

University of Louisville

ThinkIR: The University of Louisville's Institutional Repository

Electronic Theses and Dissertations

12-2023

Causal inference for the effect of continuous treatment on time-to-event outcomes and mediation analysis on health disparities in observational studies.

Triparna Poddar
University of Louisville

Follow this and additional works at: <https://ir.library.louisville.edu/etd>



Part of the [Biostatistics Commons](#), [Clinical Trials Commons](#), and the [Survival Analysis Commons](#)

Recommended Citation

Poddar, Triparna, "Causal inference for the effect of continuous treatment on time-to-event outcomes and mediation analysis on health disparities in observational studies." (2023). *Electronic Theses and Dissertations*. Paper 4229.
<https://doi.org/10.18297/etd/4229>

This Doctoral Dissertation is brought to you for free and open access by ThinkIR: The University of Louisville's Institutional Repository. It has been accepted for inclusion in Electronic Theses and Dissertations by an authorized administrator of ThinkIR: The University of Louisville's Institutional Repository. This title appears here courtesy of the author, who has retained all other copyrights. For more information, please contact thinkir@louisville.edu.

CAUSAL INFERENCE FOR THE EFFECT OF CONTINUOUS
TREATMENT ON TIME-TO-EVENT OUTCOMES AND
MEDIATION ANALYSIS ON HEALTH DISPARITIES IN
OBSERVATIONAL STUDIES

By

Triparna Poddar
B.Sc., University of Calcutta, 2015
M.Sc., University of Calcutta, 2017

A Dissertation
Submitted to the Faculty of the
School of Public Health and Information Sciences
of the University of Louisville
in Partial Fulfillment of the Requirements
for the Degree of

Doctor of Philosophy
in Biostatistics

Department of Bioinformatics and Biostatistics
University of Louisville
Louisville, Kentucky

December 2023

CAUSAL INFERENCE FOR THE EFFECT OF CONTINUOUS
TREATMENT ON TIME-TO-EVENT OUTCOMES AND
MEDIATION ANALYSIS ON HEALTH DISPARITIES IN
OBSERVATIONAL STUDIES

By

Triparna Poddar
B.Sc., University of Calcutta, 2015
M.Sc., University of Calcutta, 2017

A Dissertation Approved on

November 28, 2023

by the following Dissertation Committee:

Dr. Maiying Kong, Dissertation Director

Dr. Qi Zheng

Dr. Riten Mitra

Dr. Jeremy Gaskins

Dr. Michael Egger

DEDICATION

I dedicate this dissertation to my parents. For their teachings, blessings, love and support have helped me reach here today.

ACKNOWLEDGMENTS

I am deeply indebted to my advisor Dr. Maiying Kong for her insightful guidance and constant inspiration for my research over the last four years. I would like to thank Dr. Qi Zheng for his insightful discussion and input towards my research. I am fortunate for the support from Dr. Michael Egger, who has not only provided me with funding but also real-world data experience. I am grateful for the time and support of Dr. Riten Mitra and Dr. Jeremy Gaskins. I would also like to thank all the faculty and students of the Department of Bioinformatics and Biostatistics for their friendship and support.

I would also like to express my sincere gratitude to Dr. Bhaswati Ganguli, whose guidance and encouragement have helped me to fulfill my dream to pursue Ph.D. in the USA. Lastly I am deeply indebted to my parents and *Sagnik* who have been a constant source of support for me throughout this journey.

ABSTRACT

CAUSAL INFERENCE FOR THE EFFECT OF CONTINUOUS TREATMENT ON TIME-TO-EVENT OUTCOMES AND MEDIATION ANALYSIS ON HEALTH DISPARITIES IN OBSERVATIONAL STUDIES

Triparna Poddar

November 28, 2023

The dissertation comprises two projects related to causal inference based on observational data. In healthcare research, where abundant observational data such as claims data and electronic records are available, researchers often aim to study the treatment effect and the pathway of that effect. However, estimating treatment effects in observational data presents challenges due to confounding factors. The first project focuses on estimating continuous treatment effects for survival outcomes, while the second concentrates on mediation analysis, allowing the exploration of the pathway of the causal effect. Both projects involve addressing confounding variables.

In the first project, I investigate estimation of the average treatment effect (ATE) of continuous treatment on time to event outcome by adjusting multiple confounding factors and considering censoring observations. To adjust confounding factors, various propensity score methods such as multinomial regression and covariate balance propensity score models are used to estimate the ATE via the inverse probability of treatment weighting (IPTW) method. For continuous treatments, the IPTW is generated from covariate balancing generalized propensity score. To remedy the

possible bias in estimating ATE for time-to-event data due to censoring observations, we incorporate the censoring weights to estimate ATE. We propose using both the IPTW and the censoring weights (say, double weighting approach) to estimate ATE using the marginal structural accelerated failure time (AFT) model, where the IPTW adjusts for confounding factors and the censoring weights remedy the impact due to censored observations. Comprehensive simulation studies demonstrated our proposed method performed well. We applied our proposed method to examine if blood lead level impacts the time to death of older people in the United States, utilizing data from the NHANES III survey dataset.

In the second project, I delve into the more complex causal pathways of exposure on the outcome using mediation analyses. I begin with basic mediation analyses and progress to the more advanced four-way decomposition of causal effects from exposure to outcome. This includes the interaction between multiple mediators and the exposure. Expanding the scope of mediation analyses and four-way decomposition, I extend it to survival analysis and demonstrate the IOM-defined disparity in terms of four-way decomposition effects within the mediation analysis framework. Mediation analysis proves to be a crucial tool in unraveling the intricate pathways contributing to disparities among racial groups. Extensive simulation studies are conducted to examine the contribution of decomposition effects under various settings of mediators and outcomes. Finally, I investigate the factors influencing racial disparity among the black and white populations in the United States based on the NHANES III database.

TABLE OF CONTENTS

	PAGE
DEDICATION	iii
ACKNOWLEDGMENTS	iv
ABSTRACT	v
LIST OF TABLES	ix
LIST OF FIGURES	x
CHAPTER 1: ESTIMATION OF AVERAGE TREATMENT EFFECT FOR SURVIVAL OUTCOMES WITH CONTINUOUS TREATMENT IN OB- SERVATIONAL STUDIES	1
1.1 Introduction	1
1.2 Basic assumptions and the proposed method	5
1.2.1 Notations	5
1.2.2 The proposed method	7
1.2.3 Theoretical Properties	9
1.3 Simulation study	11
1.3.1 Simulation settings	11
1.3.2 Simulation results	14
1.4 Case study: blood lead level versus mortality based on NHANES III dataset	20
1.5 Discussion	33
CHAPTER 2: CAUSAL MEDIATION ANALYSIS FOR HEALTH RACIAL DISPARITIES	35
2.1 Introduction	35
2.2 General framework for mediation analysis	39
2.2.1 Potential outcomes and basic assumptions	39
2.2.2 Four fold decomposition of total effect	40
2.2.3 Four-way decomposition under linear models	43
2.3 Extension to survival outcomes and health racial disparity	45
2.3.1 Extension to survival outcomes	45
2.3.2 Extension to health racial disparity study	46
2.4 Simulation studies	49
2.4.1 Simulation Settings	50
2.4.2 Simulation results	52

2.5	Case study: racial disparity on all-cause mortality in the United States based on NHANES III dataset	53
2.6	Conclusion and discussion	62
	REFERENCES	64
	APPENDIX	70
	CURRICULUM VITA	77

LIST OF TABLES

TABLE		PAGE
1.1	Summary of the distribution of characteristics among participants for each categorical covariate (refer to column “Total #(%”), its association with mortality (refer to columns “Died #(%”) and “P-value”), and its association with BLL (refer to the last two columns)	26
1.2	The Mean and standard deviation (SD) for each continuous variable (refer to column “All”), stratified by mortality (refer to columns “Alive” and “Died”) and BLL (refer to columns “low BLL” and “high BLL”).	27
1.3	Estimated effect of Blood lead on time to death, Standard Error and Confidence Interval at 95%	33
2.1	Decomposition of TE of Race on Mortality Time	63

LIST OF FIGURES

FIGURE		PAGE
1.1	Causal effect of treatment A on outcome Y confounded by X	2
1.2	Boxplots of 1000 ATE estimates based on sample size $n = 1000$ at different treatment levels, with two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B). Each cluster of boxplots shows ATE estimates based on different weighting methods.	16
1.3	Boxplots of 1000 ATE estimates based on sample size $n = 5000$ at different treatment levels, with two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B). Each cluster of boxplots shows ATE estimates based on different weighting methods.	17
1.4	The scatter plot illustrates the RMSE based on three different double weighting methods and a sample size of $n = 1000$ at various treatment levels. The plot includes two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B).	18
1.5	The scatter plot illustrates the RMSE based on three different double weighting methods and a sample size of $n = 5000$ at various treatment levels. The plot includes two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B).	19
1.6	The scatter plot illustrates the true coverage rate for 95% CI based on three different double weighting methods and a sample size of $n = 1000$ at various treatment levels. The plot includes two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B). . . .	21

1.7	The scatter plot illustrates the true coverage rate for 95% CI based on three different double weighting methods and a sample size of $n = 5000$ at various treatment levels. The plot includes two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B). . . .	22
1.8	Scatterplot illustrating the relationships between continuous covariates (age, PIR, BMI, and exercise level) and Blood Lead Levels (BLL), accompanied by the fitted regression lines in both the original and the weighted sample.	28
1.9	Barplot presenting the distributions of each categorical covariate in low Blood Lead Level (BLL) and high BLL, both in the original sample (the first two bars in each panel) and the weighted sample (the last two bars in each panel).	29
1.10	Scatterplot depicting survival time (in months since survey) against Blood Lead Levels (BLL), accompanied by the fitted trend lines in both the original sample and the weighted sample.	31
1.11	Kaplan-Meier survival curves illustrating the survival outcomes for low Blood Lead Level (BLL) and high BLL in the original sample (left panel) and the weighted sample (right panel), along with the log-rank tests.	32
2.1	Illustration of directed acyclic graph (DAGs) for various study designs: (a) randomized control trials; (b) observational studies; (c) simple mediation model; and (d) a mediation model with confounding variables.	36
2.2	Illustration of mediation models with interaction of exposure and mediator: (a) one mediator and (b) multiple mediators	38
2.3	Boxplots of estimated proportion attribution by each decomposed effect at two levels of minority group distribution (35% and 50%, indicated in the first row and second row, respectively) and for two values for m^* , for single mediator and continuous outcome in the model . . .	54
2.4	Boxplots of estimated proportion attribution by each decomposed effect at two levels of minority group distribution (35% and 50%, indicated in the first row and second row, respectively) and for two values for m^* , for multiple mediators and continuous outcome in the model .	55
2.5	Boxplots of estimated proportion attribution by each decomposed effect at two levels of minority group distribution (35% and 50%, indicated in the first row and second row, respectively) and for two values for m^* , for multiple mediators and time-to-event outcome in the model	56
2.6	Plots to illustrate performance of standard error estimation: by taking ratio of SE by bootstrap method(red) or SE by Delta method(blue) and Empirical SD for decomposed effects at two levels of minority group distribution (35% and 50%, indicated in the first row and second row, respectively) and for two values for m^* , for single mediator and continuous outcome in the model	57

2.7	Plots to illustrate performance of standard error estimation: by taking ratio of SE by bootstrap method(red) or SE by Delta method(blue) and Empirical SD for decomposed effects at two levels of minority group distribution (35% and 50%, indicated in the first row and second row, respectively) and for two values for m^* , for multiple mediators and continuous outcome in the model	58
2.8	Plots to illustrate the performance of standard error (SE) estimation: each line represents the ratio of SE estimation over the empirical SE from different methods (red: bootstrap method; blue for delta method). The two rows correspond two levels of minority groups (35% and 50%), and the two columns correspond to two different values for m^*	59

CHAPTER 1

ESTIMATION OF AVERAGE TREATMENT EFFECT FOR SURVIVAL OUTCOMES WITH CONTINUOUS TREATMENT IN OBSERVATIONAL STUDIES

1.1 Introduction

In observational studies, the estimation of the Average Treatment Effect (ATE) encounters notable challenges due to the absence of random assignment, leading to potential confounding. Propensity score weighting methods are commonly employed to address this issue, although the estimation of propensity scores for continuous treatments poses distinct challenges that require careful consideration. Moreover, when dealing with time-to-event data as the outcome of interest, it is crucial to account for bias introduced by censored observations in the ATE estimation models. The primary objective of our first project is to develop a novel method for estimating ATE in situations involving continuous treatment and time-to-event outcomes.

In clinical trials, Randomized Control Trials (RCTs) stand as the gold standard for estimating the treatment effect of a new exposure or treatment on various outcomes. Within an RCT, participants are randomly allocated to different treatment groups, ensuring that the distributions of confounding factors, whether they are measured or unmeasured, are comparable among the various treatment groups. The treatment effect on outcomes can be directly estimated by calculating the difference in sample means between treated and control participants. However, in practical

terms, it is not always ethical, feasible, or cost-effective to conduct an RCT. On the contrary, observational data is often more readily available, prompting researchers to attempt to assess treatment effects using this type of data. It's important to note that in observational studies, treatment assignment, denoted as A , is frequently influenced by patients' characteristics, represented by X . Additionally, the outcome Y is influenced by both patients' characteristics X and the treatment A , as illustrated in Figure 1.1. Hence, the relationship between treatment and outcome becomes confounded by

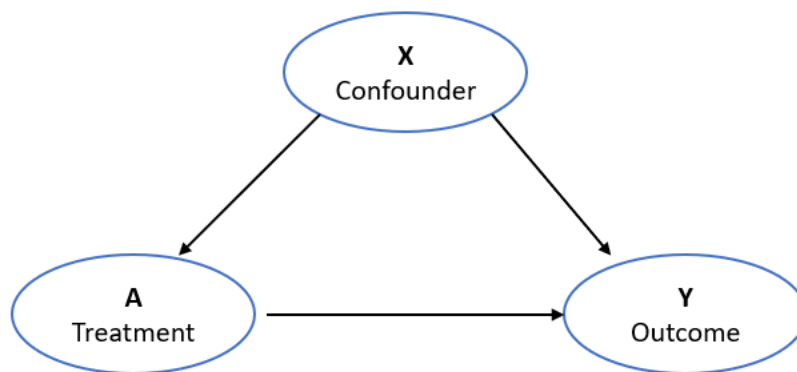


Figure 1.1: Causal effect of treatment A on outcome Y confounded by X .

patients' characteristics and various other potential confounding factors. A direct comparison of outcomes between the treated and control groups no longer serves as a valid estimator for the ATE. To estimate the ATE using observational data, it is essential to control for these confounding factors. The most commonly used methods for this purpose are propensity score-based techniques, such as matching, stratification, regression, inverse probability of treatment weighting (IPTW) (Rosenbaum, 1987), and doubly robust methods (Lunceford and Davidian, 2004). The majority of these existing propensity score estimation methods are based on binary treatment settings where the propensity scores are estimated by logistic regression. However, employing these techniques often requires researchers to dichotomize a continuous treatment, leading to the loss of crucial information and insights within the data.

In recent years, there has been an exploration of ATE estimation for multiple treatment groups using the generalized propensity score (GPS) (Imbens, 2000). However, research on ATE estimation for continuous treatments remains limited. Hirano and Imbens (2004) have extended these methods to estimate the ATE for continuous treatment when dealing with binary or continuous outcomes through the use of regression models on both treatment and propensity scores. It’s worth noting that ATE estimates relying solely on regression analysis may introduce bias (Hade and Lu, 2014), and mis-specification of GPS can also lead to biased ATE estimates. To alleviate the potential impact of mis-specification of GPS, Imai and Ratkovic (2014) introduced the concept of a covariate balancing propensity score (CBPS) for binary treatment. CBPS leverages two fundamental properties of the propensity score: estimating the likelihood of treatment assignment for each subject based on their covariates and achieving a balance in covariates across different treatment groups. CBPS effectively improves the covariate balance among various treatment groups, thereby enhancing the robustness of the estimation. Fong et al. (2018) expanded on this idea by introducing the covariate balancing generalized propensity score (CBGPS) for continuous treatment scenarios. CBGPS estimates propensity scores by minimizing the correlation between treatment and confounding covariates through weighting, thereby improving the balance of covariates across treatment groups in cases involving continuous treatments.

In this study, our primary focus is on estimating ATE when the treatment is continuous, and the outcome is represented as time-to-event data, accounting for the presence of right-censored observations. Time-to-event outcomes often involve censored observations, which occur when subjects do not experience the event outcome during the study period. Estimating ATE for time-to-event outcomes necessitates addressing the bias introduced by right censoring. To correct for this bias, the inverse probability of censoring weights (IPCW) is frequently employed, as detailed in

the works of Cain and Cole (2009); Cole and Hernán (2008); Robins and Finkelstein (2000). This approach involves estimating the probability of censoring over time, and the inverse of the probability of remaining uncensored is used as a weight for uncensored observations when modeling. Estimating the probability of censoring is typically accomplished using Kaplan-Meier curves for situations involving independent censoring and the Accelerated Failure Time (AFT) model when dealing with dependent censoring. Various methods have been developed to address right censoring and estimate ATE for time-to-event data when the treatment is binary or categorical (Andersen et al., 2017; Austin, 2010). Xie and Liu (2005) proposed an adjusted Kaplan-Meier estimator of the survival function and the log-rank test incorporating IPTW. Andersen et al. (2017) introduces the concept of creating a parallel dataset using pseudo observations to account for right censoring for survival outcomes, which can be used in traditional causal inference methodologies. Austin (2018) examines the effectiveness of the generalized propensity score in estimating the impact of continuous exposures on survival or time-to-event outcomes. To account for censoring the dose-response function was modified as the survival function which was estimated from the Cox-Proportional hazard model.

In our study, we extend these methods to estimate ATE for time-to-event outcomes when the treatment is continuous. We apply both the IPTW and IPCW to estimate ATE using a marginal structure AFT model. We also investigate the performance of the double-weighting method in estimating ATE when the generalized propensity score is estimated using the maximum likelihood method or the CBGPS method. The remainder of this paper is structured as follows. In Section 1.2, we first outline the notations and the fundamental assumptions that underpin our study. Following this, we describe the proposed double weighting method for estimating ATE for continuous treatment in the context of time-to-event outcomes and provide an in-depth exploration of the associated theoretical properties. In Section 1.3, we conduct

a series of simulated studies to assess the performance of the proposed methodology. In Section 1.4, we apply our proposed method to investigate the impact of blood lead levels on all-cause mortality among older individuals in the US population. The final section of this chapter is dedicated to an in-depth discussion.

1.2 Basic assumptions and the proposed method

In this section, we begin by defining all the terms we used for our paper and explain the identification assumptions required under the causal framework to develop our proposed method. Then we describe our proposed method with the double weighting method for estimating ATE for continuous treatment for time-to-event outcomes. In the final subsection, we delve into the theoretical properties associated with the estimators derived from our proposed method.

1.2.1 Notations

Let $\mathbf{X} \in \mathcal{X}$ denote a p -dimensional vector of covariates of a patient, and $A \in \mathcal{A}$ denote the treatment that the patient received. Here \mathcal{X} and \mathcal{A} are the support of \mathbf{X} and A respectively. In this work, \mathcal{X} is a compact set in \mathbb{R}^p and \mathcal{A} is either \mathbb{R} if A is continuous, or a set with finite many values if A is discrete. We use T to denote the actual survival time of the patient and $T^{(a)}$ to denote the potential survival time if the patient has received treatment a . As in practice, the survival time T is often subject to right censoring by C , the observed variable is $\tilde{T} = \min\{T, C\}$ and $\delta = 1\{T < C\}$ is the censoring indicator, where $1\{\cdot\}$ is an indicator function. The observed data consist of n i.i.d replicates of $\mathbf{D} := (\mathbf{X}, A, \tilde{T}, \delta)$, denoted by $\{\mathbf{D}_i = (\mathbf{X}_i, A_i, \tilde{T}_i, \delta_i), i = 1, \dots, n\}$. For convenience, we define $Y_i = \log T_i$ and $\tilde{Y}_i = \log \tilde{T}_i$. We further denote $Y^{(a)} = \log T^{(a)}$, $\tilde{T}^{(a)} = \min(T^{(a)}, C)$, and $\tilde{Y}^{(a)} = \log \tilde{T}^{(a)}$ as the notations associated with potential outcomes if the patient has received treatment a where $a \in \mathcal{A}$. For generic random variables U and V , let $f_U(\cdot)$ and $f_U(\cdot|V)$ denote the density (probability

mass) function of U and the conditional density (probability mass) function of U given V . In addition, we use $f_U(\cdot|V, \boldsymbol{\eta})$ to denote the conditional density (probability mass) function of U given V governed by parameters $\boldsymbol{\eta}$. In particular, we use $G^*(\cdot)$ to denote the survival function of C , that is, $G^*(u) = P(C > u)$.

Our target estimand is ATE which is defined as the difference of potential outcomes under two treatments (say a vs a'). That is

$$ATE(a, a') = E[Y^{(a)}] - E[Y^{(a')}] \quad (1.1)$$

Note that, not all potential outcomes are observable. Indeed only one potential outcome is observed which is the potential outcome corresponding to the received treatment. To estimate ATE based on observational data, the following assumptions are required (Brown et al., 2021; Imbens, 2000):

- (1) Weak unconfoundedness (or Ignorability): the treatment assignment A is independent of the potential outcome $Y^{(a)}$ given confounding variables \mathbf{X} . That is,

$$Y^{(a)} \perp\!\!\!\perp A | \mathbf{X} \quad \forall a \in \mathcal{A}. \quad (1.2)$$

- (2) Positivity: a subject has a non-zero probability of receiving any treatment. That is,

$$f(a|\mathbf{X}) > 0 \quad \forall a \in \mathcal{A}, \quad \mathbf{X} \in \mathcal{X}, \quad (1.3)$$

where $f(a|\mathbf{X})$ is the density function of A given \mathbf{X} .

- (3) Consistency: the observed outcome is the potential outcome corresponding to the observed treatment assignment. That is,

$$Y = \sum_{a \in \mathcal{A}} 1(A = a) Y^{(a)} \quad (1.4)$$

- (4) Correct specification of GPS model and correct specification of censoring probability model.

Given the four underlying assumptions, we present the following double weighting method to estimate ATE for time-to-event data in the presence of right-censored observations.

1.2.2 The proposed method

We have considered \mathbf{X} to be standardized with zero mean and unit variance. Let $w_p(a; \mathbf{x})$ denote the weights that balance the observed covariates \mathbf{X} across different values of the treatment variable A . We consider $w_p(a; \mathbf{x}) = f_A(a)/f_A(a|\mathbf{x})$, where the numerator $f_A(a)$ is the stabilizing factor (Robins et al., 2000) and the denominator $f_A(a|\mathbf{x})$ is the generalized propensity score (GPS) (Imbens, 2000).

In continuous treatment setup, the GPS, $f(a|\mathbf{x})$ is defined as the conditional density of receiving a treatment $A = a$ given confounding covariates $\mathbf{X} = \mathbf{x}$ (Hirano and Imbens, 2004). Following Hirano and Imbens (2004) and Imai and Van Dyk (2004), we assume that the GPS $f_A(a|\mathbf{x})$ has the conditional normal density as follows:

$$f_A(a|\mathbf{x}, \boldsymbol{\xi}^*) = \frac{1}{\sqrt{2\pi}\sigma^*} \exp\left(-\frac{(a - \mathbf{x}^\top \boldsymbol{\beta}^*)^2}{2\sigma^{*2}}\right),$$

where $\boldsymbol{\xi}^* = (\boldsymbol{\beta}^{*\top}, \sigma^*)^\top$. Then the weights that balance the confounding variables can be expressed as

$$w_p(a; \mathbf{x}) = \frac{f_A(a)}{f_A(a|\mathbf{x})} = w_p(a; \mathbf{x}, \boldsymbol{\xi}^*) = \frac{\sigma^*}{\sigma^{**}} \exp\left(\frac{(a - \mathbf{x}^\top \boldsymbol{\beta}^*)^2}{2\sigma^{*2}} - \frac{a^2}{2\sigma^{**2}}\right), \quad (1.5)$$

where $\sigma^{**2} = \boldsymbol{\beta}^T \boldsymbol{\beta} + \sigma^{*2}$. In practice, $\boldsymbol{\xi}^*$ is typically unknown but can be well estimated by $\hat{\boldsymbol{\xi}}$, a maximum likelihood estimator (MLE). According to the MLE theory (see, e.g. Le Cam, 1990), $\hat{\boldsymbol{\xi}} = \boldsymbol{\xi}^* + O_p(n^{-1/2})$. Another robust method to estimate the IPTW is the CBGPS method which estimates parameters in the propensity scores model by setting the correlation between treatment and confounding covariates in the weighted sample being zero (Fong et al., 2018). The method can be implemented using R-package CBPS (Fong et al., 2021).

With time-to-event outcomes, we consider the AFT marginal structural model

(AFT-MSM) as follows:

$$E[\log T^{(a)}] = E[Y^{(a)}] = \mathbf{Z}(a)^\top \theta^*, \quad \forall a \in \mathcal{A}. \quad (1.6)$$

In our methodology, we consider $\mathbf{Z}(a)$ as a generalized known function of treatment a . For instance, when the potential outcome $Y^{(a)}$ exhibits a linear relationship with the treatment level a , we can set $\mathbf{Z}(a) = (1, a)^\top$ to form a simple linear function, where $\mathbf{Z}(a)^\top \theta^* = \theta_0^* + a\theta_1^*$. However, if the potential outcome $Y^{(a)}$ has a piecewise linear relationship with treatment level a , we may employ different basis functions for $\mathbf{Z}(a)$. For example, $\mathbf{Z}(a) = (1, a, (a - c)_+)^\top$ helps capture relationship changes at the threshold c . $\mathbf{Z}(a)$ can be designed as a flexible vector of functions to effectively describe the relationship between the potential outcome $Y^{(a)}$ and a , characterized by the parameter θ^* . To estimate the causal parameters θ^* within the AFT-MSM, we must account for the bias arising from confounding and censoring. For simplicity, we assume that the random censoring C is independent of (T, A, \mathbf{X}) . It is worth pointing out that our proposed method can be extended into the case where C is independent of T given \mathbf{X} and A . Please see Remark 1.2.1 for more detailed discussions.

Let denote $\mathbf{Z} = \mathbf{Z}(a) = (Z_0(a), Z_1(a), \dots, Z_q(a))^\top$, let $\boldsymbol{\theta} = (\theta_0, \theta_1, \dots, \theta_q)^\top$ and $h = (\boldsymbol{\xi}, G)$ where $G(\cdot)$ is a survival function. Define $m(\mathbf{D}, \boldsymbol{\theta}, h) = w(\mathbf{D}, h)\mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta})$ and $M(\boldsymbol{\theta}, h) = E[m(\mathbf{D}, \boldsymbol{\theta}, h)]$, where $w(\mathbf{D}, h) = \delta w_p(A; \mathbf{X}, \boldsymbol{\xi})G^{-1}(Y)$. Then the marginal structure equation for the parameters in the AFT-MSM model is

$$M(\boldsymbol{\theta}^*, h^*) = 0, \quad (1.7)$$

where $h^* = (\boldsymbol{\xi}^*, G^*)$. We refer the validation of the equation to the Appendix.

Moreover, we define $M_n(\boldsymbol{\theta}, h) := n^{-1} \sum_{i=1}^n m(\mathbf{D}_i, \boldsymbol{\theta}, h)$ and consider the estimating equation as follows:

$$M_n(\boldsymbol{\theta}, \hat{h}) = \sum_{i=1}^n m(\mathbf{D}_i, \boldsymbol{\theta}, \hat{h}) = \sum_{i=1}^n w(\mathbf{D}_i, \hat{h})\mathbf{Z}_i(Y_i - \mathbf{Z}_i^\top \boldsymbol{\theta}) = 0, \quad (1.8)$$

where $\hat{h} = (\hat{\boldsymbol{\xi}}, \hat{G})$, and $\hat{G}(\cdot)$ is the Kaplan-Meier estimator of $G^*(\cdot)$. Simple algebra

yields the solution to (1.8) as

$$\hat{\boldsymbol{\theta}} = \left(\sum_{i=1}^n w(\mathbf{D}_i, \hat{h}) \mathbf{Z}_i \mathbf{Z}_i^\top \right)^{-1} \sum_{i=1}^n w(\mathbf{D}_i, \hat{h}) \mathbf{Z}_i Y_i, \quad (1.9)$$

Remark 1.2.1. *If C is independent of T given X and A , we can estimate the conditional survival probability $P(C > u|A, \mathbf{X})$ by the conditional Kaplan-Meier method (see, e.g., Dabrowska, 1989; Gonzalez-Manteiga and Cadarso-Suarez, 1994) and denote the resulting estimator by $\hat{G}(\cdot; A, \mathbf{X})$. Then we can replace $\hat{G}(\cdot)$ by $\hat{G}(\cdot; A, \mathbf{X})$ in \hat{h} and subsequently $M_n(\boldsymbol{\theta}, \hat{h})$ in Equation (1.8). It can be shown that our proposed method would still provide consistent estimates of $\boldsymbol{\theta}^*$.*

1.2.3 Theoretical Properties

In this subsection, we investigate the theoretical properties of our proposed estimators. We begin with introducing some necessary notations. We use c to represent an unspecified positive constant whose value may vary. Let $\|B\|_l$ denote the l -norm of B , where B can be a vector or a matrix. We assume that $\boldsymbol{\xi}^*$ is an interior point of Ω and $\boldsymbol{\theta}^*$ is an interior point of Θ , where Ω and Θ are two compact sets in \mathbb{R}^d . In addition, $G^* \in \mathcal{G}$ and \mathcal{G} is an infinite dimensional parameter space. Define $\Omega_\epsilon := \{\boldsymbol{\xi} \in \Omega : \|\boldsymbol{\xi} - \boldsymbol{\xi}^*\|_\infty \leq \epsilon\}$, $\Theta_\epsilon = \{\boldsymbol{\theta} \in \Theta, \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \leq \epsilon\}$, and $\mathcal{G}_\epsilon := \{h \in \mathcal{G} : \|G - G^*\|_\infty \leq \epsilon\}$. In addition, we denote $\Omega \times \mathcal{G}$ and $\Omega_\epsilon \times \mathcal{G}_\epsilon$ by \mathcal{H} and \mathcal{H}_ϵ , respectively.

For any $(\boldsymbol{\theta}, h) \in \Omega_\epsilon \times \mathcal{H}_\epsilon$, we denote the ordinary derivative of $M(\boldsymbol{\theta}, h)$ with respect to $\boldsymbol{\theta}$ as $\boldsymbol{\Gamma}_1(\boldsymbol{\theta}, h)$, which satisfies $\boldsymbol{\Gamma}_1(\boldsymbol{\theta}, h)(\boldsymbol{\theta}_1 - \boldsymbol{\theta}) = \lim_{t \rightarrow 0} t^{-1}[M(\boldsymbol{\theta} + t(\boldsymbol{\theta}_1 - \boldsymbol{\theta}), h) - M(\boldsymbol{\theta}, h)]$ for all $\boldsymbol{\theta}_1 \in \Omega$. Additionally, for any $\boldsymbol{\xi} \in \Omega_\epsilon$, we say $M(\boldsymbol{\theta}, h)$ is path-wise differentiable at $h \in \mathcal{H}_\epsilon$ in the direction $[h_1 - h]$ if $\{h + t(h_1 - h) : t \in [0, 1]\} \subset \mathcal{H}$ and $\lim_{t \rightarrow 0} t^{-1}[M(\boldsymbol{\theta}, h + t(h_1 - h)) - M(\boldsymbol{\theta}, h)]$ exists. We denote the limit by $\boldsymbol{\Gamma}_2(\boldsymbol{\theta}, h)(h_1 - h)$.

We next impose the following regularity conditions that facilitate our technical

derivations:

(C1) The study has a finite duration L such that $T \in (0, L]$ and $G^*(L) > \tau$ for some constant $\tau > 0$.

(C2) (a) Conditional exchangeability: $\{T^{(a)}, a \in \mathcal{A}\} \perp A | \mathbf{X}$; (b) Consistency: if $A = a$, $T^{(A)} = T^{(a)} = T$; (c) Positivity: if A is discrete, $f(a; \mathbf{x}, \boldsymbol{\xi}^*) > \nu > 0$ and if A is continuous, for all values of \mathbf{x} with $f(\mathbf{x}) > 0$ we have $f(a|\mathbf{x}) > 0$, where $f(\cdot)$ is the density of \mathbf{X} and $f(a|\mathbf{x})$ is the conditional density function of A given \mathbf{X} .

(C3) The fisher information matrix for $\boldsymbol{\xi}$, $\mathbf{I}(\boldsymbol{\xi}^*)$, is invertible.

Remark 1.2.2. *Without loss of generality, we consider Ω and \mathcal{G} such that for all $\boldsymbol{\xi} \in \Omega$ and $G \in \mathcal{G}$, $\|E[w_p^2(A; \mathbf{X}, \boldsymbol{\xi})]\| < c$ and $\|G - G^*\|_\infty < \tau/2$, where τ is defined in Condition (1). In fact, $\|E[w_p^2(A; \mathbf{X}, \boldsymbol{\xi})]\| < c$ is satisfied by $\sigma^2 > (2\sigma^{*2})/(2\sigma^{*2} + 1)$ for σ in $\boldsymbol{\xi}$, given that \mathcal{X} and Ω are compact.*

Theorem 1.2.1. *Under AFT-MSM (1.6) and regularity conditions (C1)–(C3), $\hat{\boldsymbol{\theta}} \rightarrow_p \boldsymbol{\theta}^*$, where \rightarrow_p denotes the convergence in probability.*

Theorem 1.2.2. *Under the same conditions as in Theorem 1.2.1,*

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \rightarrow_d N(0, (E[\mathbf{Z}\mathbf{Z}^\top])^{-1} \mathbf{V} (E[\mathbf{Z}\mathbf{Z}^\top])^{-1}),$$

where \rightarrow_d denotes the convergence in distribution, \mathbf{V} is the covariance matrix of $(m(\mathbf{D}_i, \boldsymbol{\theta}^*, h^*) + \eta(A_i, \mathbf{X}_i) + \psi(\tilde{T}_i, \delta_i))$, and $\eta(A_i, \mathbf{X}_i)$ and $\psi(\tilde{T}_i, \delta_i)$ are defined in Equations (2.27) and (2.28) respectively, in the Appendix.

The proofs of Theorem 1.2.1 and Theorem 1.2.2 are provided in the Appendix at the end of this chapter.

1.3 Simulation study

Through a comprehensive simulation study, we assess the performance of our proposed method in estimating ATE for time-to-event outcomes. Our methodology leverages double weighting, which involves both confounding and censoring weights. To gauge the effect of the weights used in estimation, we also incorporate true weights, wherein both propensity scores and censoring probabilities are generated from the underlying models. Additionally, we obtain ATE estimates without any weighting, with only weights derived from the estimated propensity scores while ignoring censoring, and with only estimated censoring weights while overlooking confounding. We also delve into the impact on ATE estimation resulting from the strength of the association between confounding and treatment assignment, the censoring rate, and the sample size.

1.3.1 Simulation settings

Let us denote \mathbf{X} as a vector of confounding variables. Let assume $\mathbf{X} \in \mathbb{R}^p$, where $p = 10$. The treatment A was generated based on the following GPS model.

$$A \sim N(\mathbf{X}^\top \beta, 2^2), \quad (1.10)$$

where $\beta = \kappa(\mathbf{1}_4, \mathbf{0}_6)$. Here $\mathbf{1}_4$ indicates a vector of 1 with four elements and $\mathbf{0}_6$ indicates a vector of 0 with six elements. κ was used to capture the degree of association between treatment and confounding variables. κ took value 0.1 to indicate a weak association between treatment and confounding variables, and κ took value 0.5 to indicate a strong association between treatment and confounding variables.

Given \mathbf{X} and A , the outcomes were generated from the following outcome model:

$$\log(T) = \theta_0 + \theta_1 A + \mathbf{X}^\top \gamma + \sigma \epsilon. \quad (1.11)$$

Here θ_0 was set as 1 and θ_1 was varied in the set $\{0, 0.25, 0.5, 0.75, 1\}$ to capture different levels of treatment effect of A on Y . $\gamma = (1, 0.5, -0.3, 0.2, \mathbf{0}_6)$, which indicated that the outcome was associated with the first four confounding variables only. We set $\sigma = 0.5$ and $\epsilon \sim N(0, 1)$ in the outcome model (1.11).

The censoring time was generated from the model:

$$\log(C) = \mu_c + \sigma_c \epsilon_c. \quad (1.12)$$

Here we set $\epsilon_c \sim \text{Gumbel}(0, 1)$ and $\mu_c = 3.6$. We set σ_c as 3.6 and 13.2 to control the censoring probability at 15% and 30% respectively.

In each of the aforementioned scenarios, we conducted simulation studies employing two distinct sample sizes: $n = 1000$ and $n = 5000$. Consequently, we have a total of 8 simulation settings, reflecting two values for κ to denote varying degrees of association between treatment and confounding variables, two different censoring probabilities, and two sample sizes. For each simulation set, we generated 1000 samples. The simulations for each setting were conducted in the following steps.

Step 1: Generated a sample with n i.i.d observations for $(\mathbf{X}, A, \tilde{T}, \delta)$, where each observation was generated through the following steps: (i) generating 10 independent covariates, denoted as $\mathbf{X} = (X_1, X_2 \dots X_{10})$, from a multivariate normal distribution $MVN(0, I)$; (ii) generating treatment A using the GPS model described in Equation (1.10), where A was only associated with the first four covariates; (iii) generating the survival time from the AFT model specified in Equation (1.11); (iv) generating the censoring outcome using the model presented in Equation (1.12); and (v) generating the observed outcome, denoted as $\tilde{T} = \min\{T, C\}$, and the censoring indicator, $\delta = 1\{T < C\}$.

Step 2: For each observation in the sample, two potential outcomes under control ($a = 0$) and under exposure at level 1 ($a = 1$) were generated. Specifically, they were generated from $\log(T^{(0)}) = \theta_0 + X^\top \gamma + \sigma \epsilon$ and $\log(T^{(1)}) = \theta_0 + \theta_1 + X^\top \gamma + \sigma \epsilon$.

Consequently, we calculated the true sample treatment effect as the difference between the sample mean of $\log(T^{(1)})$ and the sample mean of $\log(T^{(0)})$.

Step 3: Obtained the inverse probability of treatment weights for each observation based on the true propensity score model as well as the estimated propensity score models, which were obtained from MLE method and CBGPS, respectively. The weight based on the true propensity score was calculated as $W_{p.true} = \frac{f_{A.true}}{f_{A.true|X}}$, where $A.true \sim N(0, ||\beta||^2 + 2^2)$, and $A.true|(X = x) \sim N(x^\top \beta, 2^2)$. The weight based on the MLE estimated propensity score by MLE method was obtained by $W_{p.est} = \frac{f_{A.est}}{f_{A.est|X}}$, where $A.est \sim N(\mu_{sample(A)}, \hat{\sigma}_A^2)$, $A.est|(X = x) \sim N(x^\top \hat{\beta}, \hat{\sigma}_A^2)$, and $\hat{\beta}$ and $\hat{\sigma}_A$ were estimated from linear regression model. For comparison, we also estimated generalized propensity scores by CBGPS method directly using the R package.

Step 4: Obtained the censoring weights based on true censoring model and estimated censoring model. The true censoring weights were obtained from $P[\log(\tilde{T}) > t]$, where $\log(\tilde{T}) \sim Gumbel(\mu_c, \sigma_c)$ with μ_c and σ_c defined in Equation (1.12). The estimated censoring weights w_c were estimated using Kaplan-Meier estimator, which were the inverse of the survival probability of the censoring variable. We set w_c at 0.0001 if the estimated weight was 0.

Step 5: $\hat{\theta}$ was estimated by the equation (1.9) with different specifications of weights as specified in Steps 3 and 4.

Step 6: Obtained variance estimation of the treatment effects by bootstrapping procedure with 100 resampling for each specification of weighting method.

We repeat our simulation procedure 1000 times from Step 1 to Step 6 under each setting. From 1000 simulations, we obtained the mean of the 1000 estimated treatment effects, root mean square error (RMSE), mean standard error from bootstrap estimate, and mean coverage rate for 95% CI, respectively.

1.3.2 Simulation results

In Figures 1.2 and 1.3, we present boxplots illustrating the 1000 ATE estimates derived from our simulation studies. These studies aimed to explore different methods for estimating ATE under varied conditions, specifically focusing on the censoring rate and the degree of association between confounding and treatment. The simulations were conducted for two sample sizes, namely $n = 1000$ and $n = 5000$, as depicted in Figure 1.2 and Figure 1.3, respectively. For each sample size, we explored four distinct simulation scenarios. These scenarios encompassed two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively, in each figure) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B). In each cluster of boxplots, the first boxplot illustrates the ATE estimates without weighting. The second and third boxplots show ATE estimates with the inverse of the probability of treatment weighting (IPTW) and the inverse of the survival probability of censoring, respectively. These probabilities were estimated by the Maximum Likelihood Estimation (MLE) method for treatment and the Kaplan-Meier method for censoring. The last three boxplots depict ATE estimates using the double weighting method. The process starts with true propensity score and true censoring weights, followed by estimates using weights derived from the propensity score estimated by the MLE method and the estimated probability of survival of censoring by the Kaplan-Meier method. Finally, the last set of estimates uses weights derived from the propensity score estimated by the CBGPS method and the estimated probability of survival of censoring by the Kaplan-Meier method. The horizontal line for each cluster of boxplots represents the true ATE obtained from the true potential outcomes at the particular level of treatment effect on Y . The boxplots reveal that ATE estimates obtained without any weighting or with solely IPTW or censoring weighting exhibited substantial bias across all simulation settings. As the association

between confounding and treatment increased, the bias for ATE estimates by the estimated censoring weights got larger. Conversely, ATE estimates generated by the three different double weighting methods tended to show lower bias. However, as the association between confounding and treatment increased, the variance for ATE estimates for double weighting with true weights and weights estimated by MLE method and Kaplan-Meier method also increased. The smallest bias and variance were obtained by the double weighting method with weights estimated by the CBGPS method and Kaplan-Meier method in all simulation scenarios. As the sample size increases, Figure 1.3 demonstrates a reduction in overall bias and variance for all double weighting methods.

Figures 1.2 and 1.3 clearly demonstrate that the ATE estimates from double weighting methods (see the last three boxplots in each condition) were unbiased, while the unweighting method, IPTW only, and censoring only weighting methods exhibited bias. In Figure 1.4 for $n = 1000$ and Figure 1.5 for $n = 5000$, we plotted the RMSE obtained from 1000 estimates of ATE for the three double weighting methods: true weights, estimated weights by MLE and Kaplan-Meier method, and estimated weights by CBGPS and Kaplan-Meier method. From these figures, we find that the ATE estimates of the double weighting method based on CBGPS and Kaplan-Meier method had the smallest RMSE, indicating that this double weighting method works best. Comparing Panel A and Panel B, we observe that the RMSE increased as the association between confounding variables and treatment increased. Additionally, comparing Figure 1.4 for $n = 1000$ and Figure 1.5 for $n = 5000$, we conclude that as the sample size increased, the RMSE decreased.

In Figures 1.6 and 1.7, we plotted the coverage rates of the 95% CI for the true ATE obtained from 1000 simulation runs for $n = 1000$ and $n = 5000$ based on the three double weighting methods. When the association between confounding variables and treatment was small (Panel A), the coverage rates from the three double

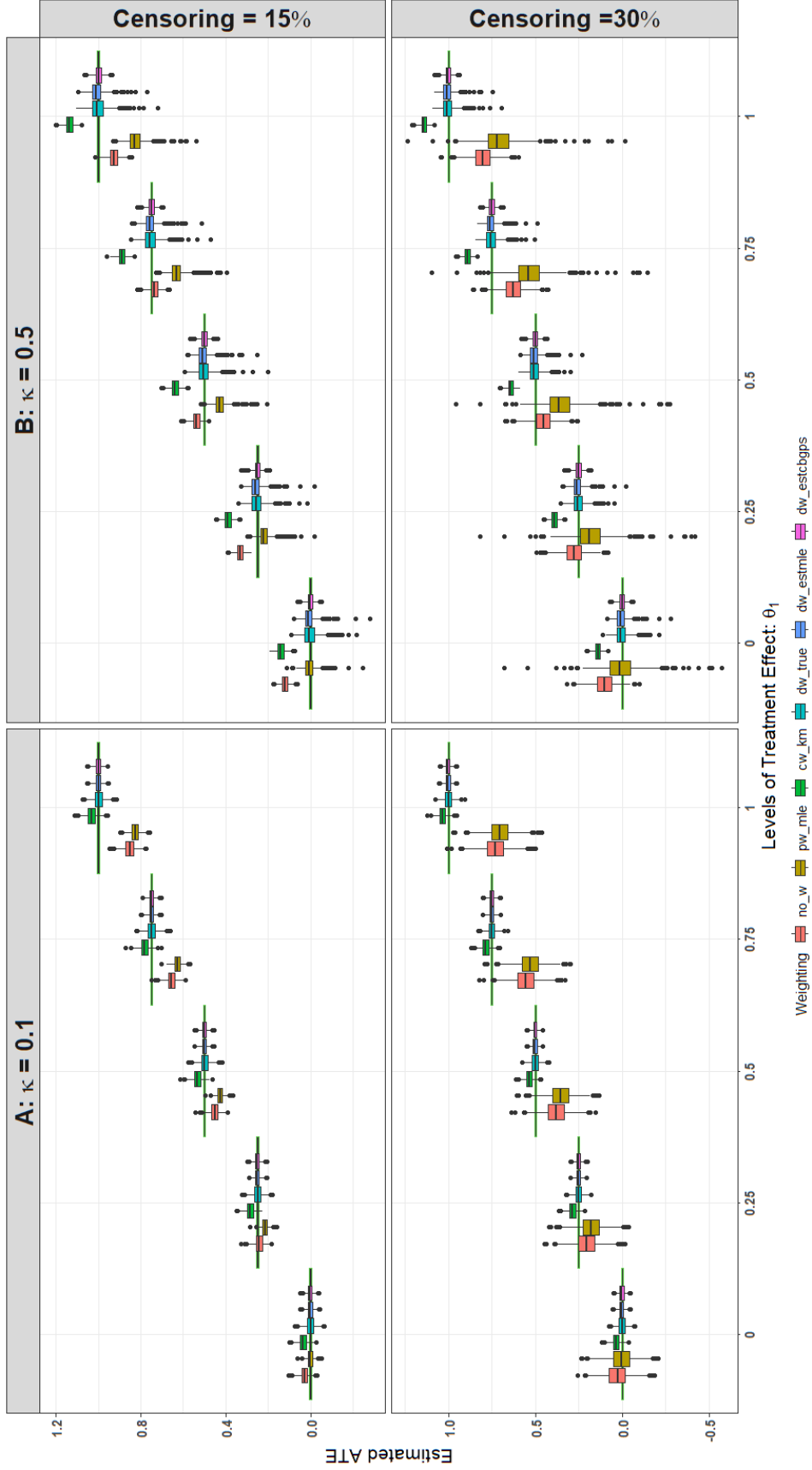


Figure 1.2: Boxplots of 1000 ATE estimates based on sample size $n = 1000$ at different treatment levels, with two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B). Each cluster of boxplots shows ATE estimates based on different weighting methods.

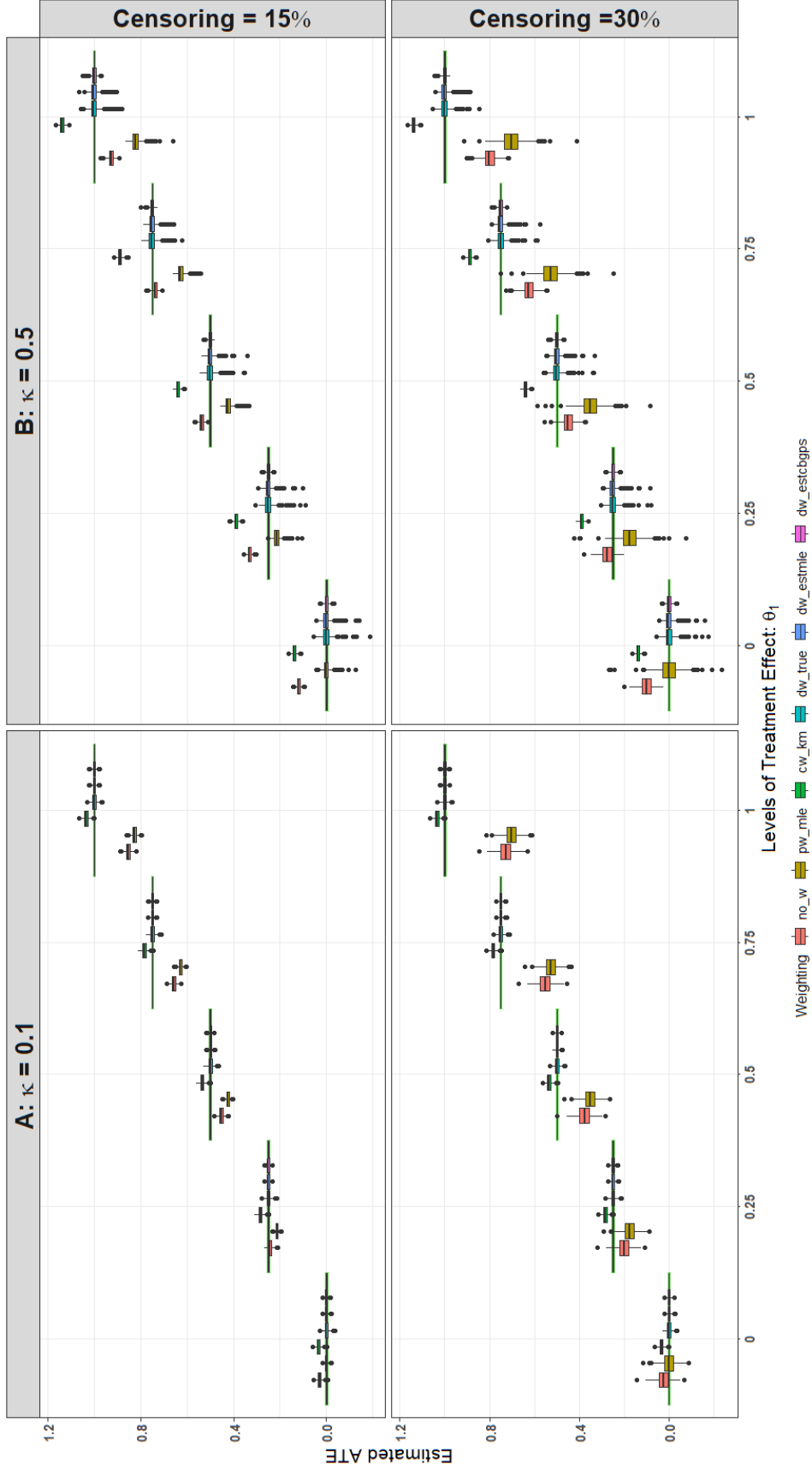


Figure 1.3: Boxplots of 1000 ATE estimates based on sample size $n = 5000$ at different treatment levels, with two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B). Each cluster of boxplots shows ATE estimates based on different weighting methods.

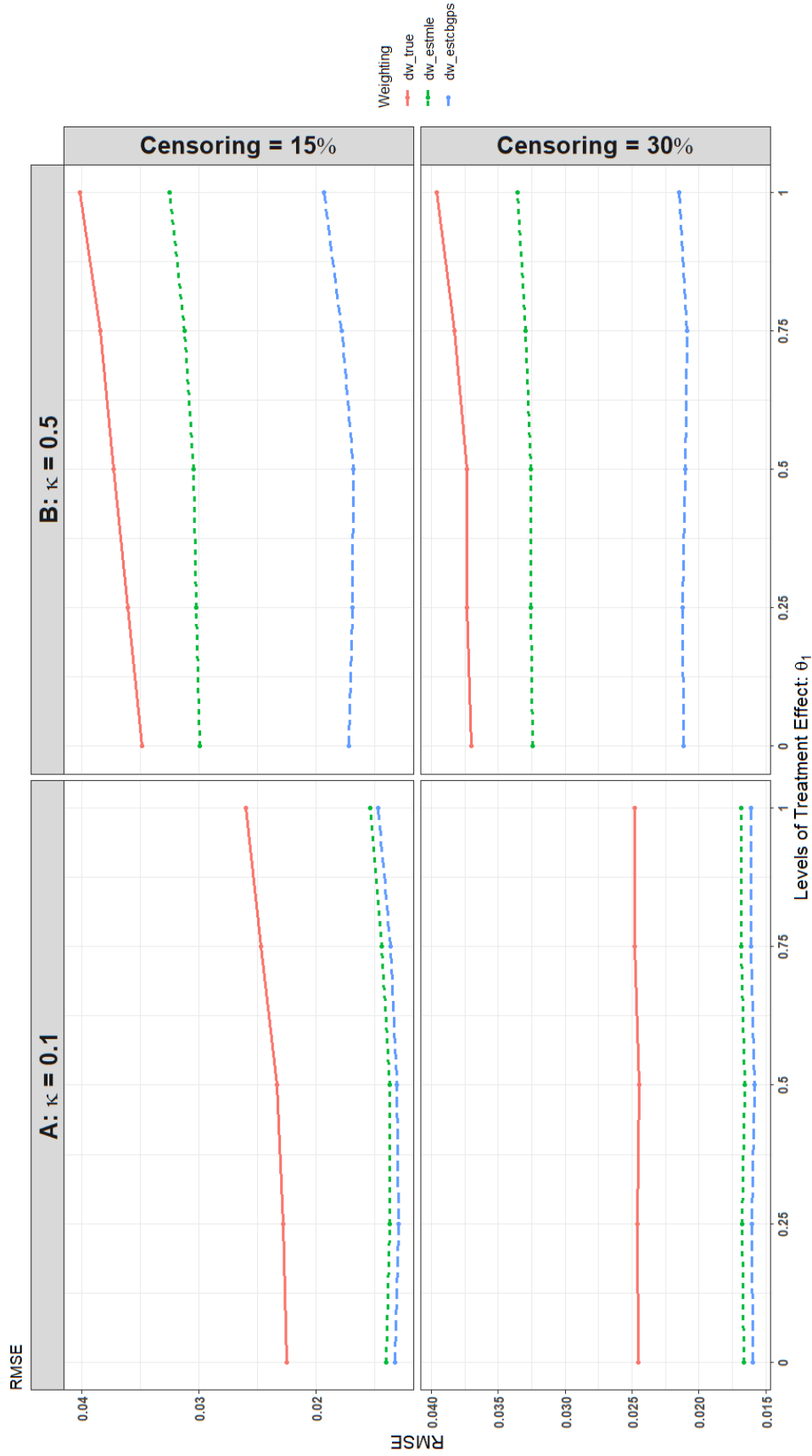


Figure 1.4: The scatter plot illustrates the RMSE based on three different double weighting methods and a sample size of $n = 1000$ at various treatment levels. The plot includes two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B).

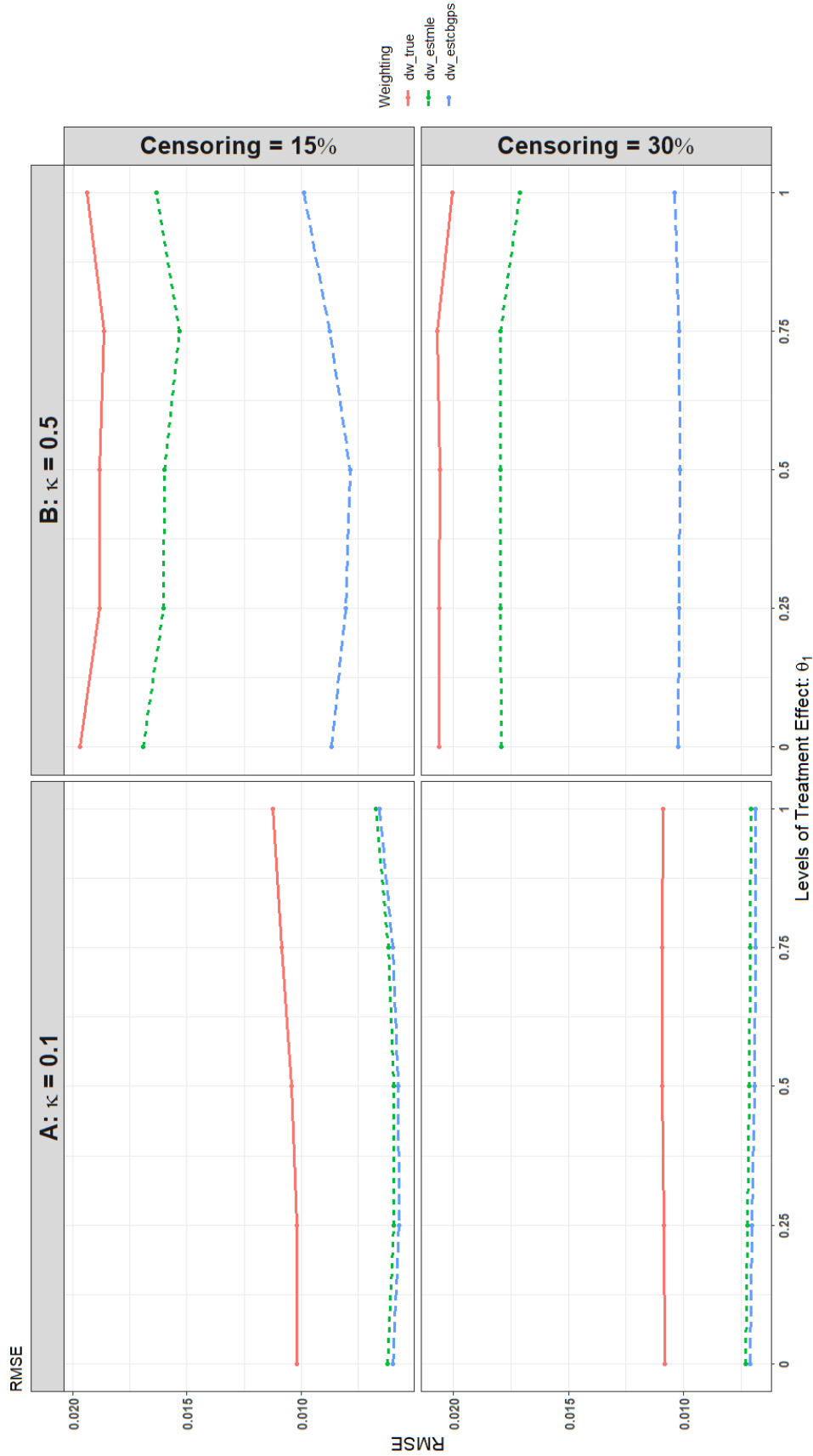


Figure 1.5: The scatter plot illustrates the RMSE based on three different double weighting methods and a sample size of $n = 5000$ at various treatment levels. The plot includes two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B).

weighting methods were quite similar and close to the nominal rate of 95%. However, when the association between confounding variables and treatment became moderate (Panel B), only the coverage rates of the double weighting method based on CBGPS were closest to the 95% nominal coverage rate.

These simulation results clearly underscore the importance of addressing confounding and censoring when estimating the ATE for time-to-event outcomes. Unbiased results can be achieved by incorporating both propensity score weights and censoring weights in the AFT-MSM for outcomes with censored observations. Increasing the sample size enhances the performance of our proposed model. Moreover, propensity score weighting by the CBGPS method proves to be more robust in estimating ATE. From all the figures, it is evident that the double weighting method by CBGPS outperforms others, exhibiting smaller bias, smaller variance, smaller RMSE, and a coverage rate closest to the 95% nominal rate.

1.4 Case study: blood lead level versus mortality based on NHANES III dataset

The relationship between blood lead levels (BLL) and mortality is a complex and multifaceted topic. Lead is a toxic metal that can affect multiple organ systems in the body, including the nervous, cardiovascular, and renal systems (Ara et al., 2015). Studies have suggested that elevated blood lead levels may be associated with an increased risk of mortality, particularly due to cardiovascular diseases. Lead exposure has been linked to hypertension, atherosclerosis, and other cardiovascular conditions that can contribute to mortality (Schober et al., 2006; Pirkle et al., 1994). However, it's essential to note that the relationship between BLL and mortality is not always straightforward, and can be confounded by socioeconomic status, lifestyle, and environmental factors. In our present study, we employed data from the third National

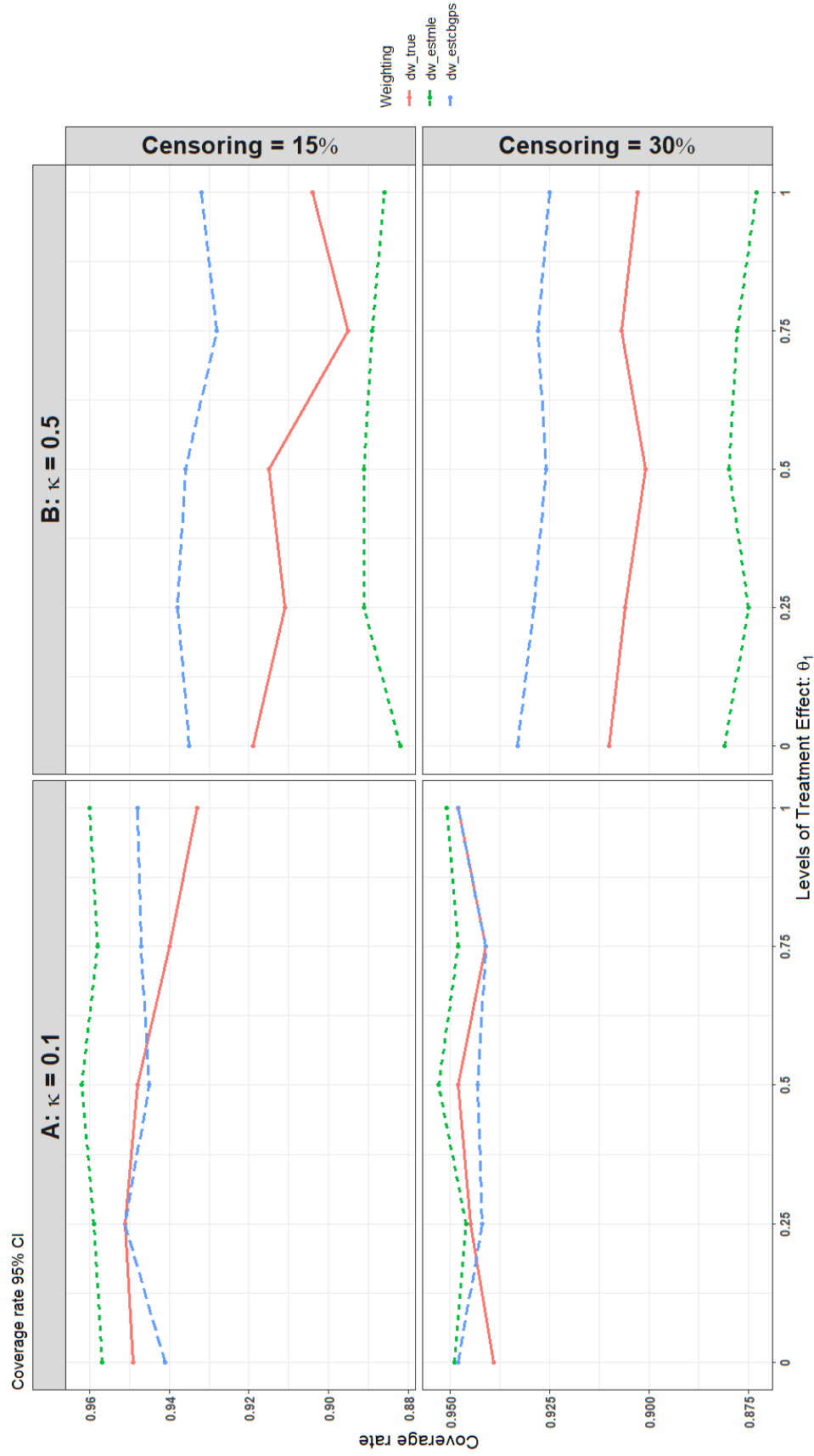


Figure 1.6: The scatter plot illustrates the true coverage rate for 95% CI based on three different double weighting methods and a sample size of $n = 1000$ at various treatment levels. The plot includes two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B).

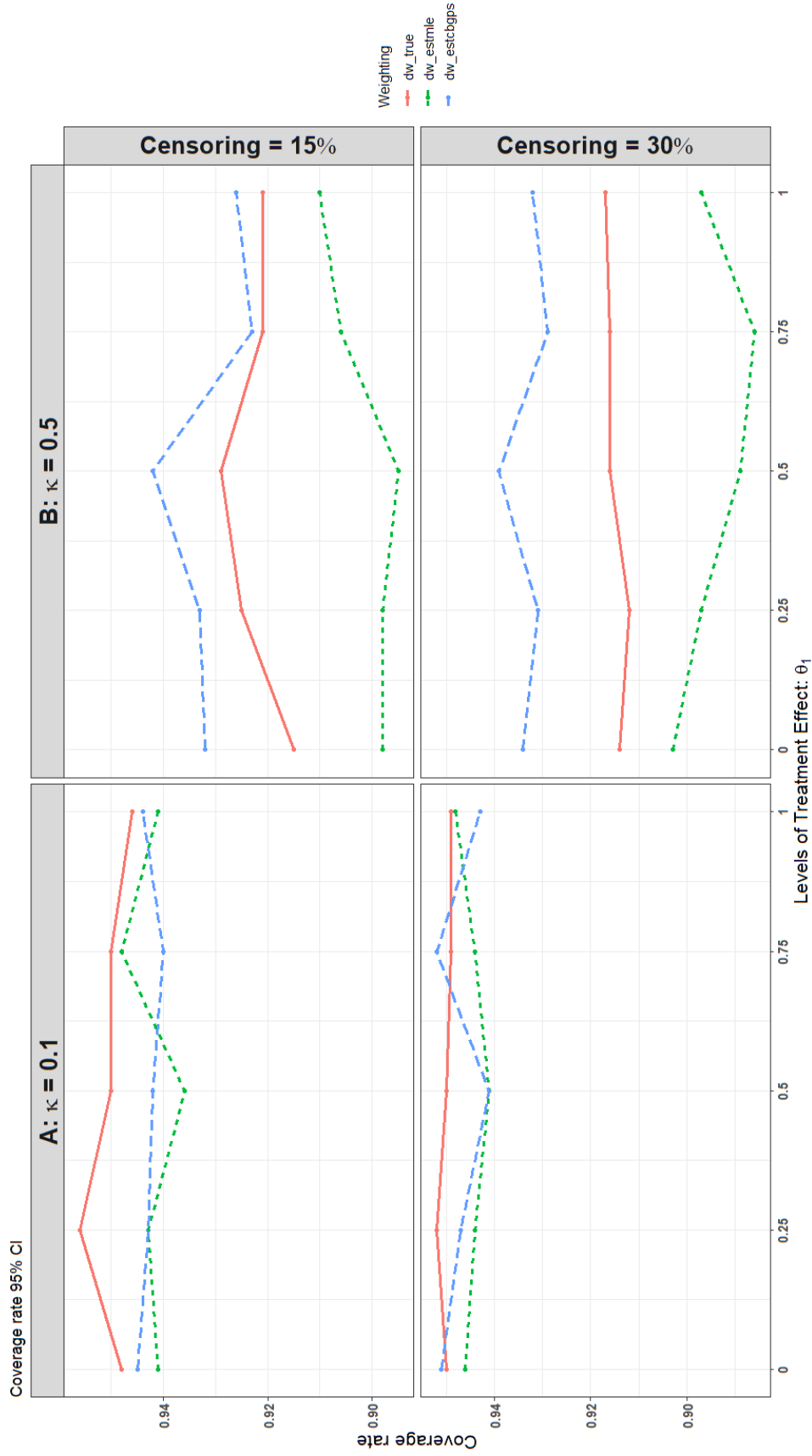


Figure 1.7: The scatter plot illustrates the true coverage rate for 95% CI based on three different weighting methods and a sample size of $n = 5000$ at various treatment levels. The plot includes two levels of censoring rates (15% and 30%, indicated in the first row and second row, respectively) and two levels of association strength between confounding variables and treatment: $\kappa = 0.1$ for low association (shown in Panel A) and $\kappa = 0.5$ for moderate association (displayed in Panel B).

Health and Nutrition Examination Survey (NHANES III), conducted in two phases from 1988 to 1994 (Data, 1994). The aim was to explore the causal relationship between BLL and mortality, taking into account confounding variables and adjusting for censored observations using our proposed method. To assess mortality, information up to December 31st, 2019, regarding NHANES III participants was gathered from the Linked Mortality Files (LMF) provided by the National Center for Health Statistics (NCHS).

The blood lead levels (BLL) for participants in NHANES III were obtained during the mobile physical examination day, and a detailed description of BLL measurements can be found in Pirkle et al. (1994). The limit of detection (LOD) for the lead level was 1.0 mg/dL, and for participants with BLL below LOD, it was imputed as $\text{LOD}/\sqrt{2}$, resulting in 0.7mg/dL . In our analysis, we focused on participants within the age range of 50 – 70 years at the time of the survey interview and excluded individuals with accidental death records. BLL was considered as the treatment variable in our analysis. To satisfy the normality assumption for the estimation of propensity score weights, we applied a log transformation to BLL.

For our covariates, we included age, sex, race-ethnicity, metro area, education level, poverty income ratio (PIR), exercise, smoking, and alcohol consumption. Additionally, we considered comorbidities such as cancer, stroke, cardiovascular disease, diabetes, and chronic kidney disease (CKD). Race and ethnicity information for participants was categorized as Non-Hispanic Black, Non-Hispanic White, Hispanic origin, and others. Education levels were classified into three categories: no education, less than high school, and College or above. The PIR (Poverty-to-Income Ratio) was determined by evaluating the family income in relation to the poverty threshold, adjusted for both family size and the annual inflation status. For the exercise variable, the frequency of any physical activity undertaken by participants was aggregated and expressed as “times per month”. Smoking status was classified into

three levels: “former” for those who had smoked more than 100+ cigarettes in their lifetime but had since quit, “current” for participants who were smoking at the time of the survey, and “never” for those who had not smoked more than 100+ cigarettes in their life time. Participants’ alcohol consumption was categorized into three groups: “former” if they had consumed at least 12 alcoholic drinks in their lifetime but had not had at least 12 drinks in the last 12 months, “current” if they had consumed at least 12 drinks in the last 12 months, and “never” otherwise. A participant was identified as having diabetes based on self-report during the interview or having glycated hemoglobin $\geq 6.5\%$ or plasma glucose $\geq 125 \text{ mg/dL}$. Participants with a history of heart disease were identified if they reported a past heart attack or had congestive heart failure. Participants were identified with CKD based on lab reports indicating estimated GFR (eGFR) $\leq 60 \text{ ml/min per } 1.73\text{m}^2$ or Urine albumin to creatinine ratio (ACR) ≥ 30 (Selvin et al., 2007). eGFR is creatinine-based estimated glomerular filtration rate. Serum creatinine was measured using a kinetic rate Jaffe method and all serum creatinine measurements were re-calibrated to standardized creatinine measurements (Coresh et al., 2007). The total number of patients within the specified age range of 50-70, considering all the mentioned covariates, exposure, and outcome, was 3621. Out of these, 2456 individuals passed away before December 31st, 2019, while the remaining 1165 participants were alive, accounting for 32% of censoring.

The study by Schober et al. (2006) suggested that individuals with BLL exceeding $\text{BLL } 5 \text{ ug/dL}$ are more susceptible to mortality. To evaluate potential confounding effects of certain covariates on the relationship between exposure to BLL and mortality, participants were classified into two groups: “low BLL” for those with $\text{BLL} \leq 5 \text{ ug/dL}$ and “high BLL” for those with $\text{BLL} > 5 \text{ ug/dL}$. The distribution of characteristics among participants for each categorical covariate is summarized in Table 1.1. The “Total #(%)”, “Died #(%)”, and “P-value” columns present the overall distribution, mortality rates, and associated p-values, respectively. Additionally, the

association between each categorical covariate and BLL exposure is detailed in the last two columns of Table 1.1. Notably, it is evident that (1) males exhibited a higher mortality rate than females (73.7% vs. 62.3%), with males also having a higher BLL exposure than females (42.6% vs. 20.7%); (2) Black participants had a higher mortality rate than white participants (73.3% vs. 67.9%), and Black individuals had higher BLL exposure than their white counterparts (47.5% vs. 23.7%). Similar trends are observed for current smokers and participants with hypertension, which are detailed in Table 1.1. Moving on to continuous confounding variables, Table 1.2 presents the mean and standard deviation (SD) values. These statistics are further stratified based on survival status and BLL exposure (low vs. high). Noteworthy findings include (1) participants who died and those exposed to high BLL had lower Poverty-to-Income Ratio (PIR); (2) individuals who died and those with high BLL exposure engaged in less physical exercise. These observations suggest a potential confounding effect of PIR and physical exercise on the relationship between BLL and mortality.

To address this confounding, we employed the proposed double weighting method with propensity score weights estimated by CBGPS. This approach aims to untangle the causal relationship between BLL and survival, accounting for the influence of confounding variables. Figure 1.8 displays a scatterplot featuring each continuous covariate plotted against Blood Lead Level (BLL), showcasing regression lines derived from both the original and weighted samples. The continuous covariates, such as age, poverty-income ratio, BMI, and physical exercise level, are presented in a logarithmic scale. Notably, the regression lines, adjusted by the CBGPS weights, reveal a shift towards zero, indicating that in the weighted sample, these covariates no longer exhibit a significant association with BLL. This adjustment underscores the effectiveness of the CBGPS weights in mitigating the covariate-exposure association.

In Figure 1.9, stacked barplots illustrate the distribution of characteristics for each categorical covariate, stratified by low and high BLL exposure levels. The initial

Table 1.1: Summary of the distribution of characteristics among participants for each categorical covariate (refer to column “Total #(%)”), its association with mortality (refer to columns “Died #(%)” and “P-value”), and its association with BLL (refer to the last two columns)

Covariates	Category	Total #(%) ^a	Died #(%) ^b	P-value	BLL High #(%) ^c	P-value (BLL)
Gender	Male	1758 (48.5%)	1296 (73.7%)	< 0.001	748 (42.5%)	< 0.001
	Female	1863 (51.4%)	1160 (62.3%)		385 (20.7%)	
Race Ethnicity	White	1754 (48.4%)	1191 (67.9%)	< 0.001	415 (23.7%)	< 0.001
	Black	849 (23.4%)	622 (73.3%)		403 (47.5%)	
	Hispanic	868 (24.0%)	564 (65.0%)		269 (31.0%)	
	Others	150 (4.1%)	79 (52.7%)		46 (30.7%)	
Education	No education	147 (4.1%)	99 (67.3%)	< 0.001	61 (41.5%)	< 0.001
	≤ High School	2626 (72.5%)	1851 (70.5%)		879 (33.5%)	
	College or higher	848 (23.4%)	506 (59.7%)		193 (22.8%)	
Metro Area	Non_Metro	1960 (54.1%)	1369 (69.8%)	0.005	543 (27.7%)	< 0.001
	Metro	1661 (45.9%)	1087 (65.4%)		590 (35.5%)	
Smoking	Never	1496 (41.3%)	871 (58.2%)	< 0.001	321 (21.5%)	< 0.001
	Current Smoker	867 (23.9%)	692 (79.8%)		432 (49.8%)	
	Former Smoker	1258 (34.7%)	893 (71.0%)		380 (30.2%)	
Alcohol Consumption	Never	691 (19.1%)	466 (67.4%)	0.02	157 (22.7%)	< 0.001
	Current Drinker	1421 (39.2%)	930 (65.4%)		563 (39.6%)	
	Former Drinker	1509 (41.7%)	1060 (70.25%)		413 (27.4%)	
Hypertension	0	2018 (55.7%)	1239 (61.4%)	< 0.001	591 (29.3%)	0.004
	1	1603 (44.3%)	1217 (75.9%)		542 (33.8%)	
Diabetes	0	2830 (78.2%)	1805 (63.8%)	< 0.001	894 (31.6%)	0.488
	1	791 (21.8%)	651 (82.3%)		239 (30.2%)	
Cardiovascular Disease	0	3271 (90.3%)	2152 (65.8%)	< 0.001	1004 (30.7%)	0.021
	1	350 (9.7%)	304 (86.9%)		129 (36.9%)	
Cancer	0	3272 (90.4%)	2191 (67.0%)	0.001	1037 (31.7%)	0.123
	1	349 (9.6%)	265 (75.9%)		96 (27.5%)	
CKD	0	406 (11.2%)	260 (64.0%)	0.093	130 (32.0%)	0.780
	1	3215 (88.8%)	2196 (68.3%)		1003 (31.2%)	
Stroke	0	3480 (96.1%)	2328 (66.9%)	< 0.001	1088 (31.3%)	0.944
	1	141 (3.9%)	128 (90.8%)		45 (31.9%)	

^a Column % for each covariate; ^b Mortality % for each level of a covariate; ^c % of people with high BLL for each level of a covariate; P-values obtained from Chi-square tests

Table 1.2: The Mean and standard deviation (SD) for each continuous variable (refer to column “All”), stratified by mortality (refer to columns “Alive” and “Died”) and BLL (refer to columns “low BLL” and “high BLL”).

Covariates	All	Alive	Died	P-val	low BLL	high BLL	P-val(BLL)
AGE	60.36 (5.97)	56.82 (5.21)	62.04 (5.57)	< 0.001	60.21 (6.02)	60.69 (5.87)	0.023
PIR	2.77 (1.96)	3.16 (2.02)	2.58 (1.9)	< 0.001	2.95 (2)	2.36 (1.79)	< 0.001
BMI	28.13 (5.6)	27.98 (5.13)	28.2 (5.81)	0.279	28.55 (5.76)	27.21 (5.11)	< 0.001
Exercise	19.76 (25.22)	21.94 (27.11)	18.73 (24.2)	< 0.001	20.4 (25.53)	18.37 (24.47)	0.025

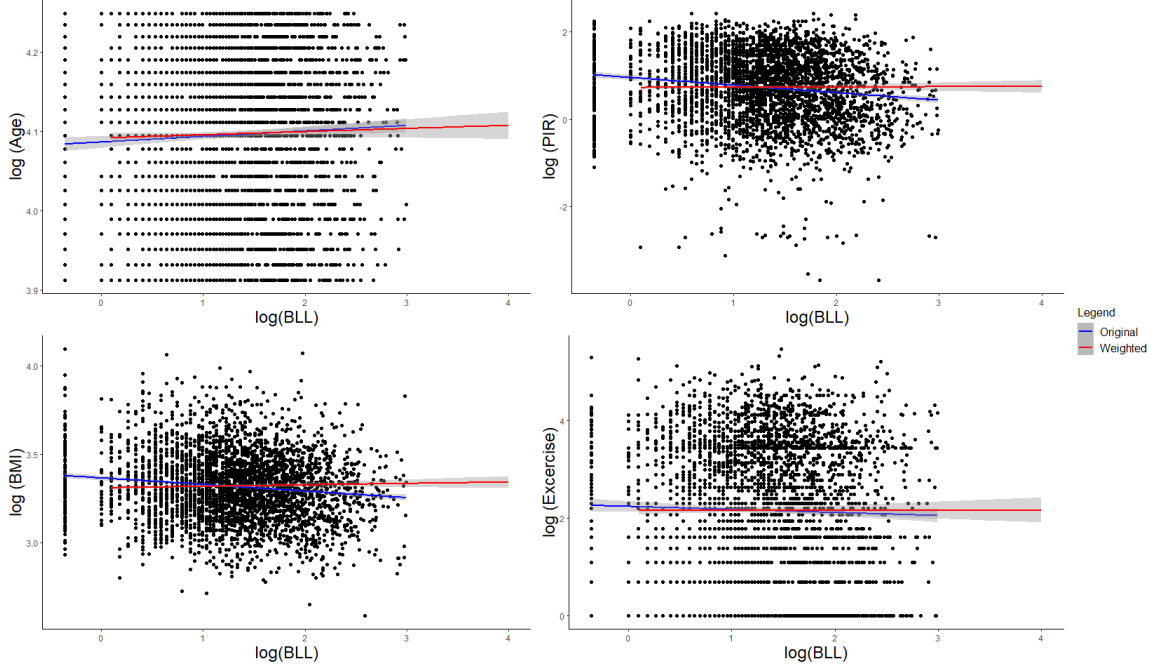


Figure 1.8: Scatterplot illustrating the relationships between continuous covariates (age, PIR, BMI, and exercise level) and Blood Lead Levels (BLL), accompanied by the fitted regression lines in both the original and the weighted sample.

observation reveals distinct covariate distributions between low and high BLL exposure. However, considering CBGPS weights in the weighted sample, covariate distributions emerge similar between low and high BLL exposure groups. For instance, in the original sample, there was a higher percentage of females in the low BLL group and a higher percentage of males in the high BLL group but in the weighted sample the gender distribution becomes remarkably similar between low and high BLL exposure groups. This further supports the effectiveness of the IPTW methodology in balancing covariate distributions, thereby enhancing the comparability between different BLL exposure levels.

Figure 1.10 illustrates the scatterplot of time-to-death in months versus BLL among the participants who died. The plot features two predicted trend lines: one derived from the original sample (solid line) and the other from the weighted sample using double weights (dashed line). In both the original sample and the weighted sample, a noticeable decrease in months to death is observed when BLL exceed 5

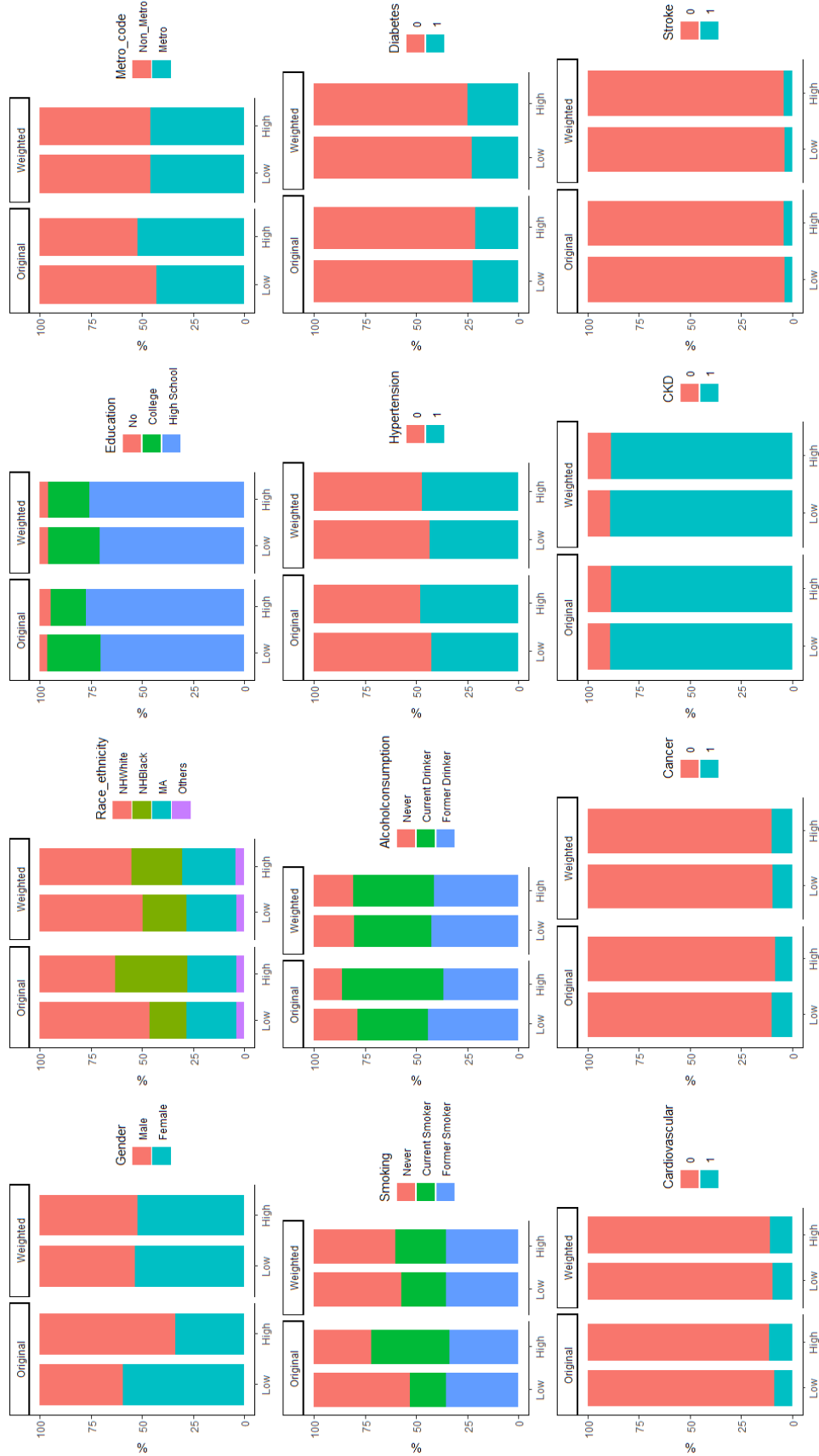


Figure 1.9: Barplot presenting the distributions of each categorical covariate in low Blood Lead Level (BLL) and high BLL, both in the original sample (the first two bars in each panel) and the weighted sample (the last two bars in each panel).

ug/dL (equivalent to 1.6 in log-scale). In other words, as BLL increases, survival time decreases, though the decline rate in the weighted sample is less pronounced than in the original sample. To formally investigate the causal relationship between BLL and survival time, considering 5 ug/dL as a change point, we employed the Accelerated Failure Time-Marginal Structural Model (AFT-MSM) with a change point at $\log(5)$:

$$\log(Y) = \theta_0 + \theta_1 \log(BLL) + \theta_2 (\log(BLL) - \log(5))_+ + \epsilon. \quad (1.13)$$

Here the function $(x)_+$ takes value x if $x \geq 0$, and 0 if $x < 0$. Table 1.3 shows the estimated parameters obtained from the model in (1.13), both without and with the application of double weights. The standard errors and 95% confidence intervals for these estimated parameters were obtained through bootstrap sampling. Based on the results from the double weighting approach, it can be concluded that when BLL is less than 5 ug/dL , there is no significant causal association with time to death. However, when BLL exceeds this threshold, it becomes significantly causally associated with mortality time. Specifically, with a 1-unit increase in BLL in log-scale, there is a 21.3% decrease in survival time in months.

In Figure 1.11, we compared the survival curves between the low BLL and high BLL in both the original sample and the weighted sample. Log-rank tests, both without and with weights, were applied to examine the survival difference between the low BLL and high BLL groups in both samples. We conclude that a significant difference exists in the survival curves between the low BLL and high BLL groups in both the original and weighted samples. Although the difference in the weighted sample is not as pronounced as in the original sample, adjusting for the effects of confounding covariates through weighting demonstrates a clear causal association between BLL and survival time.

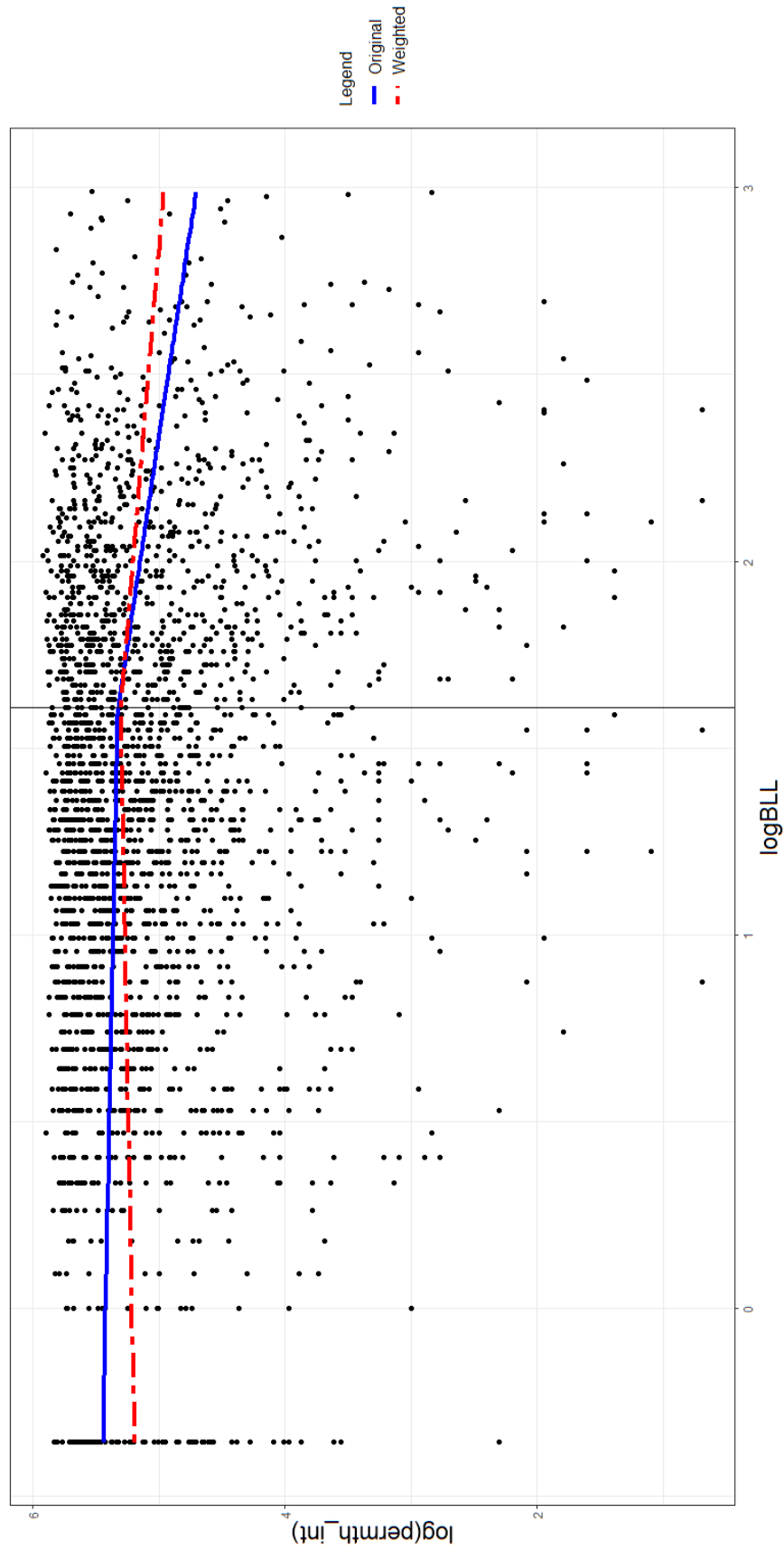


Figure 1.10: Scatterplot depicting survival time (in months since survey) against Blood Lead Levels (BLL), accompanied by the fitted trend lines in both the original sample and the weighted sample.

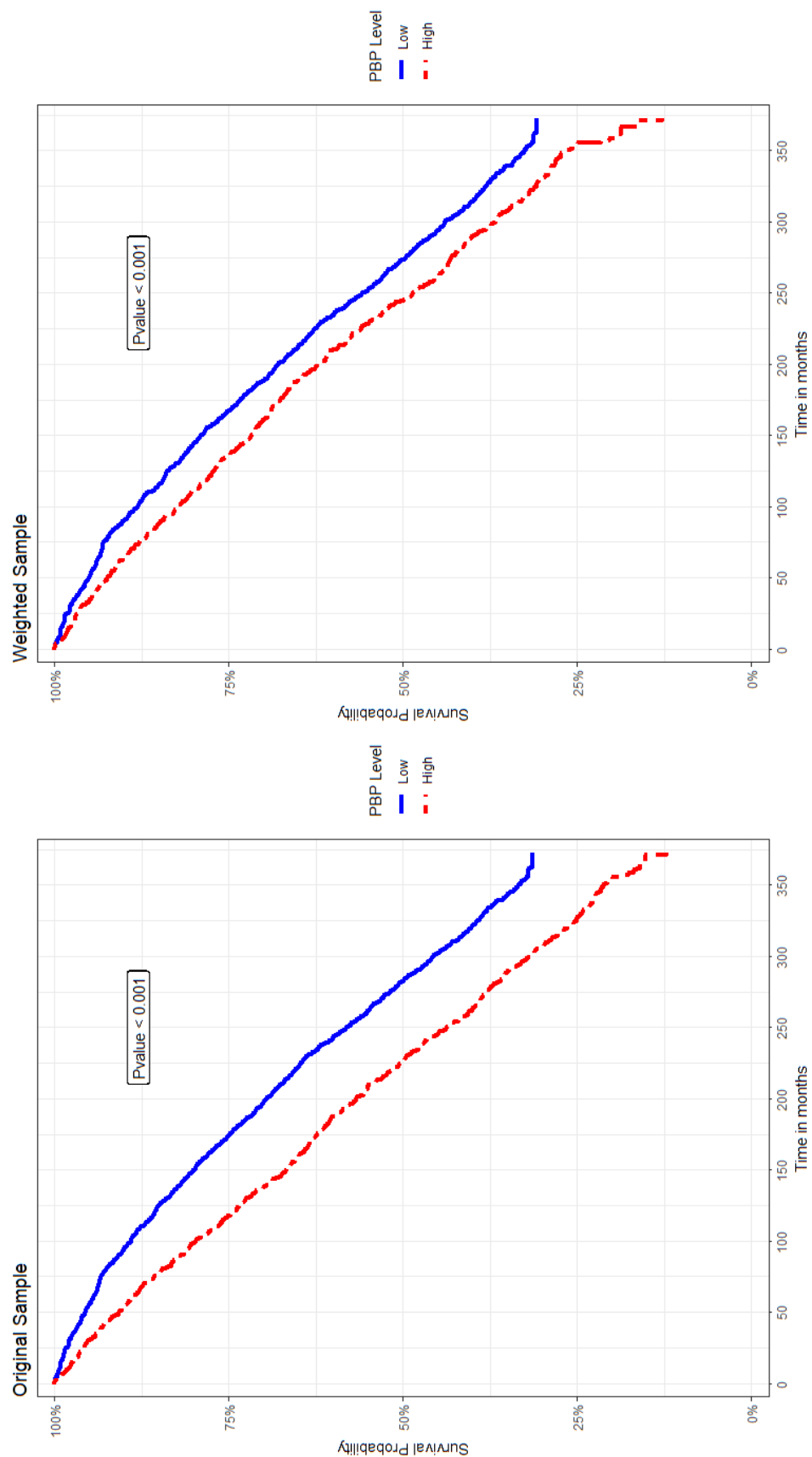


Figure 1.11: Kaplan-Meier survival curves illustrating the survival outcomes for low Blood Lead Level (BLL) and high BLL in the original sample (left panel) and the weighted sample (right panel), along with the log-rank tests.

Table 1.3: Estimated effect of Blood lead on time to death, Standard Error and Confidence Interval at 95%

No weighting			
Parameters	Estimate	Std Error	95% CI
θ_0	5.129	0.040	(5.051, 5.207)
θ_1	-0.019	0.034	(-0.085, 0.047)
θ_2	-0.368	0.093	(-0.55, -0.187)
$\theta_1 + \theta_2$	-0.388	0.074	(-0.534, -0.242)

Double Weighting: CBGPS			
Parameters	Estimate	Std Error	95% CI
θ_0	5.212	0.048	(5.118, 5.306)
θ_1	0.055	0.040	(-0.023, 0.133)
θ_2	-0.301	0.131	(-0.556, -0.045)
$\theta_1 + \theta_2$	-0.246	0.108	(-0.458, -0.034)

1.5 Discussion

In the context of abundant observational data, estimating ATE of exposure on outcomes within a causal framework becomes essential. Controlling for confounding bias arising from patient characteristics that influence both the exposure and the outcome is crucial in achieving accurate estimates. In scenarios involving continuous treatment settings, a shift from traditional propensity score weighting to the use of generalized propensity scores is employed to enhance the precision of causal inference.

In this chapter, we employed the double weighting method to estimate the average treatment effect of a continuous treatment on survival outcomes with censored observations. The double weighting methods utilized propensity score weights to control the effect of confounding and censoring weights to address bias in estimation resulting from right censoring. Based on the simulation results, we can assert the significance of incorporating censoring weights when estimating treatment effects for time-event outcomes. In this study, for obtaining propensity score weights, we compared two methods: weights estimated by MLE and the CBGPS method. Our conclusion is that the double weighting method with CBGPS performs the best, ex-

hibiting the least bias, variance, and the smallest overall variability.

In this paper, we have assumed random censoring for estimating the censoring weights for our proposed method. Our approach can be extended by using conditional Kaplan Meier estimator for censoring weights in case of dependant censoring. However, the detailed investigation will be carried out in our future work.

CHAPTER 2

CAUSAL MEDIATION ANALYSIS WITH EXTENSION TO SURVIVAL ANALYSIS AND HEALTH RACIAL DISPARITY STUDIES

2.1 Introduction

In healthcare research, randomized control trials (RCTs) are commonly regarded as the gold standard for establishing the causal effect of an exposure variable, denoted as A , on the outcome, denoted as Y . In RCTs, eligible participants are randomly assigned to either the treatment group or the control group. However, RCTs are not always feasible, especially when race is the exposure variable. On the contrary, observational studies are prevalent, and they have been utilized to investigate the causal effect of exposure with appropriate adjustments for confounding variables. In observational studies, the relationship between the exposure and the outcome is often confounded by the variables \mathbf{X} . It is crucial to consider the impact of these confounding variables and eliminate their effect. Figure 2.1 (a) illustrates a directed acyclic graph (DAG) of an RCT, where the treatment assignment of A is independent of \mathbf{X} , even though \mathbf{X} may influence the outcome Y . Figure 2.1 (b) depicts a scenario where the relationship between A and Y is confounded by \mathbf{X} , with \mathbf{X} being causally associated with both A and Y . In medicine and healthcare research, more complex pathways are often encountered. Mediation analysis aims to uncover the intricate mechanisms underlying the relationship between an exposure variable (e.g., A) and

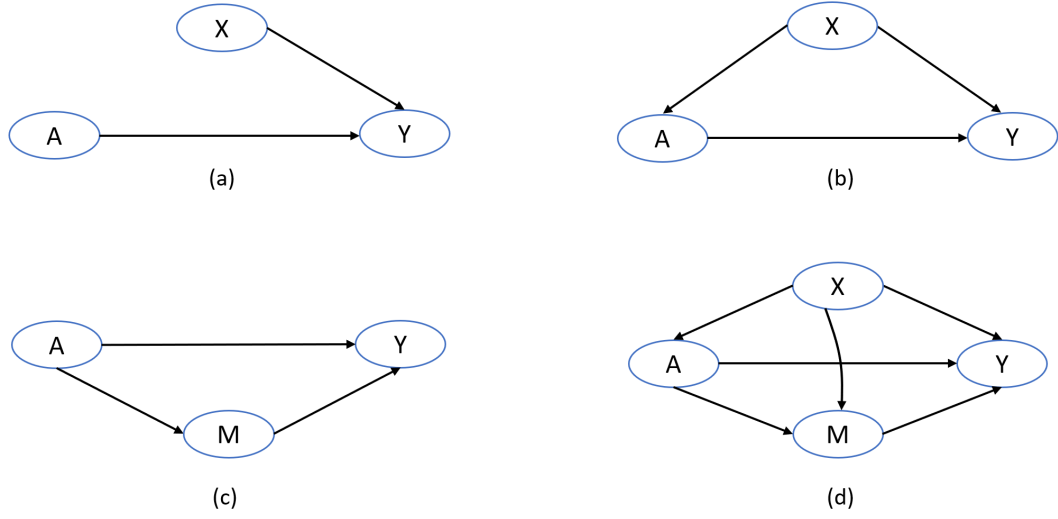


Figure 2.1: Illustration of directed acyclic graph (DAGs) for various study designs: (a) randomized control trials; (b) observational studies; (c) simple mediation model; and (d) a mediation model with confounding variables.

an outcome variable (e.g., Y) through an intermediary variable referred to as a mediator (e.g., M), as depicted in Figure 2.1 (c) (Baron and Kenny, 1986). Unlike the typical estimation of the average treatment effect of exposure variable A on outcome Y , a mediation model posits that A influences the mediator variable M , which subsequently affects Y . Subsequent research has expanded the application of mediation analysis within a causal framework, as depicted in Figure 2.1(d), representing a mediation model in the presence of confounding variables \mathbf{X} . Analyzing this pathway often provides a deeper understanding of the causal chain in complex systems. Mediation analysis decomposes the treatment effect into the mediation effect and the direct effect, where direct effect measures the effect of exposure A on outcome Y directly while indirect effect captures the effect of A on Y through mediator M . The mediator variable explains the mechanism of the relationship between the exposure and outcome variables (MacKinnon et al., 2007; Pearl, 2001).

Mediation analysis, initially introduced by Baron and Kenny (1986), has become an active research area with various methodologies and applications proposed.

For instance, VanderWeele and Vansteelandt (2009) utilized an outcome model incorporating interaction terms between the mediator and exposure, as depicted in Figure 2.2(a). This approach has been further extended to scenarios involving multiple mediators and interactions (VanderWeele, 2015). Additionally, these scenarios encompass situations where multiple mediators may operate either in parallel or in a sequential order. Figure 2.2(b) illustrates a causal mediation graph with two parallel mediators. To analyze the intricate relationships of the mediated pathways along with interactions between the mediator and exposure, the two-way decomposition has been extended to a more comprehensive 4-way decomposition (VanderWeele, 2014). This advanced decomposition of effects from the mediation model with multiple mediators facilitates the identification of the individual effect of each mediator and each of their interactions with the exposure. Thus, in the context of complex causal relationships, researchers can analyze the causal pathways and evaluate the specific contribution of each component. In this chapter, we will expand the scope of mediation analysis to encompass survival analyses and delve into its application in the context of racial disparity studies.

While mediation analysis has found applications in various research domains, it has been particularly prevalent in investigating racial disparities in healthcare, education, and medicine. When applied to the context of structural racism, mediation analysis serves as a valuable tool for researchers to explore how societal and institutional factors, rooted in historical and systemic biases, contribute to disparities among racial or ethnic groups. According to the Institute of Medicine’s (IOM) definition (Nelson, 2002), healthcare disparities related to race and ethnicity are deemed unjustifiable when they result from factors other than clinical need and patient preferences. Differences in healthcare access and use based on age, gender, health status, and patient choices are considered acceptable. However, any other sources of differences, including those influenced by a patient’s socioeconomic status and other social

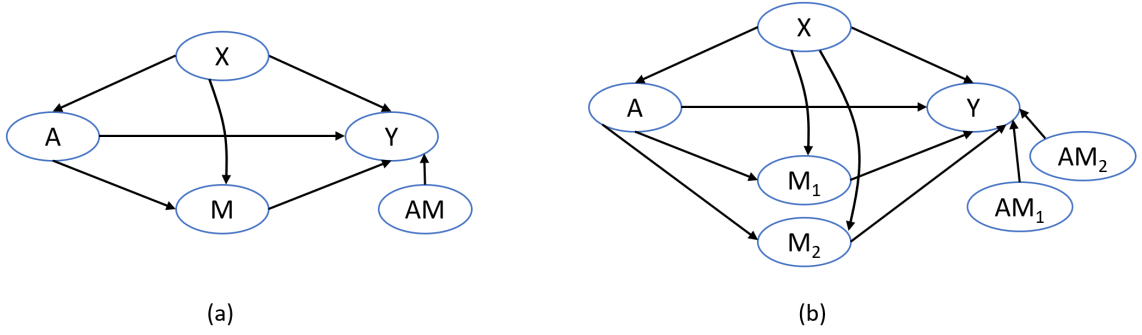


Figure 2.2: Illustration of mediation models with interaction of exposure and mediator: (a) one mediator and (b) multiple mediators

determinants of health, as well as inequalities in the functioning of the healthcare system caused by systemic racism, are regarded as unjust. In observational studies aimed at assessing health disparities, researchers often use regression models to estimate the effect of race. However, it is crucial to recognize that multiple factors can contribute to these disparities. Consequently, researchers are interested in identifying the key drivers of racial discrimination among groups—factors that can potentially be addressed through targeted interventions.

We structure the remainder of our project as follows. In Section 2.2, we provide detailed insights into expressing the total effect of the exposure on the outcome in individual effects along the pathways of multiple mediators and discuss its interpretation within a causal framework. In Section 2.3, we extend the mediation effect to survival outcomes, demonstrating that the IOM-defined disparity can be articulated in terms of the 4-way decomposition effects derived from the mediation model analysis. Section 2.4 is dedicated to extensive simulation studies designed to examine estimation accuracy. In Section 2.5, we conduct a case study to explore the factors influencing health racial disparity based on the NHANES III database. The final section is reserved for discussion and conclusion.

2.2 General framework for mediation analysis

Let us denote (A, M, \mathbf{X}, Y) as the random variables involved in a mediation model, let A denote as the exposure of interest, Y the outcome of interest, M the mediator, and \mathbf{X} confounding variables which impact all other three variables A , M , and Y .

2.2.1 Potential outcomes and basic assumptions

We start with the simple mediation model (refer to Figure 2.1 panel c), and we then proceed to elaborate the cases with confounding variables (refer to Figure 2.1 panel d) and multiple mediators (refer to Figure 2.2). For the simple mediation model, to compare between two exposure levels of A , say a and a^* , we denote the potential mediator M under exposure levels a and a^* as $M(a)$ and $M(a^*)$, respectively. The potential outcome, given that A and M are set to a and m , is defined as $Y(a, m)$. $Y(a, M(a^*))$ is the potential outcome when the exposure is set at level a and the mediator is set at the level of $M(a^*)$.

In mediation analysis, researchers often make the consistency and composition assumptions (VanderWeele and Vansteelandt, 2009). The consistency assumption states that the observed outcome Y is equal to the potential outcome of $Y(a)$ and the observed mediator M is equal to the potential mediator of $M(a)$ when the exposure A takes value of a . The composition assumption states that the potential outcome $Y(a)$ equals to the potential outcome $Y(a, M(a))$ where the exposure A is set to a and mediator is set at the value $M(a)$. Under the consistency and composition assumptions (VanderWeele and Vansteelandt, 2009), the causal effect, $E\{Y(a) - Y(a^*)\}$, which is also referred as the total effect of A on Y , can be decomposed in two components as natural direct effect (NDE) and natural indirect effect (NIE). That is,

$$E[Y(a) - Y(a^*)] = E[Y(a, M(a^*)) - Y(a^*, M(a^*))] + E[Y(a, M(a)) - Y(a, M(a^*))]. \quad (2.1)$$

Here $NDE = E[Y(a, M(a^*)) - Y(a^*, M(a^*))]$, which captures the direct causal effect of A on Y with the mediator fixed at the value of $M(a^*)$. $NIE = E[Y(a, M(a)) - Y(a, M(a^*))]$, which captures the indirect causal effect through the mediator M as mediator changes from level of $M(a)$ to $M(a^*)$. Generalising to finite number of mediators $\mathbf{M} = (M_1, \dots, M_Q), Q \in \mathbb{R}$, the identification of direct effects and indirect effects require the following four assumptions (VanderWeele and Vansteelandt, 2009; Gao et al., 2022):

- (1) There is no unmeasured exposure-outcome confounding

$$Y(a, m_1, \dots, m_Q) \perp\!\!\!\perp A \mid \mathbf{X};$$

- (2) There is no unmeasured mediator-outcome confounding

$$Y(a, m_1, \dots, m_Q) \perp\!\!\!\perp M_q \mid (A, \mathbf{X});$$

- (3) There is no unmeasured exposure-mediator confounding

$$M_q(a) \perp\!\!\!\perp A \mid \mathbf{X};$$

- (4) There is no unmeasured mediator-outcome confounding influenced by exposure

$$Y(a, m_1, \dots, m_Q) \perp\!\!\!\perp M_q(a^*) \mid \mathbf{X}.$$

2.2.2 Four fold decomposition of total effect

Considering the interaction between exposure and mediation, which can influence the outcome as illustrated in Figure 2.2 Panel (a), the total causal effect is dissected into three components. The NIE is further decomposed into mediated interaction effect

(INT_{med}) and pure indirect effect (PIE). That is,

$$\begin{aligned}
E[Y(a) - Y(a^*)] = & E[Y(a, M(a^*)) - Y(a^*, M(a^*))] \\
& + E\left[\left(Y(a, M(a)) - Y(a^*, M(a))\right) - \left(Y(a, M(a^*)) - Y(a^*, M(a^*))\right)\right] \\
& + E[Y(a^*, M(a)) - Y(a^*, M(a^*))]
\end{aligned} \tag{2.2}$$

Here, the second term in Equation (2.2) is denoted as INT_{med} , signifying its role in capturing the interaction between the mediator and the exposure. This term precisely reflects the difference in direct effects when the mediator is held at $M(a)$ compared to $M(a^*)$. To understand the complex relationships in the mediation pathways with interaction, VanderWeele (2014) has decomposed the total effect of A on Y into 4 components, where the first term in (2.2), namely, NDE is further decomposed into two components as the controlled direct effect (CDE) and interaction effect at reference level of mediator (INT_{ref}). That is,

$$\begin{aligned}
E[Y(a) - Y(a^*)] = & E[Y(a, m^*) - Y(a^*, m^*)] \\
& + E\left[\left(Y(a, M(a^*)) - Y(a^*, M(a^*))\right) - \left(Y(a, m^*) - Y(a^*, m^*)\right)\right] \\
& + E[Y(a, M(a)) - Y(a^*, M(a)) - Y(a, M(a^*)) + Y(a^*, M(a^*))] \\
& + E[Y(a^*, M(a)) - Y(a^*, M(a^*))].
\end{aligned} \tag{2.3}$$

Here m^* is a fixed value for M . CDE represents the direct effect of the exposure when M is fixed at the value m^* . $INT_{ref}(m^*)$ quantifies the variation in the direct effect of A on Y as M transitions from $M(a^*)$ to m^* while INT_{med} illustrates the change in the direct effects of A on Y as M shifts from $M(a)$ to $M(a^*)$. PIE captures the impact of A on Y through M , where M undergoes a change from $M(a^*)$ to $M(a)$. In lieu of concentrating solely on the overall influence of the mediator on the pathway from exposure to outcome, a four-way decomposition proves essential in elucidating the causal relationships inherent in the individual components through

which A influences Y . Note that the interaction effects can be expressed as:

$$INT_{ref} = \sum_m \left[Y(a, m) - Y(a^*, m) - Y(a, m^*) + Y(a^*, m^*) \right] I(M(a^*) = m)$$

$$INT_{med} = \sum_m \left[Y(a, m) - Y(a^*, m) \right] \left[I(M(a) = m) - I(M(a^*) = m) \right].$$

The additive interaction term, $\left[Y(a, m) - Y(a^*, m) - Y(a, m^*) + Y(a^*, m^*) \right] = \left[(Y(a, m) - Y(a^*, m^*)) - \{ (Y(a, m^*) - Y(a^*, m^*)) + (Y(a^*, m) - Y(a^*, m^*)) \} \right]$, indicates that the difference between $Y(a, m)$ and $Y(a^*, m^*)$ is different from the sum of the controlled direct effect at m^* and the indirect effect with M changing from m^* to m . In other words, INT_{ref} exists if M is set to $M(a^*)$ and INT_{med} exists if M varies for different levels of exposure.

In case of multiple parallel mediators \mathbf{M} , the total effect of A on Y can also be decomposed into four components. The interaction effects capture the interaction between the exposure and all the mediators and the pure indirect effect captures the effect solely due to the mediators. For simpler notation, we illustrate the decomposition by considering two mediators $\mathbf{M} = (M_1, M_2)$ but note that this can be generalized for more than 2 mediators. The total effect of exposure on the outcome, including the interaction effect with two mediators, can be written as:

$$TE = CDE(m_1^*, m_2^*) + INT_{refAM_1M_2}(m_1^*, m_2^*) + INT_{medAM_1M_2} + PIE_{M_1M_2}. \quad (2.4)$$

The total effect and all the components are defined as follows:

$$TE = Y(a) - Y(a^*) = Y(a, M_1(a), M_2(a)) - Y(a^*, M_1(a^*), M_2(a^*)),$$

$$CDE(m_1^*, m_2^*) = E[Y(a, m_1^*, m_2^*) - Y(a^*, m_1^*, m_2^*)],$$

$$INT_{refAM_1M_2}(m_1^*, m_2^*) = E[Y(a, M_1(a^*), M_2(a^*)) - Y(a^*, M_1(a^*), M_2(a^*)) - (Y(a, m_1^*, m_2^*) - Y(a^*, m_1^*, m_2^*))],$$

$$INT_{medAM_1M_2} = E[Y(a, M_1(a), M_2(a)) - Y(a^*, M_1(a), M_2(a)) \\ - (Y(a, M_1(a^*), M_2(a^*)) - Y(a^*, M_1(a^*), M_2(a^*)))] ,$$

and

$$PIE_{M_1M_2} = E[Y(a^*, M_1(a), M_2(a)) - Y(a^*, M_1(a^*), M_2(a^*))] .$$

Here m_1^* and m_2^* are the fixed values for M_1 and M_2 respectively.

2.2.3 Four-way decomposition under linear models

For a continuous outcome Y , binary exposure A , and a single continuous mediator M , we assume that the relationship between (\mathbf{X}, A, M, Y) follows the linear regression models and can be expressed as:

$$E[M|A = a, \mathbf{X} = \mathbf{x}] = \beta_0 + \beta_1 a + \beta_2' \mathbf{x}, \quad (2.5)$$

$$E[Y|A = a, M = m, \mathbf{X} = \mathbf{x}] = \theta_0 + \theta_1 a + \theta_2 m + \theta_3 am + \theta_4' \mathbf{x}. \quad (2.6)$$

Under the four causal assumptions, the decomposed effects of A on Y can be estimated by deriving the expected values of the four components defined in (2.3):

$$E[CDE|\mathbf{x}](m^*) = (\theta_1 + \theta_3 m^*)(a - a^*), \quad (2.7)$$

$$E[INT_{ref}|\mathbf{x}](m^*) = \theta_3(\beta_0 + \beta_1 a^* + \beta_2' \mathbf{x} - m^*)(a - a^*),$$

$$E[INT_{med}|\mathbf{x}] = \theta_3 \beta_1 (a - a^*)^2,$$

$$E[PIE|\mathbf{x}] = (\theta_2 \beta_1 + \theta_3 \beta_1 a^*)(a - a^*).$$

From the four equations in (2.7), we can obtain the proportion attributed to each effect to evaluate the contribution of different pathways for the causal relation of A , M and Y . If we set $m^* = 0$, then $E[CDE|\mathbf{x}](m^*)$ and $E[INT_{ref}|\mathbf{x}](m^*)$ are obtained as $\theta_1(a - a^*)$ and $\theta_3(\beta_0 + \beta_1 a^* + \beta_2' \mathbf{x})(a - a^*)$ respectively. We can obtain $E[TE|\mathbf{x}]$ by taking sum of all the 4 effects. The proportions attributed to CDE , INT_{ref} , INT_{med} and PIE can be obtained from the four ratios $\frac{E[CDE]}{E[TE]}$, $\frac{E[INT_{ref}]}{E[TE]}$,

$\frac{E[INT_{med}]}{E[TE]}$ and $\frac{E[PIE]}{E[TE]}$, respectively. The standard errors (SE) and confidence intervals (CI) for the four components and their attributable proportions can be obtained using the delta method, as demonstrated in VanderWeele and Vansteelandt (2009), or by employing the bootstrap method. This involves drawing numerous bootstrap samples from the observed data.

In case of multiple continuous mediators $M = (M_1, M_2, \dots, M_Q)$, which are not causally ordered instead are parallel, as shown in Figure 2.2 (b), we assume the relationship between \mathbf{X} , A , M , and Y are given by the following equations:

$$E[M_q|A = a, \mathbf{X} = \mathbf{x}] = \beta_{0q} + \beta_{1q}a + \beta'_{2q}\mathbf{x} \quad \text{for } q = 1, \dots, Q \quad (2.8)$$

$$E[Y|A = a, M = m, \mathbf{X} = \mathbf{x}] = \theta_0 + \theta_1a + \theta_{21}m_1 + \dots + \theta_{2Q}m_Q + \theta_{31}am_1 + \dots + \theta_{3Q}am_Q + \theta_4'\mathbf{x}. \quad (2.9)$$

Here $m = (m_1, \dots, m_Q)$. Following the definitions of the decomposed effects in (2.4), the expected values can be obtained as:

$$\begin{aligned} E[CDE|\mathbf{x}](m_1^*, \dots, m_Q^*) &= (\theta_1 + \theta_{31}m_1^* + \dots + \theta_{3Q}m_Q^*)(a - a^*) \quad (2.10) \\ E[INT_{refAM_1 \dots M_Q}|\mathbf{x}](m_1^*, \dots, m_Q^*) &= \sum_{q=1}^Q \theta_{3q}(\beta_{0q} + \beta_{1q}a^* + \beta'_{2q}\mathbf{x} - m_q^*)(a - a^*) \\ E[INT_{medAM_1, \dots, M_Q}|\mathbf{x}] &= \sum_{q=1}^Q \theta_{3q}\beta_{1q}(a - a^*)^2 \\ E[PIE_{M_1, \dots, M_Q}|\mathbf{x}] &= \sum_{q=1}^Q (\theta_{2q} + \theta_{3q}a^*)\beta_{1q}(a - a^*). \end{aligned}$$

If the j^{th} mediator is binary, while keeping the set of regression models in (2.9) the same for the outcome and all other mediators, the model for the j^{th} mediator can be formulated according to VanderWeele (2015) as:

$$\text{logit}[P(M_j = 1|A = a, \mathbf{X})] = \beta_{0j} + \beta_{1j}a + \beta'_{2j}\mathbf{x}.$$

Thus for the expected values of the decomposed effects in (2.10), the CDE would re-

main the same, and all other mediation effects would be altered by modifying the term involving the j^{th} mediator in the equations: (1) the j^{th} term in $INT_{refAM_1, \dots, M_Q}$ is replaced by $\theta_{3j} \left[\frac{\exp(\beta_{0j} + \beta_{1j}a^* + \beta'_{2j}\mathbf{x})}{1 + \exp(\beta_{0j} + \beta_{1j}a^* + \beta'_{2j}\mathbf{x})} - m_j^* \right] (a - a^*)$; (2) the j^{th} term in $INT_{medAM_1, \dots, M_Q}$ is replaced by $\theta_{3j}\phi(a - a^*)$, and (3) the j^{th} term in $PIE_{M_1 \dots M_Q}$ is replaced by $(\theta_{2j} + \theta_{3j}a^*)\phi$, where $\phi = \frac{\exp(\beta_{0j} + \beta_{1j}a + \beta'_{2j}\mathbf{x})}{1 + \exp(\beta_{0j} + \beta_{1j}a + \beta'_{2j}\mathbf{x})} - \frac{\exp(\beta_{0j} + \beta_{1j}a^* + \beta'_{2j}\mathbf{x})}{1 + \exp(\beta_{0j} + \beta_{1j}a^* + \beta'_{2j}\mathbf{x})}$. Similarly, the total effect of A on Y can be obtained by summing all the decomposed effects, and the proportion attributed to each component can be further calculated by taking the ratio of the estimated effect to the estimated total effect.

2.3 Extension to survival outcomes and health racial disparity

2.3.1 Extension to survival outcomes

Up to this point, we have focused on continuous outcome models for variable Y . However, in healthcare research, time-to-event data, often used in survival analysis, plays a crucial role. In this section, we extend the mediation analysis to survival outcomes, where the outcome is a time-to-event variable. Time-to-event outcomes are frequently subject to right censoring, occurring when the event of interest does not happen within the designated study period or when subjects are prematurely dropped from the study.

Existing literature on causal mediation analysis with time-to-event outcomes, including various survival models, has been established (VanderWeele, 2011; Lapointe-Shaw et al., 2018). In survival analysis, Cox proportional hazard (Cox PH) models are commonly used as semi-parametric models. However, applying the Cox PH model in mediation analysis can be challenging, given its assumption that the hazard ratio remains constant over time (VanderWeele, 2011). The interpretation of hazard ratios for indirect effects in mediation contexts can be complex and may not lend itself to meaningful causal insights. On the other hand, utilizing the accelerated failure

time (AFT) model for the outcome in mediation analysis offers advantages. AFT models assume a linear relationship between the logarithm of survival time and the covariates and confounding variables, making the interpretation of direct and indirect effects more straightforward. Additionally, AFT models facilitate the incorporation of exposure-mediator interaction terms. These models estimate parameters by maximizing the likelihood function, accounting for both observed and censored data (Wei, 1992). We extend the four-way decomposition to time-to-event data using the AFT model.

The AFT outcome model, considering a single mediator and interaction between exposure and mediator, can be expressed as:

$$Y = \log(T) = \theta_0 + \theta_1 a + \theta_2 m + \theta_3 am + \theta_4' \mathbf{x}. \quad (2.11)$$

Taking the counterfactual time-to-event outcome in log scale, both the direct and indirect effects remain consistent with those obtained for continuous outcomes. Thus, the direct effect can be redefined as $E[\log(T(a, M(a^*))) - \log(T(a^*, M(a^*)))]$, and the indirect effect as $E[\log(T(a, M(a))) - \log(T(a, M(a^*)))]$. The four-way decomposition can be similarly conducted, and the estimated decomposed effects will be the same as those obtained for continuous outcomes in (2.7). The extension to multiple mediators follows a similar approach.

2.3.2 Extension to health racial disparity study

In public health research, the exploration of racial disparities emerges as a substantial yet intricate issue, involving the analysis of divergent outcomes or opportunities among racial groups, particularly between black and white populations. The definition and measurement of racial disparity encompass various perspectives, with a critical consideration being whether adjustments should be made for associated factors such as socioeconomic status, education, or healthcare access. In accordance with

the Institute of Medicine’s (IOM) definition as outlined by (Nelson, 2002), racial disparity is measured by the outcome differences between racial groups influenced by intervening factors (designated as \mathbf{M}_1), including socioeconomic status, health coverage, and other social determinants of health. Non-intervening factors (referred to as \mathbf{M}_2), such as sex, age, and other specific health details of patients, are not considered in this definition. Adhering to the IOM’s definition of disparity, our focus is specifically on addressing the mediating effect of intervening variables (\mathbf{M}_1) with room for improvement. Our approach involves the use of joint regression models for both the outcome and mediators, yielding a comprehensive four-fold decomposition of effects. These effects include the controlled direct effect, interaction at reference effect, mediated interaction effect, and pure indirect effect.

In an alternative methodology, Clemans-Cope et al. (2023) employs the Kitagawa-Blinder-Oaxaca (KBO) decomposition approach (Blinder, 1973) to estimate the IOM disparity between black and white racial groups. This approach involves the utilization of two group-specific models: one for the black population and another for the white population. To assess the comparison between the black and white racial groups while considering the influence of two mediators, the two group-specific models are expressed as follows:

$$Y_B = \gamma_{0B} + \gamma_{2B}M_1 + \gamma_{3B}M_2 + \epsilon, \quad (2.12)$$

$$Y_W = \gamma_{0W} + \gamma_{2W}M_1 + \gamma_{3W}M_2 + \epsilon.$$

The group-specific models for the outcome aid in detecting disparities when measuring the difference between two groups, especially when one group is at a disadvantage compared to the other (Blinder, 1973), such as comparing black and white groups. The IOM is defined using the following terms:

- (1) Total difference = $\bar{Y}_B - \bar{Y}_W$;
- (2) Difference due to coefficients = $(\gamma_{2B} - \gamma_{2W}) * \bar{M}_{1W} + (\gamma_{3B} - \gamma_{3W}) * \bar{M}_{2W}$;

(3) Difference due to the interaction between covariates and coefficients

$$= (\gamma_{2B} - \gamma_{2W}) * (\bar{M}_{1B} - \bar{M}_{1W}) + (\gamma_{3B} - \gamma_{3W}) * (\bar{M}_{2B} - \bar{M}_{2W});$$

(4) Difference due to $M_1 = \gamma_{2W}(\bar{M}_{1B} - \bar{M}_{1W})$, and

$$\text{difference due to } M_2 = \gamma_{3W}(\bar{M}_{2B} - \bar{M}_{2W}).$$

The IOM disparity is defined as the total difference in (1) minus the difference due to M_2 in (4). The IOM disparity is equal to the summation of the difference due to coefficients in (2), the difference due to the interaction between covariates and coefficients in (3), and the difference due to M_1 . In the following, we demonstrate that the various terms in the IOM disparity definition can be connected to the mediation models and four-way decomposition.

Let's consider race as the binary exposure variable, with $a = 1$ for the black population and 0 for the white population. Let's assume we have two parallel mediators, M_1 and M_2 , without considering the confounders \mathbf{X} . In the mediation analyses, we consider the following outcome model:

$$E[Y|A = a, M_1 = m_1, M_2 = m_2] = \theta_0 + \theta_1 a + \theta_{21} m_1 + \theta_2 m_2 + \theta_{31} a m_1 + \cdots + \theta_{32} a m_2,$$

and the mediation models:

$$E[M_q|A = a] = \beta_{0q} + \beta_{1q} a \quad \text{for } q = 1, 2.$$

The total difference in the IOM disparity definition is equivalent to the total effect in the mediation model. Based on the models for mediators, we have

$$\beta_{01} = \bar{M}_{1W} = E[M_1|A = 0], \beta_{11} = \bar{M}_{1B} - \bar{M}_{1W} = E[M_1|A = 1] - E[M_1|A = 0] \text{ for } M_1,$$

$$\beta_{02} = \bar{M}_{2W} = E[M_2|A = 0], \beta_{12} = \bar{M}_{2B} - \bar{M}_{2W} = E[M_2|A = 1] - E[M_2|A = 0] \text{ for } M_2.$$

Comparing with the outcome model in the mediation analysis with the two group specific models, we have

$$\theta_0 = \gamma_{0W}, \theta_0 + \theta_1 = \gamma_{0B}, \theta_{21} = \gamma_{2W}, \theta_{21} + \theta_{31} = \gamma_{2B}, \theta_{22} = \gamma_{3W}, \theta_{22} + \theta_{32} = \gamma_{3B}.$$

Thus, the four decomposition effects defined in (2.10) can be expressed in terms of coefficients from the group-specific models:

$$CDE(m_1^*, m_2^*) = \theta_0 + \theta_{31}m_1^* + \theta_{32}m_2^* = (\gamma_{0B} - \gamma_{0W}) + (\gamma_{2B} - \gamma_{2W})m_1^* + (\gamma_{3B} - \gamma_{3W})m_2^*,$$

$$\begin{aligned} Int_{refAM_1M_2}(m_1^*, m_2^*) &= \theta_{31}(\beta_{01} - m_1^*) + \theta_{32}(\beta_{02} - m_2^*) \\ &= (\gamma_{2B} - \gamma_{2W}) * (\bar{M}_{1W} - m_1^*) + (\gamma_{3B} - \gamma_{3W}) * (\bar{M}_{2W} - m_2^*), \end{aligned}$$

$$Int_{medAM_1M_2} = \theta_{31}\beta_{11} + \theta_{32}\beta_{12} = (\gamma_{2B} - \gamma_{2W}) * (\bar{M}_{1B} - \bar{M}_{1W}) + (\gamma_{3B} - \gamma_{3W}) * (\bar{M}_{2B} - \bar{M}_{2W}),$$

$$PIE_{M_1M_2} = \theta_{21}\beta_{11} + \theta_{22}\beta_{12} = \gamma_{2W}(\bar{M}_{1B} - \bar{M}_{1W}) + \gamma_{3W}(\bar{M}_{2B} - \bar{M}_{2W}).$$

Henceforth, we can clearly see $Int_{medAM_1M_2}$ is equivalent to the difference due to interaction between mediators and coefficients while $PIE_{M_1M_2}$ can be expressed as the sum of the mediation effects due to M_1 and M_2 . As a special case, by fixing m_1^* and m_2^* to 0, we can see that $Int_{refAM_1M_2}$ is equivalent to the difference due to coefficients from the IOM definition. If we take values of m_1^* and m_2^* as the mean of control group of white, i.e., $m_1^* = \bar{M}_{1W}$ and $m_2^* = \bar{M}_{2W}$, then $Int_{refAM_1M_2}$ becomes 0 and CDE is obtained as the sum of $\gamma_{0B} - \gamma_{0W}$ and the difference due to coefficients.

2.4 Simulation studies

To assess the estimation of decomposition effects in mediation analysis we conducted a simulation study using 1000 Monte Carlo simulations. We compare the results across various configurations of mediators and outcomes. In the first setting we present the results obtained by considering a single continuous mediator in the joint mediation model. Subsequently in second setting, we extended the mediation model to incorporate multiple mediators, both continuous and binary. Additionally, we demonstrated the decomposition effects within the context of survival outcomes.

For each setting, we considered two distinct scenarios for the variation in the

exposure levels. In the first scenario, we evenly distribute 50% to each exposure level and in the second scenario we allocate 35% to the exposed group setting as later for our case study we have classified the Black race as the exposure group which is in minority. We have further compared the direct effects and the interaction at reference level for two settings of m^* , one by fixing it to 0 and another by setting it as $E[M[A = 0]]$. For comparison, we have obtained the standard error of the estimates and confidence interval both by Delta method and bootstrap method.

For each setting of mediators and outcome, we generate a set of independent covariates $\mathbf{X} \in \mathbb{R}^p, p = 3$, say with $\mathbf{X} \sim N(0, I)$ which impacts the exposure A , mediators \mathbf{M} and Y . The binary exposure A was generated from the logistic model with $p_A = \frac{\exp(\mathbf{X}^\top \delta)}{1 + \exp(\mathbf{X}^\top \delta)}$ where $\delta = (\delta_0, 1, 1, 1)$. δ_0 was varied as 0 and -0.85 to adhere to the two exposure group settings of 50% and 35% respectively. Then we generate mediators and outcome as per our different settings. The sample size for each simulation setting was set to $n = 1000$.

2.4.1 Simulation Settings

In the first setting we have considered a setting with single continuous mediator. We generate M being linearly associated with A and \mathbf{X} and Y being linearly dependant on A , M and \mathbf{X} as given in the models:

$$M = \beta_0 + \beta_1 A + \beta_2' \mathbf{X} + \epsilon_M, \quad (2.13)$$

$$Y = \theta_0 + \theta_1 A + \theta_2 M + \theta_3 AM + \theta_4' \mathbf{X} + \epsilon_Y, \quad (2.14)$$

where $\beta_0 = 0, \beta_1 = 1, \beta_2 = (0.2, 0.4, 0.6), \theta_0 = \theta_1 = 1, \theta_2 = 0.5, \theta_3 = 0.4, \theta_4 = (0.5, 1, 1.5), \epsilon_Y \perp \epsilon_M \sim N(0, 1)$.

We have obtained the estimated parameters by fitting the linear models on outcome Y and mediator M separately and from there we obtain the four estimated decomposed effects and the estimated total effect from the equations in (2.7). We calculate the

proportion attributed to the total effect by each decomposed effect taking the ratio of each effect and the total effect.

For multiple mediators scenario we have set up the second setting with three mediators. We independently generate two continuous mediators M_1 and M_2 and one binary M_3 and outcome Y from the linear models, with M_1 and M_2 being related to \mathbf{X} and A , M_3 related to A only and Y being related to A, M_1, M_2, M_3 and \mathbf{X} :

$$M_1 = \beta_{01} + \beta_{11}A + \beta'_{21}\mathbf{X} + \epsilon_{M_1} \quad (2.15)$$

$$M_2 = \beta_{02} + \beta_{12}A + \beta'_{22}\mathbf{X} + \epsilon_{M_2}, \quad (2.16)$$

$$\text{logit}(M_3) = \beta_{03} + \beta_{13}A, \quad (2.17)$$

$$Y = \theta_0 + \theta_1A + \theta_{21}M_1 + \theta_{22}M_2 + \theta_{23}M_3 \quad (2.18)$$

$$+ \theta_{31}AM_1 + \theta_{32}AM_2 + \theta_{33}AM_3 + \theta'_4\mathbf{X} + \epsilon_Y,$$

where $\beta_{01} = 0.5, \beta_{11} = 0.2, \beta_{21} = (0.1, 0.2, 0.3), \beta_{02} = 0, \beta_{12} = 0.4, \beta_{22} = (0.2, 0.4, 0.6), \beta_{03} = 0.1, \beta_{13} = 0.6$, and $\theta_0 = \theta_1 = 1, \theta_{21} = 0.5, \theta_{22} = \theta_{23} = 0.3, \theta_{31} = \theta_{32} = \theta_{33} = 0.4, \theta_4 = (0.5, 1, 1.5)$. The error terms for the models for M_1, M_2 and Y are generated from $N(0, 1)$. The generated dataset can be expressed as $D = (Y, A, M_1, M_2, M_3, \mathbf{X})$.

To illustrate performance of the mediation model for a time-to-event outcome, we generated survival time T from a lognormal distribution linked to the treatment A , \mathbf{M} and covariates \mathbf{X} . The dataset was simulated within a generalized framework that accounts for multiple mediators with the mediators being simulated under settings analogous to our multi-mediator model with a continuous outcome. We generated the survival time T from AFT model so that $Y = \log(T)$:

$$T = \exp(\theta_0 + \theta_1A + \theta_{21}M_1 + \theta_{22}M_2 + \theta_{23}M_3 + \theta_{31}AM_1 + \theta_{32}AM_2 + \theta_{33}AM_3 + \theta'_4\mathbf{X} + \epsilon_Y) \quad (2.19)$$

The parameters for the model are identical to those specified in the second setting. To address about 15% right censoring in the data, the censoring time independent of any other variables was obtained from the Gumbel distribution by $\log(C) = \mu_c + \sigma_c\epsilon_c$

where $\mu_c = 4.8, \sigma_c = 3.6$ and $\epsilon_c \sim \text{Gumbel}(0, 1)$. Then the observed outcome is generated by $\tilde{T} = \min\{T, C\}$, and the censoring indicator, $\text{Status} = 1\{T < C\}$.

2.4.2 Simulation results

In Figures 2.3, 2.4 and 2.5, we utilize boxplots to show the estimated proportion of the total effect attributed to each of the decomposition effects across 1000 simulations where each proportion was obtained by the ratio of the estimated decomposed effect and the estimated total effect. Figure 2.3 focuses on the results with single mediator in the model while figure 2.4 extends the results to cases with multiple mediators. Figure 2.5 further explores parameter settings for three mediators within the context of a time-to-event outcome. Each figure employs a 2x2 panel format to compare four distinct scenarios. Row-wise, the first row of panels demonstrates proportion attributions when the exposure group constitutes a minority, accounting for 35% of the total sample. The second row illustrates cases where both the exposed and control groups are equally distributed. The first column showcases how the effects vary when m^* is fixed at 0. In contrast, the second column depicts results when m^* depends on exposure at control group. In Figure 2.4 and Figure 2.5, the estimated proportion attributed by the interaction and the mediated effects due to three mediators are shown together.

From the three figures, it is evident if the value of m^* is held constant at 0, the contribution of the decomposed effects to the total effect remains consistent across different distributions of the exposure group. However, when m^* depends on the mean value of M taken at control group of exposure, variations in the proportion attributable to CDE and INT_{ref} are observed. However, there is no significant change in the mediated effects of INT_{med} and PIE . With single mediator in the model, there is more proportion attributed by CDE to TE when m^* is fixed at 0 compared to when $mstar$ is varied whereas with multiple mediators in the model it is the oppo-

site. The thick horizontal line in each boxplot represents the proportion attribution from the true effects obtained from the true potential outcomes implicating minimal deviation from the estimated results. This suggests robustness and reliability in the estimation process, reinforcing the validity of the mediation analysis under varying conditions and parameter settings.

Through Figures 2.6, 2.7 and 2.8 we compare the performance of SE estimation obtained through bootstrap sampling procedure and the Delta method across different settings of mediators and outcomes in the model. The red line represents the ratio of the mean SE by bootstrap method and Empirical SD for decomposed effects, while the blue line depicts the corresponding ratio for SE by Delta method. We clearly see for the decomposed effects in different settings, the ratio of SE by the bootstrap method to the empirical SD tends to converge to unity compared to the corresponding ratio obtained by the Delta method. This observation suggests that despite the computational intensity associated with the bootstrap procedure, it yields greater accuracy in variance estimation compared to the Delta method.

2.5 Case study: racial disparity on all-cause mortality in the United States based on NHANES III dataset

Although there have been a significant decline in US mortality rates over the decade due to the advancement of medical science and technologies, persistence of racial disparities in Black-white mortality rates remains a pressing concern (Benamins et al., 2021). Recent research indicates that black population experience higher age adjusted death rate in comparison to the White population in the United States of America (USA) (Haines, 2003). This kind of disparity in race is often attributed to the decades of various aspects of socioeconomic differences. By analyzing the mediated pathways and estimating the effects of these socioeconomic factors, our goal is to identify the

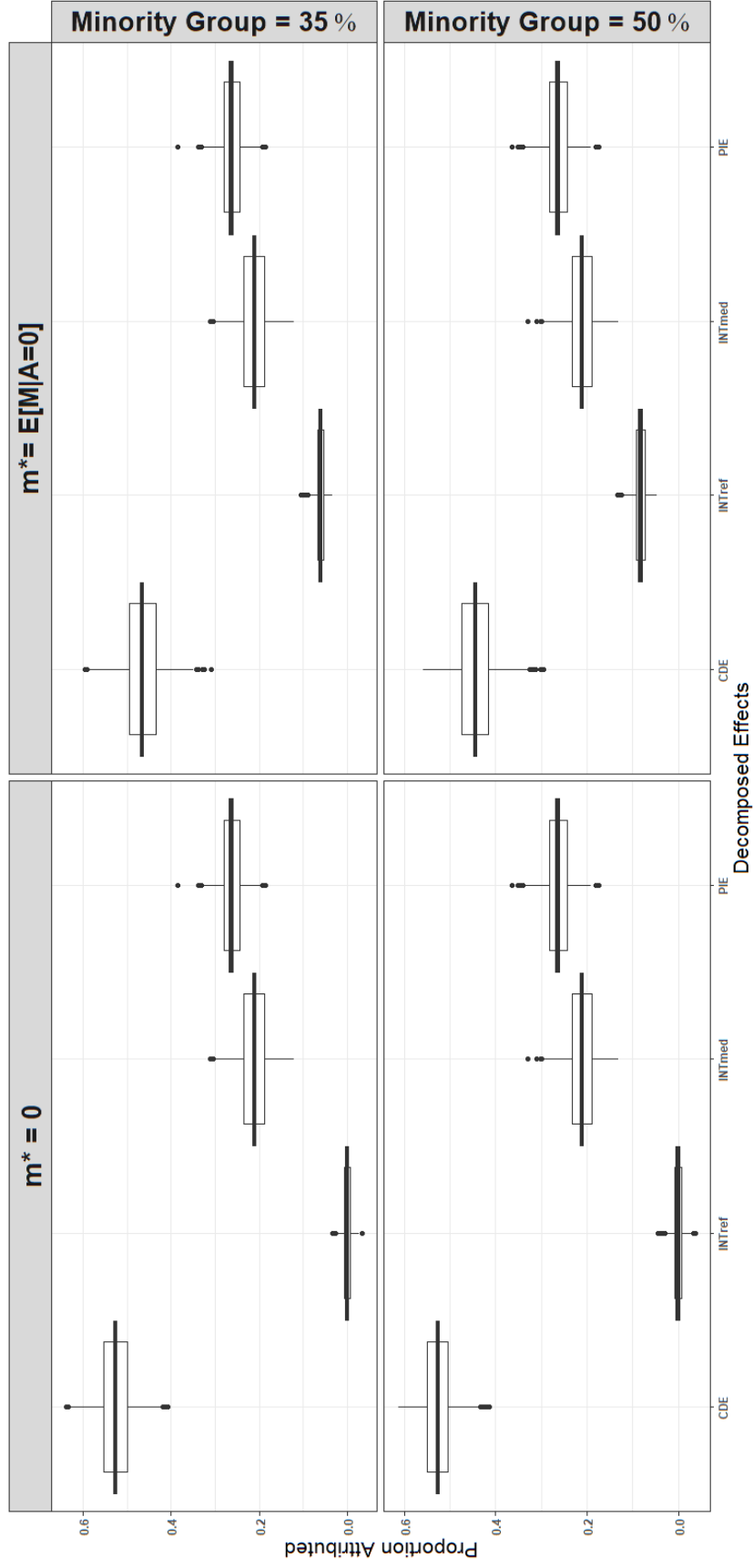


Figure 2.3: Boxplots of estimated proportion attribution by each decomposed effect at two levels of minority group distribution (35% and 50%, indicated in the first row and second row, respectively) and for two values for m^* , for single mediator and continuous outcome in the model

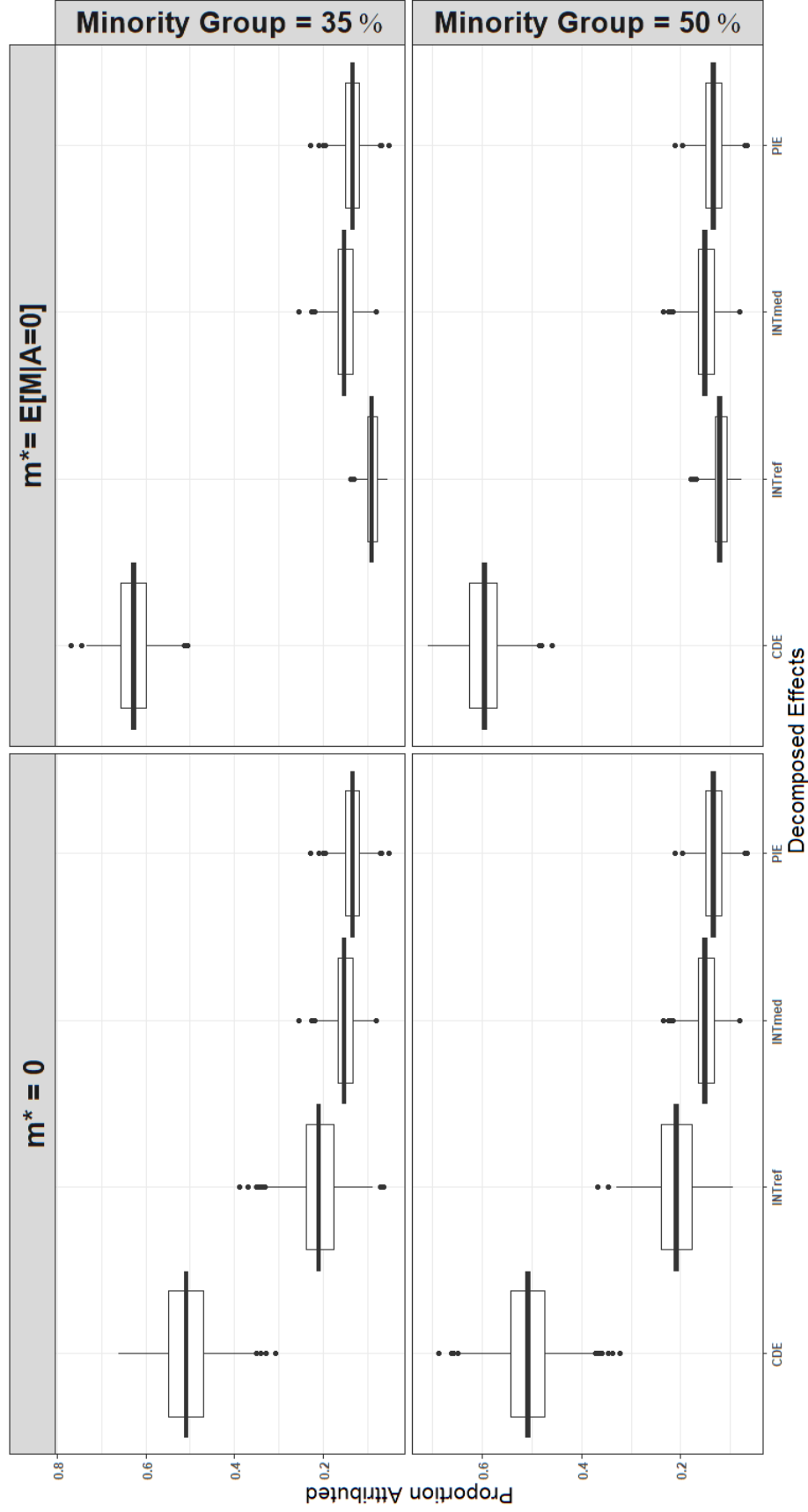


Figure 2.4: Boxplots of estimated proportion attribution by each decomposed effect at two levels of minority group distribution (35% and 50%, indicated in the first row and second row, respectively) and for two values for m^* , for multiple mediators and continuous outcome in the model

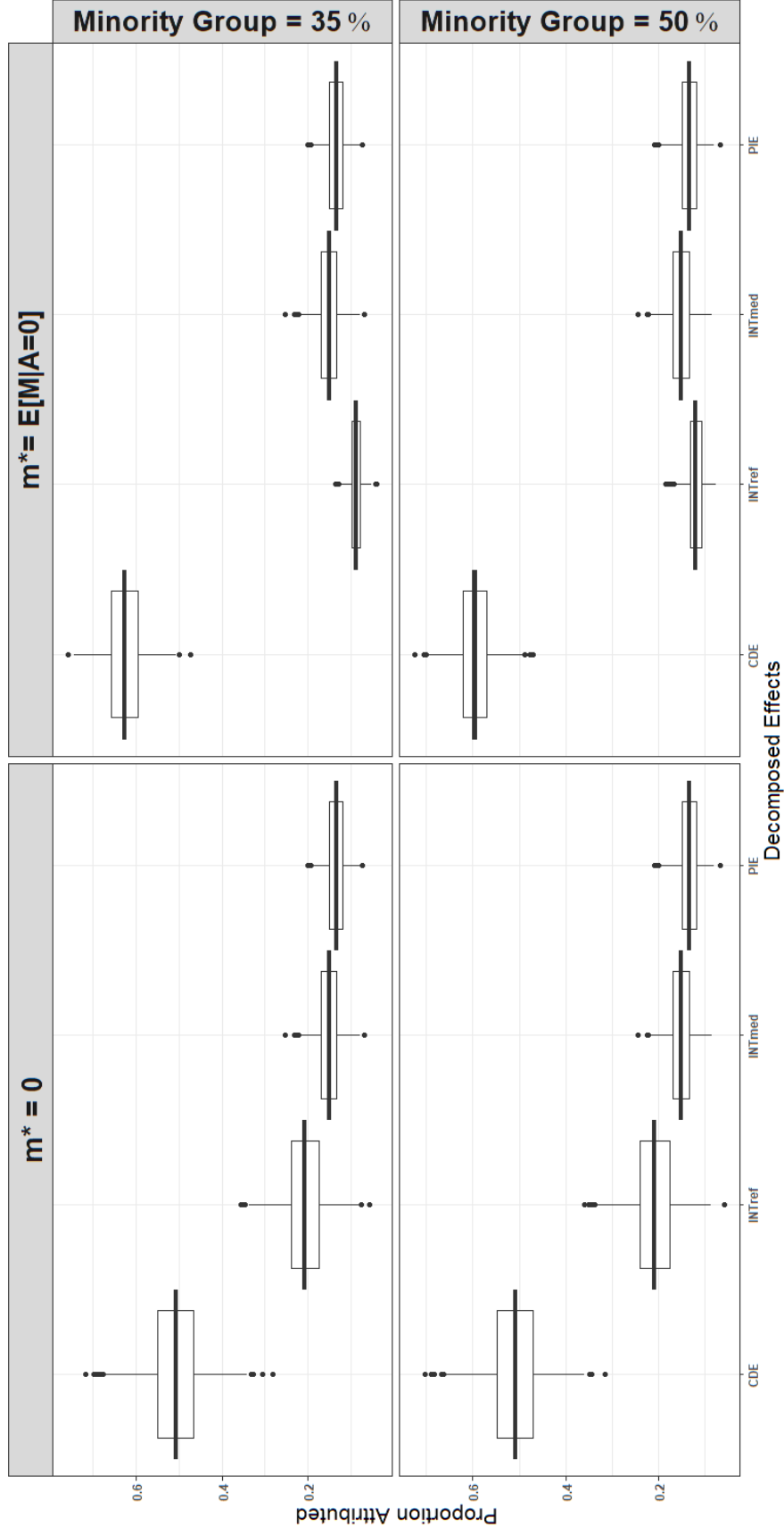


Figure 2.5: Boxplots of estimated proportion attribution by each decomposed effect at two levels of minority group distribution (35% and 50%, indicated in the first row and second row, respectively) and for two values for m^* , for multiple mediators and time-to-event outcome in the model

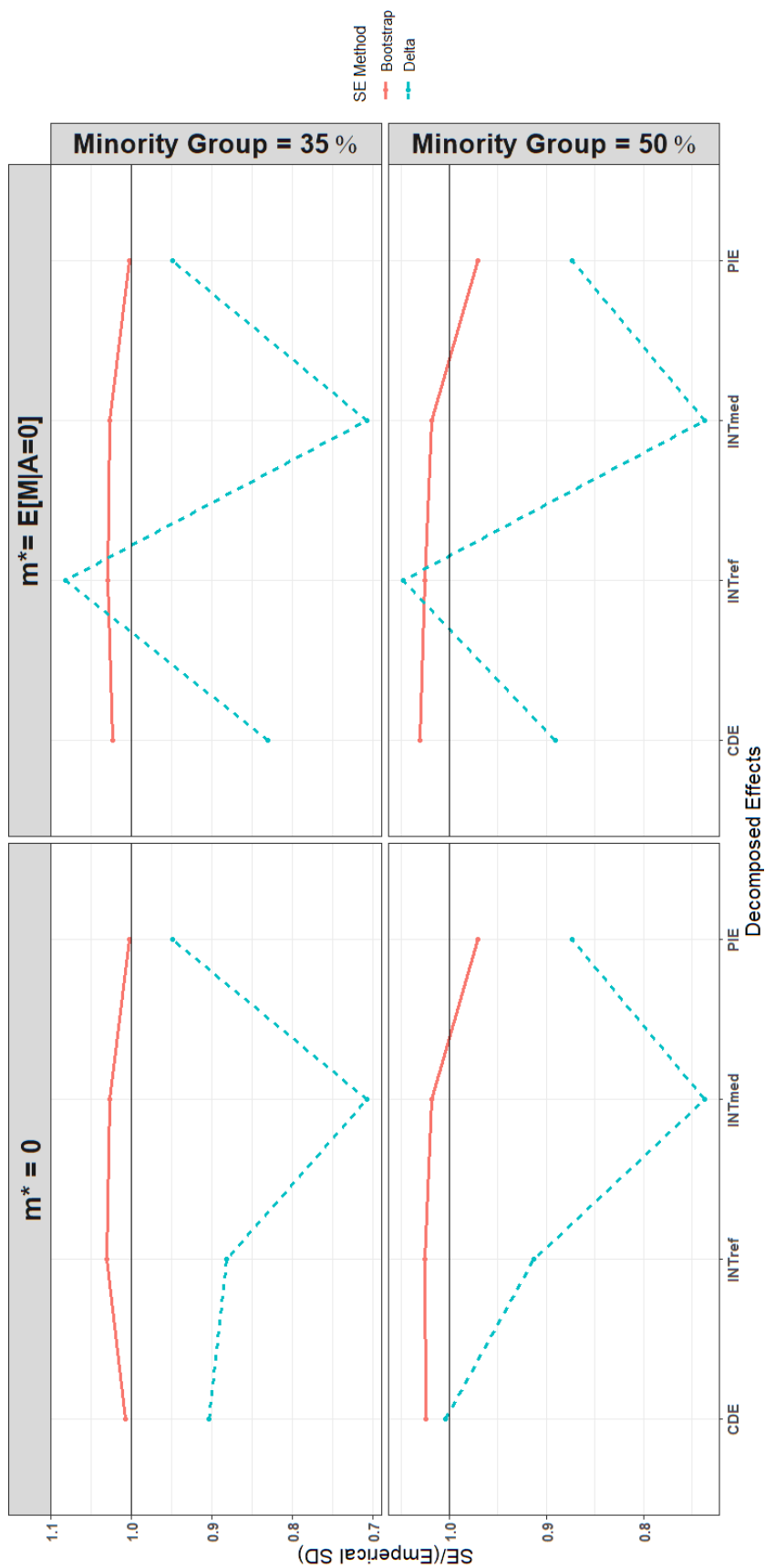


Figure 2.6: Plots to illustrate performance of standard error estimation: by taking ratio of SE by bootstrap method (red) or SE by Delta method (blue) and Empirical SD for decomposed effects at two levels of minority group distribution (35% and 50%, indicated in the first row and second row, respectively) and for two values for m^* , for single mediator and continuous outcome in the model

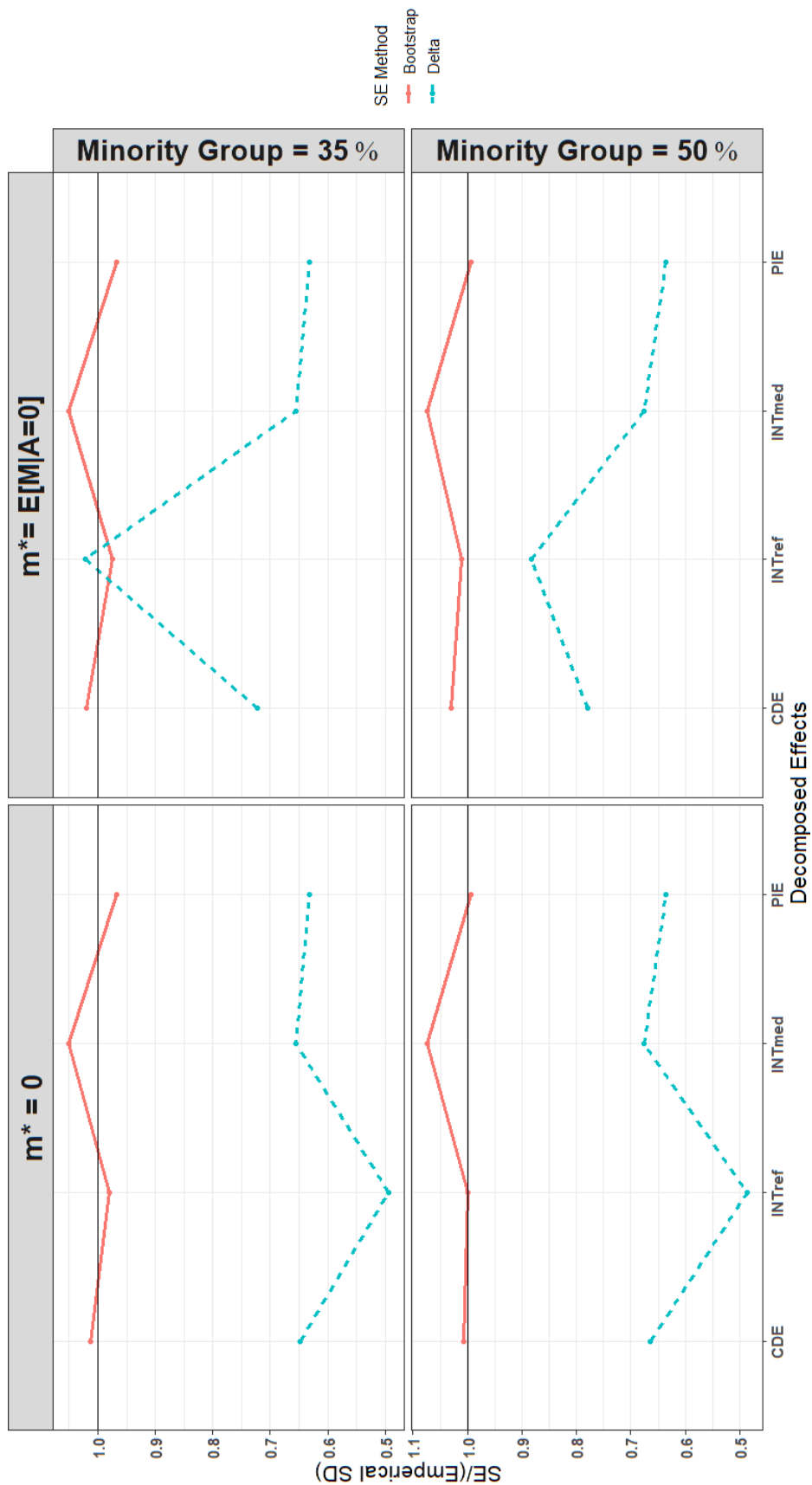


Figure 2.7: Plots to illustrate performance of standard error estimation: by taking ratio of SE by bootstrap method (red) or SE by Delta method (blue) and Empirical SD for decomposed effects at two levels of minority group distribution (35% and 50%, indicated in the first row and second row, respectively) and for two values for m^* , for multiple mediators and continuous outcome in the model

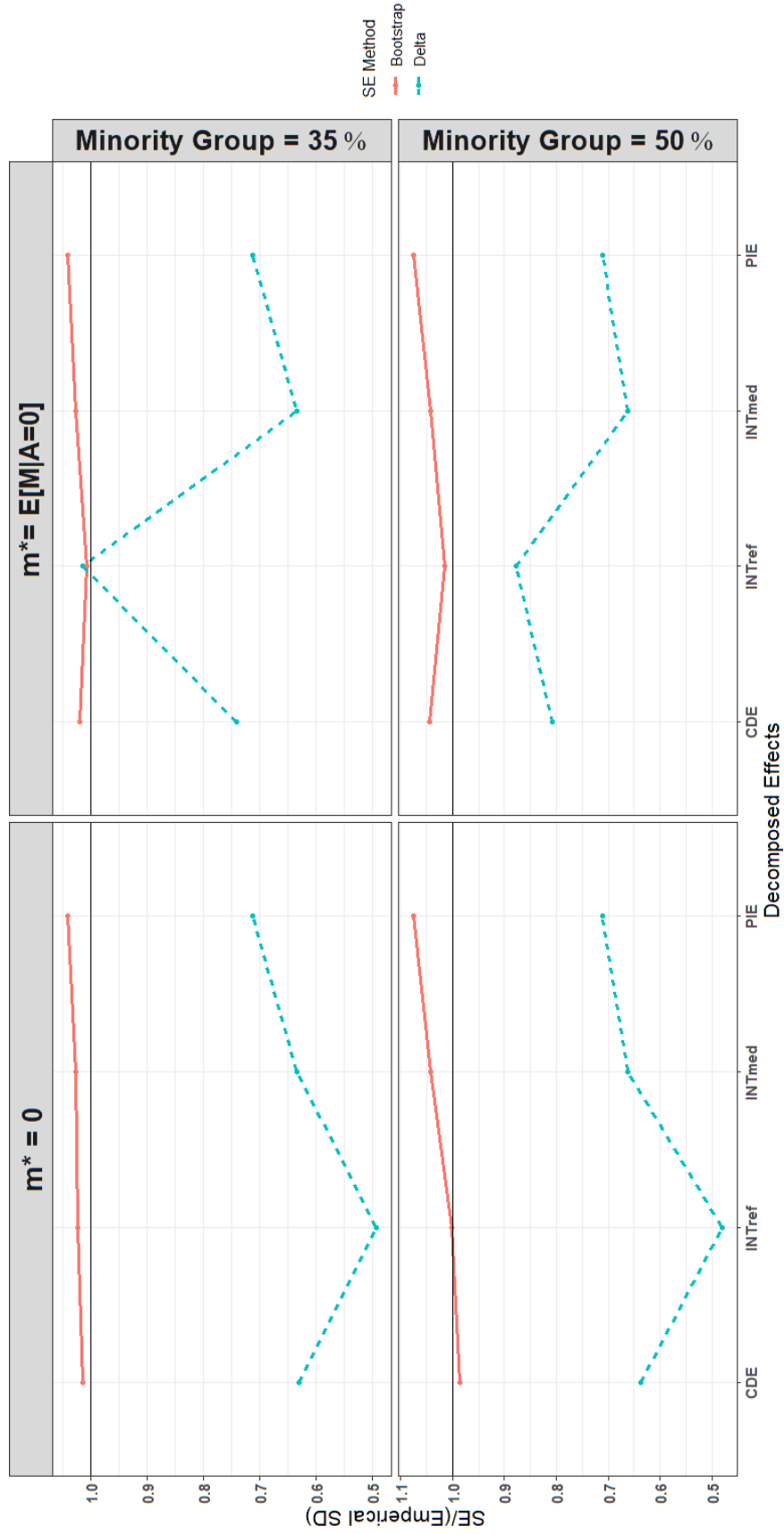


Figure 2.8: Plots to illustrate the performance of standard error (SE) estimation: each line represents the ratio of SE estimation over the empirical SE from different methods (red: bootstrap method; blue for delta method). The two rows correspond two levels of minority groups (35% and 50%), and the two columns correspond to two different values for m^* .

specific socioeconomic factors contributing to the structural racism that influences these elevated mortality rates. For our study, we have considered the third National Health and Nutrition Examination Survey (NHANES III), conducted from 1988-1994 in two phases. The 6 year survey represents the nationwide population with an overall black and white population in the USA . We are specifically interested in comparing the two racial groups (black versus white), thus we have restricted our data to the participants from black and white population only. For mortality information, we have used Linked Mortality Files (LMF) by the National Center for Health Statistics (NCHS) which have been followed up till December 31st, 2019.

In our analysis, we have considered participants within the age range of 50 – 70 years at the time of interview for the survey and excluded patients who had accidental death records. As we are primarily focused on explaining the racial disparity between black and white population, we included participants from these two races only in our study. The study cohort for our analysis include 2640 participants with 1768 white and 872 black accounting to 33% of the minority group. At the end of the study 1840 participants died and 800 were still alive leading to 30% censoring. Within the black population the mortality rate was 73.8% which is higher than the white population with 67.6% morality rate.

We have identified a set of intervening mediators as health insurance coverage, poverty income ratio (PIR) and education. Health insurance coverage was determined through participant reports of having any insurance during the last month of the interview, including Medicare/Medicaid or private. We note that approximately 5% of the white population did not have insurance coverage, whereas this percentage was about 10% in the black population. PIR was determined by evaluating the family income in relation to the poverty threshold, adjusted for both family size and the annual inflation status. Average PIR among the white population was 3.58 while that among black population was 2.1. Education was categorized into two levels as

“till high school”, and “College or high”. Education and PIR reflect the socioeconomic status of the participant. The data revealed that the proportion of participants without college education was considerably higher among the black population compared to the white population. Age, gender, metro area residence, physical exercise, smoking, and alcohol consumption, along with comorbidities such as cancer, stroke, cardiovascular disease, and diabetes are considered as the set of non-intervenable mediators. A participant’s exercise level were quantified by aggregating the frequency of any physical activity undertaken and expressing it as “times per month”. Smoking history among participants was marked as “Yes” based on if they had smoked more than 100+ cigarettes in their lifetime or smoking at the time of the survey and “No” if they had not smoked more than 100+ cigarettes in their life time. Participants’ were identified with history of alcohol consumption if they had consumed at least 12 alcoholic drinks in their lifetime or in the last 12 months. A participant was identified as having diabetes based on self-report during the interview or having glycated hemoglobin $\geq 6.5\%$ or plasma glucose $\geq 125\text{ mg/dL}$. Participants were identified with a history of heart disease if they reported a past heart attack or had congestive heart failure. Similarly history of cancer and stroke were recorded based on the self-reports of the participants. The continuous mediators PIR, BMI, age and exercise were standardized. We conducted the joint regression model including all the mediators with AFT outcome model to obtain the four decomposed effects and their proportion attributions to the total effect due to race on the survival time in months in log scale. We assumed the survival time follows a Weibull distribution. To assess the effect of the mediators collectively, we multiplied negative one for some variables so that all variables are positively related to survival time. For example, as age and BMI are negatively associated with the time-to-death, we multiplied negative one for these variables. For gender, male was considered at the reference level as they have shorter time-to-death compared to females. As participants with history of any comorbidity,

smoking or alcohol consumption had shorter survival time, positive occurrence were taken as the reference. The fixed values of m^* for each of the continuous mediators were taken at the estimated mean among the white population and for each of the binary mediators at the estimated proportion within the white population.

In Table 2.1, we present the estimated effects, its standard error, and its 95% CI obtained by bootstrap method. When we consider both the interaction effects, the combined influence of mediators and race does not exhibit a statistically significant effect on mortality. But the PIEs for both sets of mediators are significant, we conclude that the mediators alone have a significant negative impact on the time to death. The IOM score implies that there is racial disparity on all-cause mortality in the USA.

2.6 Conclusion and discussion

In this project we have illustrated the mediation analysis with parallel multiple mediators within causal framework considering the interaction between the mediation and exposure. The four-way decomposition of the total effect of exposure on outcome helps us to identify the complex pathways of mediation as it accounts for the individual effects due to neither mediation nor interaction, to just interaction between exposure and mediator, to both mediation and interaction, and to just mediation. We were able to illustrate that the four way decomposition method can be extended to survival outcome. From the findings of our case study we are able to establish the importance of considering direct, indirect as well as the interaction pathways in understanding the complex relationships between race, mediators, and mortality in the population under study. As the future aspect of our work, we intend to include mediators with multiple categories and consider continuous treatment for mediation model.

Table 2.1: Decomposition of TE of Race on Mortality Time

	Estimate	Std Error	95 % CI (Normal)	95 % CI (Percentile)	95 % CI (Basic)
CDE	-0.237	0.151	(-0.531, 0.061)	(-0.543, 0.064)	(-0.538, 0.069)
INT_{refM_2}	0.195	0.14	(-0.083, 0.468)	(-0.088, 0.48)	(-0.09, 0.478)
INT_{medM_1}	-0.025	0.03	(-0.08, 0.036)	(-0.086, 0.028)	(-0.077, 0.036)
INT_{medM_2}	0.032	0.023	(-0.013, 0.075)	(-0.013, 0.076)	(-0.013, 0.076)
PIE_{M_1}	-0.074	0.013	(-0.099, -0.05)	(-0.098, -0.049)	(-0.099, -0.05)
PIE_{M_2}	-0.046	0.02	(-0.085, -0.006)	(-0.088, -0.008)	(-0.084, -0.004)
IOM disparity	-0.109	0.03	(-0.167, -0.05)	(-0.166, -0.051)	(-0.168, -0.053)

REFERENCES

- Andersen, P. K., Syriopoulou, E., and Parner, E. T. (2017). Causal inference in survival analysis using pseudo-observations. *Statistics in medicine*, 36(17):2669–2681.
- Ara, A., Usmani, J. A., et al. (2015). Lead toxicity: a review. *Interdisciplinary toxicology*, 8(2):55–64.
- Austin, P. C. (2010). Absolute risk reductions and numbers needed to treat can be obtained from adjusted survival models for time-to-event outcomes. *Journal of clinical epidemiology*, 63(1):46–55.
- Austin, P. C. (2018). Assessing the performance of the generalized propensity score for estimating the effect of quantitative or continuous exposures on binary outcomes. *Statistics in medicine*, 37(11):1874–1894.
- Baron, R. M. and Kenny, D. A. (1986). The moderator–mediator variable distinction in social psychological research: Conceptual, strategic, and statistical considerations. *Journal of personality and social psychology*, 51(6):1173.
- Benjamins, M. R., Silva, A., Saiyed, N. S., and De Maio, F. G. (2021). Comparison of all-cause mortality rates and inequities between black and white populations across the 30 most populous us cities. *JAMA network open*, 4(1):e2032086–e2032086.
- Blinder, A. S. (1973). Wage discrimination: reduced form and structural estimates. *Journal of Human resources*, pages 436–455.

- Brown, D. W., Greene, T. J., Swartz, M. D., Wilkinson, A. V., and DeSantis, S. M. (2021). Propensity score stratification methods for continuous treatments. *Statistics in medicine*, 40(5):1189–1203.
- Cain, L. E. and Cole, S. R. (2009). Inverse probability-of-censoring weights for the correction of time-varying noncompliance in the effect of randomized highly active antiretroviral therapy on incident aids or death. *Statistics in medicine*, 28(12):1725–1738.
- Chen, X., Linton, O., and Van Keilegom, I. (2003). Estimation of semiparametric models when the criterion function is not smooth. *Econometrica*, 71(5):1591–1608.
- Clemans-Cope, L., Garrett, A. B., and McMorrow, S. (2023). How should we measure and interpret racial and ethnic disparities in health care? *Available at SSRN 4376574*.
- Cole, S. R. and Hernán, M. A. (2008). Constructing inverse probability weights for marginal structural models. *American journal of epidemiology*, 168(6):656–664.
- Coresh, J., Selvin, E., Stevens, L. A., Manzi, J., Kusek, J. W., Eggers, P., Van Lente, F., and Levey, A. S. (2007). Prevalence of chronic kidney disease in the united states. *Jama*, 298(17):2038–2047.
- Dabrowska, D. M. (1989). Uniform consistency of the kernel conditional kaplan-meier estimate. *The Annals of Statistics*, pages 1157–1167.
- Data, N. (1988-1994). Centers for disease control and prevention (cdc). national center for health statistics (nchs). national health and nutrition examination survey data. hyattsville, md: U.s. department of health and human services, centers for disease control and prevention.

- Fong, C., Hazlett, C., and Imai, K. (2018). Covariate balancing propensity score for a continuous treatment: Application to the efficacy of political advertisements. *The Annals of Applied Statistics*, 12(1):156–177.
- Fong, C., Ratkovic, M., and Imai, K. (2021). *CBPS: Covariate Balancing Propensity Score*. R package version 0.22.
- Gao, X., Li, L., and Luo, L. (2022). Decomposition of the total effect for two mediators: A natural mediated interaction effect framework. *Journal of causal inference*, 10(1):18–44.
- Gonzalez-Manteiga, W. and Cadarso-Suarez, C. (1994). Asymptotic properties of a generalized kaplan-meier estimator with some applications. *Communications in Statistics-Theory and Methods*, 4(1):65–78.
- Hade, E. M. and Lu, B. (2014). Bias associated with using the estimated propensity score as a regression covariate. *Statistics in medicine*, 33(1):74–87.
- Haines, M. R. (2003). Ethnic differences in demographic behavior in the united states has there been convergence? *Historical Methods: A Journal of Quantitative and Interdisciplinary History*, 36(4):157–195.
- Hirano, K. and Imbens, G. W. (2004). The propensity score with continuous treatments. *Applied Bayesian modeling and causal inference from incomplete-data perspectives*, 226164:73–84.
- Imai, K. and Ratkovic, M. (2014). Covariate balancing propensity score. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):243–263.
- Imai, K. and Van Dyk, D. A. (2004). Causal inference with general treatment regimes: Generalizing the propensity score. *Journal of the American Statistical Association*, 99(467):854–866.

- Imbens, G. W. (2000). The role of the propensity score in estimating dose-response functions. *Biometrika*, 87(3):706–710.
- Lapointe-Shaw, L., Bouck, Z., Howell, N. A., Lange, T., Orchanian-Cheff, A., Austin, P. C., Ivers, N. M., Redelmeier, D. A., and Bell, C. M. (2018). Mediation analysis with a time-to-event outcome: a review of use and reporting in healthcare research. *BMC medical research methodology*, 18(1):1–12.
- Le Cam, L. (1990). Maximum likelihood: an introduction. *International Statistical Review/Revue Internationale de Statistique*, pages 153–171.
- Lo, S.-H. and Singh, K. (1986). The product-limit estimator and the bootstrap: some asymptotic representations. *Probability Theory and Related Fields*, 71(3):455–465.
- Lunceford, J. K. and Davidian, M. (2004). Stratification and weighting via the propensity score in estimation of causal treatment effects: a comparative study. *Statistics in medicine*, 23(19):2937–2960.
- MacKinnon, D. P., Fairchild, A. J., and Fritz, M. S. (2007). Mediation analysis. *Annu. Rev. Psychol.*, 58:593–614.
- Nelson, A. (2002). Unequal treatment: confronting racial and ethnic disparities in health care. *Journal of the national medical association*, 94(8):666.
- Pearl, J. (2001). Direct and indirect effects. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, UAI’01, page 411–420, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Pirkle, J. L., Brody, D. J., Gunter, E. W., Kramer, R. A., Paschal, D. C., Flegal, K. M., and Matte, T. D. (1994). The decline in blood lead levels in the united states: the national health and nutrition examination surveys (nhanes). *Jama*, 272(4):284–291.

- Robins, J., Hernán, M., and Brumback, B. (2000). Marginal structural models and causal inference in epidemiology. *Epidemiology*, 11(5):550–560.
- Robins, J. M. and Finkelstein, D. M. (2000). Correcting for noncompliance and dependent censoring in an aids clinical trial with inverse probability of censoring weighted (ipcw) log-rank tests. *Biometrics*, 56(3):779–788.
- Rosenbaum, P. R. (1987). The role of a second control group in an observational study. *Statistical Science*, 2(3):292–306.
- Schober, S. E., Mirel, L. B., Graubard, B. I., Brody, D. J., and Flegal, K. M. (2006). Blood lead levels and death from all causes, cardiovascular disease, and cancer: results from the nhanes iii mortality study. *Environmental health perspectives*, 114(10):1538–1541.
- Selvin, E., Manzi, J., Stevens, L. A., Van Lente, F., Lacher, D. A., Levey, A. S., and Coresh, J. (2007). Calibration of serum creatinine in the national health and nutrition examination surveys (nhanes) 1988-1994, 1999-2004. *American journal of kidney diseases*, 50(6):918–926.
- Stute, W. (1995). The statistical analysis of kaplan-meier integrals. *Lecture Notes-Monograph Series*, pages 231–254.
- VanderWeele, T. (2015). *Explanation in causal inference: methods for mediation and interaction*. Oxford University Press.
- VanderWeele, T. J. (2011). Causal mediation analysis with survival data. *Epidemiology (Cambridge, Mass.)*, 22(4):582.
- VanderWeele, T. J. (2014). A unification of mediation and interaction: a four-way decomposition. *Epidemiology (Cambridge, Mass.)*, 25(5):749.

- VanderWeele, T. J. and Vansteelandt, S. (2009). Conceptual issues concerning mediation, interventions and composition. *Statistics and its Interface*, 2(4):457–468.
- Wei, L.-J. (1992). The accelerated failure time model: a useful alternative to the cox regression model in survival analysis. *Statistics in medicine*, 11(14-15):1871–1879.
- Xie, J. and Liu, C. (2005). Adjusted kaplan–meier estimator and log-rank test with inverse probability of treatment weighting for survival data. *Statistics in medicine*, 24(20):3089–3110.

APPENDIX

This section includes proofs of theoretical properties for Chapter 1:

Let λ_{\min} and λ_{\max} denote the smallest and largest eigenvalues of $E[\mathbf{Z}\mathbf{Z}^\top]$, respectively.

We validate the marginal structure equation (1.7). Noting that $T < C$ implies $\tilde{Y} = Y$,

$$\begin{aligned} E \left[\frac{\delta w_p(A; \mathbf{X}, \boldsymbol{\xi}^*)}{G(T)} \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}^*) \right] &= E \left[E \left[\frac{1\{T < C\} w_p(A; \mathbf{X}, \boldsymbol{\xi}^*)}{G(T)} \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}^*) \middle| T, \mathbf{Z} \right] \right] \\ &= E [w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \mathbf{Z}(Y - \mathbf{Z}^\top \boldsymbol{\theta}^*)] = E [w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \mathbf{Z}Y] - E [w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \mathbf{Z}\mathbf{Z}^\top \boldsymbol{\theta}^*]. \end{aligned}$$

For the first item, let \mathcal{Y} denote the support of $Y(a)$, $a \in \mathcal{A}$,

$$\begin{aligned} E [w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \mathbf{Z}Y] &= E [w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \mathbf{Z} E[Y|A, \mathbf{X}]] \\ &= \int_{\mathcal{A} \times \mathcal{X}} \frac{f_A(a)}{f_A(a|\mathbf{x}, \boldsymbol{\xi})} \mathbf{z} E[Y(a)|\mathbf{x}] dF_{A, \mathbf{X}}(a, \mathbf{x}) \\ &= \int_{\mathcal{A}} \int_{\mathcal{X}} \frac{f_A(a)}{f_A(a|\mathbf{x}, \boldsymbol{\xi})} \mathbf{z} \left[\int_{\mathcal{Y}} y(a) \frac{f_{Y(a), \mathbf{X}}(y(a), \mathbf{x})}{f_{\mathbf{X}}(\mathbf{x})} dy(a) \right] f_A(a|\mathbf{x}, \boldsymbol{\xi}) f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} da \\ &= \int_{\mathcal{A}} \int_{\mathcal{Y}} f_A(a) y(a) \mathbf{z} \left[\int_{\mathcal{X}} f_{Y(a), \mathbf{X}}(y(a), \mathbf{x}) d\mathbf{x} \right] dy(a) da \\ &= \int_{\mathcal{A}} f_A(a) \mathbf{z} \int_{\mathcal{Y}} y(a) f_{Y(a)}(y(a)) dy(a) da = \int_{\mathcal{A}} f_A(a) \mathbf{z} E[Y(a)] da \\ &= \int_{\mathcal{A}} f_A(a) \mathbf{z} \mathbf{z}^\top \boldsymbol{\theta}^* da = E [\mathbf{Z}\mathbf{Z}^\top \boldsymbol{\theta}^*], \end{aligned}$$

where the second last equality follows from (1.6).

For the second item,

$$\begin{aligned} E [w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \mathbf{Z}\mathbf{Z}^\top \boldsymbol{\theta}^*] &= E \left[E [w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \mathbf{Z}\mathbf{Z}^\top \boldsymbol{\theta}^* | \mathbf{X}] \right] \\ &= \int_{\mathcal{X}} \left[\int_{\mathcal{A}} \frac{f_A(a)}{f_A(a|\mathbf{x}, \boldsymbol{\xi})} \mathbf{z} \mathbf{z}^\top \boldsymbol{\theta}^* f_A(a|\mathbf{x}, \boldsymbol{\xi}) da \right] f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} = \int_{\mathcal{X}} E [\mathbf{Z}\mathbf{Z}^\top \boldsymbol{\theta}^*] f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} = E [\mathbf{Z}\mathbf{Z}^\top \boldsymbol{\theta}^*]. \end{aligned}$$

Thus, (1.6) is a valid marginal structure equation.

Proof of Theorem 1.2.1: We prove the result by invoking Theorem 1 in Chen et al. (2003) (hereafter CLK). Thus, we need to check their conditions (1.1) – (1.4) and (1.5’).

We note that condition (1.3) is trivially satisfied. Because $\hat{\boldsymbol{\xi}} = \boldsymbol{\xi}^* + O_p(n^{-1/2})$ and $\|\hat{G} - G^*\|_\infty = O_p((\log \log n/n)^{1/2})$ (Stute, 1995), $\|\hat{h} - h^*\|_\infty = o_p(n^{-1/4})$. Thus, condition (1.4) in CLK holds. We thus only need to verify conditions (1.1), (1.2) and (1.5’).

To verify condition (1.1) in CLK, we want to show that $\|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| = o_p(n^{-1/2})$, that is, $\forall u > 0 \lim_{n \rightarrow \infty} P(\sqrt{n}\|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| > u) = 0$. Let Ξ denote the event that $n^{-1} \sum_{i=1}^n w_i(\hat{h}) \mathbf{Z}_i \mathbf{Z}_i^\top$ is invertible.

$$\begin{aligned} & P(\sqrt{n}\|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| > u) \\ &= P\left(\sqrt{n}\|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| > u | \Xi\right) P(\Xi) + P\left(\sqrt{n}\|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| > u | \Xi^c\right) P(\Xi^c) \\ &\leq P\left(\sqrt{n}\|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| > u | \Xi\right) P(\Xi) + P(\Xi^c). \end{aligned} \quad (2.20)$$

Given Ξ , $\hat{\boldsymbol{\theta}} = (n^{-1} \sum_{i=1}^n w_i(\hat{h}) \mathbf{Z}_i \mathbf{Z}_i^\top)^{-1} (n^{-1} \sum_{i=1}^n w_i(\hat{h}) \mathbf{Z}_i \tilde{Y}_i)$ and $M_n(\hat{\boldsymbol{\theta}}, \hat{h}) = 0$. Thus,

$$P\left(\sqrt{n}\|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| > u | \Xi\right) P(\Xi) = 0. \quad (2.21)$$

Since $\|\hat{h} - h^*\|_\infty = o_p(1)$ and \mathcal{X} is a compact set in \mathbb{R}^p , $\sup_{a \in \mathcal{A}, \mathbf{x} \in \mathcal{X}} |w_p(a; \mathbf{x}, \hat{\boldsymbol{\xi}}) - w_p(a; \mathbf{x}, \hat{\boldsymbol{\xi}})| = o_p(1)$ and $\sup_{t \in (0, L]} |\hat{G}(t) - G(t)| = o_p(1)$. As $G(t)$ is bounded away from 0 by the regularity condition (C1), $\max_{1 \leq i \leq n} |w_i(\hat{h}) - w_i(h^*)| = o_p(1)$. Therefore, $n^{-1} \sum_{i=1}^n w_i(\hat{h}) \mathbf{Z}_i \mathbf{Z}_i^\top = n^{-1} \sum_{i=1}^n w_i(h^*) \mathbf{Z}_i \mathbf{Z}_i^\top + o_p(1) = E[w_i(h^*) \mathbf{Z}_i \mathbf{Z}_i^\top] + o_p(1) = E[\mathbf{Z}_i \mathbf{Z}_i^\top] + o_p(1)$ by the strong law of large number. Then

$$\lim_{n \rightarrow \infty} P\left(\left\|n^{-1} \sum_{i=1}^n w_i(\hat{h}) \mathbf{Z}_i \mathbf{Z}_i^\top - E[\mathbf{Z}_i \mathbf{Z}_i^\top]\right\| > \lambda_{\min}\right) = 0.$$

If $\|n^{-1} \sum_{i=1}^n w_i(\hat{h}) \mathbf{Z}_i \mathbf{Z}_i^\top - E[\mathbf{Z}_i \mathbf{Z}_i^\top]\| < \lambda_{\min}$, then $\forall \mathbf{d} \in \mathbb{R}^2$ with $\|\mathbf{d}\| = 1$,

$$\left\| n^{-1} \sum_{i=1}^n w_i(\hat{h}) \mathbf{Z}_i \mathbf{Z}_i^\top \mathbf{d} \right\| \geq \|E[\mathbf{Z}_i \mathbf{Z}_i^\top] \mathbf{d}\| - \left\| \left(n^{-1} \sum_{i=1}^n w_i(\hat{h}) \mathbf{Z}_i \mathbf{Z}_i^\top - E[\mathbf{Z}_i \mathbf{Z}_i^\top] \right) \mathbf{d} \right\|$$

$$> \lambda_{\min} - \lambda_{\min} = 0,$$

which implies that Ξ . Therefore,

$$\lim_{n \rightarrow \infty} P(\Xi^c) \leq \lim_{n \rightarrow \infty} P\left(\left\| n^{-1} \sum_{i=1}^n w_i(\hat{h}) \mathbf{Z}_i \mathbf{Z}_i^\top - E[\mathbf{Z}_i \mathbf{Z}_i^\top] \right\| > \lambda_{\min}\right) = 0. \quad (2.22)$$

Combining (2.20), (2.21), and (2.22) together yields that $\forall u > 0$,

$$\lim_{n \rightarrow \infty} P(\sqrt{n} \|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| > u) = 0 + \lim_{n \rightarrow \infty} P(\Xi^c) = 0.$$

Thus, condition (1.1) is satisfied.

To verify condition (1.2) in CLK, $\forall u > 0$, it is easy to see that

$$\inf_{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| > u} \|M(\boldsymbol{\theta}, h^*) - M(\boldsymbol{\theta}^*, h^*)\| = \|E[\mathbf{Z}\mathbf{Z}^\top](\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \geq \lambda_{\min} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \geq \lambda_{\min} u.$$

Thus, condition (1.2) is satisfied.

By Condition (C1) and Remark 1.2.2, we have $G(T) \geq \tau/2$, $\forall G \in \mathcal{G}$. Consider $\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}\| < u$ and $\|h_1 - h\|_\infty < u$, where $\boldsymbol{\theta}_1, \boldsymbol{\theta} \in \Theta$, $h_1, h \in \mathcal{H}$, and $u = o(1)$.

$$\begin{aligned} & m(\mathbf{D}, \boldsymbol{\theta}_1, h_1) - m(\mathbf{D}, \boldsymbol{\theta}, h) \\ &= \frac{\delta f_A(A)}{f_A(A|\mathbf{X}, \boldsymbol{\xi}_1)G_1(T)} \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}_1) - \frac{\delta f_A(A)}{f_A(A|\mathbf{X}, \boldsymbol{\xi})G(T)} \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \\ &= \left[\frac{\delta f_A(A)}{f_A(A|\mathbf{X}, \boldsymbol{\xi}_1)G_1(T)} - \frac{\delta f_A(A)}{f_A(A|\mathbf{X}, \boldsymbol{\xi})G(T)} \right] \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}_1) \\ &\quad + \frac{\delta f_A(A)}{f_A(A|\mathbf{X}, \boldsymbol{\xi})G(T)} \mathbf{Z}\mathbf{Z}^\top (\boldsymbol{\theta}_1 - \boldsymbol{\theta}) \\ &= \left[\frac{\delta f_A(A)}{f_A(A|\mathbf{X}, \boldsymbol{\xi}_1)G_1(T)} - \frac{\delta f_A(A)}{f_A(A|\mathbf{X}, \boldsymbol{\xi})G_1(T)} \right] \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}_1) \\ &\quad + \delta w_p(A; \mathbf{X}, \boldsymbol{\xi}) \left[\frac{1}{G_1(T)} - \frac{1}{G(T)} \right] \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}_1) + \frac{\delta w_p(A; \mathbf{X}, \boldsymbol{\xi})}{G(T)} \mathbf{Z}\mathbf{Z}^\top (\boldsymbol{\theta}_1 - \boldsymbol{\theta}) \\ &=: I_1 + I_2 + I_3. \end{aligned}$$

We first consider I_2 . As $\|h_1 - h\|_\infty < u$,

$$\begin{aligned} E[I_2^\top I_2] &\leq 2 \frac{\|G_1 - G\|_\infty^2}{\left(\frac{\tau}{2} \left(\frac{\tau}{2} - u\right)\right)^2} \left\{ (\log L)^2 E[w_p^2(A; \mathbf{X}, \boldsymbol{\xi}) \mathbf{Z} \mathbf{Z}^\top] + E[w_p^2(A; \mathbf{X}, \boldsymbol{\xi}) \mathbf{Z} \mathbf{Z}^\top \mathbf{Z} \mathbf{Z}^\top] \right\} \\ &\leq c_2 \|G_1 - G\|_\infty^2, \end{aligned} \quad (2.23)$$

for some constant c_2 , where the second inequality follows from Remark 1.2.2.

We next evaluate I_3 . By Remark 1.2.2 again,

$$E[I_3^\top I_3] \leq \frac{4}{\tau^2} (\boldsymbol{\theta}_1 - \boldsymbol{\theta})^\top E[w_p^2(A; \mathbf{X}, \boldsymbol{\xi}) \mathbf{Z} \mathbf{Z}^\top \mathbf{Z} \mathbf{Z}^\top (\boldsymbol{\theta}_1 - \boldsymbol{\theta})] \leq c_3 \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}\|^2, \quad (2.24)$$

for some constant c_3 .

We now assess I_1 . By (1.5) and the mean value theorem,

$$\begin{aligned} I_1 &= \frac{\delta}{G_1(T)} \exp\left(-\frac{A^2}{2}\right) \left[\sigma_1 \exp\left(\frac{(A - \mathbf{X}^\top \boldsymbol{\beta}_1)^2}{2\sigma_1^2}\right) - \sigma \exp\left(\frac{(A - \mathbf{X}^\top \boldsymbol{\beta})^2}{2\sigma^2}\right) \right] \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \\ &= \frac{\delta}{G_1(T)} \left[\frac{\partial}{\partial \boldsymbol{\xi}} w_p(A; \mathbf{X}, \tilde{\boldsymbol{\xi}}) \right]^\top (\boldsymbol{\xi}_1 - \boldsymbol{\xi}) \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}), \end{aligned}$$

where $\tilde{\boldsymbol{\xi}} = (\tilde{\boldsymbol{\beta}}^\top, \tilde{\sigma})^\top$ lies between $\boldsymbol{\xi}_1$ and $\boldsymbol{\xi}$ and its value depends on \mathbf{X} and A . Noting that \mathcal{X} and Ω are two compact sets in \mathbb{R}^d , according to Remark 1.2.2,

$$E[I_1^\top I_1] \leq c_1 \|\boldsymbol{\xi}_1 - \boldsymbol{\xi}\|^2, \quad (2.25)$$

for some constant c_1 . Combining (2.23), (2.24), and (2.25) together yields Equation (3.2) in CLK with $r = 2$ and $s_j = 1$.

Since Ω is a compact set in \mathbb{R}^d and $\|G - G^*\|_\infty \leq \tau/2$ by Remark 1.2.2, the covering number condition (3.3) in CLK that $\int_0^\infty \sqrt{\log N(v, \mathcal{H}, \|\cdot\|_\infty)} dv < \infty$ is satisfied, Therefore, by Theorem 3 in CLK, for all positive values $u = o(1)$,

$$\sup_{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \leq u, \|h - h^*\|_\infty \leq u} \|M_n(\boldsymbol{\theta}, h) - M(\boldsymbol{\theta}, h) - M_n(\boldsymbol{\theta}^*, h^*)\| = o_p(n^{-1/2}). \quad (2.26)$$

By the law of large number $M_n(\boldsymbol{\theta}^*, h^*) = o_p(1)$. Thus,

$$\sup_{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \leq u, \|h - h^*\|_\infty \leq u} \|M_n(\boldsymbol{\theta}, h) - M(\boldsymbol{\theta}, h)\| = o_p(1).$$

Condition (1.5) in CLK is satisfied.

By Theorem 1 in CLK, $\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^* = o_p(1)$. This completes the proof of Theorem 1.2.1. \square

Proof of Theorem 1.2.2: We prove the result by invoking Theorem 2 in CLK. Thus, we need to check their conditions (2.1) - (2.6).

In the proof of Theorem 1.2.1, we have already verified the conditions (2.4) and (2.5).

To verify the condition (2.1), as $\hat{\boldsymbol{\theta}}$ is the solution of $M_n(\boldsymbol{\theta}, \hat{h}) = 0$, $\|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| \leq \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \|M_n(\boldsymbol{\theta}, \hat{h})\|$. As $\hat{\boldsymbol{\theta}} \rightarrow_p \boldsymbol{\theta}^*$, $\forall u > 0$, $\lim_{n \rightarrow \infty} P(\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\| > u) = 0$. Since $\boldsymbol{\theta}^*$ is an interior point of $\boldsymbol{\Theta}$, we can find an u_0 such that $\boldsymbol{\Theta}_{u_0} \subset \boldsymbol{\Theta}$.

If $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\| \leq u_0$, we obtain $\|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| \geq \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \|M_n(\boldsymbol{\theta}, \hat{h})\|$ and subsequently $\|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| = \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \|M_n(\boldsymbol{\theta}, \hat{h})\|$. Thus, for any $u > 0$

$$\begin{aligned} & \lim_{n \rightarrow \infty} P \left(\sqrt{n} \left| \|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| - \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \|M_n(\boldsymbol{\theta}, \hat{h})\| \right| > u \right) \\ & \leq \lim_{n \rightarrow \infty} P \left(\left| \|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| - \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \|M_n(\boldsymbol{\theta}, \hat{h})\| \right| > 0 \right) \leq \lim_{n \rightarrow \infty} P(\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\| > u_0) = 0. \end{aligned}$$

Thus, $\|M_n(\hat{\boldsymbol{\theta}}, \hat{h})\| = \inf_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \|M_n(\boldsymbol{\theta}, \hat{h})\| + o_p(n^{-1/2})$.

To verify the condition (2.2), it is straightforward to obtain that $\boldsymbol{\Gamma}_1(\boldsymbol{\theta}, h^*) = E[w(\mathbf{D}, h^*)\mathbf{Z}\mathbf{Z}^\top]$, which is continuous at $\boldsymbol{\theta} = \boldsymbol{\theta}^*$. In addition, $\boldsymbol{\Gamma}_1(\boldsymbol{\theta}^*, h^*) = E[\mathbf{Z}\mathbf{Z}^\top]$ is of full rank. Thus, condition (2.2) is satisfied.

To verify the condition (2.3), simple algebra yields that

$$\begin{aligned}
\mathbf{\Gamma}_2(\boldsymbol{\theta}, h)(h_1 - h) &= \lim_{t \rightarrow 0} \frac{M(\boldsymbol{\beta}, h + t(h_1 - h)) - M(\boldsymbol{\beta}, h)}{t} \\
&= E \left[\lim_{t \rightarrow 0} \frac{\delta \{w_p(A; \mathbf{X}, \boldsymbol{\xi} + t(\boldsymbol{\xi}_1 - \boldsymbol{\xi})) [(G + t(G_1 - G))(T)]^{-1} - w_p(A; \mathbf{X}_i, \boldsymbol{\xi}) [G(T)]^{-1}\}}{t} \right. \\
&\quad \left. \times \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right] \\
&= E \left[\lim_{t \rightarrow 0} \frac{\delta \{w_p(A; \mathbf{X}, \boldsymbol{\xi} + t(\boldsymbol{\xi}_1 - \boldsymbol{\xi})) - w_p(A; \mathbf{X}_i, \boldsymbol{\xi})\}}{(G + t(G_1 - G))(T)t} \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right] \\
&\quad + E \left[\lim_{t \rightarrow 0} \frac{\delta w_p(A; \mathbf{X}, \boldsymbol{\xi}) \{[(G + t(G_1 - G))(T)]^{-1} - [G(T)]^{-1}\}}{t} \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right] \\
&= E \left[\left[\frac{\partial}{\partial \boldsymbol{\xi}} w_p(A; \mathbf{X}, \boldsymbol{\xi}) \right]^\top (\boldsymbol{\xi}_1 - \boldsymbol{\xi}) \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right] \\
&\quad + E \left[\frac{\delta}{G^2(T)} [(G - G_1)(T)] w_p(A; \mathbf{X}, \boldsymbol{\xi}) \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right].
\end{aligned}$$

By Condition (C1), $\forall \boldsymbol{\theta} \in \boldsymbol{\Theta}$, $\mathbf{\Gamma}_2(\boldsymbol{\theta}, h^*)(h_1 - h^*)$ exists in all directions $[h_1 - h^*] \in \mathcal{H}$.

$\forall (\boldsymbol{\theta}, h) \in \boldsymbol{\Theta}_{\epsilon_n} \times \mathcal{H}_{\epsilon_n}$ with a positive sequence $\epsilon_n = o(1)$.

(i) we show $\|M(\boldsymbol{\theta}, h) - M(\boldsymbol{\theta}, h^*) - \mathbf{\Gamma}_2(\boldsymbol{\theta}, h^*)(h - h^*)\| \leq C\|h - h^*\|_\infty^2$.

$$\begin{aligned}
&M(\boldsymbol{\theta}, h) - M(\boldsymbol{\theta}, h^*) - \mathbf{\Gamma}_2(\boldsymbol{\theta}, h^*)(h - h^*) \\
&= E \left\{ \frac{\delta}{G(T)} [w_p(A; \mathbf{X}, \boldsymbol{\xi}) - w_p(A; \mathbf{X}, \boldsymbol{\xi}^*)] \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right\} \\
&\quad + E \left\{ \delta w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \left[\frac{1}{G(T)} - \frac{1}{G^*(T)} \right] \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right\} - \mathbf{\Gamma}_2(\boldsymbol{\theta}, h^*)(h - h^*).
\end{aligned}$$

By Taylor expansion and Remark 1.2.2,

$$\begin{aligned}
&\left\| E \left\{ \frac{\delta}{G(T)} [w_p(A; \mathbf{X}, \boldsymbol{\xi}) - w_p(A; \mathbf{X}, \boldsymbol{\xi}^*)] \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right\} \right. \\
&\quad \left. - \left[\left[\frac{\partial}{\partial \boldsymbol{\xi}} w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \right]^\top (\boldsymbol{\xi} - \boldsymbol{\xi}^*) \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right] \right\| \\
&= \left\| E \left\{ (\boldsymbol{\xi} - \boldsymbol{\xi}^*)^\top \frac{\partial^2}{\partial \boldsymbol{\xi}^2} w_p(A; \mathbf{X}, \tilde{\boldsymbol{\xi}}) (\boldsymbol{\xi} - \boldsymbol{\xi}^*) \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right\} \right\| \leq C_1 \|\boldsymbol{\xi} - \boldsymbol{\xi}^*\|_\infty^2,
\end{aligned}$$

for some constant C_1 .

$$\begin{aligned}
& \left\| E \left\{ \delta w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \left[\frac{1}{G(T)} - \frac{1}{G^*(T)} \right] \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right\} \right. \\
& \quad \left. - E \left[\frac{\delta}{G^{*2}(T)} [(G^* - G)(T)] w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right] \right\| \\
&= \left\| E \left[\frac{\delta}{G^{*2}(T)G(T)} [(G^* - G)(T)]^2 w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}) \right] \right\| \leq C_2 \|G - G^*\|_\infty
\end{aligned}$$

for some constant C_2 . Combining the above two equations yields (i).

(ii) we show $\|\boldsymbol{\Gamma}_2(\boldsymbol{\theta}, h^*)(h - h^*) - \boldsymbol{\Gamma}_2(\boldsymbol{\theta}^*, h^*)(h - h^*)\| = o(\epsilon_n)$.

$$\begin{aligned}
& \|\boldsymbol{\Gamma}_2(\boldsymbol{\theta}, h^*)(h - h^*) - \boldsymbol{\Gamma}_2(\boldsymbol{\theta}^*, h^*)(h - h^*)\| \\
& \leq \left\| E \left[\exp \left(-\frac{A^2}{2} \right) \begin{pmatrix} -\exp \left(\frac{(A - \mathbf{X}^\top \boldsymbol{\beta}^*)^2}{2\sigma^{*2}} \right) \frac{(A - \mathbf{X}^\top \boldsymbol{\beta}^*) \mathbf{X}}{\sigma^*} \\ \exp \left(\frac{(A - \mathbf{X}^\top \boldsymbol{\beta}^*)^2}{2\sigma^{*2}} \right) \left(1 - \frac{(A - \mathbf{X}^\top \boldsymbol{\beta}^*)^2}{\sigma^{*2}} \right) \end{pmatrix}^\top (\boldsymbol{\xi} - \boldsymbol{\xi}^*) \mathbf{Z} \mathbf{Z}^\top \right] \right\| \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \\
& \quad + \left\| E \left[\frac{\delta}{G^{*2}(T)} [(G^* - G)(T)] w_p(A; \mathbf{X}, \boldsymbol{\xi}^*) \mathbf{Z} \mathbf{Z}^\top \right] \right\| \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \\
& = O(\epsilon_n^2),
\end{aligned}$$

where the last equality follows from $h \in H_{\epsilon_n}$ and Condition (C1). Therefore, condition (2.3) is satisfied.

To verify the condition (2.6), by Theorem 1 in Lo and Singh (1986) and Condition (C3),

$$\begin{aligned}
& M_n(\boldsymbol{\theta}^*, h^*) + \boldsymbol{\Gamma}_2(\boldsymbol{\theta}^*, h^*)(\hat{h} - h^*) \\
& = n^{-1} \sum_{i=1}^n m(\mathbf{D}_i, \boldsymbol{\theta}^*, h^*) + n^{-1} \sum_{i=1}^n \eta(A_i, \mathbf{X}_i) + n^{-1} \sum_{i=1}^n \psi(\tilde{T}_i, \delta_i) + O_p(n^{-3/4} \log^{3/4} n),
\end{aligned}$$

where

$$\begin{aligned}
\eta(A_i, \mathbf{X}_i) &= E \left[\mathbf{Z}(\tilde{Y} - \mathbf{Z}^\top \boldsymbol{\theta}^*) \left[\frac{\partial}{\partial \boldsymbol{\xi}} w_p(A; \mathbf{X}, \boldsymbol{\xi}) \right]^\top \right] \mathbf{I}(\boldsymbol{\xi}^*)^{-1} \sum_{i=1}^n \frac{\frac{\partial}{\partial \boldsymbol{\xi}} f_A(A_i; \mathbf{X}_i, \boldsymbol{\xi}^*)}{f_A(A_i; \mathbf{X}_i, \boldsymbol{\xi}^*)} \quad (2.27) \\
\psi(\tilde{T}_i, \delta_i) &= \int_{\mathcal{T} \times \mathcal{A} \times \mathbf{X}} \frac{\phi(t; \tilde{T}_i, \delta_i)}{G^*(t)} w_p(a; \mathbf{x}, \boldsymbol{\xi}^*) \mathbf{z}(\log t - \mathbf{z}^\top \boldsymbol{\theta}^*) dF_{T,A,\mathbf{X}}(t, a, \mathbf{x}), \quad (2.28) \\
\phi(t, \tilde{T}_i, \delta_i) &= G^*(t) \int_0^{\min(\tilde{T}_i, t)} \frac{f_C(s) ds}{G^{*2}(s)(1 - F_T(s))} ds + \frac{1 \left\{ \tilde{T}_i \leq t, \delta_i = 0 \right\}}{(1 - F_T(\tilde{T}_i))G^*(\tilde{T}_i)}, \\
\text{and } \mathbf{I}(\boldsymbol{\xi}^*) &= \sigma^{*-2} \begin{pmatrix} \mathbf{X}\mathbf{X}^\top & 0 \\ 0 & 2 \end{pmatrix},
\end{aligned}$$

As $m(\mathbf{D}_i, \boldsymbol{\theta}^*, h^*) + \eta(A_i, \mathbf{X}_i) + \psi(\tilde{T}_i, \delta_i)$ are independent random variables with mean zero and finite variance. Let \mathbf{V} denote the covariance matrix of $(m(\mathbf{D}_i, \boldsymbol{\theta}^*, h^*) + \eta(A_i, \mathbf{X}_i) + \psi(\tilde{T}_i, \delta_i))$. by the central limit theorem, $\sqrt{n} \left(M_n(\boldsymbol{\theta}^*, h^*) + \boldsymbol{\Gamma}_2(\boldsymbol{\theta}^*, h^*)(\hat{h} - h^*) \right) \rightarrow N(0, \mathbf{V})$. Thus, condition (2.6) is satisfied. By Theorem 2 in CLK,

$$\sqrt{n}(\sqrt{\hat{\boldsymbol{\theta}}} - \boldsymbol{\theta}^*) \rightarrow_d N(0, (E[\mathbf{Z}\mathbf{Z}^\top])^{-1} \mathbf{V} (E[\mathbf{Z}\mathbf{Z}^\top])^{-1}).$$

This completes the proof of Theorem 1.2.2. □

CURRICULUM VITA

NAME: Triparna Poddar

ADDRESS: Department of Biostatistics and Bioinformatics
University of Louisville
Louisville, KY 40202

EDUCATION: Master of Science in Statistics, 2015-2017
Department of Statistics, University of Calcutta
Kolkata, West Bengal, India

Bachelor of Science in Statistics, 2012-2015
Lady Brabourne College, University of Calcutta
Kolkata, West Bengal, India

EXPERIENCES: Graduate Research Assistant, August, 2019-Present
Department of Bioinformatics and Biostatistics,
School of Public Health and Information Sciences
University of Louisville, Louisville, KY, USA

Project Fellow, May-June, 2019
Department of Statistics, University of Calcutta,
Project: "LISA 2020: Creating Institutional Statistical

Analysis and Data Science Capacity”,
funded by USAID in collaboration with University of
Colorado Boulder

Research Fellow, March 2018-March, 2019
Department of Statistics, University of Calcutta,
Project: “Image and Imaging”

Intern as Trainee Data Analyst, July-Dec, 2017
Orbital Software Development Pvt. Ltd
West Bengal, India

PUBLICATIONS: **Poddar T.** , Zheng Q., Kong K., Estimation Of Average
Treatment Effect For Survival Outcomes With
Continuous Treatment In Observational Studies
(*Under preparation*)

Poddar T., Kong K., Causal Mediation Analysis
For Health Racial Disparities (*Under preparation*)

Egger M., Kong M., Little B., Ghosh I., **Poddar T.**,
Tyler EC., Goldsby M, Vu G. Disparities in Cancer Screening
in the Kentucky Medicaid Population (*Manuscript*)

Flaherty D., Winrich E., Eisa, M., **Poddar T.**,
Pooler A., Kong, M., Omer, E.
Colonoscopy Adenomatous Polyp Detection Rate in

Patients With Inadequate Bowel Prep
The American Journal of Gastroenterology
117(10S):e170, October 2022

Winrich E., Flaherty D., Eisa, M., **Poddar T.**, Pooler A.,
Kong, M., Omer, E. Racial And Economic Risk Factors
For Inadequate Bowel Preparation
Prior To Colonoscopy
Gastroenterology 162(7):S-469, May 2022

Flaherty D., Winrich E., Eisa, M., **Poddar T.**, Pooler A.,
Kong, M., Omer, E. Medical Risk Factors For
Inadequate Bowel Preparation Prior
To Colonoscopy *Gastroenterology 162(7):S-1029, May 2022*

Saran U., Chandrasekaran B., Kolluru V., Tyagi A.,
Nguyen KD., Valadon CL., Shaheen SP., Kong M.,
Poddar T., Ankem MK., Damodaran C.,
Diagnostic molecular markers predicting aggressive
potential in low-grade prostate cancer
Translational Research Volume 231, May 2021

Mukherjee J., **Poddar T.**, Kar M., Ganguli B.,
Chakrabarti A.,
A Feature based Automated Classification of
Subcentimeter Pulmonary Structures in Thoracic
Computed Tomography Images

PRESENTATIONS: **Poddar T.**, Zheng Q., Egger M., Kong K.,

Title: *Estimation of Impact of Blood Lead Levels on
All-cause Mortality of Older People in US Population
from NHANES III Survey,*

Poster Presentation at Research Louisville, October 2023

First Prize in Research & Practice Category

Poddar T., Zheng Q., Kong K.,

Title: *Estimation Of Average
Treatment Effect For Survival Outcomes With Continuous
Treatment In Observational Studies,*

Oral Presentation at ENAR,

Nashville, March 2023

Poddar T., Mukherjee J., Ganguli B., Kar M.,
Chakrabarti A.,

Title: *A Clinically Applicable Automated Risk
Classification model for Pulmonary Nodules,*

Oral Presentation at International

Conference on Data Management, Analytics and Innovation
at Singapore, January 2019

Poddar T., Mukherjee J.,

Title: *A CAD Tool for Automated
Detection of Subcentimeter Pulmonary*

Nodules from CT Images,

Poster Presentation at 10th International

Triennial Calcutta, Symposium on

Probability and Statistics, December, 2018

Third Prize