

*J. R. Statist. Soc. A* (2020)  
183, Part 1, pp. 311–332

# Regression-with-residuals estimation of marginal effects: a method of adjusting for treatment-induced confounders that may also be effect modifiers

Geoffrey T. Wodtke,

*University of Chicago, USA*

Zahide Alaca

*University of Toronto, Canada*

and Xiang Zhou

*Harvard University, Cambridge, USA*

[Received August 2018. Revised July 2019]

**Summary.** When making causal inferences, treatment-induced confounders complicate analyses of time-varying treatment effects. Conditioning on these variables naively to estimate marginal effects may inappropriately block causal pathways and may induce spurious associations between the treatment and the outcome, leading to bias. Although several methods for estimating marginal effects avoid these complications, including inverse probability of treatment weighted estimation of marginal structural models as well as  $g$ - and regression-with-residuals estimation of highly constrained structural nested mean models, each suffers from a set of non-trivial limitations, among them an inability to accommodate effect modification. In this study, we adapt the method of regression with residuals to estimate marginal effects with a set of moderately constrained structural nested mean models that easily accommodate several types of treatment-by-confounder interaction. With this approach, the confounders at each time point are first residualized with respect to the observed past, which involves centring them at their estimated means given prior treatments and confounders. The outcome is then regressed on all prior variables, including a set of treatment-by-confounder interaction terms, with these residuals substituted for the untransformed confounders both as ‘main effects’ and as part of any interaction terms. Through a series of simulation experiments and empirical examples, we show that this approach outperforms other methods for estimating the marginal effects of time-varying treatments.

**Keywords:**  $g$ -estimation; Inverse probability of treatment weighting; Marginal structural models; Regression with residuals; Structural nested mean models; Treatment-induced confounding

## 1. Introduction

In analyses of time-varying treatment effects, social scientists must often contend with the complications that are posed by treatment-induced confounders (e.g. Acharya *et al.* (2016), Elwert and Winship (2014) and Wodtke *et al.* (2011)). A treatment-induced confounder is a variable that is affected by a prior treatment and affects both selection into future treatment

*Address for correspondence:* Geoffrey T. Wodtke, Department of Sociology, University of Chicago, 1126 East 59th Street, Chicago, IL 60637, USA.  
E-mail: wodtke@uchicago.edu

and the outcome. For example, we consider whether living in a disadvantaged neighbourhood throughout childhood and adolescence affects academic achievement (e.g. Sampson *et al.* (2008) and Wodtke *et al.* (2011, 2016)). In studies of neighbourhood effects, parental income is likely to be affected by prior neighbourhood conditions and also likely to affect both future residential choices and child educational outcomes.

If left uncontrolled, treatment-induced confounders lead to bias in estimates of marginal effects, such as the cumulative treatment effect (CTE) in analyses of time-varying treatments. At the same time, adjusting naively for treatment-induced confounders by, for example, including them as predictors in a conventional regression model or matching on them via the propensity score also leads to bias. Specifically, conditioning on a treatment-induced confounder with conventional regression or matching methods leads to bias from overcontrol of intermediate pathways and endogenous selection (Elwert and Winship, 2014; Robins *et al.*, 2000; VanderWeele, 2015). Thus even if all relevant confounders are observed, which is a necessary condition for drawing causal inferences from any non-experimental study, treatment-induced confounders pose additional challenges for the most common approaches to covariate adjustment. Alternative methods are therefore required when estimating marginal effects in the presence of these variables.

Fortunately, there are several methods that avoid the complications that were outlined previously and that are capable of consistently estimating marginal effects, even when adjustment is required for treatment-induced confounders. These include inverse probability of treatment weighting (IPTW) estimation of marginal structural models (Robins *et al.*, 1994, 2000), *g*-estimation of highly constrained structural nested mean models (SNMMs) (Naimi *et al.*, 2017; Vansteelandt, 2009; Vansteelandt and Sjolander, 2016) and regression-with-residuals (RWR) estimation of highly constrained SNMMs (Wodtke, 2018).

Each of these methods, however, suffers from a set of non-trivial limitations. IPTW estimation is relatively inefficient, is difficult to use with continuous treatments and may suffer from finite sample bias when confounders strongly predict treatment (Lunceford and Davidian, 2004; Naimi *et al.*, 2014; Robins *et al.*, 1994). *g*- and RWR estimation of highly constrained SNMMs for marginal effects avoid the limitations of IPTW, but they are premised on the strong assumption of no effect modification (e.g. Vansteelandt (2009) and Wodtke (2018)), which is unrealistic in most social science applications. If, for example, a treatment-induced confounder also modifies the effect of a future treatment on the outcome, then these methods suffer from model misspecification bias. Because effect modification is ubiquitous in the social sciences (Morgan and Winship, 2015; Xie, 2007), this assumption may limit the utility of these methods in practice.

In this study, we adapt the method of RWR to estimate a set of moderately constrained SNMMs for marginal effects that accommodate several types of treatment-by-confounder interaction. Briefly, RWR estimation of marginal effects in a moderately constrained SNMM proceeds in two stages. First, the confounders at each time point are regressed on all prior variables and then residualized, which involves centring them at their estimated means given prior treatments and confounders. Second, the outcome is regressed on all prior variables, including a set of treatment-by-confounder interaction terms, with the residuals from the first stage substituted for the untransformed confounders both as ‘main effects’ and as part of the interaction terms. Our adaptation differs from previous implementations of RWR (e.g. Almirall *et al.* (2010), Wodtke and Almirall (2017) and Wodtke (2018)) by additionally including the residualized confounders in interaction terms with treatment, which accommodates several types of effect modification while neatly isolating the marginal effects of interest in a single, possibly vector-valued, parameter.

Under the assumptions of sequential ignorability and no model misspecification, the method proposed is consistent for marginal effects, like the CTE, even in the presence of treatment-

induced confounders. It avoids the biases that arise with naive adjustments for treatment-induced confounders because the residualized confounders are purged of their association with prior treatment and thus including them in a regression model for the outcome is unproblematic. In addition, because it does not involve weighting by a function of the conditional probability of treatment, the method avoids the limitations that are associated with IPTW estimation. Finally, because it accommodates several types of treatment-by-confounder interaction, the method also mitigates the limitations that are associated with both  $g$ - and RWR estimation of marginal effects based on a highly constrained SNMM in which effect modification is assumed away entirely.

In the sections that follow, we begin by considering the problem of estimating marginal effects for a time-varying treatment. First, we formally define the effects of interest, explain when they are identified from observed data and illustrate the problems that afflict conventional estimation methods in the presence of treatment-induced confounding. Second, we present an SNMM for the conditional, rather than marginal, effects of treatment, but we then show how these conditional effects can be additively decomposed into a set of functions that capture the marginal effects of interest and another set of functions that capture effect modification. Third, we show how to parameterize these functions appropriately and to adapt the method of RWR to estimate marginal effects with an SNMM under this alternative parameterization. Next, we briefly consider analyses of causal mediation and the problem of estimating controlled direct effects (CDEs) in the presence of mediator–outcome confounders that are affected by treatment, which we show can be accomplished with the same methods used for estimating marginal effects in the time-varying setting. Finally, with a series of simulation experiments and empirical examples, we illustrate several applications of our proposed method and show that it outperforms other common approaches.

The data that are analysed in the paper and the programs that were used to analyse them can be obtained from

<https://rss.onlinelibrary.wiley.com/hub/journal/1467985x/series-a-datasets>,

with one exception: some of the data used in this analysis were based on sensitive data files and were obtained under special contractual arrangements from the Panel Study of Income Dynamics (PSID). These data are not available on line or from the authors. Researchers interested in obtaining sensitive data from the PSID should contact PSIDhelp@umich.edu.

## 2. Background

### 2.1. Notation, estimands and identification

In this section, we formally define the marginal, or population-average, effects of interest and explain when they can be identified from observed data, drawing on the potential outcomes framework (Holland, 1986; Rubin, 1974) and directed acyclic graphs (Pearl, 2009) throughout. For expositional clarity, we focus on a simplified example with a binary treatment measured at two time points, a binary confounder measured at two time points and a continuous outcome measured at the end of follow-up, although these methods can be easily adapted for more complex analyses.

First, let  $a_t = 1$  denote exposure to treatment, and  $a_t = 0$  denote the absence of treatment, at time  $t \in \{1, 2\}$ . Second, let  $Y_i(a_1, a_2)$  denote the potential outcome for subject  $i$  if she had previously been exposed to the treatment sequence  $\{a_1, a_2\}$ . For example,  $Y_i(0, 0)$  is the potential outcome for subject  $i$  if she had never received treatment,  $Y_i(1, 0)$  is her potential outcome if she had received treatment only at time  $t = 1$ , and so on. In this framework, each subject is conceived

to have a potential outcome corresponding to each of the four possible treatment sequences, but only the single potential outcome corresponding to the treatment sequence actually received is ever observed in reality, and the others are so-called ‘counterfactuals’. Third, let  $C_{i1}$  denote the confounder for subject  $i$  measured just before treatment at time  $t = 1$ , and let  $C_{i2}(a_1)$  denote the confounder for subject  $i$  measured just before treatment at time  $t = 2$ , which is indexed by  $a_1$  as a potential outcome to reflect that it is affected by prior treatment. In other words,  $C_{i2}(a_1)$  is a treatment-induced confounder. Finally, let the set  $\{C_{i1}, A_{i1}, C_{i2}, A_{i2}, Y_i\}$  denote the observed data in temporal order.

In general, marginal effects are contrasts between different potential outcomes averaged over a population of individuals. Specifically, they give the average difference in the end-of-study outcome if everyone in the target population had received one rather than another treatment sequence. With two time points, several different marginal effects may be of interest. The first is the distal treatment effect (DTE), which can be formally defined as

$$\text{DTE}(a_1) = E\{Y_i(a_1, 0) - Y_i(0, 0)\}. \quad (1)$$

It gives the average effect of receiving treatment only at time  $t = 1$  rather than never receiving treatment. The second is the proximal treatment effect (PTE), which can be formally defined as

$$\text{PTE}(a_1, a_2) = E\{Y_i(a_1, a_2) - Y_i(a_1, 0)\}. \quad (2)$$

When  $a_1 = 0$ , it gives the average effect of receiving treatment only at time  $t = 2$  rather than never receiving treatment and, when  $a_1 = 1$ , it gives the average effect of always receiving treatment rather than receiving treatment only at time  $t = 1$ . The third is the CTE. This effect is equal to the sum of DTE(1) and PTE(1,1),

$$\begin{aligned} \text{CTE} &= \text{DTE}(1) + \text{PTE}(1, 1) \\ &= E\{Y_i(1, 0) - Y_i(0, 0)\} + E\{Y_i(1, 1) - Y_i(1, 0)\} \\ &= E\{Y_i(1, 1) - Y_i(0, 0)\}, \end{aligned} \quad (3)$$

which gives the average effect of being always *versus* never treated. Finally, the last is the interaction effect, or INE. This effect can be formally defined as

$$\begin{aligned} \text{INE} &= \text{PTE}(1, 1) - \text{PTE}(0, 1) \\ &= E\{Y_i(1, 1) - Y_i(1, 0)\} - E\{Y_i(0, 1) - Y_i(0, 0)\}, \end{aligned} \quad (4)$$

which describes how the effect of receiving treatment at time  $t = 2$  differs depending on whether an individual had previously received treatment at time  $t = 1$ .

All these effects can be non-parametrically identified from the observed data under the assumptions of stable unit treatment values, consistency, positivity and sequential ignorability (Robins *et al.*, 1994, 2000; Rubin, 1986). The stable unit treatment value assumption requires that the potential outcomes for any given individual be unaffected by the mechanisms that are used to assign treatment status or by the treatments to which other individuals are exposed. The consistency assumption requires that the observed outcome  $Y_i$  be equal to  $Y_i(a_1, a_2)$  whenever  $A_{i1} = a_1$  and  $A_{i2} = a_2$ . The positivity assumption requires that there not be any subgroups within the target population that are treated or untreated with certainty at any time point. The sequential ignorability assumption requires that the potential outcomes are independent of treatment at each time point conditional on the observed past. Formally, this assumption can be expressed as

$$\begin{aligned}
 & Y_i(a_1, a_2) \perp A_{i1} | C_{i1} && \forall (a_1, a_2), \\
 & Y_i(a_1, a_2) \perp A_{i2} | C_{i1}, A_{i1}, C_{i2} && \forall (a_1, a_2),
 \end{aligned}
 \tag{5}$$

where ‘ $\perp$ ’ denotes statistical independence. It is satisfied when there are not any unobserved variables that directly affect both selection into treatment at each time point and the outcome.

Fig. 1(a) presents a directed acyclic graph illustrating a set of causal relationships between the variables that were outlined previously in which the sequential ignorability assumption is satisfied. It shows that both treatments  $A_{i1}$  and  $A_{i2}$  directly affect the outcome  $Y_i$  and that  $A_{i1}$  also indirectly affects the outcome through  $C_{i2}$ . In addition, it shows that  $C_{i1}$  confounds the effect of  $A_{i1}$  on  $Y_i$  and that  $C_{i2}$  confounds the effect of  $A_{i2}$  on  $Y_i$ . Treatment assignment is sequentially ignorable in Fig. 1(a) because treatment at each time point is not directly affected by any unobserved variables; rather, the only unobserved variables, which are denoted by  $U_{i1}$  and  $U_{i2}$ , directly affect the observed confounders and the outcome but not either treatment. The marginal effects that were outlined previously can be consistently estimated from the observed data by appropriately adjusting for all variables that directly affect both the treatment and the outcome—in this case,  $C_{i1}$  and  $C_{i2}$ .

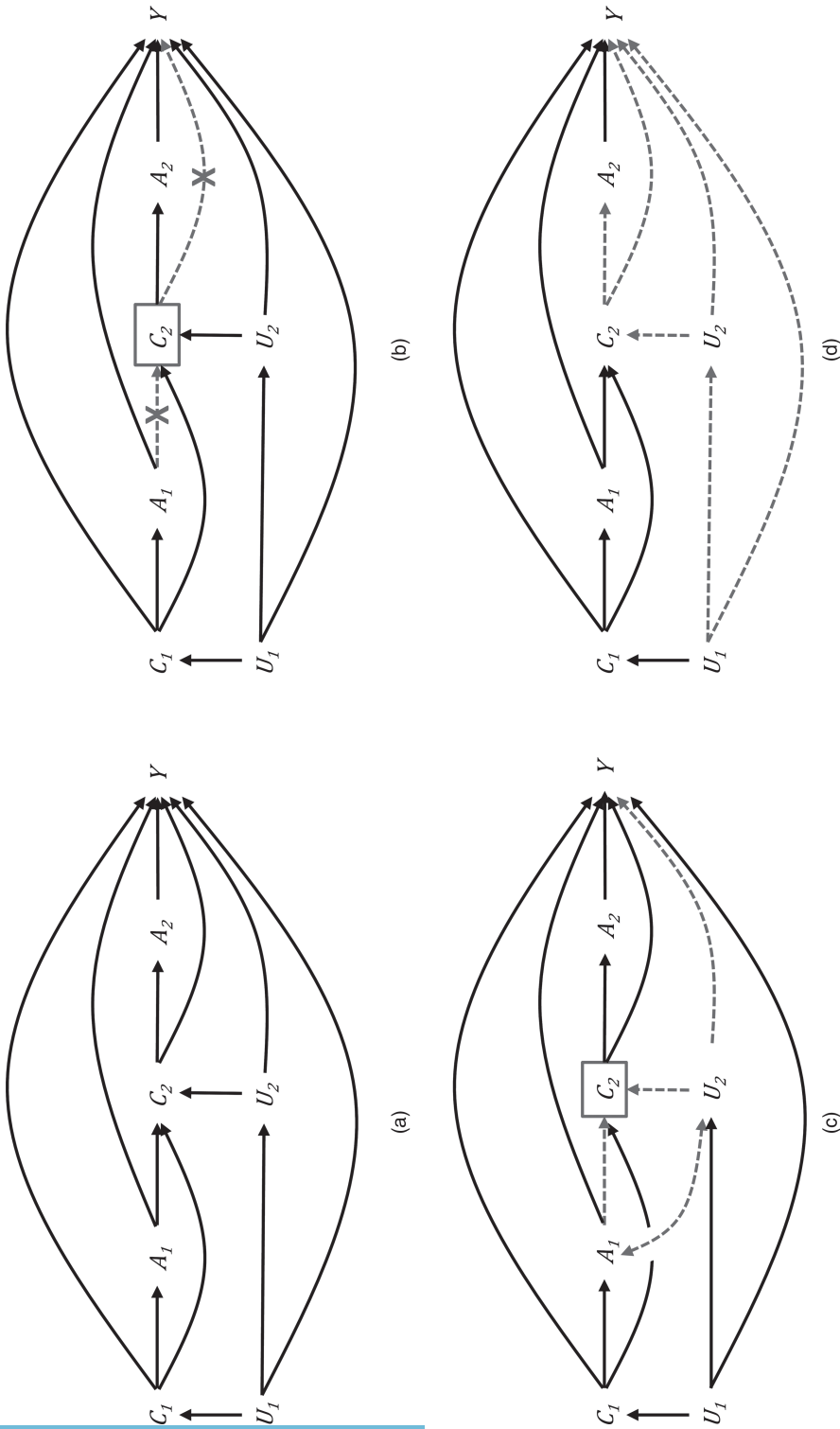
### 2.2. The problem of treatment-induced confounding

Because  $C_{i2}$  is affected by  $A_{i1}$  and confounds the effect of  $A_{i2}$  on  $Y_i$ , it is a treatment-induced confounder. Treatment-induced confounders pose several challenges for estimating marginal effects of a time-varying treatment (Daniel *et al.*, 2013). In particular, conventional methods of covariate adjustment, including conditioning, stratifying or matching directly on a treatment-induced confounder, lead to several types of bias, even when the effects of interest are identified under sequential ignorability. At the same time, failing to adjust appropriately for a treatment-induced confounder also leads to bias. Thus, treatment-induced confounders seemingly present a ‘damned if you do and damned if you don’t’ dilemma with regard to covariate adjustment.

To appreciate this, first consider the causal graph in Fig. 1(b), and recall that a path in a directed acyclic graph is ‘blocked’ when it contains

- (a) an outcome of two or more variables, known as a collider, that has not been conditioned on or
- (b) a non-collider that has been conditioned on;

otherwise, it is ‘unblocked’ (Pearl, 2009). Fig. 1(b) shows that conditioning naively on the treatment-induced confounder  $C_{i2}$  blocks the causal pathway  $A_{i1} \rightarrow C_{i2} \rightarrow Y_i$  emanating from treatment at time  $t = 1$  to the outcome, which leads to bias from overcontrol of intermediate pathways. Next consider the causal graph in Fig. 1(c). This graph shows that conditioning naively on  $C_{i2}$  also unblocks the non-causal pathway  $A_{i1} \rightarrow C_{i2} \leftarrow U_{i2} \rightarrow Y_i$  emanating from treatment at time  $t = 1$  to the outcome, which leads to bias from so-called ‘endogenous selection’ or ‘collider stratification’ (Elwert and Winship, 2014). Specifically, it shows that  $C_{i2}$  is a collider of  $A_{i1}$  and  $U_{i2}$ , and, because  $U_{i2}$  affects  $Y_i$ , conditioning on  $C_{i2}$  induces a spurious association between treatment at time  $t = 1$  and the outcome. Finally, consider the graph in Fig. 1(d). This graph shows that, when  $C_{i2}$  has not been conditioned on, the non-causal pathways emanating from treatment at time  $t = 2$  to the outcome ( $A_{i2} \leftarrow C_{i2} \rightarrow Y_i$ ,  $A_{i2} \leftarrow C_{i2} \leftarrow U_{i2} \rightarrow Y_i$  and  $A_{i2} \leftarrow C_{i2} \leftarrow U_{i2} \leftarrow U_{i1} \rightarrow Y_i$ ) remain unblocked, which leads to bias from uncontrolled confounding. Thus, conventional methods of covariate adjustment inevitably lead to bias in estimates of marginal effects when there is treatment-induced confounding, and alternative methods are required.



**Fig. 1.** Directed acyclic graphs illustrating a set of causal relationships between a time-varying treatment  $A_t$ , an observed time-varying confounder  $C_t$ , an unobserved time-varying covariate  $U_t$  and an outcome  $Y$ : (a) causal relationships; (b) overcontrol of intermediate pathways; (c) endogenous selection; (d) uncontrolled confounding

### 3. Regression with residuals for marginal effects

#### 3.1. Estimating the marginal effects of a time-varying treatment

An SNMM is a model for the conditional, or subpopulation-average, effects of a time-varying treatment given past confounder and treatment history (Almirall *et al.*, 2010; Robins, 1994; Robins *et al.*, 2007; Wodtke and Almirall, 2017). In this section, we show that conditional effects modelled with an SNMM can be additively decomposed into a set of functions that capture the marginal, or population-average, effects of interest and another set of functions that capture effect modification, i.e. whether and how the conditional effects vary around the marginal effects. We then show how to parameterize these functions appropriately and to adapt the RWR method to estimate them consistently.

An SNMM is based on the following decomposition of the conditional mean of the potential outcomes given the confounders into a set of conditional treatment effects and a set of so-called ‘nuisance’ associations:

$$E\{Y_i(a_1, a_2)|C_{i1}, C_{i2}(a_1)\} = \beta_{00} + \varepsilon_1(C_{i1}) + \mu_1(C_{i1}, a_1) + \varepsilon_2\{C_{i1}, a_1, C_{i2}(a_1)\} + \mu_2\{C_{i1}, a_1, C_{i2}(a_1), a_2\}. \quad (6)$$

Here  $\beta_{00} = E\{Y_i(0, 0)\}$  is the marginal mean of the potential outcomes under no treatment;  $\varepsilon_1(C_{i1}) = E\{Y_i(0, 0)|C_{i1}\} - E\{Y_i(0, 0)\}$  is a nuisance association that captures the relationship between the confounder at time  $t = 1$  and the outcome under no treatment;  $\mu_1(C_{i1}, a_1) = E\{Y_i(a_1, 0) - Y_i(0, 0)|C_{i1}\}$  is a causal function that captures the conditional effects of treatment at time  $t = 1$  given  $C_{i1}$ ;  $\varepsilon_2\{C_{i1}, a_1, C_{i2}(a_1)\} = E\{Y_i(a_1, 0)|C_{i1}, C_{i2}(a_1)\} - E\{Y_i(a_1, 0)|C_{i1}\}$  is another nuisance association that captures the relationship between the confounder at time  $t = 2$  and the outcome under treatment sequence  $\{a_1, 0\}$ ; and  $\mu_2\{C_{i1}, a_1, C_{i2}(a_1), a_2\} = E\{Y_i(a_1, a_2) - Y_i(a_1, 0)|C_{i1}, C_{i2}(a_1)\}$  is another causal function that captures the conditional effects of treatment at time  $t = 2$  given both prior confounders. The functions  $\varepsilon_1(C_{i1})$  and  $\varepsilon_2\{C_{i1}, a_1, C_{i2}(a_1)\}$  are called ‘nuisance’ associations because they do not contain any information about the causal effects of treatment (Wodtke and Almirall, 2017).

The first causal function,  $\mu_1(C_{i1}, a_1)$ , can be further decomposed into a marginal effect of interest and a term that captures effect modification as follows:

$$\mu_1(C_{i1}, a_1) = \mu_{11}(a_1) + \mu_{12}(C_{i1}, a_1), \quad (7)$$

where  $\mu_{11}(a_1) = E\{Y_i(a_1, 0) - Y_i(0, 0)\}$  is equal to DTE( $a_1$ ) and  $\mu_{12}(C_{i1}, a_1) = E\{Y_i(a_1, 0) - Y_i(0, 0)|C_{i1}\} - E\{Y_i(a_1, 0) - Y_i(0, 0)\}$  captures how the effect of treatment at time  $t = 1$  differs across levels of  $C_{i1}$ .

Similarly, the second causal function,  $\mu_2\{C_{i1}, a_1, C_{i2}(a_1), a_2\}$ , can also be further decomposed as follows:

$$\mu_2\{C_{i1}, a_1, C_{i2}(a_1), a_2\} = \mu_{21}(a_1, a_2) + \mu_{22}(C_{i1}, a_1, a_2) + \mu_{23}\{C_{i1}, a_1, C_{i2}(a_1), a_2\}, \quad (8)$$

where  $\mu_{21}(a_1, a_2) = E\{Y_i(a_1, a_2) - Y_i(a_1, 0)\}$  is equal to another marginal effect of interest, PTE( $a_1, a_2$ );  $\mu_{22}(C_{i1}, a_1, a_2) = E\{Y_i(a_1, a_2) - Y_i(a_1, 0)|C_{i1}\} - E\{Y_i(a_1, a_2) - Y_i(a_1, 0)\}$  captures how the effect of treatment at  $t = 2$  differs across levels of  $C_{i1}$ ; and  $\mu_{23}\{C_{i1}, a_1, C_{i2}(a_1), a_2\} = E\{Y_i(a_1, a_2) - Y_i(a_1, 0)|C_{i1}, C_{i2}(a_1)\} - E\{Y_i(a_1, a_2) - Y_i(a_1, 0)|C_{i1}\}$  captures how the effect of treatment at  $t = 2$  differs across levels of  $C_{i2}(a_1)$  within levels of  $C_{i1}$ .

Any parameterization of the marginal effects  $\mu_{11}(a_1)$  and  $\mu_{21}(a_1, a_2)$  must satisfy the constraint that they are equal to 0 when contemporaneous treatment is equal to 0. With a binary



treatment, a saturated parameterization for  $\mu_{11}(a_1)$  is

$$\mu_{11}(a_1) = \beta_{10}a_1, \quad (9)$$

and a saturated parameterization for  $\mu_{21}(a_1, a_2)$  is

$$\mu_{21}(a_1, a_2) = (\beta_{20} + \beta_{21}a_1)a_2, \quad (10)$$

where  $\beta_{10} = \text{DTE}(1)$ ,  $\beta_{20} = \text{PTE}(0, 1)$  and  $\beta_{20} + \beta_{21} = \text{PTE}(1, 1)$ . In addition, note that  $\beta_{10} + \beta_{20} + \beta_{21} = \text{CTE}$  and  $\beta_{21} = \text{INE}$ . Thus, all the marginal effects that were defined previously are given by the parameter vector  $\{\beta_{10}, \beta_{20}, \beta_{21}\}$ .

Any parameterization of  $\mu_{12}(C_{i1}, a_1)$  must satisfy the constraints that it is equal to 0 when  $a_1 = 0$  and that it has mean 0. With a treatment and confounder that are both binary, a saturated parameterization for this function is

$$\mu_{12}(C_{i1}, a_1) = \theta_{10}a_1C_{i1}^\perp, \quad (11)$$

where  $C_{i1}^\perp = C_{i1} - E(C_{i1})$  is a residual transformation of  $C_{i1}$  with respect to its marginal mean. This parameterization satisfies the zero-mean constraint because  $E(\theta_{10}a_1C_{i1}^\perp) = \theta_{10}a_1E(C_{i1}^\perp) = \theta_{10}a_1E\{C_{i1} - E(C_{i1})\} = 0$ .

Similarly, any parameterization of  $\mu_{22}(C_{i1}, a_1, a_2)$  must satisfy the constraints that it is equal to 0 when  $a_2 = 0$  and that it has mean 0. A saturated parameterization for this function is

$$\mu_{22}(C_{i1}, a_1, a_2) = (\theta_{20} + \theta_{21}a_1)a_2C_{i1}^\perp, \quad (12)$$

which has mean 0 because the expectation function is a linear operator and because  $E(C_{i1}^\perp) = 0$ , as above.

Finally, any parameterization of  $\mu_{23}\{C_{i1}, a_1, C_{i2}(a_1), a_2\}$  must satisfy the constraints that it is equal to 0 when  $a_2 = 0$  and that it has mean 0 conditionally on  $C_{i1}$ . A saturated parameterization for this function is

$$\mu_{23}\{C_{i1}, a_1, C_{i2}(a_1), a_2\} = \{\theta_{22} + \theta_{23}a_1 + (\theta_{24} + \theta_{25}a_1)C_{i1}^\perp\}a_2C_{i2}^\perp(a_1), \quad (13)$$

where  $C_{i2}^\perp(a_1) = C_{i2}(a_1) - E\{C_{i2}(a_1)|C_{i1}\}$  is a residual transformation of  $C_{i2}(a_1)$  with respect to its conditional mean given  $C_{i1}$ . This parameterization satisfies the zero-mean constraint because the expectation function is a linear operator and because  $E\{C_{i2}^\perp(a_1)|C_{i1}\} = E\{[C_{i2}(a_1) - E\{C_{i2}(a_1)|C_{i1}\}]|C_{i1}\} = E\{C_{i2}(a_1)|C_{i1}\} - E\{C_{i2}(a_1)|C_{i1}\} = 0$ . The parameter vector  $\{\theta_{10}, \theta_{20}, \theta_{21}, \theta_{22}, \theta_{23}, \theta_{24}, \theta_{25}\}$  captures how the confounders modify the effect of treatment at each time point.

The nuisance associations  $\varepsilon_1(C_{i1})$  and  $\varepsilon_2\{C_{i1}, a_1, C_{i2}(a_1)\}$  must also be parameterized under the constraint that they have mean 0 given the past, which can be accomplished by using the same residualized confounders as defined previously. Specifically, a saturated parameterization for the first nuisance association is

$$\varepsilon_1(C_{i1}) = \gamma_{10}C_{i1}^\perp, \quad (14)$$

and a saturated parameterization for the second nuisance association is

$$\varepsilon_2\{C_{i1}, a_1, C_{i2}(a_1)\} = \{\gamma_{20} + \gamma_{21}a_1 + (\gamma_{22} + \gamma_{23}a_1)C_{i1}^\perp\}C_{i2}^\perp(a_1), \quad (15)$$

where the parameter vector  $\{\gamma_{10}, \gamma_{20}, \gamma_{21}, \gamma_{22}, \gamma_{23}\}$  captures the associational (i.e. causal and possibly non-causal) effects of the confounders on the outcome.



Combining parametric expressions for the causal functions and nuisance associations yields the following saturated SNMM:

$$\begin{aligned}
 E\{Y_i(a_1, a_2)|C_{i1}, C_{i2}(a_1)\} = & \beta_{00} + \gamma_{10}C_{i1}^\perp + \beta_{10}a_1 + \theta_{10}a_1C_{i1}^\perp + \{\gamma_{20} + \gamma_{21}a_1 \\
 & + (\gamma_{22} + \gamma_{23}a_1)C_{i1}^\perp\}C_{i2}^\perp(a_1) + (\beta_{20} + \beta_{21}a_1)a_2 + (\theta_{20} + \theta_{21}a_1)a_2C_{i1}^\perp \\
 & + \{\theta_{22} + \theta_{23}a_1 + (\theta_{24} + \theta_{25}a_1)C_{i1}^\perp\}a_2C_{i2}^\perp(a_1).
 \end{aligned}
 \tag{16}$$

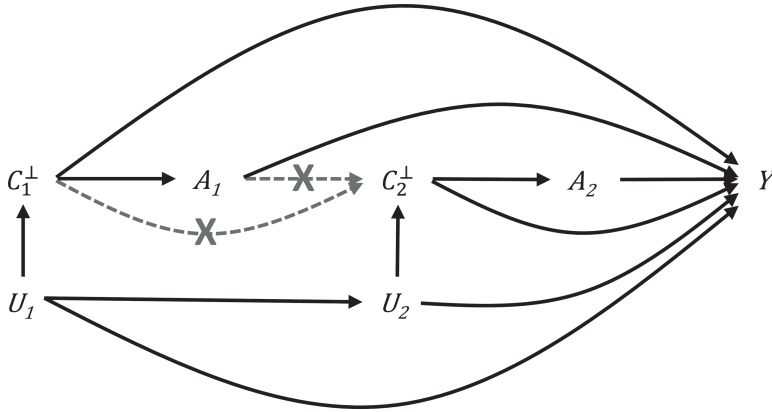
This model differs from those outlined in Almirall *et al.* (2010), Wodtke and Almirall (2017) and, for example, Vansteelandt and Sjolander’s (2016) equation (6) in that the residualized confounders, rather than the untransformed values of these variables, are included as part of interaction terms in the causal functions. It also differs from the highly constrained SNMMs that were outlined in, for example, Vansteelandt and Sjolander’s (2016) equation (3) and Wodtke (2018) in that effect modification is not assumed to be absent but rather is explicitly modelled or, in other words,  $\{\theta_{10}, \theta_{20}, \theta_{21}, \theta_{22}, \theta_{23}, \theta_{24}, \theta_{25}\}$  are free parameters that are not assumed to be zero.

An SNMM parameterized as above can be estimated by using RWR, which proceeds in two stages. In the first stage, residual terms are estimated by centring  $C_{i1}$  at its sample mean and by centring  $C_{i2}$  at its estimated conditional mean given  $C_{i1}$  and  $A_{i1}$ . Specifically,  $\hat{C}_{i1}^\perp = C_{i1} - \hat{E}(C_{i1})$  and  $\hat{C}_{i2}^\perp = C_{i2} - \hat{E}(C_{i2}|C_{i1}, A_{i1})$ , where  $\hat{E}(C_{i1}) = (1/n)\sum_i C_{i1}$  and  $\hat{E}(C_{i2}|C_{i1}, A_{i1})$  is estimated from a generalized linear model with, for example, the logit or probit link when  $C_{i2}$  is binary. Then, in the second stage, least squares estimates are computed for a linear regression of the outcome on prior treatments, the residualized confounders and interactions involving the prior treatments and residualized confounders. This regression can be expressed as

$$\begin{aligned}
 \hat{E}(Y_i|C_{i1}, A_{i1}, C_{i2}, A_{i2}) = & \hat{\beta}_{00} + \hat{\gamma}_{10}\hat{C}_{i1}^\perp + \hat{\beta}_{10}A_{i1} + \hat{\theta}_{10}A_{i1}\hat{C}_{i1}^\perp + \{\hat{\gamma}_{20} + \hat{\gamma}_{21}A_{i1} \\
 & + (\hat{\gamma}_{22} + \hat{\gamma}_{23}A_{i1})\hat{C}_{i1}^\perp\}\hat{C}_{i2}^\perp + (\hat{\beta}_{20} + \hat{\beta}_{21}A_{i1})A_{i2} + (\hat{\theta}_{20} + \hat{\theta}_{21}A_{i1})A_{i2}\hat{C}_{i1}^\perp \\
 & + \{\hat{\theta}_{22} + \hat{\theta}_{23}A_{i1} + (\hat{\theta}_{24} + \hat{\theta}_{25}A_{i1})\hat{C}_{i1}^\perp\}A_{i2}\hat{C}_{i2}^\perp,
 \end{aligned}
 \tag{17}$$

where different combinations of the estimated beta coefficients,  $\{\hat{\beta}_{00}, \hat{\beta}_{10}, \hat{\beta}_{20}, \hat{\beta}_{21}\}$ , are consistent for the marginal effects of interest under the identification assumptions that were outlined previously and under the assumption that the model is correctly specified, which is here assured by saturating it. This approach is nearly identical to conventional least squares regression except that the confounders at each time point are first residualized with respect to the observed past.

Fig. 2 displays a stylized graph that illustrates the logic of RWR estimation. It shows that residualizing the confounders at each time point with respect to the observed past purges the treatment-induced confounder  $C_{i2}$  of its association with prior treatment  $A_{i1}$ . As a result, including the residual transformation of  $C_{i2}$  in a model for the outcome avoids bias due to overcontrol and endogenous selection. In addition, because RWR adjusts for observed confounding by conditioning on residual transformations of the confounders in an outcome regression rather than by reweighting the data to balance the confounders across future treatments appropriately, it also avoids the limitations that are associated with IPTW estimation, such as the difficulty that is associated with constructing well-behaved weights for continuous treatments. Finally, by including the residualized confounders as part of interaction terms with treatment, RWR can accommodate effect modification while neatly isolating the marginal effects of interest in a single parameter vector.



**Fig. 2.** Stylized graph illustrating the logic of RWR:  $A_t$  denotes a time-varying treatment,  $C_1^\perp = C_1 - E(C_1)$  and  $C_2^\perp = C_2 - E(C_2|C_1, A_1)$  denote residualized time-varying confounders,  $U_t$  denotes an unobserved time-varying covariate and  $Y$  denotes an end-of-study outcome

**3.2. Model specification and other considerations**

In practice, estimating a saturated SNMM is often impractical, or even impossible, either because the confounders or treatments are continuous, because there are a large number of time periods or because there are many confounders. In this situation, a set of parametric constraints must be imposed on the SNMM to facilitate estimation. For example, an analyst might consider excluding all higher order interactions involving confounders at multiple time periods, in which case RWR estimation would proceed exactly as outlined previously except with the outcome regression in the second stage simplified as follows:

$$\begin{aligned} \hat{E}(Y_i|C_{i1}, A_{i1}, C_{i2}, A_{i2}) = & \hat{\beta}_{00} + \hat{\gamma}_{10}\hat{C}_{i1}^\perp + \hat{\beta}_{10}A_{i1} + \hat{\theta}_{10}A_{i1}\hat{C}_{i1}^\perp + (\hat{\gamma}_{20} + \hat{\gamma}_{21}A_{i1})\hat{C}_{i2}^\perp \\ & + (\hat{\beta}_{20} + \hat{\beta}_{21}A_{i1})A_{i2} + (\hat{\theta}_{20} + \hat{\theta}_{21}A_{i1})A_{i2}\hat{C}_{i1}^\perp \\ & + (\hat{\theta}_{22} + \hat{\theta}_{23}A_{i1})A_{i2}\hat{C}_{i2}^\perp. \end{aligned} \tag{18}$$

Of course, many other constraints are possible, but recall that RWR requires a correctly specified model for the outcome. Thus, if any of these modelling constraints are inappropriate, then RWR is biased, even when the effects of interest are identified under sequential ignorability.

Additional modelling considerations may be required of RWR when there are many time periods. Specifically, it may be necessary to impose quite stringent constraints on the SNMM to reduce its complexity. For example, an analyst might consider not merely excluding all higher order interactions but also constraining parameters to be invariant over time, in which case the outcome regression would be simplified as follows:

$$\hat{E}(Y_i|\bar{C}_i, \bar{A}_i) = \hat{\beta}_{00} + \hat{\gamma} \sum_t \hat{C}_{it}^\perp + \hat{\beta} \sum_t A_{it} + \hat{\theta} \sum_t A_{it} \hat{C}_{it}^\perp, \tag{19}$$

where overbars denote variable ‘histories’, i.e.  $\bar{C}_i = \{C_{i1}, \dots, C_{iT}\}$  and  $\bar{A}_i = \{A_{i1}, \dots, A_{iT}\}$ . Relatedly, with many time periods, it may also be necessary to constrain the models for the confounders in the first stage, and this might be accomplished by relying on some of the same simplifying constraints as considered here (e.g. no higher order interactions and time invariant coefficients). Nevertheless, RWR requires correctly specified models, which may be challenging to achieve in applications with a large number of time periods, as the potential complexity of these models increases with the dimension of the data.

Additional modelling considerations are also required with RWR when there are multiple confounders for which adjustment is necessary. First, all the different confounders must be appropriately residualized in the first stage. This is accomplished by fitting a model for each confounder at each time point by using only prior variables as predictors and then extracting their residuals. Thus, these variables are residualized exactly as outlined previously without needing to specify the causal ordering of different confounders measured concurrently in time. Second, all the residualized confounders must be included in the second-stage regression for the outcome, which may now involve additional interaction terms between treatment and the residualized confounders.

When estimating marginal effects with RWR and multiple confounders, the method can accommodate all types of treatment-by-confounder interaction except for higher order (i.e. three-way and above) interactions involving treatment and two or more different confounders measured contemporaneously. In the presence of such higher order interactions, the conditional effects of treatment cannot be conveniently decomposed and parameterized with residual terms. As an approximation, RWR could still be implemented in this situation by residualizing the various confounders and their cross-products by using separate models for each term but, with this approach, it is possible that the models could be incompatible. Thus, with multiple confounders, RWR estimation of marginal effects is most appropriate for a moderately constrained SNMM in which some especially complex forms of effect modification are assumed to be absent. Although somewhat limiting, this modelling constraint is still considerably weaker than that required of other methods for estimating marginal effects with an SNMM (e.g. Wodtke (2018)).

RWR estimation is also only appropriate for a *linear* SNMM. It is therefore best suited for use with continuous outcomes. In certain situations, it may also be suitable for use with other types of outcomes for which a linear model is sometimes a reasonable approximation, such as linear probability models for binary outcomes. Relatedly, with RWR, the outcome model is typically fitted via ordinary least squares. If the data come from a complex sampling design, RWR may also be implemented by using weighted or generalized least squares. The method therefore inherits the strengths and also some of the limitations of least squares estimators, such as sensitivity to outliers, and its performance should be assessed with regression diagnostics.

In sum, RWR estimation of a moderately constrained SNMM for marginal effects is a relatively simple adaptation of conventional least squares regression. It proceeds as follows: first, the confounders at each time point are regressed on all prior variables and then residualized and, second, the outcome is regressed on prior treatments, the residualized confounders and, to accommodate effect modification, an admissible set of interaction terms involving prior treatments and the residualized confounders. The method proposed can accommodate all types of effect modification except those involving higher order interactions between treatment and two or more confounders measured at the same point in time. RWR is consistent under the identification assumptions that were outlined previously along with the assumption of correct model specification. Because identifying correct models is challenging with many time periods and/or confounders, the method may be best suited for applications that have fewer time periods and that require adjustment for a smaller number of covariates. At the same time, however, the complexity of the modelling that is required with other methods, such as IPTW and  $g$ -estimation, also increases with the dimension of the data, as a large number of time periods and/or confounders presents a set of generic challenges for causal inference in the time-varying setting.

Valid standard errors for RWR estimates can be obtained by using the non-parametric bootstrap (Almirall *et al.*, 2014; Efron and Tibshirani, 1994). This involves repeatedly sampling  $n$  observations from the observed data with replacement to construct  $b$  replicate samples. RWR estimates are then computed from each of the  $b$  samples, and their standard deviation across

them gives the bootstrap standard error. Efron and Tibshirani (1994) suggested the use of  $b \geq 200$  replications when estimating standard errors.

### 3.3. Estimating controlled direct effects in analyses of causal mediation

In this section, we briefly demonstrate that the methods that were outlined previously can also be used to estimate CDEs in analyses of causal mediation when there are treatment-induced mediator–outcome confounders. To appreciate this, first let  $d$  denote exposure to a binary treatment, and let  $m$  denote a putative mediator that is also binary. In addition, let  $Y_i(d, m)$  denote the potential outcome for subject  $i$  if she had previously been exposed to treatment  $d$  and the mediator  $m$ . Finally, let  $X_i$  denote a treatment–outcome confounder for subject  $i$  measured at baseline, and let  $Z_i(d)$  denote a post-treatment confounder of the mediator–outcome relationship, which is indexed as a potential outcome by  $d$  to reflect that it is a treatment-induced confounder.

The CDE measures the causal relationship between treatment and the outcome when the putative mediator is fixed at the same value for all individuals. It can be formally defined as

$$\text{CDE}(d, m) = E\{Y_i(d, m) - Y_i(0, m)\}. \tag{20}$$

This quantity can be non-parametrically identified from the observed data under the assumptions of stable unit treatment values, consistency, positivity and sequential ignorability (VanderWeele, 2009, 2015). In this context, sequential ignorability is satisfied when there are no unobserved treatment–outcome or mediator–outcome confounders.

Although conventional methods are biased in the presence of treatment-induced confounders, the CDE can still be consistently estimated by using an SNMM and RWR (Zhou and Wodtke, 2019). For example, consider the following moderately constrained SNMM for the joint effects of treatment and the mediator on the outcome:

$$E\{Y_i(d, m)|X_i, Z_i(d)\} = \beta_{00} + \gamma_{10}X_i^\perp + \beta_{10}d + \theta_{10}dX_i^\perp + (\gamma_{20} + \gamma_{21}d + \gamma_{22}X_i^\perp)Z_i^\perp(d) + (\beta_{20} + \beta_{21}d)m + (\theta_{20} + \theta_{21}d)mX_i^\perp + (\theta_{22} + \theta_{23}d)mZ_i^\perp(d), \tag{21}$$

where  $X_i^\perp = X_i - E(X_i)$ ,  $Z_i^\perp(d) = Z_i(d) - E\{Z_i(d)|X_i\}$  and, for simplicity, higher order interactions involving both confounders have been excluded. With this model, the CDE is given by  $\text{CDE}(d, m) = (\beta_{10} + \beta_{21}m)d$ , any potential modification of the treatment effect by the baseline confounder  $X_i$  is captured by  $\theta_{10}$ , and any potential modification of the mediator effect by  $X_i$  or the post-treatment confounder  $Z_i(d)$  is captured by  $\{\theta_{20}, \theta_{21}, \theta_{22}, \theta_{23}\}$ . This model can be estimated via RWR by, first, centring  $X_i$  at its sample mean and centring  $Z_i$  at its estimated conditional mean given  $X_i$  and  $D_i$  and then, second, fitting a regression of the outcome on treatment, the mediator, the residualized confounders and a set of interaction terms between treatment, the mediator and the residualized confounders. As in the time-varying setting, valid standard errors can be obtained from the non-parametric bootstrap.

## 4. Simulation experiments

We use a series of simulation experiments to evaluate the performance of RWR estimation for marginal effects relative to other methods. Specifically, we use 10000 simulations of  $n = 500$  to estimate the CTE of a time-varying exposure measured at two time points. In each simulation, we generate an ‘unobserved’ continuous variable  $U_i$ , an observed continuous time-varying confounder  $\{C_{i1}, C_{i2}\}$ , a binary time-varying treatment  $\{A_{i1}, A_{i2}\}$  and a continuous end-of-study outcome  $Y_i$ . In these simulations,  $[U_i] \sim N(\mu_{U_i} = 0, \sigma_{U_i}^2 = 1)$ ,  $[C_{i1}] \sim N(\mu_{C_{i1}} = 0, \sigma_{C_{i1}}^2 = 1)$ ,  $[C_{i2}|U_i, C_{i1}, A_{i1}] \sim N(\mu_{C_{i2}|U_i, C_{i1}, A_{i1}} = 0.5U_i + 0.5C_{i1} + 0.5A_{i1}, \sigma_{C_{i2}|U_i, C_{i1}, A_{i1}}^2 = 1)$ ,  $[A_{i1}|C_{i1}] \sim$

Bernoulli $\{p_{A_{i1}|C_{i1}} = \Phi(\gamma C_{i1})\}$ ,  $[A_{i2}|C_{i1}, A_{i1}, C_{i2}] \sim \text{Bernoulli}\{p_{A_{i2}|C_{i1}, A_{i1}, C_{i2}} = \Phi(\gamma C_{i1} + 0.5A_{i1} + \gamma C_{i2})\}$  and  $[Y_i|U_i, C_{i1}, A_{i1}, C_{i2}, A_{i2}] \sim N(\mu_{Y_i|U_i, C_{i1}, A_{i1}, C_{i2}, A_{i2}} = 0.5U_i + \gamma(C_{i1} - \mu_{C_{i1}}) + A_{i1}\{0.2 + \theta(C_{i1} - \mu_{C_{i1}})\} + (C_{i2} - \mu_{C_{i2}})\{\gamma + (C_{i1} - \mu_{C_{i1}})\eta\} + A_{i2}[0.2 + 0.1A_{i1} + \theta\{(C_{i1} - \mu_{C_{i1}}) + (C_{i2} - \mu_{C_{i2}})\}], \sigma_{Y_i|U_i, C_{i1}, A_{i1}, C_{i2}, A_{i2}}^2 = 1)$ , where  $\Phi$  is the standard normal cumulative distribution function,  $\gamma$  is a parameter that is used to control the level of treatment–outcome confounding,  $\theta$  is a parameter that is used to control the level of treatment effect modification and  $\eta$  is a parameter that is used to control the degree to which the associational effect of  $C_{i2}$  on  $Y_i$  is modified by  $C_{i1}$ . In all simulations, the CTE is identified and its true value is 0.5.

We compare the performance of RWR estimation of a moderately constrained SNMM for marginal effects (henceforth ‘RWR with interactions’) with the performance of conventional least squares regression, IPTW estimation of a marginal structural model,  $g$ -estimation of a highly constrained SNMM in which effect modification is assumed to be absent and RWR estimation of the same highly constrained SNMM (henceforth ‘RWR without interactions’). To compute conventional regression estimates, we fit by least squares a linear regression of the outcome on prior treatments, the observed confounders and a treatment-by-treatment interaction. The estimated CTE is then given by the sum of the coefficients on prior treatments and the interaction term.

To compute IPTW estimates (Robins *et al.*, 1994, 2000), we fit a linear regression of the outcome on prior treatments and their interaction by using weighted least squares, with weights equal to

$$w_i = \frac{P(A_{i1} = a_{i1})}{P(A_{i1} = a_{i1}|C_{i1})} \frac{P(A_{i2} = a_{i2}|A_{i1} = a_{i1})}{P(A_{i2} = a_{i2}|C_{i1}, A_{i1} = a_{i1}, C_{i2})}, \tag{22}$$

where  $w_i$  is estimated from a series of probit models for the conditional probabilities in the numerator and denominator of the weight. At each time point, weighting by  $w_i$  balances (in expectation) prior confounders across future treatments by giving more weight to subjects with confounder histories that are underrepresented in a treatment group and less weight to subjects with confounder histories that are overrepresented in a treatment group. The estimated CTE is the sum of the coefficients on prior treatments and their interaction.

To compute  $g$ -estimates of marginal effects by using a highly constrained SNMM without any effect modification, we use the  $g$ -estimator that was proposed by Vansteelandt and Sjolander (2016). Specifically, we first fit a linear regression of the outcome on prior treatments and their interaction, estimated propensity scores for treatment at each time point, an interaction between treatment at time  $t = 1$  and the estimated propensity score for treatment at time  $t = 2$ , and the observed confounders at both time points. The coefficients on treatment at time  $t = 2$  and its interaction with treatment at time  $t = 1$  from this model provide estimates of the PTE. Then, we subtract the estimated PTE from the outcome for each respondent and regress this transformed outcome on the treatment, propensity score and the observed confounder at time  $t = 1$ . The coefficient on treatment from this model provides an estimate of the DTE and then the sum of the distal and proximal effects computed as above give the estimated CTE. Vansteelandt and Sjolander (2016) showed that this estimator is asymptotically equivalent to the doubly robust  $g$ -estimator that was considered in Robins *et al.* (1992).

To compute estimates based on RWR without interactions, we first residualize the observed confounders at each time point by regressing them on all prior variables and then centring them at their estimated conditional means. Second, we regress the outcome on prior treatments and their interaction as well as the residualized confounders. The estimated CTE is the sum of the coefficients on prior treatments and their interaction. Computing estimates based on RWR with interactions proceeds in almost exactly the same manner except that all two-way interactions

between the treatments and residualized confounders are additionally included in the outcome regression. The R code that was used to execute all the simulations outlined previously can be obtained from

<https://rss.onlinelibrary.wiley.com/hub/journal/1467985x/series-a-datasets>.

We compare the performance of these methods in terms of their bias, standard deviation SD and root-mean-squared error RMSE under various levels of treatment–outcome confounding and under various levels of effect modification. Because treatment-induced confounding is present in all simulations, we expect conventional regression to perform poorly across all scenarios. Because IPTW estimation is relatively inefficient and susceptible to finite sample bias when confounders strongly predict treatment, we expect its performance to suffer in simulations with high levels of treatment–outcome confounding. Because  $g$ - and RWR estimation of marginal effects by using a highly constrained SNMM require that the confounders must not modify the effects of treatment, we expect their performance to deteriorate in simulations with high levels of treatment effect modification. Finally, because RWR with interactions accommodates this type of effect modification, we expect it to perform well across all simulations.

Table 1 presents results from a first set of simulation experiments, wherein we varied the level of treatment–outcome confounding in the absence of effect modification. Conventional regression is badly biased at all levels of confounding, as expected. IPTW estimation performs well at lower levels of confounding but suffers from finite sample bias at higher levels and is relatively

**Table 1.** Results from simulation experiments evaluating the performance of RWR relative to other estimators under different levels of treatment–outcome confounding†

<i>Estimator and statistic</i>	<i>Magnitude of treatment–outcome confounding</i>				
	$\gamma=0.1$	$\gamma=0.2$	$\gamma=0.3$	$\gamma=0.4$	$\gamma=0.5$
Conventional regression					
Bias	−0.150	−0.200	−0.252	−0.299	−0.351
SD	0.134	0.137	0.141	0.145	0.151
RMSE	0.201	0.242	0.288	0.332	0.382
IPTW estimation					
Bias	0.000	0.002	0.005	0.020	0.057
SD	0.135	0.147	0.176	0.228	0.296
RMSE	0.135	0.147	0.176	0.229	0.302
$g$ -estimation					
Bias	0.000	0.000	−0.001	0.002	0.000
SD	0.134	0.139	0.145	0.152	0.163
RMSE	0.134	0.139	0.145	0.152	0.163
RWR without interactions					
Bias	0.000	0.000	−0.001	0.002	0.000
SD	0.134	0.139	0.145	0.151	0.161
RMSE	0.134	0.139	0.145	0.151	0.161
RWR with interactions					
Bias	0.000	0.000	−0.001	0.002	−0.001
SD	0.134	0.140	0.146	0.154	0.164
RMSE	0.134	0.140	0.146	0.154	0.164

†Results are based on 10000 simulations. See the on-line supplement for details. Across all simulations,  $\theta=0$  and  $\eta=0$ .



**Table 2.** Results from simulation experiments evaluating the performance of RWR relative to other estimators under different levels of treatment effect modification and a moderate-to-high level of treatment–outcome confounding†

Estimator and statistic	Magnitude of treatment effect modification				
	$\theta=0.1$	$\theta=0.2$	$\theta=0.3$	$\theta=0.4$	$\theta=0.5$
Conventional regression					
Bias	-0.369	-0.439	-0.508	-0.575	-0.645
SD	0.145	0.149	0.151	0.155	0.163
RMSE	0.396	0.463	0.530	0.595	0.665
IPTW estimation					
Bias	0.022	0.024	0.023	0.028	0.021
SD	0.235	0.246	0.261	0.274	0.299
RMSE	0.236	0.247	0.262	0.275	0.300
<i>g</i> -estimation					
Bias	-0.023	-0.047	-0.071	-0.094	-0.119
SD	0.155	0.161	0.164	0.168	0.177
RMSE	0.157	0.168	0.179	0.193	0.213
RWR without interactions					
Bias	-0.037	-0.076	-0.115	-0.151	-0.192
SD	0.154	0.161	0.166	0.171	0.182
RMSE	0.159	0.178	0.202	0.228	0.264
RWR with interactions					
Bias	0.001	0.001	0.001	0.000	-0.001
SD	0.156	0.161	0.164	0.167	0.175
RMSE	0.156	0.161	0.164	0.167	0.175

†Results are based on 10000 simulations. Across all simulations,  $\gamma=0.4$  and  $\eta=0$ . See the on-line supplement for details.

inefficient, also as expected. *g*- and both types of RWR estimation perform similarly in these simulations: they are all unbiased and achieve comparable gains in efficiency relative to IPTW.

Table 2 presents results from a second set of simulation experiments, wherein we varied the level of treatment effect modification after setting the level of treatment–outcome confounding at a moderate-to-high level. As expected, both conventional regression and IPTW estimation perform poorly. Although IPTW estimation accommodates effect modification, it still suffers from finite sample bias due to the high level of confounding and is relatively inefficient. Also as expected, *g*-estimation and RWR without interactions are increasingly biased as the magnitude of treatment effect modification rises, whereas RWR with interactions is unbiased and achieves the lowest RMSE across all scenarios.

Finally, because RWR requires strong modelling assumptions that may be difficult to satisfy in some applications, we evaluate the method’s performance with an SNMM that has an incorrectly specified nuisance association. Specifically, Table 3 presents results from a third set of simulation experiments, wherein we varied the degree to which the associational effect of  $C_{i2}$  on  $Y_i$  is modified by  $C_{i1}$  after setting the level of both treatment–outcome confounding and treatment effect modification at moderate-to-high levels. Because all the SNMMs that were considered in these simulations constrain the associational effects of the confounders to be invariant, they are all incorrectly specified when this type of effect modification is present.

In Table 3, both conventional regression and IPTW estimation perform poorly because, as before, these simulations involve a moderate-to-high level of treatment-induced confounding. *g*-estimation and RWR without interactions also do not perform very well because these sim-



**Table 3.** Results from simulation experiments evaluating the performance of RWR relative to other estimators under different levels of nuisance model misspecification and moderate-to-high levels of both treatment–outcome confounding and treatment effect modification†

<i>Estimator and statistic</i>	<i>Magnitude of associational effect modification</i>				
	$\eta=0.1$	$\eta=0.2$	$\eta=0.3$	$\eta=0.4$	$\eta=0.5$
Conventional regression					
Bias	−0.586	−0.597	−0.608	−0.616	−0.627
SD	0.158	0.162	0.164	0.166	0.174
RMSE	0.606	0.618	0.630	0.638	0.651
IPTW estimation					
Bias	0.030	0.031	0.029	0.033	0.027
SD	0.276	0.280	0.292	0.296	0.318
RMSE	0.277	0.282	0.294	0.297	0.319
<i>g</i> -estimation					
Bias	−0.093	−0.094	−0.095	−0.094	−0.095
SD	0.169	0.173	0.175	0.176	0.183
RMSE	0.193	0.197	0.199	0.200	0.206
RWR without interactions					
Bias	−0.160	−0.171	−0.182	−0.189	−0.201
SD	0.173	0.176	0.179	0.179	0.188
RMSE	0.236	0.246	0.255	0.261	0.275
RWR with interactions					
Bias	0.015	0.026	0.039	0.050	0.063
SD	0.169	0.172	0.174	0.176	0.185
RMSE	0.169	0.174	0.179	0.183	0.195

†Results are based on 10000 simulations. Across all simulations,  $\gamma=0.4$  and  $\theta=0.4$ . See the on-line supplement for details.

ulations involve a non-trivial level of effect modification. Note, however, that the performance of *g*-estimation is similar regardless of the degree to which the nuisance associations are incorrectly specified, which reflects its doubly robust property. The performance of RWR without interactions, by contrast, further deteriorates with greater misspecification of the nuisance associations. Similarly, the bias and RMSE of RWR with interactions also steadily increases with the degree to which the nuisance associations are incorrectly specified. Nevertheless, RWR with interactions still appears to outperform other methods, even when it is based on an SNMM with misspecified nuisance associations.

## 5. Empirical examples

### 5.1. Cumulative treatment effect of neighbourhood poverty on academic achievement

Does growing up in a disadvantaged neighbourhood inhibit academic achievement? The effects of neighbourhood composition on child development have long concerned social scientists across several disciplines (e.g. Chetty *et al.* (2016), Leventhal and Brooks-Gunn (2000), Sampson *et al.* (2008) and Wodtke *et al.* (2011)). To illustrate how the method proposed can be used with time-varying treatments, we estimate the CTE of residence in a disadvantaged neighbourhood throughout the early life course on adolescent mathematics achievement by using data from  $n=1135$  individuals in the Panel Study of Income Dynamics—Child Development Supplement (PSID—CDS) (Michigan Survey Research Center, 2014).

In these data, the outcome  $Y_i$  represents standardized scores on the Woodcock–Johnson applied problems test measured at the end of follow-up when individuals were age 13–17 years (Woodcock and Bonner Johnson, 1989). The time-varying treatment  $a_t$  is a standardized index of neighbourhood disadvantage generated via a principal component analysis of multiple census tract characteristics, such as the poverty rate, unemployment rate and median household income. Treatment is first measured during childhood when individuals were age 5–9 years and then again during adolescence when they were age 11–15 years. The vector of baseline confounders  $C_{i1}$  contains a set of time invariant factors, such as race, gender and birth cohort, as well as a set of time-varying characteristics, including equivalized family income, parental marital status and lagged achievement test scores, all measured during early childhood. Another vector of confounders  $C_{i2}$  contains the same set of time-varying characteristics only now measured just before the onset of adolescence.

Previously, Wodtke (2018) estimated the CTE of residence in a disadvantaged neighbourhood with data from the PSID—CDS by fitting a conventional regression model using least squares, a marginal structural model using IPTW and a highly constrained SNMM without any effect modification by using RWR. In that analysis, RWR estimates indicated that long-term residence in a disadvantaged neighbourhood has a severe negative effect on mathematics achievement—an effect that is obscured by bias in conventional regression models and imprecisely captured by IPTW. These estimates, however, are premised on the strong and arguably unrealistic assumption of no effect modification.

We overcome this limitation by estimating the CTE by using RWR and an SNMM that includes all two-way treatment-by-confounder interactions. Specifically, we model the distal, proximal and cumulative marginal effects of exposure to a disadvantaged neighbourhood on adolescent mathematics achievement by using the following SNMM:

$$E\{Y_i(a_1, a_2) | C_{i1}, C_{i2}(a_1)\} = \beta_{00} + \gamma_{10}^T C_{i1}^\perp + \beta_{10} a_1 + \theta_{10}^T a_1 C_{i1}^\perp + \gamma_{20}^T C_{i2}^\perp(a_1) + \beta_{20} a_2 + \theta_{20}^T a_2 C_{i1}^\perp + \theta_{21}^T a_2 C_{i2}^\perp(a_1), \quad (23)$$

where  $\beta_{10} = \text{DTE}(1)$ ,  $\beta_{20} = \text{PTE}(a_1, 1)$  and  $\beta_{10} + \beta_{20} = \text{CTE}$ . Note that this model is just a moderately constrained version of equation (16).

The first row of Table 4 reports RWR estimates for the distal, proximal and cumulative effects of living in a disadvantaged neighbourhood. We compute these estimates by, first, centring the elements of  $C_{i1}$  at their sample means and centring the elements of  $C_{i2}$  at their estimated conditional means, which come from least squares regressions of  $C_{i2}$  on the treatment and confounders measured earlier during childhood. Second, we compute the marginal effects of interest by regressing the outcome on both treatments, the residualized confounders and all two-way interactions between the treatments and residualized confounders. For comparison, the second and third rows of Table 4 report RWR and  $g$ -estimates of marginal effects from a highly constrained SNMM in which all treatment-by-confounder interactions have been excluded.

All the estimates in Table 4 indicate that the distal effect of childhood exposure to a disadvantaged neighbourhood on adolescent mathematics achievement is substantively small and fails to reach conventional significance thresholds, that the proximal effect of adolescent exposure is larger and statistically significant, and that the cumulative effect of sustained exposure is substantively large and highly significant. For example, according to these estimates, sustained exposure to a poor neighbourhood 1 standard deviation above the national mean of the disadvantage index, rather than sustained exposure to a wealthy neighbourhood 1 standard deviation below the national mean, is estimated to reduce adolescent mathematics achievement by about  $0.127 \times 2 = 0.254$  standard deviations.

**Table 4.** Estimated marginal effects of exposure to disadvantaged neighbourhoods on end-of-study mathematics achievement†

Estimator	DTE(1,0)		PTE(a <sub>1</sub> ,1)		CTE	
	Point estimate	SE	Point estimate	SE	Point estimate	SE
RWR with interactions	-0.034	(0.049)	-0.094	(0.046)‡	-0.127	(0.038)§§
RWR without interactions	-0.030	(0.044)	-0.097	(0.040)‡	-0.127	(0.038)§§
g-estimation	-0.032	(0.040)	-0.096	(0.041)‡	-0.127	(0.047)§

†The sample includes respondents who were interviewed at the 1997 wave of the PSID—CDS between age 3 and 7 years. Results are combined estimates from 100 imputations. The outcome is standardized to have mean zero and unit variance. Standard errors SE are based on the block bootstrap with 1000 replications.

‡ $p < 0.05$  for two-sided tests of no effect.

§ $p < 0.01$  for two-sided tests of no effect.

§§ $p < 0.001$  for two-sided tests of no effect.

The results in Table 4 are similar across the various methods that were employed. Consistent with this finding, a Wald test of the null hypothesis that  $\theta_{10} = \theta_{20} = \theta_{21} = 0$  provides little evidence that the effects of living in a disadvantaged neighbourhood are modified by any of the confounders ( $\chi^2 = 10.7$ ; degrees of freedom  $df = 17$ ;  $p = 0.873$ ). In this application, it therefore appears that our findings are not particularly sensitive to the presence or absence of effect modification. Nevertheless, it is the flexibility of RWR that enables an easy assessment of marginal effects under different specifications.

### 5.2. Controlled direct effect of education on mental health

Does income explain the effect of education on depression? Some prior studies have investigated the causal relationship between education and mental health (e.g. Cutler and Lleras-Muney (2006), Heckman *et al.* (2018) and Lee (2011)), but the mechanisms underlying this causal link remain unclear. Education may improve mental health by providing access to higher economic status and greater resources, or it may affect mental health through other channels—e.g. by providing greater access to health information and improving health behaviours. To illustrate the utility of RWR for analyses of causal mediation, we estimate the CDE of college completion on mental health controlling for family income as a putative mediator. In this example, a comparison of the total effect with the CDE helps to adjudicate whether family economic status explains the mental health benefits of college completion.

We use data from  $n = 2719$  individuals in the 1979 National Longitudinal Survey of Youth who were age 14–17 years when they were first interviewed in 1979. In these data, the outcome  $Y_i$  represents scores on the Center for Epidemiologic Studies depression scale when respondents were age 40 years. We standardize the depression scale scores to have mean 0 and unit variance, where a higher score implies greater depression. The treatment  $d$  is defined as completion of a 4-year college degree by age 25 years, whereas the mediator of interest,  $m$ , is the percentile rank of equivalized family income averaged over ages 36–40 years. The vector of baseline confounders  $X_i$  includes gender, race, Hispanic ethnicity, mother’s years of schooling, father’s presence in the home, number of siblings, urban residence, educational expectations and percentile scores on the armed forces qualification test. Finally, the vector of post-treatment confounders  $Z_i$  includes the

depression scale scores measured when respondents were age 27–30 years, the proportion of time married between 1990 and 1998, and the number of relationship transitions between 1990 and 1998. These variables capture mental health and family stability during young adulthood, which may be affected by treatment (college completion by age 25 years) and also affect both the mediator (family income between age 36 and 40 years) and the outcome (depression at age 40 years).

With these data, we first estimate the total effect of college completion by using the following SNMM:

$$E\{Y_i(d)|X_i\} = \beta_{00} + \gamma_{10}^T X_i^\perp + \beta_{10}d + \theta_{10}^T dX_i^\perp. \tag{24}$$

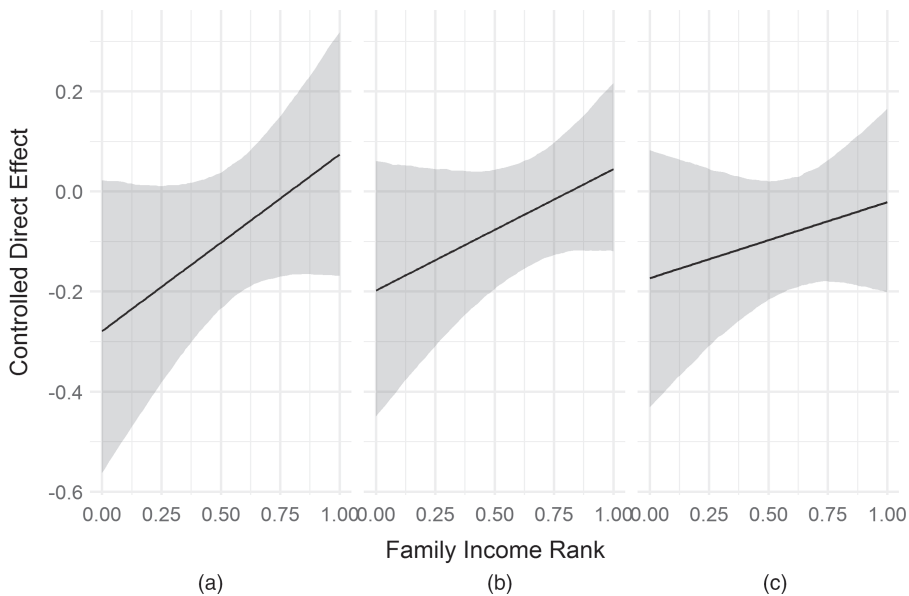
Under this specification,  $\beta_{10}$  captures the total effect of college completion on depression at age 40 years. RWR estimates of this model yield a sizable and statistically significant total effect of education on mental health, where completing college is estimated to lower depression scores by 0.165 standard deviations (standard error SE = 0.066;  $p < 0.05$ ) on average.

We then model the joint effects of college completion and family income on depression by using the following SNMM:

$$E\{Y_i(d, m)|X_i, Z_i(d)\} = \beta_{00} + \gamma_{10}^T X_i^\perp + \beta_{10}d + \theta_{10}^T dX_i^\perp + \gamma_{20}^T Z_i^\perp(d) + (\beta_{20} + \beta_{21}d)m + \theta_{20}^T mX_i^\perp + \theta_{21}^T mZ_i^\perp(d), \tag{25}$$

where the CDE is given by  $CDE(m) = \beta_{10} + \beta_{21}m$ . Note that this model is just a moderately constrained version of equation (21).

Fig. 3(a) reports estimates for the CDEs of college completion on depression computed by using RWR with interactions. These estimates are obtained by, first, computing residuals for each of the baseline confounders  $X_i$  and post-treatment confounders  $Z_i$ , which involves centring the elements of  $X_i$  at their sample means and centring the elements of  $Z_i$  at their



**Fig. 3.** Estimated CDEs of college completion on depression (controlling for family income) (the sample includes respondents to the 1979 National Longitudinal Survey of Youth who were age 13–17 years when first interviewed; confidence intervals are based on the non-parametric bootstrap with 1000 replications): (a) RWR with interactions; (b) RWR without interactions; (c) *g*-estimation

estimated conditional means given the past. Second, the CDEs are then estimated by fitting the model that was described previously by using these residual terms. Because  $CDE(m)$  may vary with  $m$ , we estimate and plot CDEs across the support of the mediator. For comparison, Figs 3(b) and 3(c) report RWR and  $g$ -estimates based on a highly constrained SNMM in which all treatment-by-confounder and mediator-by-confounder interactions are excluded.

These estimates provide some evidence that the effect of education on depression is mediated by family income. For example, point estimates of the CDE when family income is fixed at or above its sample median are substantially smaller than the estimated total effect, and the 95% confidence intervals contain 0 at every value of family income. Thus, these results suggest that at least some portion of the total effect operates through pathways involving family economic resources.

The estimates in Fig. 3 are fairly consistent across the various methods employed. Moreover, a Wald test of the null hypothesis that  $\theta_{10} = \theta_{20} = \theta_{21} = 0$  does not provide much evidence of effect modification ( $\chi^2 = 26.9$ ;  $df = 21$ ;  $p = 0.175$ ). Thus, it appears that our findings are insensitive to the inclusion of treatment-by-confounder interactions in this application as well, but recall that the flexibility of RWR is what enables us to assess easily whether estimated marginal effects are robust to different specifications.

## 6. Discussion and conclusions

In analyses of causal mediation and time-varying treatment effects, treatment-induced confounders often complicate efforts to estimate marginal effects. Several available methods avoid these complications, including marginal structural models and IPTW as well as  $g$ - and RWR estimation of highly constrained SNMMs, but they are not without limitations. Specifically, the performance of IPTW is poor with continuous treatments and/or mediators, a high degree of confounding and small samples, whereas both  $g$ - and RWR estimation of highly constrained SNMMs are biased for the marginal effects of interest when effect modification is present. To overcome these limitations, we adapt the method of RWR to estimate marginal effects with a set of moderately constrained SNMMs that easily accommodate several types of effect modification as well as continuous treatments and/or mediators. A series of simulation experiments indicate that the method proposed outperforms IPTW estimation of marginal structural models in general and that it outperforms both  $g$ - and RWR estimation of highly constrained SNMMs in the presence of effect modification. Because the method involves only simple and familiar computations, it is easily implemented with standard software, as we demonstrate across two empirical illustrations.

Nevertheless, despite its many advantages, RWR estimation of marginal effects is premised on some strong modelling assumptions. Specifically, it requires a correctly specified SNMM, which in turn requires that all the causal functions and nuisance associations that compose this model are correctly specified. It also requires the absence of more complex forms of effect modification involving two or more confounders measured contemporaneously, which complicates the decomposition and parameterization of the SNMM causal functions using residual terms. The assumption of a correctly specified SNMM may be reasonable with a relatively small number of confounders and time periods, but identifying a correct model may be challenging with high dimensional data.

In this situation, researchers might consider combining the methods that were proposed in this study with either IPTW or  $g$ -estimation to leverage their strengths while mitigating their weaknesses. For example, RWR could be used to adjust for a subset of the time-varying confounders that prove difficult to balance appropriately by using IPTW. Then, a simplified SNMM

involving only this subset of confounders and a more limited set of interaction terms could be fitted by RWR to an appropriately weighted sample in which the remaining confounders have all been balanced. Alternatively, the confounders could first be residualized with respect to the observed past and then included in interaction terms with treatment and/or a mediator at each stage of the  $g$ -estimation procedure that was outlined by Vansteelandt and Sjolander (2016). This may provide some protection against bias due to misspecification of the nuisance associations in an SNMM, as  $g$ -estimation is doubly robust, while simultaneously accommodating several types of effect modification in analyses of marginal effects. RWR might also be combined with variable-selection and regularization techniques, such as the lasso, in an effort to identify sufficiently accurate yet parsimonious models in applications with many time periods or confounders.

In sum, RWR estimation of a moderately constrained SNMM for marginal effects provides an appealing alternative to IPTW estimation of marginal structural models and to both  $g$ - and RWR estimation of highly constrained SNMMs in which effect modification is assumed away. The method proposed improves on IPTW estimation in that it is more efficient, easy to use with continuous treatments and/or mediators, and avoids finite sample bias when the magnitude of observed confounding is strong. It improves on  $g$ - and RWR estimation of highly constrained SNMMs in that it can accommodate all except highly complex forms of effect modification while still neatly isolating the marginal effects of interest in a single set of parameters. Although the method proposed is premised on some strong modelling assumptions, it can be integrated with IPTW or  $g$ -estimation in situations where these assumptions are questionable to enhance its robustness. Given their flexibility, efficiency and ease of use, we expect moderately constrained SNMMs along with the associated method of RWR to be frequently used in future studies of causal mediation and time-varying treatment effects.

## Acknowledgements

This research was supported by an insight grant from the Social Sciences and Humanities Research Council of Canada (grant 435-2018-0736).

## References

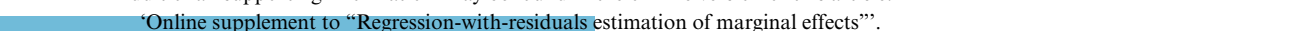
- Acharya, A., Blackwell, M. and Sen, M. (2016) Explaining causal findings without bias: detecting and assessing direct effects. *Am. Polit. Sci. Rev.*, **110**, 512–529.
- Almirall, D., Griffin, B. A., McCaffrey, D. F., Ramchand, R., Yuen, R. A. and Murphy, S. A. (2014) Time-varying effect moderation using the structural nested mean model: estimation using inverse-weighted regression-with-residuals. *Statist. Med.*, **33**, 3466–3487.
- Almirall, D., Ten Have, T. and Murphy, S. A. (2010) Structural nested mean models for assessing time-varying effect moderation. *Biometrics*, **66**, 131–139.
- Chetty, R., Hendren, N. and Katz, L. (2016) The effects of exposure to better neighborhoods on children: new evidence from the Moving to Opportunity project. *Am. Econ. Rev.*, **106**, 855–902.
- Cutler, D. M. and Lleras-Muney, A. (2006) Education and health: evaluating theories and evidence. *Working Paper 12352*. National Bureau of Economic Research, Cambridge.
- Daniel, R. M., Cousens, S. N., De Stavola, B. L., Kenward, M. G. and Sterne, J. A. C. (2013) Methods for dealing with time-dependent confounding. *Statist. Med.*, **32**, 1584–1618.
- Efron, B. and Tibshirani, R. J. (1994) *An Introduction to the Bootstrap*. New York: Chapman and Hall.
- Elwert, F. and Winship, C. (2014) Endogenous selection bias: the problem of conditioning on a collider variable. *A. Rev. Sociol.*, **40**, 31–53.
- Heckman, J. J., Humphries, J. E. and Veramendi, G. (2018) The nonmarket benefits of education and ability. *J. Hum. Cap.*, **12**, 282–304.
- Holland, P. W. (1986) Statistics and causal inference. *J. Am. Statist. Ass.*, **81**, 945–960.
- Lee, J. (2011) Pathways from education to depression. *J. Cross Cult. Geront.*, **26**, 121–135.
- Leventhal, T. and Brooks-Gunn, J. (2000) The neighborhoods they live in: the effects of neighborhood residence on child and adolescent outcomes. *Psychol. Bull.*, **126**, 309–337.



- Lunceford, J. K. and Davidian, M. (2004) Stratification and weighting via the propensity score in estimation of causal treatment effects: a comparative study. *Statist. Med.*, **23**, 2937–2960.
- Michigan Survey Research Center (2014) *Panel Study of Income Dynamics Public- and Restricted-use Datasets*. Ann Arbor: Institute for Social Research.
- Morgan, S. L. and Winship, C. (2015) *Counterfactuals and Causal Inference: Methods and Principles for Social Research*, 2nd edn. New York: Cambridge University Press.
- Naimi, A. I., Cole, S. R. and Kennedy, E. H. (2017) An introduction to G methods. *Int. J. Epidemiol.*, **46**, 756–762.
- Naimi, A. I., Moodie, E. E. M., Auger, N. and Kaufman, J. S. (2014) Constructing inverse probability weights for continuous exposures: a comparison of methods. *Epidemiology*, **25**, 292–299.
- Pearl, J. (2009) *Causality*, 2nd edn. New York: Cambridge University Press.
- Robins, J. M. (1994) Correcting for noncompliance in randomized trials using structural nested mean models. *Communs Statist. Theory Meth.*, **23**, 2379–2412.
- Robins, J. M., Hernán, M. A. and Brumback, B. (2000) Marginal structural models and causal inference in epidemiology. *Epidemiology*, **11**, 550–560.
- Robins, J. M., Hernán, M. A. and Rotnitzky, A. (2007) Effect modification by time-varying covariates. *Am. J. Epidemiol.*, **166**, 994–1002.
- Robins, J. M., Mark, S. D. and Newey, W. K. (1992) Estimating exposure effects by modelling the expectation of exposure conditional on confounders. *Biometrics*, **48**, 479–495.
- Robins, J. M., Rotnitzky, A. and Zhao, L. P. (1994) Estimation of regression coefficients when some regressors are not always observed. *J. Am. Statist. Ass.*, **89**, 846–866.
- Rubin, D. B. (1974) Estimating causal effects of treatments in randomized and nonrandomized studies. *J. Educ. Psychol.*, **66**, 688–701.
- Rubin, D. B. (1986) Which ifs have causal answers? *J. Am. Statist. Ass.*, **81**, 961–962.
- Sampson, R. J., Sharkey, P. and Raudenbush, S. W. (2008) Durable effects of concentrated disadvantage on verbal ability among African-American children. *Proc. Natn. Acad. Sci. USA*, **105**, 845–852.
- VanderWeele, T. (2009) Marginal structural models for the estimation of direct and indirect effects. *Epidemiology*, **20**, 18–26.
- VanderWeele, T. (2015) *Explanation in Causal Inference: Methods for Mediation and Interaction*. New York: Oxford University Press.
- Vansteelandt, S. (2009) Estimating direct effects in cohort and case-control studies. *Epidemiology*, **20**, 851–860.
- Vansteelandt, S. and Sjolander, A. (2016) Revisiting G-estimation of the effect of a time-varying exposure subject to time-varying confounding. *Epidem. Meth.*, **5**, 37–46.
- Wodtke, G. T. (2018) Regression-based adjustment for time-varying confounders. *Sociol. Meth. Res.*, to be published, doi 10.1177/0049124118769087.
- Wodtke, G. T. and Almirall, D. (2017) Estimating moderated causal effects with time-varying treatments and time-varying moderators: structural nested mean models and regression-with-residuals. *Sociol. Methodol.*, **47**, 212–245.
- Wodtke, G. T., Elwert, F. and Harding, D. J. (2016) Neighborhood effect heterogeneity by family income and developmental period. *Am. J. Sociol.*, **121**, 1168–1222.
- Wodtke, G. T., Harding, D. J. and Elwert, F. (2011) Neighborhood effects in temporal perspective: the impact of long-term exposure to concentrated disadvantage on high school graduation. *Am. Sociol. Rev.*, **76**, 713–736.
- Woodcock, R. W. and Bonner Johnson, M. (1989) *Tests of Achievement, Standard Battery (Form B)*. Chicago: Riverside.
- Xie, Y. (2007) Otis Dudley Duncan's legacy: the demographic approach to quantitative reasoning in social science. *Res. Soc. Strat. Mobil.*, **25**, 141–156.
- Zhou, X. and Wodtke, G. T. (2019) A regression-with-residuals method for estimating controlled direct effects. *Polit. Anal.*, **27**, 360–369.

#### Supporting information

Additional 'supporting information' may be found in the on-line version of this article:

'Online supplement to "Regression-with-residuals estimation of marginal effects"'.  




Copyright of Journal of the Royal Statistical Society: Series A (Statistics in Society) is the property of Wiley-Blackwell and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.