.	54
Figure 4.1	Architecture of a subnetwork with 5 hidden layers.	62
Figure 4.2	Network architectures of (a) SIM-Net, (b) AM-Net, and (c) AIM-Net.	64
Figure 4.3	From top to bottom, partially ordered simple to complex structures.	65
Figure 4.4	Empirical distribution and fitted density of prediction risks based on the validation dataset of Case 1.	75
Figure 4.5	Empirical distribution and fitted density of prediction risks based on the validation dataset of Case 2.	76
Figure 4.6	Empirical distribution and fitted density of prediction risks based on the validation dataset of Case 3.	76
Figure 4.7	Illustration of the proposed visualization tool.	77
Figure B.1	Boxplot of estimators under four model specification scenario: $\hat{\Delta}^{CW0}$ is worsen than $\hat{\Delta}^{CW1}$	120
Figure B.2	Boxplot of estimators for continuous outcome with population size $N = 500000$ under four model specification scenarios.	122

Figure B.3	Boxplot of estimators for binary outcome with population size $N = 500000$ under four model specification scenarios.	124
Figure C.1	Ridge functions learned from AM with polynomial splines.	126
Figure C.2	Ridge functions learned from the AM-Net.	127

1.1 Semiparametric methods

Semiparametric method is a broad term that encompasses many estimation and modeling approaches. In the simplest terms, all such approaches assume that some features of the data can be constructed parametrically (with finite-dimensional parameters), and the others can be modeled nonparametrically (with infinite-dimensional parameters). Over decades of developments, semiparametric methods and theory have become a classic branch of statistics and been applied to a wide range of scientific questions, such as missing data and causal inference. The methods popularity is due in part to the inherent balance struck in the bias-variance trade-off. Specifically, semiparametric methods fall between fully parametric methods, which are subject to misspecification (bias), and fully nonparametric methods, which lose efficiency (variance). Though semiparametric theory is well-established, several important methodology developments remain. Particularly, we study the application of semiparametric theory in the context of estimating optimal treatment regimes with missing data, causal effect generalization, and predictive modeling. These areas are of importance in the emerging fields of precision medicine,

data integration and machine learning and are the focus of this dissertation.

1.2 Overview of research work

A summary of the chapters that follow is given here as an overview of the research work included in this dissertation.

Chapter 2 is dedicated to semiparametric methods in the context of missing data. We explore this scenario by application to the subject area of dynamic treatment regimes (DTRs), an area of growing importance in the age of precision medicine. Generally speaking, dynamic treatment regimes operationalize precision medicine as a sequence of decision rules, one per stage of clinical intervention, that map up-to-date patient information to a recommended intervention. Of primary importance in this field is the identification of an "optimal" treatment regime, i.e., one that maximizes a mean utility function when applied to the population of interest. Current methods for estimating an optimal treatment regime assume that the data is fully observed, which rarely occurs in practice. A common approach to overcome this shortcoming of available data is to use multiple imputation and pool estimators across imputed (complete) datasets. However, this approach requires estimating the joint distribution of patient trajectories, which can be high-dimensional, especially when there are multiple stages of intervention. To avoid the underlying modeling step required to obtain complete data, we propose a weighted estimating equation based approach to estimate optimal treatment regimes. Our approach applies to a broad class of estimators including Q-learning and a generalization of outcome weighted learning, which are among the most popular estimators of an optimal treatment regime. In addition, we establish consistency under mild regularity conditions and demonstrate the advantages of our proposed methods in finite samples using simulation experiments and application to a schizophrenia study.

Randomized clinical trials (RCTs) are designed to and are regarded as the gold standard to evaluate treatment effects; however, due to their restricted inclusion and exclusion criteria, the findings from RCTs often lack generalizability to some real-world population of interest, which

we refer to as the target population. Another data source, called real world evidence (RWE) studies, often include large samples that are representative of a target population, but such studies are often observational and subject to complex confounding. In Chapter 3, we leverage the complementing features RCT and RWE to estimate the average treatment effect of the target population. First, we propose a calibration weighting estimator that uses only covariate information from the RWE study. Because this estimator enforces the covariate balance between the RCT and RWE study, the generalizability of the trial-based estimator is improved. We further propose an augmented calibration weighting estimator that can be applied in the event that treatment and outcome information is also available from the RWE study. This estimator achieves a semiparametric efficiency bound that we derived under the identification and outcome mean function transportability assumptions when the nuisance models are correctly specified. To resolve the misspecification issue associated with parametric approaches, a data-adaptive nonparametric sieve method is provided as an alternative. The sieve method guarantees good approximation of the nuisance models. We establish related asymptotic results under mild regularity conditions. Finite sample performances of the proposed estimators are verified in simulation studies. Finally, we apply our proposed methods to estimate the effect of adjuvant chemotherapy in early-stage resected non-small-cell lung cancer, where we utilize the data from trial CALGB 9633 and a sample from the National Cancer Database.

Chapter 4 is inspired by the use of semiparametric regression models in predictive modeling. We switch gears away from DTRs and causal inference to a statistical learning setting with two competing objectives, predictive performance and model explainability. Artificial neural networks have attracted much attention because of their high predictive power and computational scalability to large datasets. However, a key hindrance of using neural networks in practice is their lack of interpretability. We propose a new class of neural networks called the structured neural networks (sNNs). In the sNN class, we design the architecture of the feedforward neural network to represent the model structures of some commonly used semiparametric regression

models, including but not limited to the generalized additive model, the single index model, and the additive index model. The proposed sNNs leverage semiparametric regression models and neural networks to achieve a better balance between model explainability and performance. A hypothesis testing based model ranking procedure, which can be used to aid the model selection in (but not limited to) this class, is proposed and validated by simulation experiments.

CHAPTER 2

DYNAMIC TREATMENT REGIMES WITH MISSING DATA

2.1 Introduction

Dynamic treatment regimes operationalize clinical decision making as a sequence of decision rules, one per stage of intervention, that map current patient information to a recommended intervention (Murphy, 2003; Robins, 2004). An optimal treatment regime maximizes the mean utility if applied to select interventions in the patient population of interest (for alternative definitions of optimality, see Kosorok & Moodie, 2015; Linn et al., 2017; Wang et al., 2018). Optimal treatment regimes have been estimated across a wide range of application areas including anticoagulation (Henderson et al., 2010; Barrett et al., 2014; Rich et al., 2014), cancer (Wang et al., 2012b; Xu et al., 2019), mental disorders (Nahum-Shani et al., 2012; Laber et al., 2014a; Zhang et al., 2017), and HIV (Laan & Petersen, 2007; Petersen et al., 2012; Young et al., 2011). In these and nearly all other biomedical application areas, the observed data are subject to missingness, which can include missing measurements, treatments, and outcomes (Shortreed et al., 2014; Kosorok & Moodie, 2015).

There is a large body of literature on estimation of optimal treatment regimes using complete

data. This body of research includes: approximate dynamic programming methods like Q- and A-learning (Murphy, 2003; Robins, 2004; Blatt et al., 2004; Murphy, 2005a; Moodie et al., 2007; Schulte et al., 2014) and its many variants (e.g., Zhao et al., 2009; Goldberg & Kosorok, 2012; Lu et al., 2013; Moodie et al., 2014; Tian et al., 2014; Laber et al., 2014b; Zhou & Kosorok, 2017; Jeng et al., 2018; Shi et al., 2018; Kosorok & Laber, 2019, and references therein) ; direct-search methods including outcome weighted learning (Orellana et al., 2010; Zhang et al., 2012a; Zhao et al., 2012; Zhang et al., 2013; Zhao et al., 2014; Zhao et al., 2015; Zhou et al., 2017; Athey & Wager, 2017; Zhang et al., 2017; Zhang & Zhang, 2018; Liu et al., 2018; Luckett et al., 2018); and model-based planning via g-computation (see Robins, 1997; Yu & Laan, 2002; Xu et al., 2016; Xu et al., 2019; Guan et al., 2018; Laber et al., 2018, and references therein). Because these methods require complete data, it is often necessary to employ methods to address missing data. A common approach is to apply multiple imputation to complete the data, compute a given estimator of an optimal regime on each of the imputed data sets, and then aggregate these estimators. e.g., by averaging or voting (Almirall et al., 2016; Lu et al., 2016; Ertefaie et al., 2016; Nahum-Shani et al., 2017; Kilbourne et al., 2018; Kidwell et al., 2018). Though estimation of an optimal treatment regime is often but one part of a suite of secondary analyses, the requirement to develop a complete dataset is convenient as it can be used for a variety of other analyses.

Despite its appealing features, multiple imputation can be problematic with complex longitudinal data arising in the context of sequential decision problems because the data can be high-dimensional and subjects are often missing large segments of data. Further, multiple imputation requires estimating the joint distribution of patient trajectories and constructing a high-quality imputation model in this context is difficult. If a misspecified model is used to impute large amounts of missing data, the estimators may be biased and inferences inaccurate (Seaman et al., 2012). We examine a class of augmented inverse probability weighted estimators of the optimal treatment regime that applies to approximate dynamic programming and to direct-search methods when the data have a monotone missingness pattern. The application of

standard arguments of semiparametric efficiency theory establishes a double robustness property for this class of estimators (e.g., Tsiatis, 2007). We show that augmented weighting performs favorably as compared to multiple imputation and to simple inverse probability weighting in simulation examples. These results suggest that investigators should give serious consideration to using weighting methods as an alternative to multiple imputation in the context of estimating optimal treatment regimes in practice.

The remainder of this chapter is organized as follows. In Section 2.2, we set notation and define an optimal treatment regime using potential outcomes. In Section 2.3, we describe a class of estimators of an optimal treatment regimes for use with complete data; this class includes Q-learning and outcome weighted learning as special cases. In Section 2.4, we derive an augmented inverse probability weighted estimator for the proposed class of estimators that applies under a monotone missingness pattern. In Section 2.5, we present simulation examples and an application to data from a sequential multiple assignment randomized trial on schizophrenia. A brief discussion of future work is given in Section 2.6.

2.2 Setup and notation

We consider longitudinal data arising from an observational study or a sequential multiple assignment randomized trial (SMART, Lavori & Dawson, 2004; Murphy, 2005b; Kidwell, 2014). The complete data are assumed to be of the form $(X_{1i}; A_{1i}; X_{2i}; A_{2i}; \dots; X_{Ti}; A_{Ti}; Y_i)_{i=1}^n$, which comprise n independent replicates of $(X_1; A_1; X_2; A_2; \dots; X_T; A_T; Y)$; where: T is the number of treatment stages, $X_1 \in \mathbb{R}^{p_1}$ is baseline patient information and $X_t \in \mathbb{R}^{p_t}$ is information collected during stage $(t - 1)$ for $t = 2; \dots; T$, $A_t \in \mathcal{A}_t = \{0, 1\}$ is the treatment assigned during stage $t = 1; \dots; T$, and $Y \in \mathcal{Y} = \mathbb{R}$ is the terminal outcome coded so that higher values are better. The restriction to binary treatments is not necessary for approximate dynamic programming methods; however, most variants of outcome weighted learning require binary treatments (exceptions include Laber & Zhao, 2015; Chen et al., 2016) so we impose this

restriction to allow for a simple and unified notation.

Define $H_1 = X_1$, and recursively define $H_t = (H_{t-1}; A_{t-1}; X_t)$ for $t = 2; \dots; T$. Thus, H_t is the information available to the decision maker at time $t = 1; \dots; T$. Let \mathcal{H}_t denote the support of H_t , and for each $h_t \in \mathcal{H}_t$, define $\mathcal{A}_t(h_t) \subseteq A_t$ to be the set of allowable treatments for a patient presenting with history $H_t = h_t$ at time t . A treatment regime in this context is a sequence of maps, $\tau = (\tau_1; \dots; \tau_T)$, with $\tau_t : \mathcal{H}_t \rightarrow \mathcal{A}_t$ and $\tau_t(h_t) \in \mathcal{A}_t(h_t)$ for all $h_t \in \mathcal{H}_t$ and $t = 1; \dots; T$. Under τ , a patient with history $H_t = h_t$ at time t would be recommended to receive treatment $\tau_t(h_t)$. Let \mathcal{R} denote the space of all feasible regimes. An optimal treatment regime, say $\tau^{\text{opt}} \in \mathcal{R}$, maximizes the mean outcome if applied to the population of interest. We formalize this definition using potential outcomes (Rubin, 1978; Neyman, 1923).

For each $t = 1; \dots; T$, define $\bar{a}_t = (a_1; \dots; a_t)$. Let $H_t(\bar{a}_{t-1})$ denote the potential history under treatment sequence \bar{a}_{t-1} , and let $Y(\bar{a}_T)$ denote the potential outcome under treatment sequence \bar{a}_T . Therefore, the set of all potential outcomes is

$$W = \{H_t(\bar{a}_{t-1}); Y(\bar{a}_T) : a_t \in \mathcal{A}_t(H_t(\bar{a}_{t-1})); t = 1; \dots; T\};$$

where we have defined $H_1(\bar{a}_0) = H_1$. Let $1_{(\cdot)}$ be the indicator function. The potential outcome under a regime $\tau \in \mathcal{R}$ is

$$Y(\tau) = \prod_{\bar{a}_T} Y(\bar{a}_T) \prod_{v=1}^T 1_{[\tau_v(H_v(\bar{a}_{v-1})) = a_v]}.$$

Define the value of regime τ to be $V(\tau) = EY(\tau)$, i.e., the marginal mean outcome if all subjects were assigned treatment according to τ . The optimal regime, $\tau^{\text{opt}} \in \mathcal{R}$, satisfies $V(\tau^{\text{opt}}) \geq V(\tau)$ for all $\tau \in \mathcal{R}$. To identify τ^{opt} in terms of the data-generating model, we make the following assumptions: (i) consistency, $H_t = H_t(\bar{A}_{t-1})$ for $t = 2; \dots; T$ and $Y = Y(\bar{A}_T)$, (ii) strong ignorability, $A_t \perp W | H_t$ for $t = 1; \dots; T$, and (iii) positivity, $P(A_t = a_t | H_t = h_t) > 0$ for all $a_t \in \mathcal{A}_t(h_t)$ and $t = 1; \dots; T$. These assumptions are standard in the dynamic treatment regimes

literature (see Robins, 2004; Chakraborty & Moodie, 2013; Schulte et al., 2014; Kosorok & Moodie, 2015, for additional discussion). Hereafter, we implicitly assume that these conditions hold.

2.3 Estimation with complete data

In this section, we review estimation of an optimal treatment regime when the data are completely observed. We consider a class of estimators that are representable as solutions to a set of estimating equations. This class is quite broad and includes most of the estimators commonly used in practice. To illustrate this point, we show in the Appendix that Q-learning and a generalization of outcome weighted learning belong to this class.

We consider treatment regimes of the form $\tau = f_1(\cdot; \theta_1); \dots; \tau_T(\cdot; \theta_T)g$ in which the decision rules composing the regime are indexed by parameters $\theta = (\theta_1^T; \theta_2^T; \dots; \theta_T^T)^T \in B$, where B is a normed linear space with norm $\|j\|_B$. For example, one might consider linear decision rules of the form $\tau_t(h_t; \theta_t) = \text{sign}(\theta_t^T h_{t,0})$, where $h_{t,0}$ is a feature vector constructed from h_t and $\text{sign}(u)$ is 1 if u is positive and -1 otherwise. We do not exclude the case in which θ_t includes nuisance parameters so that $\tau_t(\cdot; \theta_t)$ depends only on a subvector of θ_t ; however, we do not make any special distinction for such nuisance parameters as it is not important for our purposes. We assume that an estimator $b_n = (b_{1,n}^T; \dots; b_{T,n}^T)^T$ of θ is constructed by solving the estimating equation

$$P_n m_n(H_T; A_T; Y; \cdot) = 0 \quad (2.1)$$

over $\cdot \in B$, where P_n denotes the empirical measure, and $m_n : H_T \times A_T \times Y \rightarrow \mathbb{R}^J$. The dependence of m_n on n is to allow for regularization or other factors that may vary with the sample size (Qian & Murphy, 2011; Zhao et al., 2012; Zhao et al., 2015; Jeng et al., 2018).

Many estimators of an optimal treatment regime are based on backwards recursion. For these estimators, some components of m_n will depend on the partial history S_t , $(H_t^T; A_t)^T$ for

$t = 1; \dots; T$ rather than the complete data $S_{T+1}, (H_T^T; A_T; Y)^T$. For $t = 1; \dots; T$ define J_t to be the indices of m_n such that $m_{n;j}$ depends on S_{T+1} only through S_t , i.e.,

$$J_t = \{j \in J : m_{n;j}(h_T; a_T; y) = \tilde{m}_{n;j}(h_t; a_t) \text{ for some } \tilde{m}_{n;j} : H_t \times A_t \rightarrow \mathbb{R}\};$$

and J_{T+1} are the indices that rely on the complete data. Under this representation, the estimation equation in (2.1) can be equivalently expressed as

$$P_n \tilde{m}_{n;j}(S_t; \cdot) = 0 \quad \text{for all } j \in J_t; \quad t = 1; \dots; T + 1; \quad (2.2)$$

We will exploit this representation to use more of the observed data in constructing weighted complete case estimators in Section 2.4.

Let b_n denote the population analog of b_n , i.e., the solution to (2.2) with P_n replaced by P . We say that b_n is consistent if $\|b_n - b_n\|_B$ converges to zero in probability. Because our objective is not to propose new estimators of an optimal treatment regime, we will assume that the estimating equation has been suitably constructed to ensure consistency under the data-generating model in the complete data case and avoid stating specific conditions under which such consistency holds. Giving such general conditions would be cumbersome; for example, the conditions under which Q-learning with linear models is consistent are quite different from those under which kernel-based outcome weighted learning is consistent. Before describing how to adjust the estimating equations to accommodate missing data we first briefly recount how to express Q-learning and outcome weighted learning in the form given in (2.1).

2.3.1 Q-learning with complete data

Q-learning is an approximate dynamic programming algorithm that has been applied to estimate optimal treatment regimes in a variety of biomedical and engineering applications (Murphy, 2005a; Busoniu et al., 2010; Geramifard et al., 2013; Schulte et al., 2014; Sutton & Barto,

2018). The basis for Q-learning is the dynamic programming characterization of an optimal regime. Define $Q_T(h_T; a_T) = E(Y | H_T = h_T; A_T = a_T)$, and recursively for $t = T-1; T-2; \dots; 1$ define $Q_t(h_t; a_t) = E \max_{a_{t+1} \in \mathcal{A}_{t+1}(H_{t+1})} Q_{t+1}(H_{t+1}; a_{t+1}) | H_t = h_t; A_t = a_t$. It follows from dynamic programming that the optimal regime satisfies $a_t^{\text{opt}}(h_t) = \arg \max_{a_t \in \mathcal{A}_t(h_t)} Q_t(h_t; a_t)$ (Bellman, 1957). Let $\mathcal{Q}_t(h_t; a_t; \tau)$ denote a posited class of models for $Q_t(h_t; a_t)$ indexed by $\tau \in \mathcal{B}_t$ for $t = 1; \dots; T$. The induced class of treatment regimes is thus of the form $\tau(h_t; \tau) = \arg \max_{a_t \in \mathcal{A}_t(h_t)} Q_t(h_t; a_t; \tau)$ (Zhang et al., 2012b). If $\mathcal{B}_t = \mathbb{R}^{p_t}$ and $Q(h_t; a_t; \tau)$ is differentiable in τ for all $h_t; a_t \in \mathcal{H}_t \times \mathcal{A}_t$, then $b_{T;n}$ solves

$$P_n \{f(Y - Q_T(H_T; A_T; \tau))g_T - \tau Q_T(H_T; A_T; \tau)\} = 0; \quad (2.3)$$

and for $t = T-1; T-2; \dots; 1$ the estimators $b_{t;n}$ solve

$$P_n \max_{a_{t+1} \in \mathcal{A}_{t+1}(H_{t+1})} Q_{t+1}(H_{t+1}; a_{t+1}; b_{t+1;n}) - Q_t(H_t; A_t; \tau) - \tau Q_t(H_t; A_t; \tau) = 0; \quad (2.4)$$

Thus, the estimator b_n is obtained by solving $P_n m_n(\cdot) = 0$, where m_n is constructed by stacking (2.3) and (2.4) for $T-1; \dots; 1$ so that b_n is a root of

$$P_n \begin{pmatrix} f(Y - Q_T(H_T; A_T; \tau))g_T - \tau Q_T(H_T; A_T; \tau) \\ \max_{a_{T-1} \in \mathcal{A}_{T-1}(H_{T-1})} Q_{T-1}(H_{T-1}; a_{T-1}; \tau) - Q_{T-1}(H_{T-1}; A_{T-1}; \tau) - \tau Q_{T-1}(H_{T-1}; A_{T-1}; \tau) \\ \vdots \\ \max_{a_2 \in \mathcal{A}_2(H_2)} Q_2(H_2; a_2; \tau) - Q_1(H_1; A_1; \tau) - \tau Q_1(H_1; A_1; \tau) \end{pmatrix} = 0$$

The estimated optimal decision at stage t is thus $b_{n,t}(h_t) = \arg \max_{a_t \in \mathcal{A}_t(h_t)} Q_t(h_t; a_t; b_{t;n})$.

Similar expressions can be obtained for nonparametric variants of Q-learning (Zhao et al., 2009; Moodie et al., 2014; Zhang et al., 2017). For the purpose of illustration, we briefly describe kernel ridge regression for Q-learning. Suppose that $H_t \in \mathbb{R}^{p_t}$ for all $t = 1; \dots; T$. For each t , let $K_t : \mathbb{R}^{p_t} \times \mathbb{R}^{p_t} \rightarrow \mathbb{R}$ be symmetric and positive definite and write \mathcal{H}_t to denote the corresponding reproducing kernel Hilbert space with norm $\| \cdot \|_{\mathcal{H}_t}$ (for an introduction see Cristianini & Shawe-

Taylor, 2000; Moguerza & Muñoz, 2006; Berlinet & Thomas-Agnan, 2011; Nosedal-Sanchez et al., 2012). To approximate Q_T within H_T , one solves for each $a \in \mathcal{A}$;

$$\hat{Q}_{T;n}(\cdot; a) = \arg \min_{f \in H_T} P_n 1_{A_T=a} f Y - f_a(H_T)g^2 + \lambda_{T;a,n} \sum_{i \in I_{T;a}} f_{a,i}^2; \quad (2.5)$$

where $\lambda_{T;a,n} > 0$ is a tuning parameter. For each $a \in \mathcal{A}$ define $I_{T;a} = \{i : A_{T,i} = a\}$ to be the subset of patients to receive treatment a at time T and define $Z_{T;a}^T(h_T) = \sum_{i \in I_{T;a}} K_T(H_{T,i}; h_T)g_{i,2|I_{T;a}}$.

It follows that $\hat{Q}_{T;n}(h_T; a_T) = Z_{T;a_T}^T(h_T)b_{T;a,n}$, where $b_{T;a,n}$ is a solution of

$$P_n 1_{A_T=a_T} Y - (1 + \lambda_{T;a,n})Z_{T;a_T}^T(H_T) \lambda_{T;a,n}^{-1} Z_{T;a_T}^T(H_T) = 0:$$

Notice that in the previous estimating equation, we have replaced $\lambda_{T;a,n}$ by $\lambda_{T;a,n}^{-1}$ to reflect in re-writing the estimator the penalty has been scaled by the number of subjects receiving treatment a_T . Constructing $Z_{t;a_t}(h_t)$ analogously for $t = T-1; T-2; \dots; 1$, one can construct estimators $\hat{Q}_{t;n}(h_t; a_t) = Z_{t;a_t}^T(h_t)b_{t;a,n}$ where $b_{t;a,n}$ is a solution of

$$P_n 1_{A_t=a_t} \max_{a_{t+1} \in \mathcal{A}} \hat{Q}_{t+1;n}(H_{t+1}; a_{t+1}) - (1 + \lambda_{t;a,n})Z_{t;a_t}^T(H_t) \lambda_{t;a,n}^{-1} Z_{t;a_t}^T(H_t) = 0:$$

The preceding estimating equations can be stacked to obtain a single estimating equation; see Zhang et al. (2017) for additional details.

2.3.2 Outcome weighted learning with complete data

Direct-search methods, also known as value-search methods, estimate an optimal treatment regime by maximizing an estimator of the marginal mean outcome over a pre-specified class of regimes (Orellana et al., 2010; Zhang et al., 2012a; Rubin & Laan, 2012; Zhang et al., 2013; Zhang & Zhang, 2018). Outcome weighted learning (OWL) comprises a subclass of these methods, which rely on the use of a convex surrogate to carry out the proposed maximization (see below for details). Outcome weighted learning was introduced for single-stage decisions by Zhao et al.

(2012) and has since been generalized to multi-stage decisions (Zhao et al., 2015) and undergone a number of other modifications and refinements (Zhao et al., 2014; Chen et al., 2016; Zhou et al., 2017; Liu et al., 2018; Qi & Liu, 2018a; Qi et al., 2018b).

We consider a variant of outcome weighted learning that uses a convex relaxation of the augmented inverse probability weighted estimator of the marginal mean outcome (Zhang et al., 2013; Liu et al., 2018; Zhao et al., 2019a). We use a backwards recursive procedure (see Zhao et al., 2015) to extend the single-stage procedure proposed by Zhao et al. (2019a) to multiple-stages; while this is our not main methodological contribution, it may of be interest in its own right (see Zhang & Zhang, 2018; Davidian et al., 2019, for additional details).

We describe the estimator as a sequence of models t , each indexed by their own parameters, before stacking the models and parameters into a single estimating equation. To ease bookkeeping and clarify development, we begin by using $\beta_t = (\beta_{t1}; \dots; \beta_{tT})^T$ solely to index the decision rules and use $\alpha_t = (\alpha_{t1}; \dots; \alpha_{tT})^T$ to index nuisance models; we later pool these into a single collection of parameters, $\beta = (\beta_1; \dots; \beta_T)^T$, to match the notation used in our general framework. For simplicity, we assume linear decision rules of the form $t(h_t; \beta_t) = \text{sign}(\beta_t^T h_{t,0})$, where $h_{t,0}$ is known feature vector constructed from h_t . Furthermore, let $Q_T(h_T; \alpha_T; \beta_T)$ be a posited model for $Q_T(h_T; \alpha_T) = E(Y|H_T = h_T; A_T = \alpha_T)$. For each T , define $V_T(h_T; \beta_T) = Q_T f h_T; \beta_T(h_T; \beta_T)g$ with corresponding posited model $V_T(h_T; \beta_T; \beta_T) = Q_T f h_t; \beta_T(h_T; \beta_T); \beta_T g$. Define

$$Q_{T-1}(h_{T-1}; \alpha_{T-1}; \beta_{T-1}) = E[V_T(H_T; \beta_T; \beta_{T-1}) | H_{T-1} = h_{T-1}; A_{T-1} = \alpha_{T-1}];$$

thus $Q_{T-1}(h_{T-1}; \alpha_{T-1}; \beta_{T-1})$ is the mean outcome for a patient presenting with $H_{T-1} = h_{T-1}$, treated with $A_{T-1} = \alpha_{T-1}$, and subsequently treated according to $\beta_{T-1}(\cdot; \beta_{T-1})$ assuming that the model $Q_T(h_T; \alpha_T; \beta_T)$ is correct. Let $Q_{T-1}(h_{T-1}; \alpha_{T-1}; \beta_{T-1}; \beta_{T-1})$ be a posited model for $Q_{T-1}(h_{T-1}; \alpha_{T-1}; \beta_{T-1})$. Using an underbar to denote future, e.g., $\beta_t = (\beta_{t1}; \dots; \beta_{tT})^T$, define

$$V_{T-1}(h_{T-1}; \beta_{T-1}; \beta_{T-1}) = Q_{T-1} f h_{T-1}; \beta_{T-1}(h_{T-1}; \beta_{T-1}); \beta_{T-1}; \beta_{T-1} g$$

The second set of estimating equations are based on a backwards recursive representation of an augmented inverse probability weighted estimator of the incremental regret. As noted previously, the estimated Q-functions derived above serve as augmentation terms. Define $Q_T(H_T; A_T; \tau) = Q_T(H_T; A_T; \tau)$ and $b_{T;n}(H_T; A_T) = Q_T(H_T; A_T; b_{T;n})$. Define the estimated incremental regret at stage T as

$$\begin{aligned} J_{T;n}(\tau; \tau) &= P_n \mathbb{1}_{\text{signf } W_T(H_T; A_T; \tau) g_{A_T} \mathbb{1}_{H_T,0} > 0} W_T(H_T; A_T; Y; \tau) \\ &= P_n \mathbb{1}_{\text{signf } W_T(H_T; A_T; Y; \tau) g_{A_T} \mathbb{1}_{H_T,0} < 0} W_T(H_T; A_T; Y; \tau); \end{aligned} \quad (2.7)$$

where

$$\begin{aligned} W_T(H_T; A_T; Y; \tau) \\ = \frac{Y - Q_T(H_T; A_T; \tau) - \mathbb{1}_{P(A_T | H_T) = \tau} (H_T; A_T; \tau)}{P(A_T | H_T)}. \end{aligned}$$

Define $\mathcal{J}_{T;n}(\tau) = J_{T;n}(\tau; b_{T;n})$. It can be shown that $\mathcal{J}_{T;n}(\tau)$ is, up to an additive constant that does not depend on τ , a doubly robust estimator for the difference between: (i) the marginal mean outcome for a patient receiving treatment as per protocol (i.e., under the data-generating model) for the first $T - 1$ time points, followed by treatment under an optimal regime and (ii) the marginal mean outcome for a patient receiving treatment per protocol for the first $T - 1$ time points followed by treatment under $\tau(\cdot; \tau)$ (this is a variant of the so-called estimated regret, see Davidian et al., 2019, for detailed discussion and justification of this seemingly unusual estimand). Thus, $\mathcal{J}_{T;n}(\tau)$ is a measure of the loss incurred by treatment of patients at time T, who had theretofore been treated per protocol, with $\tau(\cdot; \tau)$ rather than an optimal treatment. Because of the indicator function, direct minimization of (2.7) over τ is a mixed integer program (Wolsey & Nemhauser, 2014), which is NP-hard in general and thus typically requires the use of either specialized software or the use of heuristics (Zhang et al., 2013; Zhang & Zhang, 2018). A

distinguishing feature of outcome weighted learning is the relaxation of (2.7) by replacing the indicator function with a convex function. Given a convex function $L : \mathbb{R} \rightarrow \mathbb{R}$, called a convex surrogate, backwards recursive outcome weighted learning uses the objective

$$J_{T;n}^L(\tau; \tau) = P_n L[\text{signf } W_T(H_T; A_T; Y; \tau) g A_T^T H_{T;0}] W_T(H_T; A_T; Y; \tau);$$

which can be seen to be a convex function of τ for each T . Define $J_{T;n}^L(\tau) = J_{T;n}^L(\tau; b_{T;n})$. We define $b_{T;n}$ to be the solution to $\text{arg } \min_{\tau} J_{T;n}^L(\tau) = 0$, where the gradient can be replaced with a sub-gradient if L is not differentiable (Boyd & Vandenberghe, 2004). Thus, the estimated optimal rule at stage T is $b_T(\cdot; b_{T;n})$.

Estimators of the optimal decision rules at stages $t = T-1; T-2; \dots; 1$ solve analogous estimating equations, which are defined recursively as follows. For $t = T-1$ define

$$\begin{aligned} & J_{T-1;n}^L(H_{T-1}; A_{T-1}; \tau; \tau) \\ &= Q_{T-1} f(H_{T-1}; A_{T-1}; \tau; \tau) g Q_{T-1} f(H_{T-1}; A_{T-1}; \tau; \tau) g; \end{aligned}$$

and define $b_{T-1;n}^o(H_{T-1}; A_{T-1}) = \text{arg } \min_{\tau} J_{T-1;n}^L(H_{T-1}; A_{T-1}; \tau; \tau)$, where $b_{T-1;n}^o(b_{T;n}) = b_{T-1;n}^o(b_{T;n}; b_T)$. Subsequently, define the relaxed loss function

$$\begin{aligned} & J_{T-1;n}^L(\tau; \tau) = \\ & P_n L[\text{signf } W_{T-1}(H_T; A_T; Y; \tau; \tau) g A_{T-1}^T H_{T-1;0}] W_{T-1}(H_T; A_T; Y; \tau; \tau); \end{aligned}$$

where the weights are given by

$$W_{T-1} f_{H_T; A_T; Y; T-1; T} g = \frac{1_{A_T = T}(H_T; T) Y}{P(A_T | H_T) P(A_{T-1} | H_{T-1})}$$

$$\frac{Q_{T-1} f_{H_{T-1}; A_{T-1}; T; T-1} g}{P(A_{T-1} | H_{T-1})} \left(\frac{1_{A_T = T}(H_T; T) P(A_T | H_T)}{P(A_T | H_T)} \right)^{Q_{T-1;n}(H_{T-1}; A_{T-1}; T; T-1)}$$

$$\frac{1_{A_T = T}(H_T; T)}{P(A_T | H_T) P(A_{T-1} | H_{T-1})} \frac{1_{A_T = T}(H_T; T) P(A_T | H_T)}{P(A_T | H_T)} Q_{T;n} f_{H_T; T}(H_T; T); T g$$

De ne

$$J_{T-1;n}^L(T-1) = J_{T-1}^T(T-1; b_{T;n}; b_{T-1;n}(b_{T;n})^o$$

and let $b_{T-1;n}$ be a solution to $r_{T-1} J_{T-1;n}^L(T-1) = 0$. For a generic $t < T-1$; de ne

$$f_t(H_t; A_t; t; t+1) = Q_t f_{H_t; A_t; t+1; t+1; t} Q_t f_{H_t; A_t; t+1; t+1; t}$$

and de ne $b_{t;n}(H_t; A_t) = \begin{pmatrix} H_t; A_t; b_{t;n}(b_{t+1;n}); b_{t+1;n} \end{pmatrix}^o$, where $b_{t;n}(b_{t+1;n}) = b_{t;n}(b_{t+1;n}); b_{t+1;n}(b_{t+2;n}); \dots; b_{T;n}$. Subsequently, de ne

$$J_t^L(t; t) = P_n L \text{ sign } W_t f_{H_T; A_T; Y; t+1; t} A_t^T H_{t;0} W_t f_{H_T; A_T; Y; t+1; t}$$

where the weights are given by

$$W_t f_{H_T; A_T; Y; t+1; t} = \frac{Y^{Q_T} \prod_{s=t+1}^{Q_T} 1_{A_s = s}(H_s; s)}{\prod_{s=t}^{Q_T} P(A_s | H_s)}$$

$$\frac{Q_t f_{H_t; A_t; t+1; t+1; t}}{P(A_t | H_t)} \left(\frac{1_{A_r = r}(H_r; r) P(A_r | H_r)}{P(A_r | H_r)} \right)^{Q_r f_{H_r; r}(H_r; r); r+1; r+1; r} \frac{P(A_t | H_t)}{P(A_t | H_t)}$$

$$X^T \left(\prod_{s=t+1}^{Q_r-1} \frac{1_{A_s = s}(H_s; s)}{P(A_s | H_s)} \right)^{Q_r-1} \frac{1_{A_r = r}(H_r; r) P(A_r | H_r)}{P(A_r | H_r)} Q_r f_{H_r; r}(H_r; r); r+1; r+1; r}^o \#$$

De ne $J_{t;n}^L(t) = J_t^T(t; b_{t+1;n}; b_{t;n}(b_{t+1;n})^o$, and let $b_{t;n}$ be a solution to $r_{t;n} J_{t;n}^L(t) = 0$.

The joint estimating equations for $\beta = (\beta^T; \beta^T)^T$ are obtained by stacking (2.6) and the

estimating equations $r_{t+1} J_t(\beta_t; \theta_t) = 0$ for $t = 1; \dots; T$.

2.4 Estimation with incomplete data

2.4.1 Missingness mechanism

We assume that baseline covariate information and initial treatment assignment, $(X_1; A_1)$, are always observed. This assumption generally holds in practice because patients who do not receive an initial treatment assignment are often unenrolled from the study and excluded from subsequent analyses. We further assume that the missingness pattern is nearly monotone; i.e., any item missingness that violates this monotone pattern is sparse, and the missingness pattern has been made monotone through artificial censoring or single imputation. Because patient dropout is the primary cause for missing data in longitudinal studies, e.g. SMARTS, this assumption is common in the literature (e.g., Shortreed et al., 2014, and references therein).

Let $C \in \{1; \dots; T + 1\}$ denote the dropout time so that $C = t$ if the patient dropped out after assignment of A_t for $t = 1; \dots; T$ and $C = T + 1$ if the patient's trajectory is fully observed, i.e.,

$$C = \begin{cases} T + 1; & \text{observe } (X_1; A_1; \dots; X_T; A_T; Y) = (H_T; A_T; Y); \\ T; & \text{observe } (X_1; A_1; \dots; X_T; A_T) = (H_T; A_T); \\ \vdots & \vdots \\ t; & \text{observe } (X_1; A_1; \dots; X_t; A_t) = (H_t; A_t); \\ \vdots & \vdots \\ 1; & \text{observe } (X_1; A_1) = (H_1; A_1); \end{cases}$$

We further assume that the data are missing at random (MAR, Rubin, 1976; Little & Rubin, 2014) so that $1_{C=t} \perp (X_{t+1}; A_{t+1}; \dots; X_T; A_T; Y) | H_t; A_t$ for all $t = 1; \dots; T$.

The simplest strategy for adapting the estimating equations presented in Section 2.3 for

missing data is through inverse probability weighting of complete cases. Note that 'complete case' for a term $m_{n,j}(S_t; \cdot)$, where $j \in J_t$, is the one in which S_t is observed and not necessarily the one for which the complete trajectory, S_{T+1} , is observed.

2.4.2 Inverse probability weighted estimating equations

Inverse probability weighted complete case (IPWCC) estimators re-weight the terms in the estimating equation for an optimal regime by their respective probabilities of being observed (Tsiatis, 2007; Tsiatis et al., 2014). Define the discrete hazard of dropout at time $t = 1, \dots, T$ to be $\lambda_t(s_t) = P(C = t | C \geq t; S_t = s_t)$. Thus, $\lambda_t(s_t)$ is the probability of dropping out at stage t for a patient with covariate and treatment history $S_t = s_t$. The survivor function at time t is thus $K_t(s_t) = P(C > t | S_t = s_t) = \prod_{v=1}^t (1 - \lambda_v(s_v))$. Under the MAR assumption, we can model the hazards using a binary regression model. For concreteness, we use a logistic regression model so that

$$\lambda_t(s_t; \beta) = \text{expit}(\beta' g_t(s_t; \beta));$$

where $\text{expit}(u) = \frac{\exp(u)}{1 + \exp(u)}$, $\beta \in \mathbb{R}^q$ is a vector of parameters, and $g_t(s_t; \beta)$ is continuously differentiable in β for all s_t . Define $\beta = (\beta_1^T; \dots; \beta_T^T)^T \in \mathbb{R}^q$; where $q = q_1 + \dots + q_T$. Define $\psi(C; S_t; \beta)$, $1_{C \geq t} \lambda_t(s_t; \beta) - 1_{C < t} \lambda_t(s_t; \beta) \text{expit}(\beta' g_t(s_t; \beta))$ to be the score function of the posited logistic regression model, and let $\hat{\beta}_{t,n}$ be a solution to $P_n \psi(C; S_t; \beta) = 0$. Let $\bar{\beta}_t = (\bar{\beta}_{t,1}^T; \bar{\beta}_{t,2}^T; \dots; \bar{\beta}_{t,T}^T)^T$ so that $\hat{\beta}_{t,n} = (\hat{\beta}_{t,1,n}^T; \hat{\beta}_{t,2,n}^T; \dots; \hat{\beta}_{t,n}^T)^T$. The estimated survivor function is $K_t(s_t; \hat{\beta}_{t,n}) = \prod_{v=1}^t (1 - \lambda_v(s_v; \hat{\beta}_{v,n}))$.

Define the complete case weights at level $t = 2, \dots, T+1$ and $j \in J_t$ under parameters β to be $w_j^{cc}(C; S_{t-1}; \bar{\beta}_{t-1}) = 1_{C > (t-1)} / K_{t-1}(S_{t-1}; \bar{\beta}_{t-1})$. The IPWCC estimator of an optimal treatment regime based on (2.2) solves

$$P_n w_j^{cc}(C; S_{t-1}; \hat{\beta}_{t-1,n}) m_{n,j}(S_t; \cdot) = 0; \quad j \in J_t; \quad t = 1, \dots, T+1;$$

with the understanding that $w_j^{cc}(c; s_0; \bar{0}) = 1$ for $j \in J_1$. The preceding equations along with those for $\bar{0}$ could be expressed as a single stacked estimating equation by concatenating the score equation for the logistic regression models onto m_n . Let $P_n m_n(S_{T+1}; \bar{0}) = 0$ denote this joint estimating equation. It follows from the derivations given in the Appendix that both the Q-learning and outcome weighted learning estimators can thus be constructed under a monotone missingness pattern using IPWCC by means of the preceding estimating equation. The following result can be used to establish consistency of the IPWCC estimator when combined with standard conditions for Z-estimators, e.g., the estimating equation has a unique isolated minimizer (see Kosorok, 2007).

Theorem 1 Assume that the survivor function is correctly specified so that $K_t(s_t) = K_t(s_t; \bar{t})$ for all t and s_t , for some $\bar{t} \in K$ and that the $b_n \rightarrow \bar{t}$ in probability. If m_n satisfies $\|P_n m_n(S_{T+1}; b_n)\| = o_p(1)$, then $\|P_n m_n(S_{T+1}; b_n)\| = o_p(1)$.

2.4.3 Augmented inverse probability weighted estimating equations

Using semi-parametric efficiency theory for monotone coarsening, of which monotone missingness is a special case, we derive an augmented inverse probability weighted complete case (AIPWCC) estimator (Robins & Rotnitzky, 1995; Tsiatis, 2007). Define for each $t = 1, \dots, T$, the conditional mean $d_{n,t}(s_t) = E[m_n(S_{T+1}; \bar{0}) | S_t = s_t]$ for which we posit parametric models $d_{n,t}(s_t; \bar{t})$ indexed by $\bar{t} \in R^{V_t}$. Let $\bar{t} = (\bar{t}_1, \dots, \bar{t}_T)^T$. An estimator of $d_{n,t}$ can thus be obtained by regressing $m_n(S_{T+1}; b_n)$ on $d_{n,t}(S_t; \bar{t})$ restricted to patients with $C = T + 1$, which gives $b_{t,n}$ and subsequently $\hat{d}_{n,t}(s_t) = d_{n,t}(s_t; b_{t,n})$ for each $t = 1, \dots, T$. Let $\bar{t}(s_t; \bar{t})$ and $K_t(s_t; \bar{t})$ be as defined in the previous section. For each $t = 2, \dots, T + 1$, $j \in J_t$, and $r = 1, \dots, t - 1$, define the augmentation weights

$$w_{r,j}^{aug}(c; s_r; \bar{r}) = \frac{1_{c=r} \bar{r}(s_r; \bar{r})}{K_r(s_r; \bar{r})}$$

The AIPWCC estimating equations are

$$P_n \sum_{j=1}^J w_j^{cc}(C; S_{t-1}; \mathbf{b}_{t-1;n}) \mathbb{E}_{n;j}(S_t; \cdot) + \sum_{r=1}^{J-1} w_{r;j}^{aug}(C; S_r; \mathbf{b}_{r;n}) d_{n;r;j}(S_r; \mathbf{b}_{r;n}) = 0; \quad \text{for all } j \in J_t; t = 1, \dots, T+1:$$

To obtain a single set of estimating equations one could concatenate the estimating equations for $\mathbb{E}_{n;j}(S_t; \cdot)$ and $d_{n;r;j}(S_r; \mathbf{b}_{r;n})$ to those given above. Let $P_n m_n(S_{T+1}; \cdot; \cdot) = 0$ denote the joint estimating equations. Explicit forms of these estimating equations for both Q-learning and outcome weighted learning are provided in the Appendix.

Theorem 2 Assume that the hazard functions for dropout are correctly specified so that $\mathbb{E}_{n;t}(s_t) = \mathbb{E}_{n;t}(s_t; \cdot)$ for all s_t for some $\epsilon_t \geq K_t$ and $b_n \rightarrow 0$ in probability or that the regression functions are correctly specified so that $d_{n;t}(s_t) = d_{n;t}(s_t; \cdot)$ for some $\epsilon_t \geq a_t$ and $b_n \rightarrow 0$ in probability. If b_n satisfies $\|P_n m_n(S_{T+1}; b_n)\| = o_p(1)$, then $\|P_n m_n(S_{T+1}; b_n; b_n)\| = o_p(1)$. Thus, the AIPWCC estimator is doubly robust.

Remark 1 When the dimension of the trajectory space is high, the solutions to the estimating equations may be unstable. In these cases, regularization may be necessary to stabilize the solution and to reduce overfitting. (see Fu, 2003; Johnson et al., 2008, for additional discussion of regularization with estimating equations). In our simulation experiments, we regularize the estimating equation using an adaptive ridge penalty; i.e., given a preliminary estimator b_n^0 of $\mathbb{E}_{n;j}(S_t; \cdot)$ based on the un-penalized estimating equation, we compute b_n^n as the solution to $P_n m_n(S_{T+1}; \cdot; b_n; b_n) + \lambda_n \mathbf{b}_n = 0$, where the division is taken elementwise and $\lambda_n \rightarrow 0$ is a tuning parameter.

2.5 Simulation and Data Application

2.5.1 Simulation Studies

We compare the performances of IPWCC, AIPWCC, and multiple imputation (MI) in terms of the value of the estimated optimal regime; these methods are applied with monotone missing data and the estimation is done using either Q-learning or outcome weighted learning. Data are simulated to mimic a two-stage SMART with binary treatments at each stage. The complete data are generated as follows:

$$X_1 = (X_{11}; \dots; X_{1p})^T; X_{1k} \sim \text{Bernoulli}(0.5); k = 1; \dots; p; \quad (2.8)$$

$$A_1 \sim \text{Uniform}(f_1; 1g);$$

$$X_2 | X_1 = x_1; A_1 = a_1 \sim N_p(\mu_0 + \mu_1 a_1) x_1; \Sigma_1;$$

$$A_2 \sim \text{Uniform}(f_1; 1g);$$

$$Y | X_1 = x_1; A_1 = a_1; X_2 = x_2; A_2 = a_2 \sim N(\gamma(x_1; a_1; x_2; a_2)); \Sigma_2 g;$$

where $\gamma(x_1; a_1; x_2; a_2) = \mu_0 + \mu_1 a_1 + \mu_2 x_1 a_1 + \mu_3 x_2 + (\mu_0 + \mu_1 a_1 + \mu_2 x_2) a_2$. Thus, the model is indexed by the matrices $\mu_0; \mu_1 \in \mathbb{R}^{p \times p}$, coefficients $\mu_0; \mu_1; \mu_2; \mu_3; \mu_0; \mu_1; \mu_2$, and variance components $\Sigma_1; \Sigma_2 > 0$.

For the missingness mechanism, we consider hazard functions of the form

$$P(C = t | C > t; H_t; A_t) = \lambda_t(H_t; A_t) = \expit(f(1; X_{t1}; A_t; X_{t2})^T \theta_t g); \quad t = 1; 2;$$

We vary the parameters θ_1 and θ_2 to obtain 35% and 65% missingness. The actual parameter values are provided in the Appendix. All simulation experiments use training sets of size $n = 1000$ and 500 Monte Carlo replications.

In our implementation of Q-learning, we use linear models of the form $seQ_t(H_t; A_t; \theta_t) =$

$\mathbb{1}^\top \mathbf{B}_{t;0}$, where $\mathbf{B}_{1;0} = (1; \mathbf{X}_1^\top; \mathbf{A}_1; \mathbf{X}_1^\top \mathbf{A}_1)^\top$ and $\mathbf{B}_{2;0} = (1; \mathbf{X}_1^\top; \mathbf{A}_1; \mathbf{X}_1^\top \mathbf{A}_1; \mathbf{A}_2; \mathbf{A}_1 \mathbf{A}_2; \mathbf{X}_2^\top \mathbf{A}_2)^\top$. For outcome weighted learning, we use linear decision rules of the form $\eta_t(\mathbf{H}_t; \mathbf{t}) = \mathbb{1}^\top \mathbf{H}_{t;0}$, where $\mathbf{H}_{1;0} = (1; \mathbf{X}_1^\top)^\top$ and $\mathbf{H}_{2;0} = (1; \mathbf{X}_1^\top; \mathbf{A}_1; \mathbf{X}_1^\top \mathbf{A}_1; \mathbf{X}_2^\top)$. These rules are estimated using logistic loss, also known as the entropy loss, as the convex surrogate (Zhao et al., 2019a; Jiang et al., 2019).

To illustrate the double-robustness property of the AIPWCC estimator, we consider both correctly and incorrectly specified models for the hazard functions. In the correctly specified case, we fit a logistic regression model at each stage with the correct features, i.e., $(\mathbf{X}_{t1}; \mathbf{A}_t \mathbf{X}_{t2})$ for $t = 1; 2$. For the incorrectly specified model, we fit a logistic regression model with features $(1; \mathbf{X}_1; \mathbf{X}_{11} \mathbf{X}_{12})$ at stage 1 and $(1; \mathbf{X}_{21}^2; \mathbf{X}_{22}^2)$ at stage 2. The conditional mean model is not correctly specified throughout.

We evaluate the performance of an estimated optimal regime \hat{b} , in terms of its relative value, which is defined as

$$RV(\hat{b}) = \frac{V(\hat{b}) - V(\pi_0)}{V(\hat{b}^{\text{opt}}) - V(\pi_0)};$$

where π_0 is the stochastic policy that assigns treatments at each stage using a fair-sided coin flip. The reason for using the relative value, for instance, instead of the raw value, is to allow for comparison across a range of generative models, e.g., different problem dimensions. The requisite values are estimated using Monte Carlo methods with 40,000 simulated patients (see Section A.6. of Schulte et al., 2014, for additional details on simulation-based value estimation).

We approximate the roots of the estimating equations using R package `leqslv` with multiple starts. We implement MI using R package `MICE` with default settings and 10 imputed datasets (Buuren & Groothuis-Oudshoorn, 2011). The conditional expectation model in the AIPWCC estimator is fitted using ridge regression and tuned using the 5-fold cross-validation estimator of the value under the optimal regime.

The results for the correctly specified hazard models with Q-learning and outcome weighted learning are displayed in Figures 2.1 and 2.2, respectively. It can be seen that both the IPWCC and AIPWCC estimators generally outperform MI when the dimension of the covariates at

each stage, p , is large. The results for the incorrectly specified hazard models estimated using Q-learning and outcome weighted learning are displayed in Figures 2.3 and 2.4, respectively. The results are qualitatively similar though the robustness of the AIPWCC shows improved performance relative to the IPWCC in some scenarios.

2.5.2 CATIE Trial Analysis

We use data from the CATIE schizophrenia study (Stroup et al., 2003) to illustrate the proposed methods. The CATIE study is a SMART, which enrolled 1460 schizophrenia patients. This dataset was chosen in part because it was used as an illustrative case study with MI by others (Shortreed et al., 2011; Shortreed & Moodie, 2012; Laber et al., 2014c; Shortreed et al., 2014).

As done elsewhere in the literature, we compare two treatments of primary clinical interest at each stage: Perphenazine (coded -1) and Olanzapine (coded 1). The dataset consisting of 506 patients receiving these treatments, 46% of whom followed the entire course (i.e., they are complete cases), 34% dropped out after stage 1, and 20% dropped out after stage 2. The missingness pattern is shown in Figure 2.5. Item missingness was sparse (less than 2% of the observed data) and singly imputed using mean imputation.

The Positive and Negative Syndrome Scale (PANSS) score is the standard medical scale for measuring symptom severity in schizophrenia. This score is a time-varying variable and was measured at each stage: baseline (PANSS0); stage 1 (PANSS1); and stage 2 (PANSS2). A higher PANSS score is associated with more severe symptoms, so we use $100 - \text{PANSS2}$ as the final outcome to match our convention of higher values representing better clinical outcomes. We include four baseline covariates in our analyses: PANNS0, baseline PANSS; EXACER, an indicator that the patient has been recently hospitalized; SEX; and TD, an indicator that the patient has Tardive Dyskinesia, a serious movement disorder associated with some antipsychotic medications. In addition, we include PANNS1, first stage PANSS, in our second stage models. As in the simulation study, we used linear models for the Q-functions of Q-learning and linear

Figure 2.1 Relative value of Q-learning with IPWCC estimator, AIPWCC estimator and MI when the missingness model is correctly specified. The length of the corresponding vertical bar is the Monte Carlo standard deviation of the relative value estimates

Figure 2.2 Relative value of outcome weighted learning with IPWCC estimator, AIPWCC estimator and MI when the missingness model is correctly specified. The length of the corresponding vertical bar is the Monte Carlo standard deviation of the relative value estimates.

Figure 2.3 Relative value of Q-learning with IPWCC estimator, AIPWCC estimator and MI when the missingness model is misspecified. The length of the corresponding vertical bar is the Monte Carlo standard deviation of the relative value estimates.

Figure 2.4 Relative value of outcome weighted learning with IPWCC estimator, AIPWCC estimator and MI when the missingness model is misspecified. The length of the corresponding vertical bar is the Monte Carlo standard deviation of the relative value estimates.

Table 2.1 Cross-validated value estimates of the optimal regimes estimated using different methods.

	Q-MI	Q-IPW	Q-AIPW	OWL-MI	OWL-IPW	OWL-AIPW
$\hat{v}(b^{\text{opt}})$	35.406	40.684	41.147	48.220	49.366	48.164

Figure 2.5 Missingness pattern of CATIE study after cleaning.

decision rules for outcome weighted learning.

We compared the following six approaches: Q-learning with MI (Q-MI), Q-learning with IPWCC (Q-IPW), Q-learning with AIPWCC (Q-AIPW), outcome weighted learning with MI (OWL-MI), outcome weighted learning with IPWCC (OWL-IPW) and outcome weighted learning with AIPWCC (OWL-AIPW). The cross-validated value for each approach is reported in Table 2.1. With the CATIE data, outcome weighted learning generally performed favorably to Q-learning. In terms of adjustment for missing data, MI appears to be worse than AIPWCC/IPWCC with Q-learning but about the same with outcome weighted learning.

2.6 Discussion

Missing data is essentially unavoidable with SMARTS and other longitudinal study designs commonly used to estimate optimal treatment regimes. Multiple Imputation has been shown to be an effective tool for accommodating missing data in such studies. However, in some settings, imputation can involve modeling complex processes, which may be prone to misspecification and high variance. We examined the use of inverse probability weighted methods and showed such methods are consistent for a broad class of estimators of an optimal treatment regime. Furthermore, in empirical experiments, these methods outperformed imputation with the gap in performance widening with the increasing trajectory dimension as well as with the increasing amount of missingness. Thus, we recommend that such weighting methods be given serious consideration by researchers estimating optimal treatment regimes from randomized or observational studies. Alternatively, it may be beneficial to forgo choosing between imputation and weighting and instead combine them (e.g., Shortreed et al., 2014). We leave such a hybrid approach to future work.

CHAPTER 3

INTEGRATIVE ANALYSIS OF RANDOMIZED CLINICAL TRIALS WITH REAL WORLD EVIDENCE STUDIES

3.1 Introduction

Randomized clinical trials (RCTs) are the gold standard for evaluation of treatment effects. However, due to restrictive inclusion and exclusion criteria for patient eligibility, the trial sample is narrowly defined and can be systematically different from the real-world patient population to which the new treatment is supposed to be given. Therefore, the findings from RCTs often lack generalizability to the target population of interest. Real world evidence (RWE) studies often include large samples that are representative of real-world patient populations; however, there are concerns about whether or not confounding has been addressed adequately in the analyses of RWE studies. In cancer research, there is an ongoing discussion on the strengths and limitations of utilizing data from RCT and RWE in comparative effectiveness analyses (Hahn & Schilsky, 2012; Hershman & Wright, 2012; Korn & Freidlin, 2012). The cancer research community has reached a consensus on the need to incorporate evidence from real world data, but how to effectively integrate them in pooled analysis of individual patients data from RCTs

remains under-developed.

There is considerable interest in bridging the findings from a RCT to the target population. This problem has been termed as generalizability (Cole & Stuart, 2010; Stuart et al., 2011; Hernan & VanderWeele, 2011; Tipton, 2013; O'Muircheartaigh & Hedges, 2014; Stuart et al., 2015; Keiding & Louis, 2016; Buchanan et al., 2018), external validity (Rothwell, 2005) or transportability (Pearl & Bareinboim, 2011; Rudolph & Laan, 2017) in the statistics literature and has connections to the covariate shift problem in machine learning (Sugiyama & Kawanabe, 2012). Most of the existing methods rely on direct modeling of the sampling score, which is the sampling analog of the propensity score (Rosenbaum & Rubin, 1983a). The subsequent sampling score adjustments include inverse probability of sampling weighting (IPSW, Cole & Stuart, 2010; Buchanan et al., 2018), stratification (Tipton, 2013; O'Muircheartaigh & Hedges, 2014), and augmented IPSW (Dahabreh et al., 2018). There are two major drawbacks in the sampling score adjustment approaches. Firstly, they require correct model specification of the sampling score. The IPSW estimators are unstable or inconsistent if the sampling score is too extreme or misspecified. Secondly, they assume the RWE sample to be a simple random sample from the target population, and implicitly require either the population size or all the baseline information of the population to be available. For example, Dahabreh et al. (2018) set the target population to be all trial-eligible individuals and assumed that all population baseline covariates are known, which is rarely the case in practice.

In this chapter, we consider combining a RCT sample and a RWE sample to estimate the average treatment effect (ATE) of the target population, where the RCT sample is subject to selection bias and the RWE sample is representative of the target population with a known sampling mechanism. It is worth noting that we allow the RWE sample to be a general random sample from the target population. This relaxation is particularly useful when real-world cohort studies are based on stratified random sampling in order to sample sufficient representations of some subgroups.

To address the selection bias of the RCT sample, we estimate the sampling score weights directly by calibrating covariate balance between the RCT sample and the design-weighted RWE sample, in contrast to the dominant approaches that focus on predicting sample selection probabilities. Calibration weighting (CW) is widely used to integrate auxiliary information in survey sampling (Wu & Sitter, 2001; Chen et al., 2002; Kott, 2006; Chang & Kott, 2008; Kim et al., 2016), and causal inference, such as in Constrained Empirical Likelihood (Qin & Zhang, 2007), Entropy Balancing (Hainmueller, 2012), Inverse Probability Tilting (Graham et al., 2012), and Covariate Balance Propensity Score (Imai & Ratkovic, 2014; Fan et al., 2016). Chan et al. (2015) showed that estimating ATE by empirical balancing calibration weighting can achieve global efficiency. To the best of our knowledge, our paper is the first to use calibration in combining a RCT sample and a RWE sample to construct ATE estimators. We show that calibration weighting has several advantages in the data integration problem. First, its estimation does not require the population baseline covariates or the population size to be available, and it allows the RWE sample to be a general random sample, not necessarily a simple random sample, from the target population. Second, the weights are estimated directly from an optimization problem instead of inverting the estimated sampling probability, which requires careful monitoring to avoid extreme weights that result in highly variable estimates. Third, similar to Zhao & Percival (2017), we show that the calibration weighting estimators are doubly robust in the sense that the estimators are consistent if the parameterization for either the outcome or the sampling score model is correctly specified.

The efficiency of the CW estimators can be further improved if we have additional treatment and outcome information from the RWE sample. We derive the first semiparametric efficiency bound for the ATE under identification assumptions and the outcome mean function transportability assumption. We propose an augmented CW estimator that is doubly robust and achieves the semiparametric efficiency bound if both nuisance models are correctly specified. In the presence of complex confounding, we adopt the method of sieves (Shen, 1997; Chen, 2007),

which allows flexible data-adaptive estimation of the nuisance functions while retaining usual root-n consistency under mild regularity conditions.

The rest of this chapter is organized as follows. We formalize the causal framework and assumptions in Section 3.2. The CW estimators are introduced in Section 3.3. We provide the semiparametric efficiency bound and propose the augmented CW estimator in Section 3.4. The finite sample performances of the proposed estimators are evaluated and compared in simulation studies in Section 3.5. We apply the proposed methods to estimate the effect of adjuvant chemotherapy by integrating the RCT data for early-stage resected non-small-cell lung cancer with the real-world data from the National Cancer Database in Section 3.6. Section 3.7 concludes with future directions.

3.2 Basic setup

3.2.1 Notation: causal effect and two data sources

We let X be the p -dimensional vector of covariates; let A be the treatment assignment with two levels $\{0, 1\}$, where 0 and 1 are the labels for control and active treatments, respectively; and let Y be the outcome of interest, which can be either continuous or binary. We use the potential outcomes framework (Neyman, 1923; Rubin, 1974) to formulate the causal problem. Under the Stable Unit Treatment Value Assumption (Rubin, 1980), for each level of treatment a , we assume that each subject in the target population has a potential outcome $Y(a)$, representing the outcome had the subject, possibly counterfactual, been given treatment a . The conditional average treatment effect (CATE) is defined as $\tau(X) = E\{Y(1) - Y(0) \mid X\}$. We are interested in estimating the population ATE $\tau_0 = E\{\tau(X)\}$, where the expectation is taken with respect to the target population. If the outcome is binary, the ATE is referred to as the average causal risk difference $\tau_0 = P\{Y(1) = 1\} - P\{Y(0) = 1\}$.

We consider a scenario where we have access to two samples: one sample is from a RCT

that compares the two treatments, and the other sample is from an observational RWE study that can be used to characterize the target population. The target population of size N consists of all patients with certain diseases to whom the new treatment is intended to be given. In practice, it is often difficult to identify the entire population, and the population size, N , is not necessarily known. Let $e = 1$ denote RCT participation, and let $e = 0$ denote the RWE study participation. The RWE sample is assumed to be a random sample drawn from the target population with a known sampling mechanism. Let $d = 1 = P(e = 1 | X)$ be the design weight in the RWE sample. For example, in health research, many cohort studies utilized stratified sampling to over-represent some subgroups. Suppose the population consists of L strata with sizes N_1, \dots, N_L , and suppose that the RWE study selects a fixed number of subjects n_l from the N_l subjects in the l th stratum. Then, the design weight for subject j from the l th stratum is $d_j = N_l/n_l$. We denote data from the RCT of size n to be $(X_i; A_i; Y_i; e_i = 1) : i = 1, \dots, n$; the data from the RWE study of size m to be either $(X_j; e_j = 1) : j = n+1, \dots, n+m$ if we only have covariate information, or $(X_j; A_j; Y_j; e_j = 1) : j = n+1, \dots, n+m$ if treatment and outcome information are also available in the study. The data structure is demonstrated in Figure 3.1.

3.2.2 Identification assumptions

A fundamental problem in causal inference is that we can observe at most one of the potential outcomes for an individual subject. To identify the ATE from the observed data, we make the following assumptions throughout this chapter.

Assumption 1 (Consistency) The observed outcome is the potential outcome under the actual received treatment: $Y = AY(1) + (1 - A)Y(0)$.

Assumption 2 (Randomization and Transportability) (i) $Y(a) \perp A_j | (X_j, e_j = 1)$ for $a = 0, 1$; and (ii) $E\{Y(1) - Y(0) | X_j, e_j = 1\} = \tau(X)$.

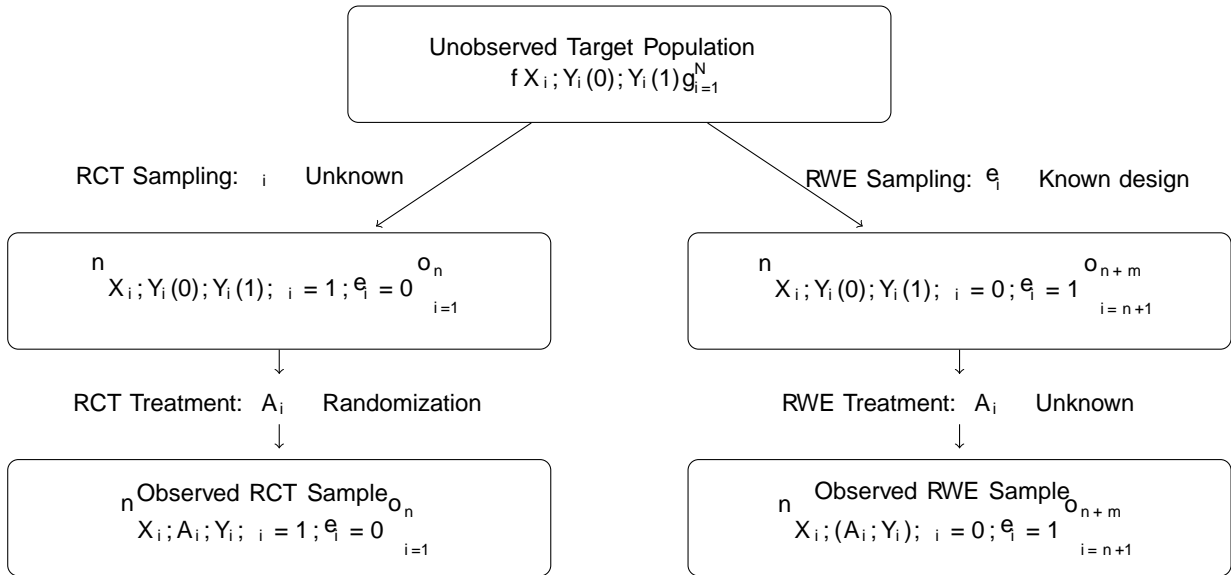


Figure 3.1 Demonstration of the sampling and treatment assignment regimes for the RCT and RWE samples within the target population.

Assumption 2 (i) holds for the RCT by default. Assumption 2 (ii) requires that the CATE function is transportable from the RCT to the target population. This assumption relaxes the ignorability assumption on trial participation (Stuart et al., 2011; Buchanan et al., 2018), i.e., $f(Y(0); Y(1))_{j=X}$, and the mean exchangeability assumption over treatment assignment and trial participation (Dahabreh et al., 2018), i.e. $E\{f(Y(a))_{j=X; \epsilon=1; A=a}\} = E\{f(Y(a))_{j=X; \epsilon=1}\}$ as well as $E\{f(Y(a))_{j=X; \epsilon=1}\} = E\{f(Y(a))_{j=X}\}$ for $a = 0; 1$.

Moreover, we require adequate overlap of the covariate distribution between the trial sample and the target population, and also the treatment groups over the trial sample, formalized by the following assumption. Define the sampling score as $\epsilon(X) = P(\epsilon = 1 | X)$.

Assumption 3 (Positivity) There exists a constant c such that with probability 1, $\epsilon(X) > c$; and $0 < P(A = 1 | X = x; \epsilon = 1) < 1$ for all x such that $P(X = x | \epsilon = 1) > 0$.

Under Assumptions 1-3, the ATE satisfies

$$\begin{aligned} \tau_0 &= E[f(X)]g = E[E\{f(Y(1) - Y(0) | X; A=1)\}] \\ &= E\left[\frac{f(X)}{h(X)}\{E(Y | X; A=1) - E(Y | X; A=0)\}\right] \\ &= E\left[\frac{1}{h(X)}\{E(Y | X; A=1) - E(Y | X; A=0)\}\right]; \end{aligned}$$

and thus is identifiable.

3.2.3 Existing estimation methods

Because the RCT assigns treatments randomly to the participants, $f(X)$ is identifiable and can be estimated by standard estimators solely from the RCT. However, the covariate distribution of the RCT sample $f(X | A=1)$ is different from that of the target population $f(X)$ in general; therefore, $E\{f(X) | A=1\}g$ is different from τ_0 , and the ATE estimator using trial data only is biased of τ_0 generally.

A widely-used approach is the IPSW estimator that predicts the sampling score $f(X)$ and uses the inverse of the estimated sampling score to account for the shift of the covariate distribution from the RCT sample to the target population. Specifically, most of the empirical literature assumes that $f(X)$ follows a logistic regression model $f(X; \beta)$ and can be estimated by $f(X; \hat{\beta})$. The IPSW estimator of the ATE is

$$\hat{\tau}^{\text{IPSW}} = \frac{\sum_{i=1}^n \frac{(X_i; \hat{\beta})^{-1} A_i Y_i}{(X_i; \hat{\beta})^{-1} A_i}}{\sum_{i=1}^n \frac{(X_i; \hat{\beta})^{-1} (1 - A_i) Y_i}{(X_i; \hat{\beta})^{-1} (1 - A_i)}}. \quad (3.1)$$

The IPSW estimator has several drawbacks as discussed in Section 3.1. In the next two sections, we propose approaches to (i) calibrate the covariate distribution of the trial sample to that of the design-weighted RWE sample so that the estimated treatment effects can be generalized to the target population, and (ii) leverage the predictive power of the RWE sample to improve the precision for the generalized ATE estimator.

3.3 Calibration weighting estimators

We propose to use calibration originated in survey sampling to eliminate the selection bias in the trial-based ATE estimator. The calibration weighting approach is similar to the idea of entropy balancing weights introduced by Hainmueller (2012). We calibrate subjects in the RCT sample so that after calibration, the covariate distribution of the RCT sample empirically matches the target population. Our insight is based on the observation that for any vector-valued function $g(X)$,

$$E \frac{g(X)}{p(X)} = E \sum_{i=1}^n \frac{d_i}{n} g(X_i) = E f g(X) g;$$

Here, $g(X)$ contains the covariate functions to be calibrated, which could be moment functions of the original covariate X or any sensible transformations of X .

To this end, we assign a weight q_i to each subject i in the RCT sample so that

$$\sum_{i=1}^N q_i g(X_i) = \mathbf{g}; \tag{3.2}$$

where $\mathbf{g} = \sum_{i=1}^N \frac{d_i}{n} g(X_i) = \sum_{i=1}^N \frac{d_i}{n} g_i$ is a design-weighted estimate of $E f g(X) g$ from the RWE sample. Constraint (3.2) is referred to as the balancing constraint, and weights $Q = \{q_i : i = 1, \dots, n\}$ are the calibration weights. The balancing constraint calibrates the covariate distribution of the RCT sample to the target population in terms of $g(X)$. The choice of $g(X)$ is important for both bias and variance considerations, which we will discuss in Section 3.4.2.

We estimate Q by solving the following optimization problem:

$$\min_Q \sum_{i=1}^n q_i \log q_i; \tag{3.3}$$

subject to $q_i \geq 0$; for all i ; $\sum_{i=1}^n q_i = 1$, and the balancing constraint (3.2).

The objective function in (3.3) is the entropy of the calibration weights; thus, minimizing

this criteria ensures that the empirical distribution of calibration weights are not too far away from the uniform, such that it minimizes the variability due to heterogeneous weights. This optimization problem can be solved using convex optimization with Lagrange multiplier. By introducing Lagrange multiplier θ , the objective function becomes

$$L(\mathbf{q}; \mathbf{Q}) = \sum_{i=1}^n q_i \log q_i + \theta \left(\sum_{i=1}^n q_i g(X_i) - \mathbf{1} \right) \quad (3.4)$$

Thus by minimizing (3.4), the estimated weights are

$$\mathbf{q} = \mathbf{q}(X_i; \mathbf{b}) = \frac{\exp(\mathbf{b}^T g(X_i))}{\sum_{i=1}^n \exp(\mathbf{b}^T g(X_i))} \mathbf{1};$$

and \mathbf{b} solves the equation

$$U(\mathbf{b}) = \sum_{i=1}^n \exp(\mathbf{b}^T g(X_i)) f g(X_i) - \theta \mathbf{1} = 0; \quad (3.5)$$

which is the dual problem to the optimization problem (3.3).

There are two general types of weighted estimators for population means, namely the Horvitz-Thompson estimator (HT, Horvitz & Thompson, 1952) and the Hajek estimator (Hajek, 1971). We propose both estimators with the superscript 0 for the HT estimator and the superscript 1 for the Hajek estimator. Let $A_i = P(A_i = 1 | X_i; i = 1)$ be the treatment propensity score for subject i . For RCTs, it is common that the propensity score is known with $A_i = 0.5$, for all $i = 1; \dots; n$.

Based on the calibration weights, we propose the HT CW estimator

$$\hat{\mu}^{CW0} = \sum_{i=1}^n \mathbf{q}_i \frac{A_i Y_i}{A_i} - \frac{(1 - A_i) Y_i}{1 - A_i}; \quad (3.6)$$

and the Hajek CW estimator

$$\hat{\tau}^{CW1} = \frac{\sum_{i=1}^n \mathbf{q}_i Y_i}{\sum_{i=1}^n \mathbf{q}_i} - \frac{\sum_{i=1}^n (1 - A_i) Y_i}{\sum_{i=1}^n (1 - A_i) \mathbf{q}_i} \quad (3.7)$$

To investigate the properties of the proposed CW estimators, we impose the following regularity conditions on the sampling designs for both the RWE and the RCT samples.

Assumption 4 Let $\tau_0 = E\{g(X)\}$. The design weighted estimator $\hat{\tau}_g = N^{-1} \sum_{i=1}^N e_i d_i g(X_i)$ satisfies $V(\hat{\tau}_g) = O(m^{-1})$, and $\sqrt{N}(\hat{\tau}_g - \tau_0) \xrightarrow{d} N(0, \tau_0)$ in distribution, as $m \rightarrow \infty$.

Assumption 5 The sampling score of RCT participation follows a loglinear model, i.e. $\pi(X) = \exp\{\beta_0 + g(X)\}$ for some β_0 .

Note that if the sampling score follows a logistic regression model $\pi(X; \beta) = \frac{\exp\{\beta_0 + g(X)\}}{1 + \exp\{\beta_0 + g(X)\}}$ and the sampling fraction of the RCT is small, the loglinear model in Assumption 5 is close to the logistic regression model. Our simulation studies demonstrate that the proposed estimators perform well under a logistic regression model.

In addition, in the estimation of calibration weights we only require specifying $\pi(X)$. Thus, CW evades explicitly modeling either the sampling score model or the outcome mean models. Under Assumption 5, we show that there is a direct correspondence between calibration weight $q(X_i; \hat{\beta})$ and the estimated sampling score $\pi(X_i; \hat{\beta})$, i.e. $q(X_i; \hat{\beta}) = \frac{\pi(X_i; \hat{\beta})}{N^{-1} + \pi(X_i; \hat{\beta})}$; see the proof of Theorem 3 in the Appendix.

The following assumption is on the linearity of the CATE in $g(X)$.

Assumption 6 $\tau(X) = \tau_0 + g(X)$.

Based on the above assumptions, we establish the double robustness property of the CW estimators in the following theorem and relegate all proofs to the Appendix.

Theorem 3 (Double robustness of the CW estimators) Under Assumptions 1-4, if either Assumption 5 or Assumption 6 holds, not necessarily both, $\hat{\mu}^{CW0}$ and $\hat{\mu}^{CW1}$ in (3.6) and (3.7) are consistent for μ_0 .

The Hajek-type estimator often improves the HT-type estimator empirically (Sarndal et al., 2003). However, this improvement is not insured in our context. We compare their asymptotic variances in the proof of Theorem 3 to show that either one could outperform the other in certain scenarios; see Section B.2.1 for details.

We provide a simple bootstrap procedure to estimate the variance of the CW estimators. The first step is to draw B bootstrap samples from both the RCT sample and the RWE sample respectively, which results in B pairs of bootstrapped samples. Then for each resampled pair, we can obtain a replicate of the CW estimator. The variance of the CW estimator is given by the sample variance of the B bootstrapped estimators.

3.4 Semiparametric efficient estimator when Y and A are available in RWE

3.4.1 Augmented calibration weighting estimator

We have utilized the design weighted covariate distribution of the RWE sample to adjust for selection bias of the RCT sample. Now we consider the setting where we have access to additional treatment and outcome information $(Y; A)$ from the RWE sample.

Define the trial conditional outcome mean function as $\mu_a(X) = E(Y | X; A = a; \tau = 1)$ for $a = 0; 1$. To leverage the predictive power of the RWE sample, we assume the transportability of $\mu_a(X)$ from the RCT sample to the RWE sample, as formulated below.

Assumption 7 For $a = 0; 1$, $E(Y | X; A = a; \tau = 1) = \mu_a(X)$.

Assumption 7 is plausible if X includes a rich set of covariates that influence the outcome, and is also testable because it is based only on the observed data.

The following theorem gives the semiparametric efficiency bound for θ_0 in our data integration setting.

Theorem 4 (Semiparametric efficiency bound) Under Assumptions 1 - 4 and 7, the semiparametric efficiency score for θ_0 is

$$S(X; A; Y; \theta; \theta_0) = \frac{A f_Y(1|X)g}{(X)} - \frac{(1-A) f_Y(0|X)g}{1-A} + \theta f(X) \theta_0 g$$

The semiparametric efficiency bound for θ_0 is

$$V_e = E \left[\frac{V f_Y(1|X)g}{A} + \frac{V f_Y(0|X)g}{1-A} + \theta^2 f(X) \theta_0^2 g^2 \right]$$

The semiparametric efficiency bound is decomposed naturally into two components as shown in its form. The first term accounts for the variance of potential outcomes adjusted for treatment assignment and sample selection; the second term quantifies the treatment heterogeneity in the RWE sample. Rudolph & Laan (2017) also provided a semiparametric efficiency score for transporting the ATE from one study site to another, where one site is regarded as a population without the sampling assumption.

The result in Theorem 4 serves as a foundation to derive efficient estimators combining two data sources. The score $S(X; A; Y; \theta; \theta_0)$ has unknown nuisance functions $f(X)$ and $g_a(X)$, ($a = 0; 1$). Therefore, to estimate θ_0 , we posit models for the nuisance functions, denoted by $f(X; \beta)$ and $g_a(X; \alpha_a)$. For example, we assume $f(X)$ is a loglinear model as in Assumption 5. By the correspondence between the loglinear model and the calibration weighing algorithm, we can estimate θ_0 following the optimization algorithm in (3.3). We also posit models $g_a(X; \alpha_a)$; $a = 0; 1$. By Assumption 7, we are able to obtain a consistent estimator $\hat{\beta}_a$ based on the RWE sample.

Based on the semiparametric efficiency score, we propose a new estimator for the ATE. As the

outcome mean models in the semiparametric efficiency score can be viewed as an augmentation to the CW estimator, we refer to the proposed estimator as the augmented calibration weighting (ACW) estimator, and it is given by

$$\hat{\tau}^{ACW} = \frac{\sum_{i=1}^n \frac{A_i Y_i - \tau_1(X_i; b_1)}{A_i} + \sum_{i=1}^n \frac{(1 - A_i) Y_i - \tau_0(X_i; b_0)}{1 - A_i}}{\sum_{i=1}^n \frac{A_i}{1 - A_i}} + \frac{\sum_{i=1}^n e_i d_i}{\sum_{i=1}^n d_i} \tau_1(X_i; b_1) - \tau_0(X_i; b_0) : (3.8)$$

We show in the following theorem that $\hat{\tau}^{ACW}$ achieves double robustness and local efficiency. For a vector v , we use $\|v\|_2 = (\sum v^2)^{1/2}$ to denote its Euclidean norm. For a function $f(V)$, where V is a generic random variable, we define its L_2 -norm as $\|f(V)\| = \int f(v)^2 dP(v)^{1/2}$.

Theorem 5 (Double robustness and local efficiency of the ACW estimator) Under Assumptions 1{4 and 7, if either Assumption 5 or Assumption 6 holds, not necessarily both, $\hat{\tau}^{ACW}$ is consistent for τ_0 . When both Assumption 5 and Assumption 6 hold $\sqrt{n}(\hat{\tau}^{ACW} - \tau_0) \rightarrow N(0; V_e)$ in distribution, as $n \rightarrow \infty$, where V_e is defined in Theorem 4, i.e. $\hat{\tau}^{ACW}$ is locally efficient with its asymptotic variance achieves the semiparametric efficiency bound.

By the empirical processes theory, the effect of nuisance parameter estimation in $\hat{\tau}^{ACW} - \tau_0$ is bounded by $\|k(X; b) - \tau(X)\| \sum_{a=0}^1 \|k_a(X; b_a) - \tau_a(X)\|$; see Section B.2.4 in the Appendix for details. If this bound is of rate $o_p(n^{-1/2})$, then it is asymptotically negligible, and thus $\hat{\tau}^{ACW}$ is semiparametric efficient. This rate can occur in various scenarios. For example, if all nuisance models are correctly specified parametric models, both $\|k(X; b) - \tau(X)\|$ and $\sum_{a=0}^1 \|k_a(X; b_a) - \tau_a(X)\|$ are $o_p(n^{-1/2})$ and their product is $o_p(n^{-1})$, then $\hat{\tau}^{ACW}$ is locally efficient. Another instance is to approximate the nuisance functions by flexible machine learning methods. In such case, if both $\|k(X; b) - \tau(X)\|$ and $\sum_{a=0}^1 \|k_a(X; b_a) - \tau_a(X)\|$ reach $o_p(n^{-1/4})$, then their product still achieves $o_p(n^{-1/2})$, and therefore is asymptotically negligible. In general, there exist different combinations of convergence rates of $\|k(X; b) - \tau(X)\|$ and $\sum_{a=0}^1 \|k_a(X; b_a) - \tau_a(X)\|$ ($a = 0, 1$)

that result in a negligible error bound accommodating different smoothness conditions of the underlying true nuisance functions. The following theorem formalizes the above statement.

Theorem 6 Let $(X; b)$ and $a(X; b_a)$ ($a = 0; 1$) be general semiparametric models for (X) and $a(X)$ ($a = 0; 1$), respectively. Assume the following regularity conditions hold:

Condition 1 $k(X; b) - (X)k = o_p(1)$, and $k(a(X; b_a) - a(X)k = o_p(1)$, for $a = 0; 1$;

Condition 2 $k(X; b) - (X)k \sum_{a=0}^1 k(a(X; b_a) - a(X)k = o_p(n^{-1/2})$.

Then $\hat{\theta}^{ACW}$ is consistent for θ_0 and achieves the semiparametric efficiency bound.

The semiparametric efficiency bound is attained as long as either b or $(b_0; b_1)$ approximate the underlying sampling score model or the outcome models well. Condition 1 states that we require that the posited models to be consistent. Condition 2 states that the combined rate of convergence of the posited models is $o_p(n^{-1/2})$. In Section 3.4.2, we construct such estimators using the method of sieves, which satisfies Condition 1 and 2 under mild regularity conditions.

For locally efficient estimator $\hat{\theta}^{ACW}$, the variance estimator can be calculated empirically as

$$\hat{v}^{ACW} = \frac{1}{n} \sum_{i=1}^n \frac{1}{(X_i; b)} \left(\frac{1}{A} \psi(Y_{ij}X_i; A_i = 1) + \frac{1}{1-A} \psi(Y_{ij}X_i; A_i = 0) \right) + \frac{1}{n} \sum_{i=1}^n \frac{1}{(X_i; b)} \left(\frac{1}{A} \psi(Y_{ij}X_i; A_i = 1) + \frac{1}{1-A} \psi(Y_{ij}X_i; A_i = 0) \right)^2 \quad (3.9)$$

The bootstrap variance estimator introduced in Section 3.3 can be applied here as well. The bootstrap variance estimator is more straightforward than (3.9), and it can accommodate situations where either one of the nuisance models is misspecified. Thus, the bootstrap variance estimator is recommended in practice and adopted in the simulation study.

3.4.2 Semiparametric models by the method of sieves

To overcome the model misspecification issue inherent to parametric models, we consider the method of sieves (Geman & Hwang, 1982), which allows flexible models for (X) and $\mu_a(X)$; ($a = 0; 1$) and also achieves Conditions 1 and 2 under mild regularity conditions. In comparison with other nonparametric methods such as the kernel method, the sieve method is particularly well-suited to calibration weighting. Although general sieve basis functions such as Fourier series, splines, wavelets, and artificial neural networks (see Chen, 2007, for a comprehensive review) are applicable, the power series is the most common class of functions. For a vector of non-negative integers $\mathbf{j} = (j_1; \dots; j_p)$, let $|\mathbf{j}| = \sum_{l=1}^p j_l$ and $X^{\mathbf{j}} = \prod_{l=1}^p X_l^{j_l}$. Define a series $f(\mathbf{j}) : \mathbf{j} = 1; 2; \dots; g$ for all distinct vectors of \mathbf{j} such that $|\mathbf{j}(\mathbf{k})| \leq |\mathbf{j}(\mathbf{k}+1)|$. Based on this series, we consider \mathbb{R}^g -vector $\mathbf{g}(X) = (g_1(X); \dots; g_g(X))$ and $\mu_a(X) = \sum_{\mathbf{j}} f(\mathbf{j}) X^{\mathbf{j}}$.

To accommodate different type of variables, we approximate (X) and $\mu_a(X)$ by the generalized sieves functions

$$\mu_a(X; \eta) = \sum_{\mathbf{j}} \eta_{\mathbf{j}} g_{\mathbf{j}}(X); \quad \mu_a(X; \eta) = m_a \left(\sum_{\mathbf{j}} \eta_{\mathbf{j}} g_{\mathbf{j}}(X) \right) \quad (a = 0; 1);$$

where for a continuous outcome, $m_a(\cdot)$ is the identity link function, and for a binary outcome, $m_a(\cdot)$ is the expit (inverse logit) link function with

$$\eta_{\mathbf{j}} = \arg \min_{\eta} E \left[\mu_a(X) - \sum_{\mathbf{j}} \eta_{\mathbf{j}} g_{\mathbf{j}}(X) \right]^2;$$

$$\eta_{\mathbf{j}} = \arg \min_{\eta} E \left[\mu_a(X) - m_a \left(\sum_{\mathbf{j}} \eta_{\mathbf{j}} g_{\mathbf{j}}(X) \right) \right]^2; \quad (a = 0; 1);$$

We assume that (X) and $\mu_a(X)$; ($a = 0; 1$) are sufficiently smooth, with (x) s -times continuously differentiable, and $\mu_a(x)$ s_a -times continuously differentiable for any x in the support of X with a condition of $\min(s_0; s_1) > 2p$. Under standard regularity conditions specified in Section B.3 in the Appendix, the deterministic differences between the true functions

and the sieves approximations are bounded (Newey, 1997)

$$\sup_{x \in X} |j(x) - \exp\{ \int_{\mathcal{X}} g(x) g_j \}| = O_p(K^{-1} n^{-2p}) ; \quad (3.10)$$

$$\sup_{x \in X} |j_a(x) - m_a \int_{\mathcal{X}} g(x) g_j| = O_p(K^{-1} n^{-2p}) ; \quad (a = 0; 1):$$

Thus, the approximation errors can be made sufficiently small by choosing a large K . However, it is also necessary that K increases slowly with the sample size n to control the variance of the sieves estimators. Formally, for the sieves estimators to be consistent, K should satisfy that $1/K + K/n \rightarrow 0$ as K and n grow (Chen, 2007). Therefore, in the presence of data with complex confounding, choosing K becomes important. In spite of the importance, there has limited study of the selection of K in sieves estimation. Imbens et al. (2005) proposed a selection method through minimizing the mean-squared-error (MSE) of the ATE over a pre-defined candidate set of K , which requires a rather complicated estimation of the population MSE and is not automatic.

To cope with this issue, we propose a new basis (variable) selection procedure for sieves estimation of $j(X)$ and $j_a(X)$. The number of basis functions controls the smoothness of sieves estimators. From this viewpoint, we can specify a sufficiently large K as an initial number and apply the penalization technique to regularize the variability of the estimators.

For penalized sieves estimation of $j(X)$, the key insight is that the dual problem of calibration leads to solving a system of estimating equations given by equation (3.5). Therefore, we adopt the penalized estimating equation approach (Johnson et al., 2008; Wang et al., 2012a) to facilitate basis selection. We consider solving the penalized estimating equations

$$U(\beta) = U(\beta) - \lambda \sum_{j=1}^K q(\beta_j) \text{sign}(\beta_j);$$

for $\beta = (\beta_1, \dots, \beta_K)^T$, where $q(\beta_j) = \lambda [q(\beta_{j-1}), \dots, q(\beta_{j+1})]^T$ is some continuous function, $q(\beta_j) \text{sign}(\beta_j)$ is the element-wise product of $q(\beta_j)$ and $\text{sign}(\beta_j)$. We let $q(x) = d p(x) = dx$,

where $p(x)$ is some penalization function. In this paper, we specify $p(x)$ to be a folded-concave SCAD penalty function (Fan & Li, 2001), although the same discussion applies to different penalty functions, such as adaptive LASSO (Zou, 2006). Accordingly, for the SCAD penalty, we have

$$q(j_{kj}) = \begin{cases} 0 & \text{if } |j_{kj}| < \lambda \\ \frac{(b - |j_{kj}|)_+}{(b - \lambda)} & \text{if } |j_{kj}| \geq \lambda \end{cases}; k = 1, \dots, K;$$

and we specify $b = 3.7$ following the suggestion of the literature.

Implicitly, we turn the problem of choosing the number of basis functions K to the problem of choosing the tuning parameter λ . The tuning parameter λ controls the magnitude of the regularization. To help understand the penalized estimating equation, we discuss two scenarios. If $|j_{kj}|$ is large, then $q(j_{kj})$ is zero, and the k th component of the estimating equation $U_k(\beta)$ is not penalized. Whereas, if $|j_{kj}|$ is small but not zero, then $q(j_{kj})$ is nonzero and $U_k(\beta)$ is penalized. Consequently, the penalty term forces β_k to be zero and excludes the k th element in $g(X)$ from the final selected set of variables.

In the simulation studies, we use 5-fold cross validation to select tuning parameter λ . The loss function is specified to be $l(\beta) = \frac{1}{n} \sum_{i=1}^n q(X_i; \beta) + \sum_{k=2}^K g(X_i; \beta)$, which quantifies the average degree of covariate balancing with calibration weights $q(X_i; \beta); i = 1, \dots, n$. Under standard regularity conditions, the estimator that solves the penalized estimating equation with the SCAD penalty would possess the oracle property (see Wang et al., 2012a, for details). Therefore, with $b = 3.7$, we can obtain $\|\beta\|_2 = O_p(n^{-1/2})$.

For penalized sieves estimation of $a(X)$, we can apply the standard penalization technique for regression models with the pre-specified basis functions based on the RWE sample. Specifically, let

$$\hat{b}_a = \arg \min_{\beta \in \mathbb{R}^K} \sum_{i=1}^n \epsilon_i |A_i - a| - \sum_{i=1}^n Y_i \beta + \sum_{j=1}^K p_a(j_{jj}) A_j;$$

where $p_a(\cdot)$ is the SCAD penalty function, for $a = 0, 1$. Under certain regularity conditions given in Fan & Li (2001), \hat{b}_a satisfies the selection consistency and oracle properties under

penalized likelihood for both linear regression or logistic regression, i.e. $\| \hat{\beta}_a - \beta_a \|_2 = O_p(n^{-1/2})$, for $a = 0; 1$.

Assuming $\min(s_0; s_1) = (2p) > 3=4$ in (3.10), the penalized sieves estimators of (X) and $\hat{\beta}_a(X)$ satisfies the two conditions in Theorem 6. Therefore, $\hat{\tau}^{ACW}$ with flexible approximation of the two nuisance functions still achieves rootn consistency and is semiparametric efficient.

3.5 Simulation

In this section, we evaluate the finite sample performances of the CW estimators and the ACW estimator via a set of simulation experiments. We first generate a target population of size $N = 50000$. Covariate $X \in \mathbb{R}^4$ is generated from $X_j \sim N(1; 1)$ for each $j = 1; \dots; 4$. For continuous outcomes, the potential outcome model is

$$Y(a)jX = 100 + 27:4aX_1 + 13:7X_2 + 13:7X_3 + 13:7X_4 + \epsilon;$$

where $\epsilon \sim N(0; 1)$ for $a = 0; 1$. Therefore, the true ATE under this setting is $\tau_0 = 27:4E(X_1) = 27:4$. For binary outcomes, we generate potential outcome according to

$$Y(a)jX \sim \text{Bernoulli}(p_a(X));$$

where

$$\text{logit}(p_a(X)) = 1 - 2aX_1 - X_2 - X_3 + X_4; \tag{3.11}$$

under which the average causal risk difference is $\tau_0 = 0:24$. We generate the indicator of selection into the RCT sample according to $I_j \sim \text{Bernoulli}(p(X))$; where $\text{logit}(p(X)) = 2:5 - 0:5X_1 - 0:3X_2 - 0:5X_3 - 0:4X_4$. By this design, the RCT selection rate is around 2%, which results in a roughly $n = 1000$ subjects in the RCT sample. The treatment assignment in the RCT sample is $A \sim \text{Bernoulli}(0:5)$. For the remaining subjects in the population, we take a random

sample of size $n = 5000$ to form a RWE sample. For patients in the RWE sample, treatment assignment is $A_j | X \sim \text{Bernoulli}(e_A(X))$; where $\text{logit}(e_A(X)) = X_1 + 0.5X_2 - 0.25X_3 - 0.1X_4$. The actual observed outcome Y in both samples is generated by $Y = AY(1) + (1 - A)Y(0)$.

To study the impact of model misspecification, following Kang & Schafer (2007), we define a nonlinear transformation of X to be

$$X = (\exp(X_1/3); X_2 = 1 + \exp(X_1)/10; X_1X_3 = 25 + 0.6; X_1 + X_4 + 20)^T;$$

and further scale and center X such that $E(X_j) = 1$ and $V(X_j) = 1$, for $j = 1, \dots, 4$. Throughout, we use X for fitting models. We assume X to be unobserved, but can be used in the true generative models, in which cases the fitted models are misspecified. We compare the following four model specification scenarios:

- ^ Scenario 1 (O:C/S:C): both outcome and sampling score models are correctly specified;
- ^ Scenario 2 (O:C/S:W): the outcome model is correctly specified; the sampling score model is incorrectly specified by using X in the generative model;
- ^ Scenario 3 (O:W/S:C): the outcome model is incorrectly specified by using X in the generative model; the sampling score model is correctly specified;
- ^ Scenario 4 (O:W/S:W): both outcome model and sampling score models are incorrectly specified by using X in the generative model.

To demonstrate the double robustness of the ACW estimator against parametric model misspecification, we consider the calibration variables $g_1(X) = (X_1; X_2; X_3; X_4)^T$ in all four scenarios. Moreover, we consider the ACW estimator using sieves estimation, with an initial large number of basis functions. Specifically, we extend the basis functions $\text{irg}_1(X)$ to its second order power series that includes all two-way interaction terms and quadratic terms, i.e. $g_2(X) = (X_1; \dots; X_p; X_1X_2; \dots; X_{p-1}X_p; X_1^2; \dots; X_p^2)^T$. For $p = 4$, $g_2(X)$ contains 14 basis functions.

We compare the following estimators for ATE:

1. Naive: the difference in sample means of the two treatment groups in the RCT sample to demonstrate the degree of selection bias;
2. IPSW: the inverse probability of sampling weighting estimator defined by (3.1), where the sampling weights estimated by logistic regression;
3. CW0: the HT calibration weighting estimator defined by (3.6) with $g(X) = g_1(X)$;
4. CW1: the Hajek calibration weighting estimator defined by (3.7) with $g(X) = g_1(X)$;
5. ACW: the augmented calibration weighting estimator defined by (3.8) with $g(X) = g_1(X)$;
6. ACW(S): the penalized augmented calibration weighting estimator using the method of sieves with $g(X) = g_2(X)$.

We use bootstrap variance estimation for all estimators with $B = 50$. All simulations are all based on 1000 Monte Carlo replications.

Table 3.1 and Figure 3.2 summarize the results for continuous outcome; Table 3.2 and Figure 3.3 summarize the results for binary outcome. It can be seen that the CW estimators are doubly robust compared to the IPSW estimator for continuous outcome. The CW estimators for binary outcome demonstrate some bias in Scenario 2. This phenomenon is because of Assumption 6 does not hold under the data generating mechanism (3.11). With additional information on $(A; Y)$ in the RWE sample, the ACW estimator is shown to be doubly robust and more efficient than the IPSW estimator and the CW estimators. In Scenario 4, where both outcome and sampling score models are misspecified, we show that the ACW estimator using the method of sieves is still unbiased and efficient. Moreover, the empirical coverage rates for the unbiased ACW estimator are close to the nominal level. As the sample size increases, which we show in the Appendix, they become closer to the nominal level.

Table 3.1 Simulation results for continuous outcome bias of point estimates, Monte Carlo variance, relative bias of bootstrap variance estimate, and the coverage rate of 95% Wald confidence interval.

	Naive	IPSW	CW0	CW1	ACW	ACW(S)
	Bias					
1. O:C/S:C	13:03	0:00	0:15	0:20	0:16	0:15
2. O:C/S:W	11:45	0:93	0:10	0:17	0:12	0:12
3. O:W/S:C	10:99	0:08	0:19	0:13	0:18	0:21
4. O:W/S:W	10:30	1:92	1:53	1:61	1:55	0:00
	Monte Carlo variance					
1. O:C/S:C	3:45	9:41	1:81	8:13	0:15	0:15
2. O:C/S:W	2:94	5:89	1:54	5:32	0:15	0:15
3. O:W/S:C	3:34	14:25	6:07	11:98	4:78	0:72
4. O:W/S:W	3:00	5:94	1:69	5:42	0:98	0:31
	Relative bias (%) of bootstrap variance estimate					
1. O:C/S:C	3:7	0:2	2:6	7:7	0:1	1:4
2. O:C/S:W	4:2	3:8	0:3	1:0	2:3	3:2
3. O:W/S:C	6:9	4:1	3:0	4:5	5:6	0:3
4. O:W/S:W	1:0	2:4	8:2	2:9	5:9	10:6
	95% Wald CI coverage rate					
1. O:C/S:C	0:0	93:9	94:1	93:2	94:1	93:6
2. O:C/S:W	0:0	93:4	95:6	95:5	94:3	95:7
3. O:W/S:C	0:0	94:7	94:3	93:6	94:3	92:8
4. O:W/S:W	0:0	85:3	80:1	87:4	59:5	92:1

Figure 3.2 Boxplot of estimators for continuous outcome under four model specification scenarios.

Table 3.2 Simulation results for binary outcome: bias of point estimates, Monte Carlo variance, relative bias of bootstrap variance estimate, and the coverage rate of 95% Wald confidence interval.

	Naive	IPSW	CW0	CW1	ACW	ACW(S)
	Bias 1000					
1. O:C/S:C	91:6	0:4	1:4	1:9	1:9	1:7
2. O:C/S:W	83:4	7:1	17:1	16:4	1:5	2:0
3. O:W/S:C	77:0	0:1	1:4	1:9	2:1	2:0
4. O:W/S:W	72:7	-25:1	29:8	29:1	19:8	4:4
	Monte Carlo variance 1000					
1. O:C/S:C	0:92	1:51	1:42	1:44	0:92	0:97
2. O:C/S:W	0:81	1:19	1:23	1:22	0:85	1:02
3. O:W/S:C	0:90	1:56	1:86	1:51	1:25	1:24
4. O:W/S:W	0:90	1:31	1:78	1:32	1:19	1:21
	Relative bias (%) of bootstrap variance estimate					
1. O:C/S:C	2:6	1:7	6:9	0:3	1:1	2:4
2. O:C/S:W	12:0	10:9	5:4	5:9	0:7	2:6
3. O:W/S:C	0:6	0:1	0:2	1:7	1:5	0:2
4. O:W/S:W	2:4	3:8	1:6	0:7	1:5	2:4
	95% Wald CI coverage rate					
1. O:C/S:C	15:4	94:0	95:1	93:2	94:5	94:6
2. O:C/S:W	20:6	94:5	92:6	90:8	94:6	93:1
3. O:W/S:C	26:6	93:8	93:8	93:3	94:8	94:5
4. O:W/S:W	31:7	89:6	87:3	85:9	89:3	95:3

Figure 3.3 Boxplot of estimators for binary outcome under four model specification scenarios.

Table 3.3 Number of patients by treatment of the CALGB 9633 trial sample and the NCDB sample.

	Observation (A = 0)	Adjuvant chemotherapy (A = 1)	Total
RCT: CALGB 9633	163	156	319
RWE: NCDB	10936	4271	15207

3.6 Real data application

We apply the proposed estimators to evaluate the effect of adjuvant chemotherapy for early-stage resected non-small cell lung cancer (NSCLC). Adjuvant chemotherapy for resected NSCLC is shown to be effective in stages II and IIIA disease on the basis of RCTs (Massarelli et al., 2003); however, its utility in the early-stage disease remains unclear. Cancer and Leukemia Group B (CALGB) 9633 is the only trial designed specifically for stage IB NSCLC (Strauss et al., 2008). Nonetheless, it consists of only 319 patients, which is undersized to detect clinically meaningful improvements. National Cancer Database (NCDB), as another data source, is a large clinical oncology database that contains information of more than 70% of the newly diagnosed cancer patients in the US with more than 34 million historical records.

The comparable sample from the NCDB includes 15207 patients diagnosed with NSCLC between years 2004 - 2016 with stage IB disease who first had surgery and then received either adjuvant chemotherapy or on observation (i.e. no chemotherapy) and with age greater than 20; these patients also did not receive any of the neoadjuvant chemotherapy, radiation therapy, induction therapy, immunotherapy, hormone therapy, transplant/endocrine procedures, or systemic treatment before their surgery. Thus, the patients in both samples are stage IB NSCLC patients who either received adjuvant chemotherapy or were on observation after their surgery without any other interventions.

In our analysis, the treatment indicator A is coded as 1 for adjuvant chemotherapy and 0 for on observation. A summary of the sample sizes in both samples by treatment groups is given in Table 3.3. We include four covariates in the analysis: X_1 is gender (1 = male, 0 = female); X_2 is

age; X_3 is the indicator for histology (1 = squamous, 0 = non-squamous); X_4 is the tumor size in cm. The outcome is the indicator of cancer recurrence within 3 years after the surgery, i.e. $Y = 1$ if recurrence occurred and $Y = 0$ otherwise. Table 3.4 reports the covariate and outcome means of the two samples. It can be seen that the patients in the CALGB 9633 trial are healthier compared to the population represented by the NCDB sample, i.e. the trial patients are younger and have smaller tumor sizes.

We compare the same six methods as in the simulation studies with the same set of basis functions. Bootstrap variance estimation is applied to estimate the standard errors. In the estimation we standardize age by its mean 67 and standard deviation 102, and tumor size by its mean 4.9 and standard deviation 3.0. Table 3.5 reports the results.

The results indicate that in the RCT sample there is a 8.3% decrease in the risk of recurrence for adjuvant chemotherapy over observation. The IPSW, CW, and ACW estimators, which utilized RWE sample information, show a 9%-11% decrease in the risk of recurrence. However, the causal effect is not significant according to the 95% confidence interval. The ACW estimator with the method of sieves gives an estimate of 17% risk decrease; the risk difference is significant at 0.05 level. The difference between the nonparametric ACW estimator and other adjusted estimators is because the former has second order polynomials selected and used for both calibration and outcome regression, which also indicates that the sampling score or/and outcome mean functions used in IPSW, CW, and ACW estimators are likely to be misspecified. The bootstrap variance estimates are not distinctive from one another, which is consistent with our simulation results that there is no clear improvement in efficiency for binary outcomes. All of the RWE-adjusted estimators have deeper decline in recurrence risk compared to the trial-based Naive estimator, which suggests that the causal risk difference in the target population is larger than the one of the RCT sample, i.e. the effect of adjuvant chemotherapy is more profound in the real world population.

Table 3.4 Covariate and outcome means comparison of the CALGB 9633 trial sample and the NCDB sample.

	Gender (X_1)	Age (X_2)	Histology(X_3)	Tumor size(X_4)	Recurrence (Y)
CALGB 9633	0:64	60:83	0:40	4:60	0:25
NCDB	0:55	67:87	0:39	4:94	0:33

Table 3.5 Point estimate, standard error and 95% Wald confidence interval of the causal risk difference between adjuvant chemotherapy and observation based on the CALGB 9633 trial sample and the NCDB sample.

	Est.	S.E	95% Wald confidence interval
Naive	0:083	0:048	(0:177; 0:011)
IPSW	0:088	0:060	(0:205; 0:029)
CW0	0:107	0:062	(0:227; 0:014)
CW1	0:105	0:061	(0:225; 0:014)
ACW	0:110	0:065	(0:236; 0:017)
ACW(S)	0:172	0:082	(0:333; 0:011)

3.7 Concluding remarks

In this chapter, we have developed a new semiparametric framework to evaluate the average treatment effects integrating the complementary features of the RCT and RWE study. There are several directions for future work: (i) we will extend this framework to the setting with survival outcomes. Additional challenges arise due to possible right censoring. We will follow the technique in Bai et al. (2013) to establish the corresponding semiparametric efficiency theory and estimators in the context of data integration; (ii) RCTs are often underpowered for the heterogeneity of treatment effect (HTE); on the contrary, RWE studies provide rich information on individual treatment responses (see Wu et al., 2018; Zhao et al., 2019b, for such developments in precision medicine). We will develop a semiparametric framework for integrative analyses of the HTE utilizing randomization in the RCT and power of treatment response prediction in the RWE study.

CHAPTER 4

MODEL BUILDING WITH STRUCTURED NEURAL NETWORKS

4.1 Introduction

Neural networks (NNs) have attracted much attention in the areas of machine learning and artificial intelligence because of their high predictive power and computational scalability to large datasets; however, often being referred to as black-box models, a key hindrance of using NNs in practice is their lack of interpretability. Model explainability/interpretability is crucial and is even mandated by law in some settings (e.g. EU law on explainable AI, Koszegi, 2019). The restriction of using machine learning models is especially reinforced in areas with strict regulations, such as the banking industry. The US Federal Reserve has been advocating the development of explainable machine learning tools for the use in financial services (Brainard, 2018).

To balance the two competing objectives, model performance and interpretability, it is favorable to use semiparametric regression models that contain both parametric and nonparametric components in their model structures. The parametric components help in model interpretation, and the nonparametric components add flexibility to the model so that it can capture the

underlying complex features of data. We survey a list of popular semiparametric regression models in econometrics and machine learning: the single index model (SIM, Ichimura, 1993; Hardle et al., 1993; Xia & Li, 1999); the (generalized) additive model (AM, Hastie & Tibshirani, 1987; Hastie, 2017) and its regularized versions (Lin & Zhang, 2006; Ravikumar et al., 2009); and the additive index model (AIM, Xia et al., 2002; Ruan & Yuan, 2010), which is also known as the projection pursuit regression (PPR, Friedman & Stuetzle, 1981). Because of their straightforward model structures, those models have similar level of interpretability as linear models.

In this chapter, we represent a continuum of models spanning linear model, AM, SIM and AIM using a new neural network architecture. We design the architecture of feedforward neural networks (FNN) to represent the model structures of the aforementioned semiparametric regression models. We name the proposed NN class as the structured neural network (sNN). The additive-index model based explainable neural network in Vaughan et al. (2018) is included as a special case in the sNN.

The sNN class is designed to leverage the NNs and classical semiparametric regression models to strike a better balance between model explainability and predictive performance. Though we essentially use the same model structures, there are several advantages of using the sNN representation and thus the NN training procedure over semiparametric regression methods. First of all, the universal approximation theorem of neural networks justifies its usage in functional approximation (Hornik et al., 1989). In addition, NN can be classified as nonlinear parametric models from a statistical view, which means that the gradient of its cost function can be easily computed through backpropagation (Rumelhart et al., 1985) and therefore gradient-based optimization algorithms can be directly applied. On the contrary, the existing estimation methods for semiparametric regression models mostly based on backfitting algorithm or its modifications, which are of higher cost computationally. The backfitting algorithm iterates between the estimation functional forms and the parametric components of the model. The functional forms are estimated by nonparametric smoothing operators, including smoothing

splines or kernel regression (Wahba, 1990; Gu, 2013), which are subject to curse of dimensionality and thus are computationally burdensome for large and high-dimensional datasets. Lastly, using standard deep learning software, the training process of the sNNs are unified and convenient, whereas the fitting methods for semiparametric regression models have many variants in the literature.

Within the sNN class, a remaining issue lies in how to determine the optimal model structure from the observed data. The existing selection criteria include information criteria (AIC, BIC and NIC Akaike, 1973; Akaike, 1977; Murata, 1991) and cross-validated errors (see Anders & Korn, 1999, for details). These methods do not have satisfactory performance in neural network models (Curry & Morgan, 2006). More importantly, they are not designed to suit in the context of balancing performance and explainability of the models. To this end, we propose a model ranking procedure that can be applied (but not limited) to compare the predictive performances of NN models. We provide a partial ordering on the structural complexity of models within the sNN class. Then by sample split, we construct hypothesis tests on the predictive errors that leads to valid comparative inference for models along this path.

The rest of this chapter is organized as the following. In Section 4.2 we set up the problem and notations. In Section 4.3 we propose the sNN model class. Section 4.4 presents a model ranking procedure based on hypothesis testing. Section 4.5 presents the simulation studies that compare the sNNs with semiparametric regression models and evaluate of the proposed model ranking procedure. Finally, Section 4.6 concludes and provides future directions.

4.2 Set-up and notations

We consider the predictive modeling or supervised learning setting. Denote $X \in \mathbb{R}^p$ as a p -dimensional covariate/predictors. We focus on a continuous outcome $Y \in \mathbb{R}$, though generalization to other outcome types is feasible. As in a typical supervised learning setting, the data observed are $\{(Y_i; X_i)\}_{i=1}^n$; and $(Y_i; X_i) \stackrel{i.i.d.}{\sim} P_0$, where P_0 is the underlying joint distribution of $(Y_i; X_i)$.

Particularly, we consider the setting of large sample size with moderate number of predictors (< 30), which is common in credit models or credit scoring models in the banking industry. We assume that the conditional expectation of outcome satisfies $E(Y|X = x) = \mu(x)$.

The goal is to model $\mu(x)$ such that the estimated model can be used to predict future observation with high accuracy. In the next section, we propose a novel class of NN architectures in Section 4.3 to approximate the unknown function $\mu(x)$ that combines the predictive power of NN and the interpretability of semiparametric regression models.

4.3 Structured neural networks

4.3.1 Architecture of structured neural networks

To start with, we briefly review the model structures of some commonly used semiparametric regression models. For the unified notation, we refer to the linear coefficients as indices and refer to the unknown functions as ridge functions in all the following models. The single index model has the form of $\mu(x) = f(\beta^T x)$. It has one index, β , and one ridge function $f(\cdot)$. The additive model has the form $\mu(x) = \sum_{j=1}^p \alpha_j f_j(x_j)$. It assumes that each predictor X_j has its own ridge function $f_j(\cdot)$ and is additive to each other. The additive index model, which generalizes both the single index model and the additive model, has the form $\mu(x) = \sum_{k=1}^q \alpha_k f_k(\beta_k^T x)$. It has q ridge functions, additive to each other, and within each there is linear combination of the predictors.

To take advantage of the explainability of semiparametric regression models inherited to their structures, we design the architecture of the FNNs to represent these models. We refer to such FNN models as the structured neural networks (sNN) to highlight the special design in their structures. We list the semiparametric models and their sNN counterparts in Table 4.1. Linear model is also listed as the simplest model form and a special case.

The key idea in sNN is to learn both the indices and the ridge functions by NNs. The

Table 4.1 Summary of structured neural networks and their model structures.

(Semi)parametric Regression Model	sNN	Model Structure
Linear Model		$(x) = \mathbf{w}^T \mathbf{x}$
Single Index Model	SIM- Net	$(x) = f(\mathbf{w}^T \mathbf{x})$
Additive Model	AM-Net	$(x) = \sum_{j=1}^p c_j f_j(x_j)$
Additive Index Model	AIM-Net	$(x) = \sum_{k=1}^q c_k f_k(\mathbf{w}_k^T \mathbf{x})$

Figure 4.1 Architecture of a subnetwork with 5 hidden layers.

indices are represented by the weights of a projection layer. The ridge functions are learned by subnetworks. A subnetwork serves as a nonlinear function approximator. The structure of a layer subnetwork is a fully-connected FNN that takes a single neural input and a single neural output. In between the first and the last hidden layer, the subnetwork consists of 2 fully-connected layers of arbitrary number of hidden units in each layer. We denote the architecture of a layer fully-connected FNN by $[n_1; n_2; \dots; n_l]$. For example, $[30; 10; 5]$ is a 3-layer fully-connected FNN with 30, 10 and 5 hidden units within each layer, respectively. Figure 4.1 shows a 5-layer subnetwork with architecture $[1; 5; 5; 3; 1]$.

In spirit, subnetwork can be related to the modular neural networks first appeared in Jacobs et al. (1991), where the network's architecture is designed to contain multiple split neural networks. Each split neural network serves as an expert of a sub-task (see Chapter 16, Rojas, 2013, for a comprehensive review). Watanabe et al. (2018) takes the other way by first training

the FNN and then decomposing it into multiple small networks to gain insights of FNN. In either case, the idea is to explain the model behavior by capsule the structures into several divided pieces.

To be more specific, Figure 4.2 illustrates the architectures of the three sNN classes in Table 4.1, i.e. SIM-Net, AM-Net and AIM-Net. The predictors first enter a projection layer, which then produces univariate inputs to be fed into the subnetwork(s). The weights of the projection layer can be either trained or assigned, depending on the class of models. For example, AM-Net has the projection layer weights fixed to map each predictor to a separate subnetwork. The subnetwork learns a nonlinear functional form. The final output of the sNN is produced from a combination layer, which could be a direct output from the subnetwork (as in SIM-Net), or linear combination of the outputs of several subnetworks.

We also give a unified representation of the sNNs as

$$\underset{\text{input}}{X} \rightarrow \underset{\text{projection layer}}{BX} \rightarrow \underset{\text{subnetworks}}{g(BX)} \rightarrow \underset{\text{combination layer}}{c^T g(BX)} \rightarrow \underset{\text{output}}{\hat{Y}}; \quad (4.1)$$

where $B \in \mathbb{R}^{q \times p}$ is the weight matrix of the projection layer that maps X to a lower dimension: $\mathbb{R}^p \rightarrow \mathbb{R}^q$ where in general $q < p$. g is a vector of nonlinear transformation(s) learned by the subnetwork(s), and each applied to one element of $BX \in \mathbb{R}^q$. $c \in \mathbb{R}^q$ is the weight vector of the combination layer that combines the output of the learned subnetworks. Explicit forms of each sNN class under this representation is given in Appendix C.1.

Given the structural design, sNN combines the strengths of NN and semiparametric regression models. We show that the sNNs and the semiparametric regression model have similar predictive performance in Section 4.5.1.

We note that there is a partial ordering in terms of structural complexity among these models, as shown in Figure 4.3. From top to bottom, the models become more flexible but also more complex in structure. Specifically, linear model is the simplest model; AM and SIM are two generalizations of linear model; AIM model generalizes both SIM and AM and Fully-connected

(a)

(b)

(c)

Figure 4.2 Network architectures of (a) SIM-Net, (b) AM-Net, and (c) AIM-Net.

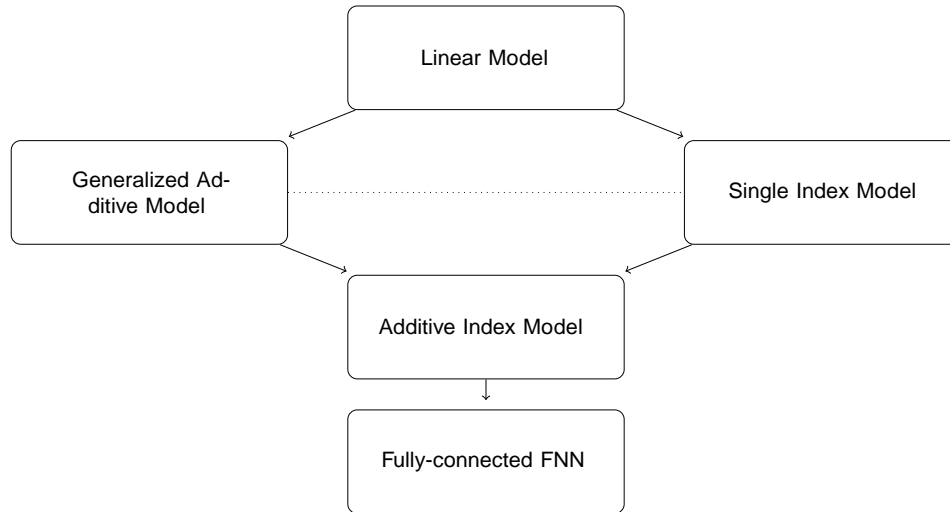


Figure 4.3 From top to bottom, partially ordered simple to complex structures.

FNN is the most flexible model under consideration. The partially ordered hierarchy helps in model selection, which we discuss in Section 4.4.

4.3.2 Regularization

Though we do not focus on high-dimensional X , we point out that penalty can be added to the projection layer and the combination layer to achieve variable/ridge function selection. Based on representation (4.1), the regularized sNN can be written as

$$\underset{\text{input}}{X} \rightarrow \underset{\text{projection layer}}{BX + p(B)} \rightarrow \underset{\text{subnetworks}}{B^T X} \rightarrow \underset{\text{combination layer}}{c^T g(B^T X) + p(c)} \rightarrow \underset{\text{output}}{c^T g(B^T X)} = \hat{Y}$$

Here, $p(B)$ is a q -dimensional vector of penalty functions added to the projection layer, one for each row of B , and $p(c)$ is the penalty added to the combination layer. The tuning parameters λ and γ control the magnitude of the penalties. The choice of penalty functions is flexible, e.g. LASSO, adaptive LASSO (Tibshirani, 1996; Zou, 2006) can be applied.

4.3.3 Computation

We use the python deep learning library Keras with TensorFlow backend to implement the sNN. We use the tanh function (i.e. hyperbolic tangent) as the activation function for the subnetwork. In simulation study, we set batch size to be 1000 and number of epochs to be 500 for both training and testing. The optimizer used is ADAM algorithm (Kingma & Ba, 2014).

A remarkable advantage of the sNN over semiparametric regression algorithm is that the training and validating process of the sNN can be paralleled, same as other NN models. Therefore, with the power of GPU computing, the speed of training and tuning over sNN models for large scale datasets can be escalated. This parallelization is readily applicable using the existing software to all sNN models.

4.4 Prediction error based model ranking

4.4.1 Prediction error

Model selection is a very practical issue in predictive modeling. In particular, for modeling with sNN that contains multiple competing structures, the issue remains in which structure to use in practice. We propose a prediction error based model ranking approach to facilitate the model selection process.

We start by constructing a set of candidate models $M = \{f_j(x) : j = 1, \dots, J; j \geq 2\}$; parameter j indexes model j and \mathcal{X}_j is its parameter space, which can be of either finite or infinite dimensional. The number of candidate models, J is assumed to be fixed. Among the models in M , we want to find the best performing model. We define "best" with respect to the prediction squared error, also called the prediction risk. For a model indexed by j , the prediction risk is

$$R_j = \int_{\mathcal{X}_j} (f_j(x) - y)^2 dP_0 \quad (4.2)$$

The benefits for choosing this criterion are twofold: (i). it measures the expected performance of the model on future (unobserved) data, representing our main interest in predictive modeling; and (ii). it can be easily estimated by the means of sample split, which we will describe soon.

To construct valid inference, we use sample split in training and validating models. We split observed data into two datasets: training set and testing set. Normally, we make an 80%/20% training and testing split as the sample size is assumed to be large. The training set D_{train} with size m is used to build a candidate set: $M_{m,j} = \{f_{m,j}(x) : j = 1, \dots, J; m, j \geq 2\}$. Denote the testing set to be D_{test} and its size to be $n - m$. An estimate of the prediction risk is the mean squared error (MSE) evaluated on D_{test} ,

$$b_{r,j} = \frac{1}{r} \sum_{i \in D_{test}} (f_{m,j}(x_i) - y_i)^2$$

Based on MSE, we introduce a simple yet valid hypothesis testing procedure to facilitate model selection.

4.4.2 Hypothesis test

Let $b_{r,j} = (b_{r,1}, \dots, b_{r,J})^T$ be the vector of MSEs of all the models in the candidate set evaluated on D_{test} , and $\theta = (\theta_1, \dots, \theta_J)^T$ be the vector of the true predictive risks under the data generating distribution P_0 , as defined in (4.2).

For a fixed candidate model set indexed by j , by the central limit theorem, we have an asymptotic multivariate normal distribution

$$P_{\bar{r}}(b_{r,j}) \stackrel{d}{\rightarrow} N_J(\theta; \Sigma) \quad (4.3)$$

The covariance matrix Σ can be estimated empirically as

$$\hat{\Sigma} = \frac{1}{r} \sum_{i \in D_{test}} \begin{pmatrix} z_{1,i} \\ \vdots \\ z_{J,i} \end{pmatrix} \begin{pmatrix} z_{1,i}^T \\ \vdots \\ z_{J,i}^T \end{pmatrix}$$

where $\sigma_{ij}^2 = \int y_i^2(x) g^2(x) dx; \dots; \int y_j^2(x) g^2(x) dx$;

As we favor the model $m_i(x)$ with the smallest σ_{ij}^2 in M , the challenge is to identify such "best" model from limited observed data. We utilize the asymptotic result (4.3) to construct a hypothesis testing procedure to compare the model performance and consequently select the best model in the candidate set. The procedure is able to provide (i) pointwise confidence intervals for each σ_{ij}^2 ; and (ii) pairwise comparisons between models with family-wise error control.

The hypothesis test compares the prediction risk of model $m_i(x)$ to the prediction risk of model $m_j(x)$ for $i < j$, i.e. test the null hypothesis

$$H_{0;ij} : \sigma_{ij}^2 \geq 0;$$

against the alternative

$$H_{A;ij} : \sigma_{ij}^2 < 0;$$

for all $j > i; i = 1; \dots; J - 1$. For each pairwise test, the test statistic is

$$t_{ij} = \frac{b_i - b_j}{\sqrt{b_{ii} + b_{jj} - 2b_{ij}}}; \tag{4.4}$$

which is a one-sided t-test comparing the difference of two means. If the null hypothesis is rejected, we conclude that $m_j(x)$ has a significant smaller prediction risk compared to $m_i(x)$ and thus is a better model for prediction purposes.

When comparing J models, there are $K = (J - 1)J/2$ number of pairwise comparisons. Therefore, we need to control family-wise error rate (FWER). We use Holm step-down procedure to control FWER (Lehmann & Romano, 2006). Let α be the target FWER level. A brief review of Holm step-down procedure is the following:

Step 1: Order p-values of the K hypothesis tests from small to large

$$p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(K)}$$

and the corresponding null hypotheses as $H_{(1)}; H_{(2)}; \dots; H_{(K)}$;

Step 2: For a given target significant level α , let l be the minimal index such that $p_{(l)} > \frac{\alpha}{K+1}$;

Step 3: Reject $H_{(1)}; H_{(2)}; \dots; H_{(l-1)}$, do not reject $H_{(l)}; H_{(l+1)}; \dots; H_{(K)}$;

Based on the pairwise comparison results, we can provide a partially ordered list of models. It is partially ordered when one or some of the models are not significantly different from some of the others; and then those models are grouped as a set of similar models. On the cases where one model stands out, it is the selected "best" model. When no such distinction exists, we can still identify a "best" set of models. We can then choose among them incorporating domain knowledge and explainability concerns, i.e. the most parsimonious model is preferred within the set of equally performing models.

Remark 2 We make an important remark on the model ranking procedure that we do not aim to uncover the true data generating model since in practice it is mostly unknown; instead we aim to select the best performing model within a pre-defined candidate set. As we only specify model configurations and train the models by the training data, the actual trained models and thus their performance rankings are data dependent, which means that whether each H_0 holds is conditional on the candidate set built with one specific dataset.

In simulation study, we are able to identify which hypothesis holds given a fitted candidate model set by the aid of a large hold-out dataset D_{val} generated under the same mechanism that generates the training and testing set. The true prediction risk is approximated by the MSE on D_{val} . As the sample size of D_{val} grows arbitrarily large, it will be very close to the true prediction risk inherited to model $f(x)$. Therefore, in simulation study we are able to detect whether H_0

holds in each run, and thus the power and type I error rate regarding the whole sample-splitting and hypothesis testing procedure are accessible. The simulation study for such examinations is in Section 4.5.2.

4.5 Simulation studies

4.5.1 Comparative study of sNN and semiparametric regression models

In this section, we benchmark the sNN models against their semiparametric counterparts to demonstrate the equivalence. Specifically, SIM-net is compared with projection pursuit regression with one ridge function using python package `projection-pursuit` (Komarov, 2018); AM-Net is compared with generalized additive model using python package `pyGAM` (Seven & Brummitt, 2018); AIM-Net is compared with projection pursuit regression with multiple ridge functions. The smooth function used in semiparametric models is polynomials of degree 2. The subnetwork architecture is [1; 30; 10; 5; 1] throughout.

We simulate $X \in \mathbb{R}^p$ in the following fashion. First we generate Z_k , $k = 1; \dots; p$ and an auxiliary variable U i.i.d. from $\text{Uniform}(0; 1)$; then the covariate is generated by $X_k = (Z_k + tU)/(1 + t)$, for $k = 1; \dots; p$. By this formulation, the pairwise correlation between two covariate is $\text{Corr}(X_i; X_j) = t^2/(1 + t^2)$ for $i \neq j$.

The metric used for comparison is MSE on the testing set, as predictive performance is our primary goal. We compare results based on two levels of pairwise correlation, i.e. $t = 0$ ($\text{Corr} = 0$) and $t = 1$ ($\text{Corr} = 0.5$), and two sample sizes of $n = 10;000$ and $n = 80;000$. We make an 80%/20% sampling splitting for training and testing set respectively. Within each class of the sNN model, we generate data according to two outcome models. Results are based on 500 Monte Carlo replications.

SIM-Net : we generate $X \in \mathbb{R}^8$ and use the following two SIM models to generate outcomes:

$$\text{SIM1 } Y = \sin(2 \pi X_1) - 2 \sin(2 \pi X_2) + e;$$

$$\text{SIM2 } Y = \exp(3 \frac{1}{2} X) + e.$$

$$\beta_1 = (1; 2; 3; 4; 0; 0; 0; 0)^T, \quad \beta_2 = (1; 1; 1; 1; 1; 1; 1; 1)^T, \text{ and } e \sim N(0; 1).$$

AM-Net : we use the following two additive models to generative outcomes:

$$\text{AM1 } Y = 2X_1 + 3(2X_2 - 1)^2 + 3\sin(2X_3) + 2\sin(2X_3)g + 3\log(X_4 + 0.5) + e;$$

$$\text{AM2 } Y = X_1 + 2X_2 + X_3 + 2X_4 + 3X_5 + (2X_6 - 1)^2 + \sin(2X_7) + 2\sin(2X_7)g + \exp(3X_8) + \log(2X_9 + 0.5) + \cos(X_{10}) + e;$$

We generate $X \sim R^8$ for AM1, $X \sim R^{20}$ for AM2 and $e \sim N(0; 1; 5^2)$.

AIM-Net : we generate $X \sim R^8$. Outcomes are generated from

$$\text{AIM1 } Y = \frac{1}{3}X + (\frac{1}{4}X)^2 + \sin(2\beta_5^T X) + 2\sin(2\frac{1}{5}X)g + e;$$

$$\text{AIM2 } Y = \exp(3\frac{1}{3}X) + \log(2\frac{1}{4}X + 0.5) + \cos(\frac{1}{5}X) + e.$$

$$\beta_3 = (1; 0; 1; 0; \dots; 0)^T, \quad \beta_4 = (0; 1; 0; 1; \dots; 0)^T, \quad \beta_5 = (0; 0; 1; 0; 1; 0; \dots; 0)^T \text{ and } e \sim N(0; 2^2).$$

Results are summarized in Table 4.2. In addition to MSE comparison, we also check the shape of learned ridge functions, the results are regulated to Appendix C.2. From the results, sNN and semiparametric regression produce similar results though one outperforms the other in different scenarios. However, we made an observation that the performance of sNN exhibits more prominent improvement when the sample size increases. This verifies that sNNs are more suitable for large dataset as compared to the semiparametric regression methods.

4.5.2 Simulation studies on model ranking

In this section, we examine the performance of our proposed model ranking procedure.

We generate covariate $X \sim R^8$ and $X_i \stackrel{i.i.d}{\sim} U[0; 1]$ for $i = 1; \dots; 8$. Throughout this section, we set the total sample size to be $n = 10;000$. Again, we make 80% / 20% training-testing splitting, which results in $m = 8;000$ training data D_{train} and $r = 2;000$ testing data D_{test} . In each simulation, we use D_{train} to train the J models in the pre-defined candidate set. Each model

takes training data as input, and outputs a fitted model $\hat{m}_{mj}(x)$ for $j = 1, \dots, J$. Then we evaluate $\hat{m}_{mj}(x)$ using D_{test} , which results in an estimated prediction risk $\hat{b}_{r, mj}$. An enormous validation dataset D_{val} of size 1,000,000 is created to assess the true prediction risk m_{mj} for $j = 1, \dots, J$. Results in this section are all based on 1,000 Monte Carlo simulations. The nominal level for hypothesis test is $\alpha = 0.05$.

The metrics for the pointwise estimate of prediction risk are: $\text{MSE} \hat{b}_{r, mj}$ and its standard error; 95% coverage of the Wald-type confidence interval for $\hat{b}_{r, mj}$. The metrics for pairwise comparisons are: MC mean of p-values of each pairwise test; p-value; the number of times H_0 is true $\# H_0$ and the number of rejections given H_0 is true $\# \text{Rej}_0$; the number of times H_A is true $\# H_A$ and the number of rejections given H_A is true $\# \text{Rej}_A$; the empirical type 1 error rate $\hat{\alpha} = \# \text{Rej}_0 / \# H_0$, and the empirical power $1 - \hat{\beta} = \# \text{Rej}_A / \# H_A$ with and without Holm step-down FWER adjustment. We also provide a sanity check with respect to the percentage that the estimated best model set contains the best model and the average length (number of models) of the estimated best model set.

We provide three cases with different candidate model sets and data generating mechanisms. The candidate model sets are all formed from the following models:

Model 1: LM - ordinary least square;

Model 2: SIM-Net - subnetwork layer architecture [1; 30; 10; 5; 1];

Model 3: AM-Net - subnetwork layer architecture [1; 30; 10; 5; 1];

Model 4: AIM-Net - subnetwork layer architecture [1; 30; 10; 5; 1];

Model 5: FFNN1 - fully-connect FNN with hidden layer architecture [5; 5];

Model 6: FFNN2 - fully-connect FNN with hidden layer architecture [30; 10; 5].

Case 1: we start with a simple case, where only three models are compared. The candidate

model set is f 1.AIM-Net; 2.FFNN1; 3.FFNN2 g. Outcomes are generated from

$$Y = \frac{q}{\exp(3 \beta_1^T X + 10) \cos(\frac{1}{2} X^T \beta_2)} + 2 X_2 X_8 + e;$$

where $\beta_1 = (0.5; 0.5; 0.5; 0; \dots; 0)^T$, $\beta_2 = (0.1; 0; 0.1; 0; 0.1; 0; 0)^T$ and $e \sim N(0; 1)$.

Case 2: we expand the candidate set to contain five models: f 1.SIM-Net; 2.AM-Net; 3.AIM-Net; 4.FFNN1; 5.FFNN2 g. Therefore, there are 10 hypothesis tests to be conducted simultaneously. Outcomes are generated by the same model as in Case 1. I

Case 3: we generate outcomes from a more complex model

$$Y = \frac{q}{\exp(3 \beta_1^T X + 10) \cos(\frac{1}{2} X^T \beta_2)} + 2 X_2 X_8 + 2 X_3 \sin(\frac{1}{3} X^T \beta_3) + e;$$

where $e \sim N(0; 1)$, $\beta_1 = (0.5; 0.5; 0.5; 0; \dots; 0)^T$, $\beta_2 = (0; 0.1; 0; 0.1; 0; 0.1; 0; 0)$ and $\beta_3 = (0; 0; 0; 0; 0.2; 0.2; 0; 0)^T$. The candidate model set is f 1.LM; 2.SIM-Net; 3.AM-Net; 4.AIM-Net; 5.FFNN2 g. Again, we conduct 10 hypothesis tests simultaneously. The purpose of this case is to see the performance when the performance among models are more distinct.

The density plot of the prediction risk for models in the candidate set are shown in Figure 4.4 - 4.6 for Case 1 - 3 respectively. It shows the distribution of the prediction risk evaluated using D_{val} for each model configuration.

Simulation results of each case are summarized in Table 4.3 - 4.5. The results show that the coverage rate of the Wald confidence intervals is very close to the nominal level, and the estimated best set can almost always contain the best model. Besides, there are additional insightful observations from the results. The empirical type I error rate is always smaller than the nominal level, which indicates that the test is conservative. As known, the power of the test is contingent on the effect size, which in our case is the difference between two prediction risks. Note that the metric p -value quantifies the effect size. It can be seen from the results that, the power of the test is higher when p -value is either close to 0 or 1, which is consistent with the fact

that the larger the effect size, the higher the power.

Recall that the actual fitted candidate set is data-dependent. Even though the model configurations in the candidate set are fixed, it is possible to observe either event " H_0 holds" or event " H_A holds" in each run of simulation. Notice that the larger the p-value, the more likely the occurrence of " H_0 holds"; the smaller the p-value, the more likely the occurrence of " H_A holds". To see this better, Figure 4.4 shows that FFNN2 performs the best and FFNN1 performs the worst on average in Case 1. This is consistent with the result that the third test in Case 1 has p-value close to 0. The effect size in the third test is the largest, thus it has the highest power. Also notice that the event " H_A holds" dominates with 9853 occurrences in 1,000 Monte Carlo runs. Other results can be interpreted in the same manner.

Comparing the results in Case 2 and 3, we show the impact of effect size on the model selection. In Case 2, the distributions of prediction risk are similar among models, except for the AM-Net, as shown in Figure 4.5. On the contrary, the distributions are more separated in Case 3. It is thus more easily to select the best model from the candidate set in Case 3. Therefore, the average length of the estimated best model set in Case 3 (3:08) is shorter than that of Case 2 (3:31).

4.6 Concluding remarks and future work

We have presented a new model framework that uses neural networks to represent the semiparametric regression models, along with a model ranking procedure. The proposed sNN class is highly practical to predictive modeling objectives and in the meanwhile maintains the model explainability. The sNN class can be generalized to other semiparametric model of interest easily.

For future work, we propose to construct a visualization tool to aid the decision making in model selection for practitioners. An illustration of such tool is shown in Figure 4.7. From the top to bottom, we have models that have four to one ridge function(s). The left panel shows the heat maps that correspond to the estimated indices; the right panel shows the solution paths for each

Figure 4.4 Empirical distribution and fitted density of prediction risks based on the validation dataset of Case 1.

of the ridge functions for each model. With the aid of such tool, we can enable an interactive model building process that helps explain the model behaviour to the practitioners.

Figure 4.5 Empirical distribution and fitted density of prediction risks based on the validation dataset of Case 2.

Figure 4.6 Empirical distribution and fitted density of prediction risks based on the validation dataset of Case 3.

Figure 4.7 Illustration of the proposed visualization tool.

Table 4.2 Comparison on MSE (S.E.) of sNN with semiparametric regressions: (a). SIM-Net verses PPR with one ridge function; (b). AM-Net verses GAM; (c). AIM-Net verses PPR with three ridge functions.

(a) SIM-Net vs PPR(1)

	Corr = 0		Corr = 0.5	
	SIM-Net	PPR(M = 1)	SIM-Net	PPR(M = 1)
n = 10;000				
SIM1	1:058(0:048)	1:122(0:034)	1:054(0:044)	1:095(0:035)
SIM2	1:026(0:032)	1:002(0:029)	1:020(0:035)	1:000(0:033)
n = 80;000				
SIM1	1:011(0:036)	1:123(0:012)	1:002(0:014)	1:094(0:012)
SIM2	1:005(0:013)	1:001(0:011)	1:001(0:010)	1:002(0:010)

(b) AM-Net vs GAM

	AM-Net	GAM	AM-Net	GAM
	n = 10;000			
AM1	2:275(0:083)	2:305(0:045)	2:267(0:089)	2:264(0:055)
AM2	2:278(0:332)	2:277(0:055)	2:271(0:114)	2:269(0:054)
n = 80;000				
AM1	2:256(0:045)	2:298(0:063)	2:258(0:071)	2:254(0:020)
AM2	2:263(0:064)	2:261(0:022)	2:257(0:070)	2:252(0:026)

(c) AIM-Net vs PPR(3)

	AIM-Net	PPR	AIM-Net	PPR
	n = 10;000			
AIM1	4:158(0:078)	4:140(0:007)	4:089(0:023)	4:073(0:008)
AIM2	4:137(0:176)	4:030(0:128)	4:053(0:159)	4:024(0:133)
n = 80;000				
AIM1	4:056(0:076)	4:099(0:002)	4:023(0:007)	4:034(0:001)
AIM2	4:022(0:071)	4:029(0:044)	4:005(0:045)	4:009(0:046)

Table 4.3 Summary of Case 1: point estimate of prediction risk, 95% coverage of the Wald-type confidence intervals for each model; mean of p-values, empirical type I error rate and power for each pairwise comparison w/o Holm adjustment; percentage of the estimated best set contains the true best model w/o Holm adjustment and the average length of the the estimated best set.

Model			1. AIM-Net	2. FFNN1	3. FFNN2
MSE (S.E.)			1:012(0.033)	1:023(0.033)	1:007(0.032)
95% Wald CI Coverage			0.951	0.951	0.950

H ₀	p-value	# H ₀ (# Rej ₀)	b	b ^{adj}	# H _A (# Rej _A)	1	b	1	b ^{adj}
2	1	0	0:796	8717(31)	0004	0001	1283(505)	0894	0337
3	1	0	0:281	1538(28)	0018	0007	8462(3249)	0884	0286
3	2	0	0:079	147(4)	0027	0007	9853(7252)	0736	0616

Percentage of the estimated best set contains the true best model			
Unadjusted	(Ave. length)	Holm Adjusted	(Ave. length)
99:64%	(208)	99:88%	(225)

Table 4.4 Summary of Case 2: point estimate of prediction risk and 95% coverage of the Wald-type confidence interval; mean of p-value, empirical type I error rate and power w/o Holm adjustment; percentage of the estimated best set contains the true best model w/o Holm adjustment and the average length of the the estimated best set.

Model			1. SIM-Net	2.AM-Net	3.AIM-Net	4.FFNN1	5.FFNN2
MSE (S.E.)			1:03(0.06)	1:03(0.03)	1:01(0.03)	1:01(0.03)	1:01(0.03)
95% Wald CI Coverage			0.949	0.951	0.950	0.951	0.951

H ₀	p-value	# H ₀ (# Rej ₀)	b	b ^{adj}	# H _A (# Rej _A)	1	b	1	b ^{adj}
2	1	0	0:818	8471(7)	0001	0	1529(1242)	0812	0761
3	1	0	0:260	1120(44)	0039	0003	8880(3639)	0410	0272
3	2	0	0:027	0(0)	-	-	10000(8641)	0864	0630
4	1	0	0:387	3629(52)	0014	0002	6371(2970)	0466	0336
4	2	0	0:032	92(3)	0033	0	9908(8870)	0885	0673
4	3	0	0:650	7511(98)	0013	0001	2489(428)	0172	0047
5	1	0	0:222	321(8)	0025	0003	9679(4066)	0420	0264
5	2	0	0:021	0(0)	-	-	10000(8884)	0888	0662
5	3	0	0:435	3587(87)	0024	0004	6413(1289)	0201	0062
5	4	0	0:297	1555(37)	0024	0006	8445(2853)	0338	0157

Percentage of the estimated best set contains the true best model			
Unadjusted	(Ave. length)	Holm Adjusted	(Ave. length)
98:67%	(331)	99:77%	(385)

Table 4.5 Summary of Case 3: point estimate of prediction risk and 95% coverage of the Wald-type confidence interval; mean of p-value, empirical type I error rate and power w/o Holm adjustment; percentage of the estimated best set contains the true best model w/o Holm adjustment and the average length of the the estimated best set.

Model			1.LM	2.SIM-Net	3.AM-Net	4.AIM-Net	5.FFNN2			
MSE (S.E.)			1:04(0:03)	1:04(0:04)	1:05(0:03)	1:02(0:03)	1:01(0:03)			
95% Wald CI Coverage			0.951	0.951	0.950	0.948	0.950			
H ₀			p-value	# H ₀ (# Rej ₀)	b	b ^{adj}	# H _A (# Rej _A)	1 b 1	b ^{adj}	
2	1	0	0:199	1080(5)	0005	0	8920(4896)	0549	0322	
3	1	0	0:672	9846(132)	0013	0003	154(13)	0084	0060	
3	2	0	0:810	9144(21)	0002	0	856(408)	0477	0380	
4	1	0	0:033	11(0)	0	0:	9989(8513)	0466	0652	
4	2	0	0:098	204(2)	001	0	9796(6669)	0852	0468	
4	3	0	0:028	6(0)	0	0:	9994(8745)	0681	0688	
5	1	0	0:010	0(0)	-	-	10 000(9502)	0950	0803	
5	2	0	0:027	0(0)	-	-	10 000(8702)	0870	0678	
5	3	0	0:007	0(0)	-	-	10 000(9671)	0967	0843	
5	4	0	0:183	535(24)	0045	007	9465(4612)	0487	0291	
Percentage of the estimated best set contains the true best model										
Unadjusted			(Ave. length)			Holm Adjusted			(Ave. length)	
99:76%			(1:68)			99:96%			(230)	

BIBLIOGRAPHY

- Akaike, H. (1973). "Information theory and an extension of the maximum likelihood principle". In: 2nd International Symposium on Information Theory, 1973. Akademiai Kiado, pp. 267-281.
- Al (1977). "On entropy maximization principle". *Application of Statistics*, pp. 27-41.
- Almirall, D. et al. (2016). "Longitudinal effects of adaptive interventions with a speech-generating device in minimally verbal children with ASD". *Journal of Clinical Child & Adolescent Psychology* 45.4, pp. 442-456.
- Anders, U. & O. Korn (1999). "Model selection in neural networks". *Neural networks* 12.2, pp. 309-323.
- Athey, S. & S. Wager (2017). "Efficient policy learning". arXiv preprint arXiv:1702.02896.
- Bai, X. et al. (2013). "Doubly-Robust Estimators of Treatment-Specific Survival Distributions in Observational Studies with Stratified Sampling". *Biometrics* 69, pp. 830-839.
- Barrett, J. K. et al. (2014). "Doubly robust estimation of optimal dynamic treatment regimes". *Statistics in Biosciences* 6.2, pp. 244-260.
- Bellman, R. (1957). "Dynamic programming". Princeton University Press, p. 151.
- Berlinet, A. & C. Thomas-Agnan (2011). *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Springer Science & Business Media.
- Bickel, P. J. et al. (1993). *Efficient and Adaptive Inference in Semiparametric Models*. Johns Hopkins University Press Baltimore.
- Blatt, D et al. (2004). "A-learning for approximate planning". Tech. rep. University of Michigan, Ann Arbor, MI, pp. 48109-2122.
- Boos, D. D. & L. A. Stefanski (2013). *Essential Statistical Inference: Theory and Methods* Vol. 120. Springer Science & Business Media.
- Boyd, S. & L. Vandenberghe (2004). *Convex Optimization*. Cambridge university press.
- Brainard, L. (2018). "What Are We Learning about Artificial Intelligence in Financial Services?" <https://www.federalreserve.gov/newsevents/speech/brainard20181113a.htm>.
- Buchanan, A. L. et al. (2018). "Generalizing evidence from randomized trials using inverse probability of sampling weights". *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 181, pp. 1193-1209.

- Busoniu, L. et al. (2010). Reinforcement Learning and Dynamic Programming using Function Approximators. CRC press.
- Buuren, S. & K. Groothuis-Oudshoorn (2011). \mice: Multivariate imputation by chained equations in R". Journal of Statistical Software 45.3.
- Chakraborty, B. & E. Moodie (2013). Statistical Methods for Dynamic Treatment Regimes. Springer.
- Chan, K. C. G. et al. (2015). \Globally efficient non-parametric inference of average treatment effects by empirical balancing calibration weighting". Journal of the Royal Statistical Society: Series B (Statistical Methodology) 78, pp. 673{700.
- Chang, T. & P. S. Kott (2008). \Using calibration weighting to adjust for nonresponse under a plausible model". Biometrika 95, pp. 555{571.
- Chen, G. et al. (2016). \Personalized dose finding using outcome weighted learning" Journal of the American Statistical Association 111.516, pp. 1509{1521.
- Chen, J et al. (2002). \Using empirical likelihood methods to obtain range restricted weights in regression estimators for surveys" Biometrika 89, pp. 230{237.
- Chen, X. (2007). \Large sample sieve estimation of semi-nonparametric models" Handbook of Econometrics 6, pp. 5549{5632.
- Cole, S. R. & E. A. Stuart (2010). \Generalizing evidence from randomized clinical trials to target populations: The ACTG 320 trial". American Journal of Epidemiology 172, pp. 107{115.
- Cristianini, N., J. Shawe-Taylor, et al. (2000). An Introduction to Support Vector Machines and other Kernel-based Learning Methods Cambridge university press.
- Curry, B. & P. H. Morgan (2006). \Model selection in neural networks: some difficulties". European Journal of Operational Research 170.2, pp. 567{577.
- Dahabreh, I. J. et al. (2018). \Generalizing causal inferences from individuals in randomized trials to all trial-eligible individuals". Biometrics 0. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/biom.13009> .
- Davidian, M. et al. (2019). Introduction to Treatment Regimes. Chapman Hall (forthcoming).
- Ertefaie, A. et al. (2016). \Q-learning residual analysis: application to the effectiveness of sequences of antipsychotic medications for patients with schizophrenia" Statistics in Medicine 35.13, pp. 2221{2234.

- Fan, J. & R. Li (2001). "Variable selection via nonconcave penalized likelihood and its oracle properties". *Journal of the American Statistical Association* 96, pp. 1348{1360.
- Fan, J. et al. (2016). Improving covariate balancing propensity score: A doubly robust and efficient approach . Tech. rep. Technical report, Princeton University.
- Friedman, J. H. & W. Stuetzle (1981). "Projection pursuit regression". *J Am Stat Assoc* 76, pp. 817{823.
- Fu, W. J. (2003). "Penalized estimating equations". *Biometrics* 59.1, pp. 126{132.
- Geman, S. & C.-R. Hwang (1982). "Nonparametric maximum likelihood estimation by the method of sieves". *The Annals of Statistics* 10, pp. 401{414.
- Geramifard, A. et al. (2013). "A tutorial on linear function approximators for dynamic programming and reinforcement learning". *Foundations and Trends in Machine Learning* 6.4, pp. 375{451.
- Goldberg, Y. & M. R. Kosorok (2012). "Q-learning with censored data". *Annals of statistics* 40.1, p. 529.
- Graham, B. S. et al. (2012). "Inverse probability tilting for moment condition models with missing data". *The Review of Economic Studies* 79, pp. 1053{1079.
- Gu, C. (2013). *Smoothing Spline ANOVA Models*. Vol. 297. Springer Science & Business Media.
- Guan, Q. et al. (2018). "Bayesian Nonparametric Policy Search with Application to Periodontal Recall Intervals". *arXiv preprint arXiv:1810.04338*.
- Hahn, J. (1998). "On the role of the propensity score in efficient semiparametric estimation of average treatment effects". *Econometrica* 66, pp. 315{331.
- Hahn, O. M. & R. L. Schilsky (2012). "Randomized controlled trials and comparative effectiveness research". *Journal of Clinical Oncology* 30, pp. 4194{4201.
- Hainmueller, J. (2012). "Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies". *Political Analysis* 20, pp. 25{46.
- Hájek, J (1971). "Comment on a paper by D. Basu". *Foundations of Statistical Inference* 236.
- Hardle, W. et al. (1993). "Optimal smoothing in single-index models". *The Annals of Statistics* 21, pp. 157{178.
- Hastie, T. & R. Tibshirani (1987). "Generalized additive models: some applications". *Journal of the American Statistical Association* 82.398, pp. 371{386.

- Hastie, T. J. (2017). "Generalized additive models". In: *Statistical models* in S. Routledge, pp. 249{307.
- Henderson, R. et al. (2010). "Regret-regression for optimal dynamic treatment regimes"*Biometrics* 66.4, pp. 1192{1201.
- Hernan, M. A. & T. J. VanderWeele (2011). "Compound treatments and transportability of causal inference".*Epidemiology* 22, p. 368.
- Hershman, D. L. & J. D. Wright (2012). "Comparative effectiveness research in oncology methodology: observational data". *Journal of Clinical Oncology* 30, pp. 4215{4222.
- Hirano, K. et al. (2003). "Efficient estimation of average treatment effects using the estimated propensity score". *Econometrica* 71, pp. 1161{1189.
- Hornik, K. et al. (1989). "Multilayer feedforward networks are universal approximators". *Neural networks* 2.5, pp. 359{366.
- Horvitz, D. G. & D. J. Thompson (1952). "A generalization of sampling without replacement from a finite universe". *Journal of the American Statistical Association* 47, pp. 663{685.
- Ichimura, H. (1993). "Semiparametric least squares (SLS) and weighted SLS estimation of single-index models".*Journal of Econometrics* 58, pp. 71{120.
- Imai, K. & M. Ratkovic (2014). "Covariate balancing propensity score". *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 76, pp. 243{263.
- Imbens, G. W. et al. (2005). "Mean-square-error calculations for average treatment effects". Department of Economics, UC Berkeley, unpublished manuscript
- Jacobs, R. A. et al. (1991). "Adaptive mixtures of local experts." *Neural computation* 3.1, pp. 79{87.
- Jeng, X. J. et al. (2018). "High-dimensional inference for personalized treatment decision". *Electronic Journal of Statistics* 12.1, pp. 2074{2089.
- Jiang, B. et al. (2019). "Entropy Learning for Dynamic Treatment Regimes". *Statistica Sinica*.
- Johnson, B. A. et al. (2008). "Penalized estimating functions and variable selection in semiparametric regression models".*Journal of the American Statistical Association* 103, pp. 672{680.
- Kang, J. D. Y. & J. L. Schafer (2007). "Demystifying Double Robustness: A Comparison of Alternative Strategies for Estimating a Population Mean from Incomplete Data". *Statistical Science* 22, pp. 523{539.

- Keiding, N. & T. A. Louis (2016). "Perils and potentials of self-selected entry to epidemiological studies and surveys". *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 179, pp. 319-376.
- Kennedy, E. H. (2016). "Semiparametric theory and empirical processes in causal inference". In: *Statistical Causal Inferences and Their Applications in Public Health Research* Springer, pp. 141-167.
- Kidwell, K. M. (2014). "SMART designs in cancer research: Past, present, and future". *Clinical trials* 11.4, pp. 445-456.
- Kidwell, K. M. et al. (2018). "Design and analysis considerations for comparing dynamic treatment regimens with binary outcomes from sequential multiple assignment randomized trials". *Journal of Applied Statistics* 45.9, pp. 1628-1651.
- Kilbourne, A. M. et al. (2018). "Adaptive School-based Implementation of CBT (ASIC): clustered-SMART for building an optimized adaptive implementation intervention to improve uptake of mental health interventions in schools". *Implementation Science* 13.1, p. 119.
- Kim, J. K. et al. (2016). "Calibrated propensity score method for survey nonresponse in cluster sampling". *Biometrika* 103, pp. 461-473.
- Kingma, D. P. & J. Ba (2014). "Adam: A method for stochastic optimization". arXiv preprint arXiv:1412.6980
- Komarov, P. (2018). "Projection-pursuit: An implementation of multivariate projection pursuit regression and univariate classification" <https://github.com/pavelkomarov/projection-pursuit> .
- Korn, E. L. & B. Freidlin (2012). "Methodology for Comparative Effectiveness Research: Potential and Limitations". *Journal of Clinical Oncology* 30, pp. 4185-4187.
- Kosorok, M. R. (2007). *Introduction to Empirical Processes and Semiparametric Inference* Springer Science & Business Media.
- Kosorok, M. R. & E. E. Moodie (2015). *Adaptive Treatment Strategies in Practice: Planning Trials and Analyzing Data for Personalized Medicine* Vol. 21. SIAM.
- Kosorok, M. & E. Laber (2019). "Precision Medicine". *Annual Review of Statistics and Its Application* In press .
- Koszegi, S. T. (2019). *High-Level Expert Group on Artificial Intelligence* .
- Kott, P. S. (2006). "Using calibration weighting to adjust for nonresponse and coverage errors". *Survey Methodology* 32, pp. 133-142.

- Laan, M. J. van der & M. L. Petersen (2007). "Causal effect models for realistic individualized treatment and intention to treat rules". *The International Journal of Biostatistics* 3.1.
- Laber, E. & Y. Zhao (2015). "Tree-based methods for individualized treatment regimes". *Biometrika* 102.3, pp. 501{514.
- Laber, E. B. et al. (2014a). "Dynamic treatment regimes: Technical challenges and applications". *Electronic journal of statistics* 8.1, p. 1225.
- Laber, E. B. et al. (2014b). "Interactive model building for Q-learning". *Biometrika* 101.4, pp. 831{847.
- Laber, E. B. et al. (2014c). "Set-valued dynamic treatment regimes for competing outcomes". *Biometrics* 70.1, pp. 53{61.
- Laber, E. B. et al. (2018). "Optimal treatment allocations in space and time for on-line control of an emerging infectious disease". *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 67.4, pp. 743{789.
- Lavori, P. W. & R. Dawson (2004). "Dynamic treatment regimes: practical design considerations". *Clinical trials* 1.1, pp. 9{20.
- Lehmann, E. L. & J. P. Romano (2006). *Testing Statistical Hypotheses* Springer Science & Business Media.
- Lin, Y., H. H. Zhang, et al. (2006). "Component selection and smoothing in multivariate nonparametric regression". *The Annals of Statistics* 34.5, pp. 2272{2297.
- Linn, K. A. et al. (2017). "Interactive Q-learning for Quantiles". *Journal of the American Statistical Association* 112.518, pp. 638{649.
- Little, R. J. & D. B. Rubin (2014). *Statistical Analysis with Missing Data*. Vol. 333. John Wiley & Sons.
- Liu, Y. et al. (2018). "Augmented outcome-weighted learning for estimating optimal dynamic treatment regimens". *Statistics in Medicine*.
- Lu, W. et al. (2013). "Variable selection for optimal treatment decision". *Statistical Methods in Medical Research* 22.5, pp. 493{504.
- Lu, X. et al. (2016). "Comparing dynamic treatment regimes using repeated-measures outcomes: modeling considerations in SMART studies". *Statistics in Medicine* 35.10, pp. 1595{1615.
- Luckett, D. J. et al. (2018). "Estimating Dynamic Treatment Regimes in Mobile Health Using V-learning". *Journal of the American Statistical Association* just-accepted, pp. 1{39.

- Massarelli, E et al. (2003). "A retrospective analysis of the outcome of patients who have received two prior chemotherapy regimens including platinum and docetaxel for recurrent non-small-cell lung cancer". *Lung Cancer* 39, pp. 55{61.
- Moguerza, J. M., A. Muñoz, et al. (2006). "Support vector machines with applications". *Statistical Science* 21.3, pp. 322{336.
- Moodie, E. E. et al. (2007). "Demystifying optimal dynamic treatment regimes". *Biometrics* 63.2, pp. 447{455.
- Moodie, E. E. et al. (2014). "Q-learning: Flexible learning about useful utilities". *Statistics in Biosciences* 6.2, pp. 223{243.
- Murata, N. (1991). "A criterion for determining the number of parameters in an artificial neural network model". *Artificial Neural Networks* , pp. 9{14.
- Murphy, S. A. (2003). "Optimal dynamic treatment regimes". *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 65.2, pp. 331{355.
- | (2005a). "A generalization error for Q-learning". *Journal of Machine Learning Research* 6.Jul, pp. 1073{1097.
- | (2005b). "An experimental design for the development of adaptive treatment strategies". *Statistics in Medicine* 24.10, pp. 1455{1481.
- Nahum-Shani, I. et al. (2012). "Q-learning: A data analysis method for constructing adaptive interventions." *Psychological methods* 17.4, p. 478.
- Nahum-Shani, I. et al. (2017). "A SMART data analysis method for constructing adaptive treatment strategies for substance use disorders" *Addiction* 112.5, pp. 901{909.
- Newey, W. K. (1997). "Convergence rates and asymptotic normality for series estimators". *Journal of Econometrics* 79, pp. 147{168.
- Neyman, J. (1923). "Sur Les applications de la thar des probabilités aux expériences Agraricales: Essay de principe. English translation of excerpts by Dabrowska, D. and Speed, T." *Statistical Science* 5, pp. 465{472.
- Nosedal-Sanchez, A. et al. (2012). "Reproducing kernel Hilbert spaces for penalized regression: A tutorial". *The American Statistician* 66.1, pp. 50{60.
- O'Muircheartaigh, C. & L. V. Hedges (2014). "Generalizing from unrepresentative experiments: a stratified propensity score approach". *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 63, pp. 195{210.

- Orellana, L. et al. (2010). \Dynamic regime marginal structural mean models for estimation of optimal dynamic treatment regimes, part I: main content". The International Journal of Biostatistics 6.2.
- Pearl, J. & E. Bareinboim (2011). \Transportability of causal and statistical relations: A formal approach". In: Data Mining Workshops (ICDMW), 2011 IEEE 11th International Conference on. IEEE, pp. 540{547.
- Petersen, M. L. et al. (2012). \Diagnosing and responding to violations in the positivity assumption". Statistical Methods in Medical Research 21.1, pp. 31{54.
- Qi, Z., Y. Liu, et al. (2018a). \D-learning to estimate optimal individual treatment rules". Electronic Journal of Statistics 12.2, pp. 3601{3638.
- Qi, Z. et al. (2018b). \Multi-armed Angle-based Direct Learning for Estimating Optimal Individualized Treatment Rules with Various Outcomes". Journal of the American Statistical Association just-accepted, pp. 1{35.
- Qian, M. & S. A. Murphy (2011). \Performance guarantees for individualized treatment rules". Annals of Statistics 39.2, p. 1180.
- Qin, J. & B. Zhang (2007). \Empirical-likelihood-based inference in missing response problems and its application in observational studies". Journal of the Royal Statistical Society: Series B (Statistical Methodology) 69, pp. 101{122.
- Ravikumar, P. et al. (2009). \Sparse additive models". Journal of the Royal Statistical Society: Series B (Statistical Methodology) 71.5, pp. 1009{1030.
- Rich, B. et al. (2014). \Simulating sequential multiple assignment randomized trials to generate optimal personalized warfarin dosing strategies". Clinical Trials 11.4, pp. 435{444.
- Robins, J. M. (1997). \Causal inference from complex longitudinal data". In: Latent variable modeling and applications to causality. Springer, pp. 69{117.
- | (2004). \Optimal structural nested models for optimal sequential decisions". In: Proceedings of the second seattle Symposium in Biostatistics Springer, pp. 189{326.
- Robins, J. M. & A. Rotnitzky (1995). \Semiparametric efficiency in multivariate regression models with missing data". Journal of the American Statistical Association 90.429, pp. 122{129.
- Rojas, R. (2013). Neural Networks: A Systematic Introduction. Springer Science & Business Media.
- Rosenbaum, P. R. & D. B. Rubin (1983a). \The Central Role of the Propensity Score in Observational Studies for Causal Effects". Biometrika 70, pp. 41{55.

- Rothwell, P. M. (2005). "External validity of randomised controlled trials: to whom do the results of this trial apply?" *The Lancet* 365, pp. 82{93.
- Ruan, L. & M. Yuan (2010). "Dimension reduction and parameter estimation for additive index models". *Statistics and Its Interface* 3.4, pp. 493{499.
- Rubin, D. B. & M. J. van der Laan (2012). "Statistical issues and limitations in personalized medicine research with clinical trials". *The International Journal of Biostatistics* 8.1.
- Rubin, D. B. (1974). "Estimating causal effects of treatments in randomized and nonrandomized studies". *Journal of Educational Psychology* 66, pp. 688{701.
- | (1976). "Inference and missing data". *Biometrika* 63.3, pp. 581{592.
- | (1978). "Bayesian inference for causal effects: The role of randomization" *Annals of statistics*, pp. 34{58.
- | (1980). "Comment on "Randomization analysis of experimental data: The Fisher randomization test" by D. Basu". *Journal of the American Statistical Association* 75, pp. 591{593.
- Rudolph, K. E. & M. J. van der Laan (2017). "Robust estimation of encouragement design intervention effects transported across sites". *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 79, pp. 1509{1525.
- Rumelhart, D. E. et al. (1985). "Learning internal representations by error propagation" Tech. rep. California Univ San Diego La Jolla Inst for Cognitive Science.
- Saunders, C.-E. et al. (2003). "Model assisted survey sampling" Springer Science & Business Media.
- Schulte, P. J. et al. (2014). "Q-and A-learning methods for estimating optimal dynamic treatment regimes". *Statistical Science* 29.4, p. 640.
- Seaman, S. R. et al. (2012). "Combining multiple imputation and inverse-probability weighting". *Biometrics* 68.1, pp. 129{137.
- Servén, D. & C. Brummitt (2018). "pyGAM: Generalized Additive Models in Python". <https://github.com/dswah/pyGAM> .
- Shen, X. (1997). "On methods of sieves and penalization". *The Annals of Statistics* 25, pp. 2555{2591.
- Shi, C. et al. (2018). "High-dimensional A-learning for optimal dynamic treatment regimes". *The Annals of Statistics* 46.3, pp. 925{957.

- Shortreed, S. M. & E. E. Moodie (2012). \Estimating the optimal dynamic antipsychotic treatment regime: evidence from the sequential multiple-assignment randomized Clinical Antipsychotic Trials of Intervention and Effectiveness schizophrenia study". *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 61.4, pp. 577{599.
- Shortreed, S. M. et al. (2011). \Informing sequential clinical decision-making through reinforcement learning: an empirical study". *Machine Learning* 84.1-2, pp. 109{136.
- Shortreed, S. M. et al. (2014). \A multiple imputation strategy for sequential multiple assignment randomized trials". *Statistics in Medicine* 33.24, pp. 4202{4214.
- Strauss, G. M. et al. (2008). \Adjuvant paclitaxel plus carboplatin compared with observation in stage IB non-small-cell lung cancer: CALGB 9633 with the Cancer and Leukemia Group B, Radiation Therapy Oncology Group, and North Central Cancer Treatment Group Study Groups". *Journal of Clinical Oncology* 26, p. 5043.
- Stroup, T. S. et al. (2003). \The National Institute of Mental Health Clinical Antipsychotic Trials of Intervention Effectiveness (CATIE) project: schizophrenia trial design and protocol development." *Schizophrenia Bulletin* 29.1, p. 15.
- Stuart, E. A. et al. (2011). \The use of propensity scores to assess the generalizability of results from randomized trials". *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 174, pp. 369{386.
- Stuart, E. A. et al. (2015). \Assessing the generalizability of randomized trial results to target populations". *Prevention Science* 16, pp. 475{485.
- Sugiyama, M. & M. Kawanabe (2012). *Machine Learning in Non-stationary Environments: Introduction to Covariate Shift Adaptation* . MIT press.
- Sutton, R. S. & A. G. Barto (2018). *Reinforcement Learning: An Introduction* . MIT press.
- Taylor, J. M. et al. (2015). \Reader reaction to A robust method for estimating optimal treatment regimes by Zhang et al.(2012)". *Biometrics* 71.1, pp. 267{273.
- Tian, L. et al. (2014). \A simple method for estimating interactions between a treatment and a large number of covariates". *Journal of the American Statistical Association* 109.508, pp. 1517{1532.
- Tibshirani, R. (1996). \Regression shrinkage and selection via the lasso". *Journal of the Royal Statistical Society: Series B (Methodological)* 58, pp. 267{288.
- Tipton, E. (2013). \Improving generalizations from experiments using propensity score subclassification: Assumptions, properties, and contexts". *Journal of Educational and Behavioral Statistics* 38, pp. 239{266.

- Tsiatis, A. (2007). *Semiparametric Theory and Missing Data* Springer Science & Business Media.
- Tsiatis, A. A. et al. (2014). *Handbook of Missing Data Methodology* Chapman and Hall/CRC.
- Vaart, A. W. van der & J. A. Wellner (1996). *Weak Convergence and Empirical Processes: With Applications to Statistics*. New York: Springer.
- Vaughan, J. et al. (2018). "Explainable Neural Networks based on Additive Index Models". arXiv e-prints, arXiv:1806.01933.
- Wahba, G. (1990). *Spline models for observational data* Vol. 59. Siam.
- Wang, L. et al. (2012a). "Penalized generalized estimating equations for high-dimensional longitudinal data analysis". *Biometrics* 68, pp. 353{360.
- Wang, L. et al. (2018). "Quantile-optimal treatment regimes". *Journal of the American Statistical Association* 113.523, pp. 1243{1254.
- Wang, L. et al. (2012b). "Evaluation of viable dynamic treatment regimes in a sequentially randomized trial of advanced prostate cancer". *Journal of the American Statistical Association* 107.498, pp. 493{508.
- Watanabe, C. et al. (2018). "Modular representation of layered neural networks". *Neural Networks* 97, pp. 62{73.
- Wolsey, L. A. & G. L. Nemhauser (2014). *Integer and Combinatorial Optimization*. John Wiley & Sons.
- Wu, C. & R. R. Sitter (2001). "A model-calibration approach to using complete auxiliary information from survey data". *Journal of the American Statistical Association* 96, pp. 185{193.
- Wu, P. et al. (2018). "Matched Learning for Optimizing Individualized Treatment Strategies Using Electronic Health Records". *Journal of the American Statistical Association* just-accepted, pp. 1{35.
- Xia, Y. & W. Li (1999). "On single-index coefficient regression models". *Journal of the American Statistical Association* 94.448, pp. 1275{1285.
- Xia, Y. et al. (2002). "An adaptive estimation of dimension reduction space". *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 64.3, pp. 363{410.
- Xu, Y. et al. (2016). "Bayesian nonparametric estimation for dynamic treatment regimes with sequential transition times". *Journal of the American Statistical Association* 111.515, pp. 921{950.

- Xu, Y. et al. (2019). \Bayesian non-parametric survival regression for optimizing precision dosing of intravenous busulfan in allogeneic stem cell transplantation". *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 68.3, pp. 809{828.
- Young, J. G. et al. (2011). \Comparative effectiveness of dynamic treatment regimes: an application of the parametric g-formula". *Statistics in Biosciences* 3.1, p. 119.
- Yu, Z. & M. J. van der Laan (2002). \Construction of counterfactuals and the G-computation formula".
- Zhang, B. & M. Zhang (2018). \C-learning: A new classification framework to estimate optimal dynamic treatment regimes". *Biometrics* 74.3, pp. 891{899.
- Zhang, B. et al. (2012a). \A robust method for estimating optimal treatment regimes". *Biometrics* 68.4, pp. 1010{1018.
- Zhang, B. et al. (2012b). \Estimating optimal treatment regimes from a classification perspective". *Stat* 1.1, pp. 103{114.
- Zhang, B. et al. (2013). \Robust estimation of optimal dynamic treatment regimes for sequential treatment decisions". *Biometrika* 100.3, pp. 681{694.
- Zhang, Y. et al. (2017). \Estimation of Optimal Treatment Regimes Using Lists". *Journal of the American Statistical Association* just-accepted.
- Zhao, Q. & D. Percival (2017). \Entropy balancing is doubly robust". *Journal of Causal Inference* 5, DOI: <https://doi.org/10.1515/jci-2016-0010>.
- Zhao, Y.-Q. et al. (2014). \Doubly robust learning for estimating individualized treatment with censored data". *Biometrika* 102.1, pp. 151{168.
- Zhao, Y.-Q. et al. (2015). \New statistical learning methods for estimating optimal dynamic treatment regimes". *Journal of the American Statistical Association* 110.510, pp. 583{598.
- Zhao, Y.-Q. et al. (2019a). \Efficient augmentation and relaxation learning for individualized treatment rules using observational data". *Journal of Machine Learning Research* To appear .
- Zhao, Y.-Q. et al. (2019b). \Robustifying trial-derived optimal treatment rules for a target population". *Electronic Journal of Statistics* 13, pp. 1717{1743.
- Zhao, Y. et al. (2012). \Estimating individualized treatment rules using outcome weighted learning". *Journal of the American Statistical Association* 107.499, pp. 1106{1118.
- Zhao, Y. et al. (2009). \Reinforcement learning design for cancer clinical trials". *Statistics in Medicine* 28.26, pp. 3294{3315.

Zhou, X. & M. R. Kosorok (2017). "Causal nearest neighbor rules for optimal treatment regimes". arXiv preprint arXiv:1711.08451.

Zhou, X. et al. (2017). "Residual weighted learning for estimating individualized treatment rules". Journal of the American Statistical Association 112.517, pp. 169-187.

Zou, H. (2006). "The adaptive lasso and its oracle properties". Journal of the American Statistical Association 101, pp. 1418-1429.

APPENDICES

APPENDIX A

APPENDIX FOR CHAPTER 2

A.1 IPWCC estimating equations

We give the explicit forms of the IPWCC estimating equations for Q-learning and OWL as the following. Denote

$$V_t = \max_{a_{t+1}} Q_{t+1}(H_{t+1}; a_{t+1}; \pi_{t+1}); t = 1; \dots; T-1$$

For Q-learning, the stacked IPWCC estimating equation for ψ is

$$P_n \begin{pmatrix} \psi_T \\ \vdots \\ \psi_1 \end{pmatrix} = \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n \left(Y_{T-1} - Q_T(S_{T-1}; \pi_{T-1}) \right) \left(r_{T-1} + Q_T(S_{T-1}; \pi_{T-1}) - Q_T(S_{T-1}; \pi_{T-1}) \right) \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n \left(Y_1 - Q_1(S_1; \pi_1) \right) \left(r_1 + Q_1(S_1; \pi_1) - Q_1(S_1; \pi_1) \right) \end{pmatrix} = 0$$

For OWL, the stacked IPWCC estimating equation for $\beta = (\beta_T; \beta_{T-1})^T$ is

$$P_n \begin{pmatrix} Y - Q_T(H_T; A_T; \beta_T) - r_T Q_T(H_T; A_T; \beta_T) 1_{C>T} = K_T(S_T; \beta_T) \\ \vdots \\ V_2(H_2; \beta_2) - Q_1(H_1; A_1; \beta_2; \beta_1) - r_1 Q_1(H_1; A_1; \beta_2; \beta_1) 1_{C>1} = K_1(S_1; \beta_1) \\ \vdots \\ r_T J_T(\beta_T; \beta_T) 1_{C=T+1} = K_T(S_T; \beta_T) \\ \vdots \\ r_1 J_1(\beta_1; \beta_1) 1_{C=T+1} = K_T(S_T; \beta_T) \end{pmatrix} = 0: \quad (A.1)$$

A.2 AIPWCC estimating equations

For Q-learning, at stage $t = T; \dots; 1$, the AIPWCC estimating equation for β_t is

$$P_n \frac{1_{C>t}}{K_t(S_t; \beta_t)} \left[Y_t - Q_t(H_t; A_t; \beta_t) - r_t Q_t(H_t; A_t; \beta_t) + \sum_{r=1}^T \frac{1_{C=r}}{K_r(S_r; \beta_r)} 1_{C=r} E_{r_t}^{h_n} \left[Q_t(H_t; A_t; \beta_t) - r_t Q_t(H_t; A_t; \beta_t) \mid S_r; r_t \right] \right] = 0: \quad (A.1)$$

A remark is that for Q-learning, the estimating equation at stage t , i.e. (A.1) belongs to J_{t+1} . The joint estimating equation is formed by stacking (A.1) for $t = 1; \dots; T$ with the estimating equations for β_T and β_{T-1} .

For OWL, at stage $t = T; \dots; 1$, the AIPWCC estimating equation for β_t is

$$P_n \frac{1_{C=T+1}}{K_T(S_T; \beta_T)} \left[r_{t+1} J_t^L(\beta_t; \beta_t) + \sum_{r=1}^T \frac{1_{C=r}}{K_r(S_r; \beta_r)} 1_{C=r} E_{r_t} \left[r_{t+1} J_t^L(\beta_t; \beta_t) \mid S_r; r_t \right] \right] = 0: \quad (A.2)$$

It can be seen that (A.2) belongs to J_{T+1} because outcome Y is required to get the outcome-based weights $W_t(H_t; A_T; Y; \beta_{t+1}; \beta_t)$ in $J_t(\beta_t; \beta_t)$ at every stage $T; \dots; 1$. Again, the joint estimating equation is formed by stacking (A.2) with the ones for β_T , β_{T-1} , and β_{T-2} .

A.3 Simulation settings

Let $\tau = (\tau_{20}; \tau_{21}; \tau_{22}; \tau_{23})^T$ and $\theta = (\theta_{20}; \theta_{21}; \theta_{22})^T$.

For $p = 2$:

$$\tau = (1; 0; 0.2; 0.2; 0.2; 0; 0; 0.1; 0.1)^T, \quad \theta = (0; 1; 0; 1; 0; 1; 0; 1)^T,$$

$$0 = \begin{matrix} B & 1 & 0.2 & C \\ @ & A & & A \end{matrix} \text{ and } 1 = \begin{matrix} B & 0.2 & 0.0 & C \\ @ & A & 0.0 & A \end{matrix}.$$

For $p = 10$; $p = 30$ and $p = 50$,

$$\tau = (1; 0.2; 0.2; 0.2; \underbrace{0; \dots; 0}_{p-2 \text{ zeros}}; 0; 0.1; 0.1; \underbrace{0; \dots; 0}_{p-2 \text{ zeros}})^T;$$

$$\theta = (0; 1; 0; 1; 0; 1; 0; 1; \underbrace{0; \dots; 0}_{p-2 \text{ zeros}})^T, \quad 0 = \begin{matrix} 1 & 0.2 & 0 & \dots & 1 \\ @ & A & & & A \\ 0.3 & 1 & 0 & \dots & 0 \\ \vdots & & & \ddots & \vdots \\ 0 & \dots & 0 & & 1 \end{matrix}, \text{ and}$$

$$1 = \text{diag}(0.2; \dots; 0.2)$$

For missing rate 65%: $\tau = (0.5; 0.3; 0.3)^T$, $\theta = (0.4; 0.2; 0.2)^T$.

For missing rate 35%: $\tau = (1.2; 0.6; 0.2)^T$, $\theta = (1.2; 0.6; 0.3)^T$.

A.4 Proof of Theorem 1

By stacking the estimating equations for τ and IPWCC estimating equations for θ , we construct a Z-estimator $(b_n; b_n)$ that solves

$$P_n \begin{matrix} 2 \\ \vdots \\ 4 \end{matrix} \begin{matrix} 1_{C=T} \tau_{T^c} g_T(S_T; \tau) \\ \vdots \\ 1_{C=1} \tau_{1^c} g_1(S_1; \tau) \end{matrix} \begin{matrix} 1_{C=T} \tau_{T^c} g_T(S_T; \tau) \text{expit } f g_T(S_T; \tau) g \\ \vdots \\ 1_{C=1} \tau_{1^c} g_1(S_1; \tau) \text{expit } f g_1(S_1; \tau) g \end{matrix} \begin{matrix} 3 \\ \vdots \\ 5 \end{matrix} = P_n m_n(S_T; Y; \tau; \theta) = 0;$$

$$w^{cc}(C; S_T; \tau) m_n(H_T; A_T; Y; \theta)$$

where $w_j^{cc}(C; S_T; \cdot)$ and $m_n(H_T; A_T; Y; \cdot)$ are constructed by aggregating $w_j^{cc}(C; S_{t-1}; \bar{S}_{t-1})$ and $\mathbb{m}_{n,j}(S_t; \cdot)$ in the order of $j \in J_{T+1} \cup \dots \cup J_T \cup \dots \cup J_1$ respectively; the \odot notation is the elementwise product operator for vectors.

Assume that the following standard regularity conditions for Z-estimator hold.

1. $P m_n(S_T; Y; \cdot; \cdot)$ exists for all $\cdot \in K \times B$; there exists $\cdot_n \in K \times B$ such that $\|P m_n(S_T; Y; \cdot; \cdot_n)\| = o(1)$, and $\|P m_n(S_T; Y; \cdot; \cdot)\| \leq o(1)$ for $\cdot \in \cdot_n$.
2. Each function in $m_n(\cdot)$ is continuous in $K \times B$ and bounded by an integrable function of the data that does not depend on $(\cdot; \cdot)$.

We show the consistency of b_n given that $t(s_t) = t(s_t; \bar{S}_t)$. Assuming the correctness of the hazard models, it follows that $K_t(s_t) = K_t(s_t; \bar{S}_t)$ for all t and s_t , for some $\cdot \in K$ and that the $b_n \rightarrow \cdot$ in probability. We only need to show

$$E \{ w_j^{cc}(C; S_{t-1}; \bar{S}_{t-1}) \mathbb{m}_{n,j}(S_t; \cdot) g \} = E \{ \mathbb{m}_{n,j}(S_t; \cdot) g \}; j \in J_t; t = 1; \dots; T+1:$$

Take the double expectation with the inner expectation condition on S_t and by MAR assumption, we have

$$\begin{aligned} E \left\{ \frac{1_{C>t-1}}{K_{t-1}(S_{t-1})} \mathbb{m}_{n,j}(S_t; \cdot) \right\} &= E \left\{ E \left\{ \frac{1_{C>t-1}}{K_{t-1}(S_{t-1})} \mathbb{m}_{n,j}(S_t; \cdot) \mid S_t \right\} \right\} \\ &= E \left\{ \mathbb{m}_{n,j}(S_t; \cdot) \frac{E(1_{C>t-1} \mid S_t)}{K_{t-1}(S_{t-1})} \right\} \\ &= E \{ \mathbb{m}_{n,j}(S_t; \cdot) g \} \end{aligned}$$

Therefore, if \cdot_n satisfies $\|P m_n(S_{T+1}; \cdot_n)\| = o(1)$, we have $\|P m_n(S_{T+1}; \cdot; b_n)\| = o_p(1)$.

A.5 Proof of Theorem 2

We prove Theorem 2 by showing the consistency in either one of two scenarios: 1. the missingness model is correctly specified, i.e. $q_t(s_t) = q_t(s_t; \theta_t)$ for all s_t for some $\theta_t \in K_t$ and $b_n \rightarrow \theta_t$ in probability; 2. the conditional expectation is correctly specified, i.e. $d_{n,t}(s_t) = d_{n,t}(s_t; \beta_t)$ for some $\beta_t \in \mathcal{A}_t$.

By stacking together the AIPWCC estimation equation for θ_t , estimating equation for β_t and τ_t , we form a unified Z-estimator representation $P_n m_n(S_{T+1}; \theta, \beta, \tau) = 0$. Under same regularity conditions as in proof of Theorem 1 for general Z-estimator, it is sufficient to show for $t = 1, \dots, T+1$, one has

$$E \left[w_j^{cc}(C; S_{t-1}; \theta_{t-1;n}) \tau_{n,j}(S_t; \theta) + \sum_{r=1}^{X-1} w_{r,j}^{aug}(C; S_r; \beta_{r;n}) d_{n,r,j}(S_r; \beta_{r;n}) \right] = E \left[\tau_{n,j}(S_t; \theta) \right] \quad \text{for all } j \in J_t \quad (\text{A.3})$$

By Lemma 10.4 of (Tsiatis, 2007), we have

$$1 - \frac{1_{C>t}}{K_t(S_t)} = \sum_{r=1}^{X-t} \frac{1_{C=r} \tau_r(S_r; \theta_r) 1_{C=r}}{K_r(S_r; \theta_r)}; \quad \text{for all } t = 1, \dots, T \quad (\text{A.4})$$

Therefore

$$\frac{1_{C>t-1}}{K_{t-1}(S_{t-1}; \theta_{t-1})} \tau_{n,j}(S_t; \theta) = \tau_{n,j}(S_t; \theta) + \left(\frac{1_{C>t-1}}{K_{t-1}(S_{t-1}; \theta_{t-1})} - 1 \right) \tau_{n,j}(S_t; \theta);$$

and the AIPWCC estimation equation on the left of (A.3) can be rearranged as

$$\tau_{n,j}(S_t; \theta) \left(\frac{1_{C>t-1}}{K_{t-1}(S_{t-1}; \theta_{t-1})} - 1 \right) + \sum_{r=1}^{X-1} \frac{1_{C=r} \tau_r(S_r; \beta_{r;n}) 1_{C=r}}{K_r(S_r; \beta_{r;n})} \tau_{n,j}(S_t; \theta) - d_{n,r,j}(S_r; \beta_{r;n}) g$$

As we assume that the first term $\mathbb{E} [f_{n,j}(S_t; \cdot)] = o(1)$, now it is sufficient to show that

$$\mathbb{E} \left[\frac{1_{C=r} \int_{\mathcal{S}_r} f_{n,j}(S_t; \cdot) d_{n,r,j}(S_r; \mathbf{b}_{r,n}) g}{K_r(S_r; \mathbf{b}_{r,n})} \right] = 0; \quad (A.5)$$

for $r = 1, \dots, t-1$ for a specific $t = 1, \dots, T+1$

We first prove for scenario 1 where the hazards models are correctly specified so that $\beta_n \rightarrow \beta$. Define a sequence of random vectors $\mathbf{G}_r = (S_t^T; 1_{C=1}; \dots; 1_{C=r-1})^T$ for $r = 1, \dots, t-1$. By the law of iterated expectations, we take conditional expectation on \mathbf{G}_r . Noting that the only unknown random variable is $1_{C=r}$ given \mathbf{G}_r , we have

$$\mathbb{E} \left[\frac{\mathbb{E}(1_{C=r} | \mathbf{G}_r) \int_{\mathcal{S}_r} f_{n,j}(S_t; \cdot) d_{n,r,j}(S_r; \mathbf{b}_{r,n}) g}{K_r(S_r)} \right]$$

Because of MAR assumption, we have

$$\mathbb{E} [1_{C=r} | \mathbf{G}_r] = P(C = r | C \leq r; S_t) 1_{C=r} = P(C = r | C \leq r; S_r) 1_{C=r} = \int_{\mathcal{S}_r} f_{n,j}(S_t; \cdot) d_{n,r,j}(S_r; \mathbf{b}_{r,n}) g$$

Thus, we prove that (A.5) holds.

Now we consider scenario 2 - the conditional expectation models are correctly specified so that $d_{n,t}(S_t) = d_{n,t}(S_t; \cdot)$. Rewrite the left hand side of (A.5) as

$$\mathbb{E} \left[\frac{1_{C=r} \int_{\mathcal{S}_r} f_{n,j}(S_t; \cdot) d_{n,t}(S_t) g}{K_r(S_r; \mathbf{b}_{r,n})} \right] = \mathbb{E} \left[\frac{\int_{\mathcal{S}_r} f_{n,j}(S_t; \cdot) d_{n,t}(S_t) g}{K_r(S_r; \mathbf{b}_{r,n})} \right]$$

Applying the law of iterated expectation again, we take conditional expectation of the first term in the above equation on $(1_{C=r}; S_r^T)^T$. This leads to

$$\mathbb{E} \left[\frac{1_{C=r}}{K_r(S_r; \mathbf{b}_{r,n})} [\mathbb{E} f_{n,j}(S_t; \cdot) | 1_{C=r}; S_r] - \mathbb{E} f_{n,j}(S_t; \cdot) | S_r \right]$$

By MAR, we have $P(S_t \perp 1_{C=r}; S_r) = P(S_t \perp S_r)$, i.e. $Ef_{\theta_{n;j}}(S_t; \cdot) \perp 1_{C=r}; S_r = Ef_{\theta_{n;j}}(S_t; \cdot) \perp S_r$. This implies the first term in is zero. Similarly, we take conditional expectation of the second term on $(1_{C=r}; S_r^T)^T$ and obtain

$$E \frac{r(S_r; \cdot) 1_{C=r}}{K_r(S_r; \cdot)} [Ef_{\theta_{n;j}}(S_t; \cdot) \perp 1_{C=r}; S_r - Ef_{\theta_{n;j}}(S_t; \cdot) \perp S_r] = 0$$

By MAR again, we have $P(S_t \perp 1_{C=r}; S_r) = P(S_t \perp S_r)$, which implies the second term in is zero.

Therefore we show that (A.5) holds under the two scenarios and the doubly robust conclusion follows.

APPENDIX B

APPENDIX FOR CHAPTER 3

B.1 Additional treatment calibration

We have developed methods that use calibration to enforce the covariate balance between the RCT sample and the target population. We can also apply the calibration technique to balance the covariate distributions between the two treatment groups in the RCT. Toward this end, we further assign weight w_{1i} to each subject i in the treatment group with $A_i = 1$ and assign weight w_{0i} to subject i in the treatment group with $A_i = 0$. The weights w_{1i} and w_{0i} are obtained by solving the following optimization problem:

$$\begin{aligned} \min_{w_{1i}; w_{0i}; \delta_i} & \sum_{i \in 2A} I(A_i = 1) w_{1i} \log w_{1i} + I(A_i = 0) w_{0i} \log w_{0i} \\ \text{subject to } w_{ai} & \geq 0; \sum_{i \in 2A} I(A_i = a) \frac{w_{ai}}{a^{A_i} + (1-a)(1-A_i)} = 1 \quad (a = 0; 1); \\ & \sum_{i \in 2A} I(A_i = 1) \frac{w_{1i}}{A_i} g(X_i) = \sum_{i \in 2A} I(A_i = 0) \frac{w_{0i}}{1-A_i} g(X_i); \end{aligned}$$

Similar as in the estimation of Q in (3.3), we can estimate w_{ai} using Lagrangian multiplier as

well and obtain w_{ai} for $a = 0; 1$ and $i = 1; \dots; n$. We refer to the estimators that calibrate both the sampling and the treatment covariate balance as the Double Covariate Balancing (DCB) estimators. The HT-type DCB estimator is

$$\hat{\Delta}^{DCB0} = \sum_{i=1}^n q \frac{w_{1i} A_i Y_i}{A_i} - \frac{w_{0i} (1 - A_i) Y_i}{1 - A_i} :$$

The Hajek-type DCB estimator is

$$\hat{\Delta}^{DCB1} = \sum_{i=1}^n q \frac{p \frac{w_{1i} A_i Y_i}{\sum_{i=1}^n q w_{1i} A_i}}{p \frac{w_{0i} (1 - A_i) Y_i}{\sum_{i=1}^n q w_{0i} (1 - A_i)}} :$$

Similar weights can be added to the ACW estimator as well. It is straightforward, so we omit the formula here.

We examined the performances of the DCB estimators through simulation study, where we found that their improvement over the CW estimators are not evident. This is likely due to that the RCT sample already has a good balance between the two treatment groups. Therefore, we exclude the DCB estimators from the main text. However, we recommend the practitioners to consider using the DCB estimators if there is a belief on the unbalance of covariate distributions between the two treatment groups.

B.2 Proofs

B.2.1 Proof of Theorem 3

Proof of the double robustness of the CW estimators

Let $g_0 = E[f g(X) | g]$, $g_0 = g(X) - g_0$. To use the M-estimator theory (Boos & Stefanski, 2013), we write (3.5) as the following estimating equations

$$\frac{1}{N} \sum_{i=1}^N C(X_i; e_i; g) = \frac{1}{N} \sum_{i=1}^N e_i d_i f g(X_i) - g g = 0; \quad (\text{B.1})$$

$$\frac{1}{N} \sum_{i=1}^N (X_i; i; ; g) = \frac{1}{N} \sum_{i=1}^N i \exp^{-\eta} g(X_i) - f g(X_i) - g g = 0; \quad (\text{B.2})$$

First consider the case where Assumption 5 holds, we have $(X) = \exp^{-\eta} g(X) | g$. Notice that g_0 is the solution to $E[f C(X; g)] = 0$. Taking expectation on the left hand side of (B.2) with $g = g_0$ leads to

$$E[f (X; ; ; g_0) | g_0] = E[(X) \exp^{-\eta} g(X) | g_0] - E[f g(X) | g_0] :$$

For the above conditional expectation to be zero, one needs $(X) \exp^{-\eta} g(X) | g$ to be constant. As $(X) = \exp^{-\eta} g(X) | g$, we have $(X) \exp^{-\eta} g(X) | g = \exp^{-\eta} g(X) | g$. Thus $= 0$ makes (B.2) a system of unbiased estimating equations. We point out that denominator in η is

an estimator of population size N , i.e.

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \exp^{b > g(X_i)} g &= \frac{1}{N} \sum_{i=1}^N \exp^{b > g(X_i)} g \\ &= \frac{1}{N} \sum_{i=1}^N \exp^{n \geq 0} g(X_i) + O_p(n^{-1/2} N^{-1}) \\ &= 1 + O_p(N^{-1/2}) + O_p(n^{-1/2} N^{-1}) \\ &= 1 + o_p(1): \end{aligned}$$

Therefore,

$$\hat{q} = q(X_i; b) = \frac{\sum_{i=1}^n \exp^{b > g(X_i)} g}{\sum_{i=1}^n \exp^{b > g(X_i)}} = \frac{1}{N} \frac{1}{(X_i; 0)} + O_p(n^{-1/2} N^{-1}); \quad (B.3)$$

i.e. $\hat{q} \xrightarrow{P} (X_i; 0)g^{-1}$ as $n \rightarrow \infty$. Based on (B.3), we have

$$\begin{aligned} \hat{\Delta}^{CW0} &= \sum_{i=1}^N \hat{q}_i \frac{A_i Y_i}{A_i} \frac{(1 - A_i) Y_i}{1 - A_i} \\ &= \frac{1}{N} \sum_{i=1}^N \frac{1}{(X_i; 0)} \frac{A_i Y_i}{A_i} \frac{(1 - A_i) Y_i}{1 - A_i} = 0 + O_p(N^{-1/2}) + O_p(n^{-1/2}) \\ &= 0 + o_p(1): \end{aligned} \quad (B.4)$$

Therefore, $\hat{\Delta}^{CW0}$ is consistent for 0 .

The Hajek-type CW estimator assumes that $A_i = 1$; for all i in the RCT. Under this assumption, based on (B.3), we have

$$\sum_{i=1}^N \hat{q}_i A_i = \frac{1}{N} \sum_{i=1}^N \frac{1 \cdot A_i}{(X_i; 0)} + O_p(n^{-1/2}) = 1 + O_p(N^{-1/2}) + O_p(n^{-1/2}) = 1 + o_p(1);$$

and thus $\sum_{i=1}^N \hat{q}_i (1 - A_i) = 1 - 1 + o_p(1)$. Based on (B.4), it follows that $\hat{\Delta}^{CW1}$ is consistent

for θ_0 .

Now consider the case where Assumption 6 holds. Then we have

$$\begin{aligned}
 E \sum_{i=1}^N \mathbf{q} \frac{A_i Y_i}{A_i} \frac{(1 - A_i) Y_i}{1 - A_i} &= E \sum_{i=1}^N \mathbf{q} E \left[\frac{A_i Y_i}{A_i} \frac{(1 - A_i) Y_i}{1 - A_i} \mid X_i \right] \\
 &= E \sum_{i=1}^N \mathbf{q} E \left[Y(1) - Y(0) \mid X_i \right] \\
 &= E \sum_{i=1}^N \mathbf{q} g(X_i) = E \frac{1}{N} \sum_{i=1}^N \mathbf{q} g(X_i) \\
 &= E \sum_{i=1}^N \mathbf{q} g(X_i) = 0;
 \end{aligned}$$

where the equation on the third line is obtained by the balancing constraint (3.2). Under mild regularity conditions for unbiased M-estimators, $\hat{\theta}^{CW0}$ is consistent for θ_0 .

Again, for the Hajek CW estimator,

$$E \left(\sum_{i=1}^N \mathbf{q} A_i \right) = E \left(\sum_{i=1}^N \mathbf{q} E(A_i \mid X_i) \right) = E \left(\sum_{i=1}^N \mathbf{q} E(A_i \mid X_i) \right) = 1;$$

And similarly, $E \sum_{i=1}^N \mathbf{q} (1 - A_i) = 1 - 1 = 0$. Under regularity conditions for unbiased M-estimators, we have $\sum_{i=1}^N \mathbf{q} A_i \rightarrow 1$ as $n \rightarrow \infty$. Combining this with the consistency results for $\hat{\theta}^{CW0}$, $\hat{\theta}^{CW1}$ is consistent for θ_0 . We thus conclude the double robustness of both $\hat{\theta}^{CW0}$ and $\hat{\theta}^{CW1}$.

Proof of the asymptotic variances for the CW estimators

We compare the asymptotic variances of $\hat{\theta}^{CW0}$ and $\hat{\theta}^{CW1}$ under Assumption (5) and (6), because the comparison under the doubly robust condition (i.e. either Assumption (5) or Assumption (6) holds) is cumbersome and unavailing.

Let $\theta = (\beta; \gamma; \alpha; \lambda)^T$ to denote all parameters. The estimating function for θ is

$$\begin{aligned}
 \psi(X; A; Y; \theta; e) = & \begin{pmatrix} C(X; e; \gamma) \\ (X; \beta; \gamma) \\ h(X; A; \alpha; \lambda) \\ t(X; A; Y; \alpha; \lambda) \end{pmatrix};
 \end{aligned}$$

where the first two functions are (B.1), (B.2) and

$$\begin{aligned}
 h(X; A; \alpha; \lambda) &= \exp\{\beta g(X)g(A; \lambda)\}; \\
 t(A; X; Y; \alpha; \lambda) &= \exp\{\beta g(X)g\} \frac{AY}{A} - \frac{(1-A)Y}{1-A} :
 \end{aligned}$$

Then $\hat{\theta} = (\hat{\beta}; \hat{\gamma}; \hat{\alpha}; \hat{\lambda}^{CW1})^T$ solves the joint estimating equation

$$\frac{1}{N} \sum_{i=1}^N \psi(X_i; A_i; Y_i; \hat{\theta}) = 0;$$

Under standard regularity conditions in the M-estimator theory, we have

$$\hat{\theta} - \theta_0 = \frac{1}{N} \sum_{i=1}^N \psi(X_i; A_i; Y_i; \theta_0) + o_p(N^{-1/2});$$

where $A(\theta_0) = E \psi(\theta_0)$, and $\theta_0 = (\beta_0; \gamma_0; \alpha_0; \lambda_0)^T$. The asymptotic variance of $N^{-1/2}(\hat{\theta} - \theta_0)$ is $A^{-1}(\theta_0)B(\theta_0)A^{-1}(\theta_0)$, where $B(\theta_0) = E \psi(\theta_0) \psi(\theta_0)^T$.

To further express the asymptotic variance, we denote $\alpha_0 = \alpha_0(X) = \exp\{\beta_0 g(X)g\}$ and $(Y; A) = f AY = A(1-A)Y = (1-A) \alpha_0 g$. Note that $E f(Y; A) | X = 1 g = (X) \alpha_0$ and $E \psi = 1$. Under Assumption 5 and 6, $(X) \alpha_0(X) = 1$ and $(X) \alpha_0 = \beta_0 g_0$. In the following derivation we use $\stackrel{a}{=}$ to indicate equality when both Assumption 5 and 6 hold.

Using iterated expectation, we have

$$\begin{aligned}
 A(\theta_0) &= E \left[f'(\theta_0) g(\theta_0) \right] \\
 &= E \begin{pmatrix} 0 & 0_{1 \times K} & 0_{1 \times K} & 0_{1 \times K} \\ 0_{1 \times K} & q_0 I_K & q_0 g_0 g_0' & 0_{1 \times K} \\ 0_{1 \times K} & q_0 (A - A) g_0' & q_0 & 0 \\ 0_{1 \times K} & q_0 (Y; A) g_0' & q_0 \left(\frac{AY}{2} + \frac{(1-A)Y}{(1-A)^2} \right) & q_0 \end{pmatrix} \\
 &= E \begin{pmatrix} I_K & 0_{K \times K} & 0_{K \times 1} & 0_{K \times 1} \\ E(q_0) I_K & E f q_0 (g_0 g_0') g_0' & 0_{K \times 1} & 0_{K \times 1} \\ 0_{1 \times K} & 0_{1 \times K} & E(q_0) & 0 \\ 0_{1 \times K} & E(X) q_0 f'(X) g_0 g_0' & E q_0 \left(\frac{AY}{2} + \frac{(1-A)Y}{(1-A)^2} \right) & E(q_0) \end{pmatrix}
 \end{aligned}$$

By block matrix inversion,

$$A(\theta_0)^{-1} = \begin{pmatrix} 0 & 0_{1 \times K} & 0_{1 \times K} & 0_{1 \times K} \\ E(q_0) I_K & E f q_0 (g_0 g_0') g_0' & 0_{K \times 1} & 0_{K \times 1} \\ 0_{1 \times K} & 0_{1 \times K} & E(q_0)^{-1} & 0 \\ A_{41} & A_{42} & A_{43} & A_{44} \end{pmatrix}^{-1}$$

where

$$\begin{aligned}
 A_{41} &= E \left[q_0(X) f'(X) g_0 g_0' \right] E \left[q_0(X) g_0 g_0' \right]^{-1} I_K; \\
 A_{42} &= E \left[f q_0(X) g_0' \right] E \left[(X) q_0(X) f'(X) g_0 g_0' \right] E \left[q_0(X) g_0 g_0' \right]^{-1} I_K; \\
 A_{43} &= E \left[f q_0(X) g_0' \right] E \left[q_0(X) \left(\frac{AY}{2} + \frac{(1-A)Y}{(1-A)^2} \right) \right] E \left[q_0(X) \left(\frac{Y(1)}{A} + \frac{Y(0)}{1-A} \right) \right]; \\
 A_{44} &= E \left[q_0(X) g_0' \right]^{-1} I_K;
 \end{aligned}$$

Under Assumption 5 and 6, we have

$$A_{41}B_{11}A_{41}^> = E \left[\frac{h}{n} \frac{f(X)}{q_0(X)} \right] \frac{1}{n} ;$$

$$(A_{42}B_{22} + A_{44}B_{24}^>)A_{42}^> = 0 ;$$

$$A_{42}B_{24} = E \left[\frac{h}{n} q_0(X) f(X) \right] \frac{1}{n} ;$$

$$A_{43}B_{34} = E \left[\frac{f(Y(1) - A + Y(0) = (1 - A)g)}{q_0^2(X)} f((1 - A)Y(1) + AY(0))g \right] ;$$

$$A_{43}^2B_{33} = [E \left[\frac{f(Y(1) - A + Y(0) = (1 - A)g)}{q_0(X)} \right]]^2 \frac{1}{n} E \left[\frac{f(Y(1) - A + Y(0) = (1 - A)g)}{q_0(X)} \right] ;$$

Therefore, V^{CW0} can be simplified as

$$V^{CW0} = E \left[\frac{h}{n} \frac{f(X)}{q_0(X)} \right] \frac{1}{n} + E \left[\frac{q_0^2(X)}{q_0(X)} \left(\frac{Y(1)^2}{A} + \frac{Y(0)^2}{1 - A} \right) \right] \frac{1}{n} ;$$

The difference in the asymptotic variance of the two CW estimators is given by

$$V^{CW1} - V^{CW0} = E \left[\frac{Y(1)}{A} + \frac{Y(0)}{1 - A} \right]^2 \frac{1}{n} E \left[\frac{f(Y(1) - A + Y(0) = (1 - A)g)}{q_0(X)} \right] - E \left[\frac{h}{n} \frac{f(X)}{q_0(X)} \right] \frac{1}{n} \\ - 2E \left[\frac{Y(1)}{A} + \frac{Y(0)}{1 - A} \right] E \left[\frac{q_0^2(X)}{q_0(X)} \left(\frac{Y(1)}{A} + \frac{Y(0)}{1 - A} \right) \right] \frac{1}{n} ;$$

We now compare the asymptotic variances under the special case where $\alpha = 0.5$. Define the sum of two potential outcomes to be $Y^+ = Y(1) + Y(0)$. The asymptotic variances are

$$V^{CW0} = E \left[\frac{h}{n} \frac{f(X)}{q_0(X)} \right] \frac{1}{n} + 2E \left[\frac{q_0^2(X)}{q_0(X)} \left(\frac{Y(1)^2}{A} + \frac{Y(0)^2}{1 - A} \right) \right] \frac{1}{n} ;$$

$$V^{CW1} = V^{CW0} + E \left[\frac{Y^+}{2} \right]^2 E \left[\frac{f(Y(1) - A + Y(0) = (1 - A)g)}{q_0(X)} \right] - E \left[\frac{h}{n} \frac{f(X)}{q_0(X)} \right] \frac{1}{n} - 2E \left[\frac{Y^+}{2} \right] E \left[\frac{q_0^2(X)}{q_0(X)} Y^+ \right] \frac{1}{n} ;$$

and their differences is

$$V^{CW1} - V^{CW0} = E \left[\frac{Y^+}{2} \right]^2 E \left[\frac{f(Y(1) - A + Y(0) = (1 - A)g)}{q_0(X)} \right] - E \left[\frac{h}{n} \frac{f(X)}{q_0(X)} \right] \frac{1}{n} - 2E \left[\frac{Y^+}{2} \right] E \left[\frac{q_0^2(X)}{q_0(X)} Y^+ \right] \frac{1}{n} ;$$

It is important to note that the difference can be negative, zero, or positive. Therefore, there is no guarantee that which one of the two estimators is better than the other. An observation is that the expected value of Y^+ , which is the sum of the two potential outcomes, plays an important role in the asymptotic variance. In the simulation studies provided in the Section B.4, we demonstrate that one outperforms the other in different scenarios, which confirms our findings.

B.2.2 Proof of Theorem 4

Let $Z = (X; A; Y; ; \theta)$ be a vector of random variables. Assumptions 1 - 3 and 7 constitute the semiparametric model. The semiparametric likelihood based on a single Z is

$$f(Z) = f(X) f(A; Y|X; \theta) g(X) f(A; Y|X; \theta = 1) e(X)^{\theta};$$

where $f(\cdot)$ is a density function for a continuous random variable and is a probability mass function for a discrete random variable.

Assuming that $\theta = 0$, the score function (Hahn, 1998) satisfies

$$S(X; A; Y; ; \theta) = S(X; A; Y;) + S(X; A; Y; \theta);$$

We first list four identities that are used in the following derivation of the efficiency bound:

1. For any function $h(X; A;)$, we have $E f h(X; A;) S(Y|X; A;) g = 0$;
2. For any function $h(X; A;)$, we have $E [h(X; A;) f Y - E(Y|X; A;) g] = 0$;
3. For any $h(X; A; Y)$, if $E \frac{\partial}{\partial \theta} h(X; A; Y) = 0$, we have $E f \frac{\partial}{\partial \theta} h(X; A; Y) S(X; A; Y;) g = 0$;
4. For any $h(X; A; Y)$, if $E f h(X; A; Y) g = 0$, we have $E \frac{\partial}{\partial \theta} h(X; A; Y) S(X; A; Y; \theta) = 0$;

To derive the semiparametric efficiency score, we use the method of parametric submodel (Bickel et al., 1993). Let $f_t(Z) : t \in \mathbb{R}$ be a regular parametric submodel which contains the truth at $t = 0$, i.e. $f_t(Z)|_{t=0} = f(Z)$:

Note that $\theta_t = E(Y|X; A = 1; t) - E(Y|X; A = 0; t)$ and $\theta_0 = E(Y|X)$. Let $\dot{\theta}_t = E \left[\frac{\partial}{\partial t} \theta_t(X) \right]$ denote the parameter evaluated with respect to the regular parametric submodel $f_t(Z)$. Following Bickel et al. (1993), the semiparametric efficiency score $S(Z)$ is the pathwise derivative of the target parameter in the sense that

$$\frac{\partial}{\partial t} \theta_t \Big|_{t=0} = E f S(Z);$$

where $S(Z) = \frac{\partial}{\partial t} \log f_t(Z) \Big|_{t=0}$. Toward this end, we express

$$\frac{\partial}{\partial t} \theta_t \Big|_{t=0} = E \left[\frac{\partial}{\partial t} \theta_t(X) S(X) \right] + E \left[\frac{\partial \theta_t(X)}{\partial t} \Big|_{t=0} \right] \quad (B.5)$$

To express (B.5) further,

$$\begin{aligned} \frac{\partial \theta_t(X)}{\partial t} \Big|_{t=0} &= \int_Z y \frac{\partial}{\partial t} f_t(y|X; A=1) dy - \int_Z y \frac{\partial}{\partial t} f_t(y|X; A=0) dy \\ &= \int_Z y S(y|X; A=1) f_t(y|X; A=1) dy - \int_Z y S(y|X; A=0) f_t(y|X; A=0) dy \\ &= E \left[\frac{AY}{(X)_A} S(Y|A; X) \right] - E \left[\frac{(1-A)Y}{(X)(1-A)} S(Y|A; X) \right] \\ &= E \left[\frac{AY}{(X)_A} - \frac{(1-A)Y}{(X)(1-A)} \right] S(Y|A; X) \end{aligned}$$

Therefore,

$$\begin{aligned}
 E \frac{\partial_t(X)}{\partial t} \Big|_{t=0} &= E \frac{AY}{(X)} \frac{(1-A)Y}{(1-A)} S(Y|A; X; \cdot) \\
 &= E \frac{AfY}{(X)} \frac{1(X)g}{A} \frac{(1-A)fY}{1-A} \frac{o(x)g}{A} S(Y|A; X; \cdot) \quad (B.6) \\
 &= E \frac{AfY}{(X)} \frac{1(X)g}{A} \frac{(1-A)fY}{1-A} \frac{o(x)g}{A} S(X; A; Y; \cdot) \quad (B.7) \\
 &= E \frac{AfY}{(X)} \frac{1(X)g}{A} \frac{(1-A)fY}{1-A} \frac{o(x)g}{A} S(X; A; Y; \cdot; \epsilon) : \\
 & \hspace{20em} (B.8)
 \end{aligned}$$

In the above derivation, (B.6) follows by identity 1, (B.7) follows by identity 2, and (B.8) follows by identity 3 and 4. Therefore,

$$\begin{aligned}
 E \frac{\partial^n}{\partial t^n} (X) S(X) &= E \frac{\partial^n}{\partial t^n} f(X) \quad \text{og} S(X) \\
 &= E \frac{\partial^n}{\partial t^n} f(X) \quad \text{og} S(X) + S(X; A; Y; \cdot) + S(X; A; Y; \epsilon) \quad (B.9) \\
 &= E \frac{\partial^n}{\partial t^n} f(X) \quad \text{og} S(X; A; Y; \cdot; \epsilon) ;
 \end{aligned}$$

where (B.9) holds because

$$\begin{aligned}
 E \frac{\partial^n}{\partial t^n} f(X) \quad \text{og} S(X; A; Y; \epsilon) &= E \frac{\partial^n}{\partial t^n} f(X) \quad \text{og} S(X; A; Y; \epsilon = 1) \\
 &= E \frac{\partial^n}{\partial t^n} f(X) \quad \text{og} E S(X; A; Y; \epsilon = 1) X; \epsilon \\
 &= E \frac{\partial^n}{\partial t^n} f(X) \quad \text{og} E S(X; A; Y; \epsilon = 1) X; \epsilon = 1 \\
 &= 0:
 \end{aligned}$$

Substituting back to (B.5), we have

$$\frac{\partial}{\partial t} \Big|_{t=0} = E \left[\psi f(X) \right] + \frac{A f(Y=1|X)g}{(X)} - \frac{(1-A) f(Y=0|X)g}{1-A} S(X; A; Y; ; \epsilon)$$

Thus, the semiparametric efficiency score is

$$= \psi f(X) + \frac{A f(Y=1|X)g}{(X)} - \frac{(1-A) f(Y=0|X)g}{(1-A)}$$

It follows that the semiparametric efficiency bound is

$$E(\sigma^2) = E \left[\psi^2 f(X) \right] + \frac{V f(Y=1|X)g}{A} + \frac{V f(Y=0|X)g}{1-A}$$

B.2.3 Proof of Theorem 5

Proof of the double robustness of the ACW estimator

Let $\theta = (\alpha; \beta; \gamma)$ denote the vector of nuisance parameters. Note that θ^{ACW} is the solution to the estimating equation $N^{-1} \sum_{i=1}^N \psi(X_i; A_i; Y_i; \theta; \epsilon; \gamma) = 0$, where

$$\psi(X; A; Y; ; \epsilon; ; \gamma) = \frac{A f(Y=1|X; \gamma)g}{(X; \gamma)} - \frac{(1-A) f(Y=0|X; \gamma)g}{1-A} + \psi f_1(X; \gamma) - \psi f_0(X; \gamma)$$

Let θ_n be the probability limits of θ_n . It suffices to show that $E \left[\psi(X; A; Y; ; \epsilon; \theta_n; \gamma) \right] = 0$ if either $\psi(X; \gamma)$ or $\psi_a(X; \gamma) (a = 0; 1)$ is correctly specified. Under standard regularity conditions

for M-estimators, $\hat{\theta}^{ACW}$ is consistent for θ_0 . Use iterated expectation, we can write

$$E \left[\frac{1}{n} \sum_{i=1}^n \psi(X_i; A_i, Y_i; \theta; \theta_0) \right] = E \left[\frac{AY}{(X; \theta)} - \frac{(1-A)Y}{1-A} \theta_0 + E \left[\frac{\psi(X; \theta)}{(X; \theta)} \mid f_1(X; \theta_1) \theta_0(X; \theta_0)g \right] \right] : (B.10)$$

The first term on the left-hand-side of (B.10) is 0 either one of the $(X; \theta)$ or $\theta_0(X; \theta_0)$ ($a = 0; 1$) is correctly specified, as shown in the proof of consistency in the CW estimators. Now consider the second term on the left-hand-side of (B.10).

Firstly, if $(X; \theta)$ is correctly specified, we have $(X; \theta) = (X)$. Take iterated expectation conditional on X , we have the second term on the left-hand-side of (B.10)

$$E \left[\frac{\psi(X; \theta)}{(X; \theta)} \mid f_1(X; \theta_1) \theta_0(X; \theta_0)g \right] = E \left[f_1(X; \theta_1) \theta_0(X; \theta_0)g E \left[\frac{\psi(X; \theta)}{(X; \theta)} \mid X \right] \right] = 0;$$

as $E \left[\frac{\psi(X; \theta)}{(X; \theta)} \mid X \right] = 0$. Thus, (B.10) equals to zero.

Secondly, if outcome model $\theta_0(X; \theta_0)$ ($a = 0; 1$) is correctly specified, we have $f_1(X; \theta_1) \theta_0(X; \theta_0) = \theta_0(X)$. Then the second term on the left-hand-side of (B.10) satisfies

$$E \left[\frac{\psi(X; \theta)}{(X; \theta)} \mid f_1(X; \theta_1) \theta_0(X; \theta_0)g \right] = E \left[\frac{\psi(X; \theta)}{\theta_0(X)} \mid \theta_0 \right] = E \left[\frac{\psi(X; \theta)}{\theta_0(X)} \mid \theta_0 \right] = 0$$

by the balancing constraint (3.2). Thus, (B.10) equals to zero under this scenario as well. This completes the proof of the double robustness of $\hat{\theta}^{ACW}$.

B.2.4 Proof of Theorem 5 and Theorem 6

Proof of local efficiency

Following the empirical process literature, let P_N denote the empirical measure. For a random variable V , $Pf(V)g = \int f(v)g dP$ is the expectation of $f(V)$ under the true data-generating process. Recall that $Z = (X; A; Y; ; \Theta)$, $\theta = (\theta_1; \theta_0; \theta_1)$, θ_0 is the probability limits of $\hat{\theta}_0$ and θ_1 is the corresponding true parameter value. Let

$$\begin{aligned} \mathbb{E}(Z; \theta) &= \frac{AfY}{(X; \theta)} \frac{1(X; \theta_1)g}{A} - \frac{(1-A)fY}{1-A} \frac{\theta_0(X; \theta_0)g}{A} \\ &\quad + \mathbb{E}f(1(X; \theta_1) - \theta_0(X; \theta_0))g \\ &= \frac{AfY}{(X; \theta)} \frac{1(X; \theta_1)g}{A} + \mathbb{E}f(1(X; \theta_1) \\ &\quad - \frac{(1-A)fY}{(X; \theta)} \frac{\theta_0(X; \theta_0)g}{1-A} - \mathbb{E}f\theta_0(X; \theta_0))g \\ &=: \theta_1(Z; \theta) - \theta_0(Z; \theta): \end{aligned}$$

Under the conditions specified in Theorem 5 or the conditions specified in Theorem 6, and assume that $(Z; \theta)$ belongs to Donsker classes (Vaart & Wellner, 1996; Kennedy, 2016), $P\theta_1(Z; \theta) = \theta_1$, $P\theta_0(Z; \theta) = \theta_0$ and $P(Z; \theta) = \theta_1 - \theta_0 = \theta$. Thus,

$$\begin{aligned} \hat{\theta}_0^{ACW} &= P_N(Z; \hat{\theta}) - P(Z; \hat{\theta}) \\ &= (P_N - P)(Z; \hat{\theta}) + Pf(Z; \hat{\theta}) - (Z; \hat{\theta})g \\ &= (P_N - P)(Z; \hat{\theta}) + Pf(Z; \hat{\theta}) - (Z; \hat{\theta})g + o_p(N^{-1/2}): \end{aligned} \quad (B.11)$$

We now show that

$$Pf(Z; \hat{\theta}) - (Z; \hat{\theta})g = Pf(\theta_1(Z; \hat{\theta}) - \theta_1(Z; \hat{\theta}))g - Pf(\theta_0(Z; \hat{\theta}) - \theta_0(Z; \hat{\theta}))g$$

is a small order term under conditions in Theorem 5 or Theorem 6. We write

$$\begin{aligned}
 \text{Pf}_1(Z; b) - \text{Pf}_1(Z; \hat{b}) &= P \frac{1}{n} \frac{\sum_{i=1}^n Y_i \phi_1(X_i; b_1)}{\phi_1(X; b)} + o_p(n^{-1/2}) \\
 &= P \frac{1}{n} \frac{\sum_{i=1}^n Y_i \phi_1(X_i; b_1)}{\phi_1(X; b)} + o_p(n^{-1/2}) \\
 &= P \frac{\int Y \phi_1(X; b_1) dP}{\int \phi_1(X; b) dP} + o_p(n^{-1/2}) :
 \end{aligned}$$

Similarly, we have

$$\text{Pf}_0(Z; b) - \text{Pf}_0(Z; \hat{b}) = P \frac{\int \phi_0(X) dP}{\int \phi_0(X; b_0) dP} + o_p(n^{-1/2}) :$$

Therefore, by Cauchy-Schwarz inequality and the positivity of $\phi_1(X; b)$, $|\text{Pf}_1(Z; \hat{b}) - \text{Pf}_1(Z; b)|$ is bounded above by

$$k(X) - \int \phi_1(X; b) dP \leq \sqrt{k_a(X) - \int \phi_1(X; b_a) dP} : \quad (\text{B.12})$$

Under the conditions in Theorem 5, if $\phi_1(X; \cdot)$ is a correctly specified parametric model for $\phi_1(X)$, then $k(X) - \int \phi_1(X; b) dP = O_p(n^{-1/2})$; and if $\phi_a(X; \cdot)$ is a correctly specified parametric model for $\phi_a(X)$, then $k_a(X) - \int \phi_a(X; b_a) dP = O_p(n^{-1/2})$. Therefore, the product (B.12) is $O_p(n^{-1})$, which makes $|\text{Pf}_1(Z; \hat{b}) - \text{Pf}_1(Z; b)|$ in (B.11) asymptotically negligible. Under the conditions in Theorem 6, the product (B.12) is $o_p(n^{-1/2})$ and therefore the term $|\text{Pf}_1(Z; \hat{b}) - \text{Pf}_1(Z; b)|$ in (B.11) is asymptotically negligible. The result follows.

B.3 Conditions for the sieves estimator

Following Hirano et al. (2003), we assume the following regularity conditions on the data generating process and the nuisance functions.

Condition 3 (Distribution of X) Let $X \subseteq \mathbb{R}^p$ be the support of X . Assume that X is a Cartesian product of compact intervals, i.e. $X = \prod_{j=1}^p [l_j; u_j]$, $l_j; u_j \in \mathbb{R}$. The density of X , $f(X)$, is bounded above and below away from 0 on X .

Condition 4 (Basis functions) There exist constant l and u such that

$$l \min_{\mathbf{x}} \frac{\mathbf{1}^T \mathbf{g}(\mathbf{x})}{\mathbf{1}^T \mathbf{g}(\mathbf{x})} \leq \frac{\mathbf{1}^T \mathbf{g}(\mathbf{x})}{\mathbf{1}^T \mathbf{g}(\mathbf{x})} \leq u \max_{\mathbf{x}} \frac{\mathbf{1}^T \mathbf{g}(\mathbf{x})}{\mathbf{1}^T \mathbf{g}(\mathbf{x})}$$

almost surely where λ_{\min} and λ_{\max} denote the minimum and maximum eigenvalues of a matrix.

Condition 5 (Potential outcomes) The second moment of the potential outcomes are finite. i.e. $E[Y(a)^2] < \infty$, for $a = 0; 1$.

Condition 6 (Smoothness) The log sampling score function $\log \pi(\mathbf{x})$ is s -times continuously differentiable and the outcome mean function $\mu_a(\mathbf{x})$ is s_a -times continuously differentiable, $\mathbf{x} \in X$; $a = 0; 1$; The sieve estimators of $\log \pi(X)$ and $\mu_a(X)$ use a power series and the smoothness condition is $s = 2p + 1$, for $s = s$ and $s = s_a$ respectively.

B.4 Additional simulation study

We conduct additional simulation studies to evaluate the effect of potential outcome model and sample sizes on the performance of our proposed estimators.

B.4.1 Comparison of the CW estimators

In this simulation study, we compare the performance of $\hat{\mu}^{CW0}$ and $\hat{\mu}^{CW1}$. Under the data generating mechanism in Section 3.5, $\hat{\mu}^{CW0}$ is more efficient than $\hat{\mu}^{CW1}$. In this section, the simulation setting is the same as that in Section 3.5, except that we generate the potential

outcomes according to

$$Y(a) | X = 100 + 27.4X_1 + a + 13.7X_2 + 13.7X_3 + 13.7X_4 + \epsilon;$$

where $\epsilon \sim N(0, 1)$ for $a = 0, 1$. Figure B.1 reports the simulation results. Under this setting, it can be seen that the variances of $\hat{\tau}^{CW0}$ become dramatically larger in all four scenarios than the variances of $\hat{\tau}^{CW1}$. Combining with the simulation results in Section 3.5, we show that $\hat{\tau}^{CW0}$ and $\hat{\tau}^{CW1}$ can outperform each other in different scenarios.

B.4.2 Increased sample sizes

In this simulation study, we increase the population size to $N = 500000$, the RWE sample size to $m = 10000$, and the RCT selection rate to 0.45% that leads to $n = 2250$ to study the effect of sample size. Results are shown in Figure B.2 and Table B.1 for continuous outcome and . It can be seen that when the sample sizes are larger, the advantages of our proposed estimators become more obvious. The coverage rate for both continuous outcome and binary outcome are closer to the nominal rate.

Figure B.1 Boxplot of estimators under four model specification scenario: $\hat{\Lambda}^{CW0}$ is worsen than $\hat{\Lambda}^{CW1}$.

Table B.1 Simulation results for continuous outcome with population size $N = 500000$: bias of point estimates, Monte Carlo variance, relative bias of bootstrap variance estimate, and the coverage rate of 95% Wald confidence interval.

	Naive	IPSW	CW0	CW1	ACW	ACW(S)
	Bias					
1. O:C/S:C	13:64	0:09	0:07	0:00	0:05	0:05
2. O:C/S:W	11:80	1:20	0:06	0:08	0:06	0:06
3. O:W/S:C	11:39	0:10	0:03	0:04	0:02	0:00
4. O:W/S:W	11:01	2:68	1:59	1:56	1:58	0:02
	Monte Carlo variance					
1. O:C/S:C	1:53	5:23	1:02	4:42	0:08	0:08
2. O:C/S:W	1:38	2:98	0:74	2:79	0:08	0:08
3. O:W/S:C	1:51	8:87	4:10	7:34	3:23	0:09
4. O:W/S:W	1:42	2:98	0:88	2:84	0:48	0:10
	Relative bias (%) of bootstrap variance estimate					
1. O:C/S:C	2:3	0:6	4:4	4:3	1:2	0:1
2. O:C/S:W	4:3	1:0	8:2	3:1	2:1	2:5
3. O:W/S:C	1:1	0:3	11:0	9:5	13:7	1:2
4. O:W/S:W	1:9	0:8	12:6	4:4	2:8	1:4
	95% Wald CI coverage					
1. O:C/S:C	0:0	94:0	93:6	93:2	94:6	94:6
2. O:C/S:W	0:0	88:0	95:1	94:5	93:9	93:8
3. O:W/S:C	0:0	92:3	94:1	91:8	91:7	94:9
4. O:W/S:W	0:0	62:3	63:2	81:6	33:7	94:9

Figure B.2 Boxplot of estimators for continuous outcome with population size $N = 500000$ under four model specification scenarios.

Table B.2 Simulation results for continuous outcome with population size $N = 500000$: bias of point estimates, Monte Carlo variance, relative bias of bootstrap variance estimate, and the coverage rate of 95% Wald confidence interval.

	Naive	IPSW	CW0	CW1	ACW	ACW(S)
	Bias 1000					
1. O:C/S:C	94:75	0:49	1:08	0:22	0:22	0:28
2. O:C/S:W	87:52	8:73	16:20	16:37	1:59	1:58
3. O:W/S:C	80:35	1:93	0:50	1:35	1:23	0:93
4. O:W/S:W	76:53	24:52	26:57	26:72	17:82	2:95
	Monte Carlo variance 1000					
1. O:C/S:C	0:40	0:70	0:69	0:68	0:43	0:44
2. O:C/S:W	0:38	0:55	0:54	0:54	0:36	0:44
3. O:W/S:C	0:38	0:71	0:85	0:70	0:61	0:56
4. O:W/S:W	0:38	0:57	0:75	0:55	0:50	0:52
	Relative bias (%) of bootstrap variance estimate					
1. O:C/S:C	1:8	4:8	2:5	6:0	3:4	4:2
2. O:C/S:W	2:7	5:2	4:6	3:9	2:8	4:3
3. O:W/S:C	1:1	2:8	5:6	4:5	10:0	7:3
4. O:W/S:W	1:9	8:3	3:4	8:3	2:9	0:9
	95% Wald CI coverage					
1. O:C/S:C	0:4	94:3	94:5	93:9	94:4	93:8
2. O:C/S:W	1:0	94:0	90:1	89:1	95:3	94:2
3. O:W/S:C	2:3	94:0	93:6	94:8	93:2	93:4
4. O:W/S:W	2:1	82:9	83:2	79:6	86:7	95:3

Figure B.3 Boxplot of estimators for binary outcome with population size $N = 500000$ under four model specification scenarios.

Figure C.1 Ridge functions learned from AM with polynomial splines.

Nonlinear transformation $g(\cdot)$ contains q subnetworks, one for each row of \mathbf{B} .

C.2 Ridge functions

We take the AM2 model in Section 4.5.1 as a showcase. Both AM and AM-Net can recover the true ridge functions, as shown in Figure C.1 and C.2.

